
Variational Gaussian Process

Dustin Tran
Harvard University
dtran@g.harvard.edu

Rajesh Ranganath
Princeton University
rajeshr@cs.princeton.edu

David M. Blei
Columbia University
blei@cs.columbia.edu

Abstract

Representations offered by deep generative models are fundamentally tied to their inference method from data. Variational inference methods require a rich family of approximating distributions. We construct the **variational Gaussian process (vGP)**, a Bayesian nonparametric model which adapts its shape to match complex posterior distributions. The **vGP** generates approximate posterior samples by generating latent inputs and warping them through random non-linear mappings; the distribution over random mappings is learned during inference, enabling the transformed outputs to adapt to varying complexity. We prove a universal approximation theorem for the **vGP**, demonstrating its representative power for learning any model. For inference we present a variational objective inspired by auto-encoders and perform black box inference over a wide class of models. The **vGP** achieves new state-of-the-art results for unsupervised learning, inferring models such as the deep latent Gaussian model and the recently proposed DRAW.

1 Introduction

Originally developed in the 1990s (Hinton and Van Camp, 1993; Waterhouse et al., 1996; Jordan et al., 1999), variational inference has enjoyed renewed interest around developing scalable optimization for large datasets (Hoffman et al., 2013), deriving generic strategies for easily fitting many models (Ranganath et al., 2014), and applying neural networks as a flexible parametric family of approximations (Kingma and Welling, 2014; Rezende et al., 2014). This research has been particularly successful for computing with deep Bayesian models (Neal, 1990; Hinton et al., 2006), which require inference of a complex posterior distribution.

Despite these advances, research has been constrained by the lack of an expressive family of approximating distributions, which can generalize across many models. Newer research aims toward richer families that allow the latent variables to be dependent. One way to introduce dependence is to consider the variational family itself as a model of the hidden variables (Lawrence, 2000; Salimans et al., 2015; Ranganath et al., 2015), and then to use models that go beyond the simple mean field. These *variational models* naturally extend to Bayesian hierarchies, which retain the mean-field “likelihood” but introduce dependence through variational latent variables.

In this paper we develop a powerful new variational model—the **variational Gaussian process (vGP)**. The **vGP** is a Bayesian nonparametric variational model; its complexity grows efficiently and towards *any* distribution, adapting to the inference problem at hand. We present a variational lower bound inspired by auto-encoders and, for black box inference, derive an efficient stochastic optimization algorithm. We report new state-of-the-art results on the binarized MNIST data set.

2 Variational model

Let $p(\mathbf{z} | \mathbf{x})$ denote a posterior distribution over d latent variables $\mathbf{z} = (z_1, \dots, z_d)$ conditioned on a data set \mathbf{x} . Assume a choice of mean-field distribution $\prod_{i=1}^d q(z_i; \lambda_i)$. For example, each factor can simply be a Gaussian with parameters $\lambda_i = (\mu_i, \sigma_i^2)$. The **vGP** generates \mathbf{z} by generating latent

inputs, warping them with random non-linear mappings, and using the warped inputs as parameters to a mean-field distribution. The random mappings are drawn conditional on “variational data,” which is itself learned as part of variational inference. We will show that the **vGP** enables samples from the mean-field to follow arbitrarily complex posteriors.

Let $\mathcal{D} = \{(\mathbf{s}_n, \mathbf{t}_n)\}_{n=1}^m$ be variational data, comprising input-output pairs that are parameters to the variational distribution. (This idea appears in a different context in [Blei and Lafferty \(2006\)](#).) The **vGP** specifies the following generative process for posterior latent variables \mathbf{z} :

1. Draw latent input $\xi \in \mathbb{R}^c$: $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
2. Draw non-linear mapping $f : \mathbb{R}^c \rightarrow \mathbb{R}^d$ conditioned on \mathcal{D} : $f \sim \prod_{i=1}^d \mathcal{GP}(\mathbf{0}, \mathbf{K}_{\xi\xi}) \mid \mathcal{D}$.
3. Draw approximate posterior samples $\mathbf{z} \in \text{supp}(p)$: $\mathbf{z} = (z_1, \dots, z_d) \sim \prod_{i=1}^d q(f_i(\xi))$.

[Figure 1](#) displays a graphical model for the **vGP**. Marginalizing over all non-linear mappings and latent inputs, the **vGP** is

$$q_{\text{vGP}}(\mathbf{z}; \boldsymbol{\theta}, \mathcal{D}) = \iint \left[\prod_{i=1}^d q(z_i \mid f_i(\xi)) \right] \left[\prod_{i=1}^d \mathcal{GP}(f_i; \mathbf{0}, \mathbf{K}_{\xi\xi}) \mid \mathcal{D} \right] \mathcal{N}(\xi; \mathbf{0}, \mathbf{I}) \, d\mathbf{f} \, d\xi, \quad (1)$$

which is parameterized by kernel hyperparameters $\boldsymbol{\theta}$ and variational data.

As a variational model, the **vGP** forms an infinite ensemble of mean-field distributions. A mean-field distribution is specified *conditional* on a fixed function $f(\cdot)$ and input ξ ; the d outputs $f_i(\xi) = \lambda_i$ are the mean-field’s parameters. The **vGP** is a form of a hierarchical variational model ([Ranganath et al., 2015](#)); it places a continuous Bayesian nonparametric prior over mean-field parameters.

We emphasize that the **vGP** needs variational data because—unlike typical **Gaussian process (GP)** regression—there is no observed data available to learn a distribution over non-linear mappings. The variational data appear in the conditional distribution of f , anchoring the random non-linear mappings at certain input-output pairs. Thus, when we optimize the **vGP**, the learned variational data enables a complex distribution of variational parameters $f(\xi)$.

We show that the **vGP** is a universal approximator.

Theorem 1 (Universal approximation). *Let $q(\mathbf{z}; \boldsymbol{\theta}, \mathcal{D})$ denote the **variational Gaussian process**. For any posterior distribution $p(\mathbf{z} \mid \mathbf{x})$ with a finite number of latent variables and continuous quantile function (inverse CDF), there exist a set of parameters $(\boldsymbol{\theta}, \mathcal{D})$ such that*

$$\text{KL}(q(\mathbf{z}; \boldsymbol{\theta}, \mathcal{D}) \parallel p(\mathbf{z} \mid \mathbf{x})) = 0.$$

[Theorem 1](#) states that any posterior distribution with strictly positive density can be represented by a **vGP**. Thus the **vGP** is a flexible model for learning posterior distributions.

3 Black box inference

We derive an algorithm for performing black box inference over a wide class of generative models. The original **evidence lower bound (ELBO)** is analytically intractable due to the log density $\log q_{\text{vGP}}(\mathbf{z})$ ([Eq.1](#)). Inference on variational models with latent variables typically requires computation of an auxiliary model $r(\xi, f \mid \mathbf{z})$ ([Salimans et al., 2015](#); [Ranganath et al., 2015](#)). Unlike previous approaches, we present the variational objective in terms of auto-encoders:

$$\tilde{\mathcal{L}} = \mathbb{E}_{q_{\text{vGP}}}[\log p(\mathbf{x} \mid \mathbf{z})] - \mathbb{E}_{q_{\text{vGP}}} \left[\text{KL}(q(\mathbf{z} \mid f(\xi)) \parallel p(\mathbf{z})) + \text{KL}(q(\xi, f) \parallel r(\xi, f \mid \mathbf{z})) \right].$$

In auto-encoder parlance, we maximize the expected negative reconstruction error, regularized by an expected divergence between the variational model and the original model’s prior, and an expected divergence between the auxiliary model and the variational model’s prior. This is simply a nested instantiation of the variational auto-encoder bound ([Kingma and Welling, 2014](#)): a KL divergence between the inference model and a prior is taken as regularizers on both the posterior and variational spaces. This interpretation justifies the previously proposed bound for variational models; as we shall see, it also enables the use of several analytic expectations.

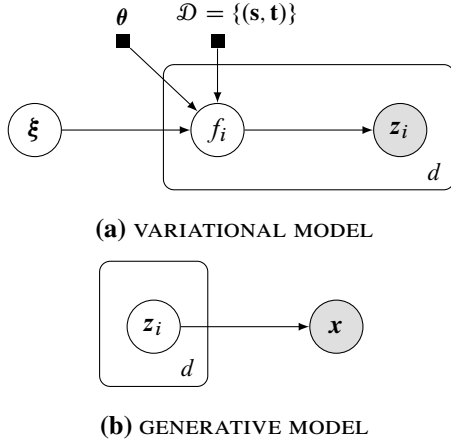


Figure 1: (a) Graphical model of the **variational Gaussian process**. Approximate posterior samples \mathbf{z} for the posterior distribution of a generative model (b).

Algorithm 1: Black box inference with a **vGP**

Input: Model $p(\mathbf{x}, \mathbf{z})$,
Mean-field family $\prod_i q(\mathbf{z}_i | f_i(\xi))$.

Output: Variational and auxiliary parameters (θ, ϕ) .

Initialize (θ, ϕ) randomly.

while not converged **do**

 Draw noise samples $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $\epsilon \sim w$.

 Parameterize variational samples
 $\mathbf{z} = \mathbf{z}(\epsilon; f(\xi))$, $f(\xi) = \mathbf{f}(\xi; \theta)$.

 Update (θ, ϕ) with stochastic gradients $\nabla_{\theta} \tilde{\mathcal{L}}$,
 $\nabla_{\phi} \tilde{\mathcal{L}}$.

end

3.1 Auto-encoding variational models

Inference networks provide a flexible parameterization of approximating distributions as used in Helmholtz machines (Hinton and Zemel, 1994), deep Boltzmann machines (Salakhutdinov and Larochelle, 2010), and variational auto-encoders (Kingma and Welling, 2014; Rezende et al., 2014). To auto-encode the **vGP** we specify inference networks to parameterize both the variational and auxiliary models:

$$\mathbf{x}_n \mapsto q(\mathbf{z}_n | \mathbf{x}_n; \theta_n), \quad \mathbf{x}_n, \mathbf{z}_n \mapsto r(\xi_n, f_n | \mathbf{x}_n, \mathbf{z}_n; \phi_n),$$

where q has local variational parameters given by the variational data \mathcal{D}_n , and r is specified as a fully factorized Gaussian with local variational parameters $\phi_n = (\mu_n \in \mathbb{R}^{c+d}, \sigma_n^2 \in \mathbb{R}^{c+d})$.¹

3.2 Stochastic optimization

We maximize the variational objective $\tilde{\mathcal{L}}(\theta, \phi)$ over both θ and ϕ , where θ newly denotes both the kernel hyperparameters and the inference network’s parameters for the **vGP**, and ϕ denotes the inference network’s parameters for the auxiliary model. Following the standard procedure in black box methods, we write the gradient as an expectation and apply stochastic approximations (Robbins and Monro, 1951), sampling from the variational model and evaluating stochastic gradients.

First, we simplify the expression for the stochastic gradients by analytically deriving any tractable expectations. The KL divergence between $r(\xi, f | \mathbf{z})$ and $q(\xi, f)$ is analytic as we’ve specified both joint distributions to be Gaussian. The KL divergence between $q(\mathbf{z} | f(\xi))$ and $p(\mathbf{z})$ is standard and used to reduce variance in traditional variational auto-encoders: it is analytic for widely used deep generative models such as the deep latent Gaussian model (Rezende et al., 2014) and deep recurrent attentive writer (Gregor et al., 2015). See Appendix B for these calculations.

To derive black box gradients, we can first reparameterize the **vGP**, separating noise generation of samples from the parameters in its generative process (Kingma and Welling, 2014; Rezende et al., 2014). The **GP** easily enables reparameterization: for latent inputs $\xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, the transformation $\mathbf{f}(\xi; \theta) = \mathbf{L}\xi + \mathbf{K}_{\xi s} \mathbf{K}_{ss}^{-1} \mathbf{t}_s$ is a location-scale transform, where $\mathbf{L}\mathbf{L}^\top = \mathbf{K}_{\xi\xi} - \mathbf{K}_{\xi s} \mathbf{K}_{ss}^{-1} \mathbf{K}_{\xi s}^\top$. This is equivalent to evaluating ξ with a random mapping from the **GP**. Suppose the mean-field $q(\mathbf{z} | f(\xi))$ is also reparameterizable, and let $\epsilon \sim w$ such that $\mathbf{z}(\epsilon; \mathbf{f}) \sim q(\mathbf{z} | f(\xi))$. This two-level reparameterization is equivalent to the generative process for \mathbf{z} outlined in Section 2.

We now rewrite the variational objective as

$$\tilde{\mathcal{L}}(\theta, \phi) = \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\log p(\mathbf{x} | \mathbf{z}(\epsilon; \mathbf{f}(\xi; \theta))) \right] \right] \quad (2)$$

¹We let the kernel hyperparameters of the **vGP** be fixed across data points.

Model	$-\log p(\mathbf{x})$	\leq
DLGM + VAE [1]		86.76
DLGM + HVI (8 leapfrog steps) [2]	85.51	88.30
DLGM + NF ($k = 80$) [3]		85.10
EoNADE-5 2hl (128 orderings) [4]	84.68	
DBN 2hl [5]	84.55	
DARN 1hl [6]	84.13	
Convolutional VAE + HVI [2]	81.94	83.49
DLGM 2hl + IWAE ($k = 50$) [1]		82.90
DRAW [7]		80.97
DLGM 1hl + VGP		83.64
DLGM 2hl + VGP		81.90
DRAW + VGP		80.11

Table 1: Negative predictive log-likelihood for binarized MNIST. Previous best results are [1] (Burda et al., 2015), [2] (Salimans et al., 2015), [3] (Rezende and Mohamed, 2015), [4] (Raiko et al., 2014), [5] (Murray and Salakhutdinov, 2009), [6] (Gregor et al., 2014), [7] (Gregor et al., 2015).

$$- \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\text{KL}(q(z | \mathbf{f}(\xi; \theta)) \| p(z)) + \text{KL}(q(\xi, f; \theta) \| r(\xi, f | z(\epsilon; \mathbf{f}(\xi; \theta)); \phi)) \right] \right].$$

Eq.2 enables gradients to move inside the expectations and backpropagate over the nested reparameterization. Thus we can take unbiased stochastic gradients, which exhibit low variance due to both the analytic KL terms and reparameterization.

An outline is given in Algorithm 1. For gradients of the model log-likelihood, we employ convenient differentiation tools such as those in Stan and Theano (Carpenter et al., 2015; Bergstra et al., 2010). For non-differentiable latent variables \mathbf{z} , we apply the score function estimator for gradients of mean-field expectations (Ranganath et al., 2014). Complexity analysis is available in Appendix D.

4 Experiments

The binarized MNIST data set (Salakhutdinov and Murray, 2008) consists of 28x28 pixel images with binary-valued outcomes. Training a deep latent Gaussian model (DLGM), we apply two stochastic layers of 100 random variables and 50 random variables respectively, and inbetween each stochastic layer is a deterministic layer with 100 units using tanh nonlinearities. We apply mean-field Gaussian distributions for the stochastic layers and a Bernoulli likelihood. For DRAW (Gregor et al., 2015), we augment the mean-field Gaussian distribution originally used to generate the latent samples at each time step with the VGP, as it places a complex variational prior over its parameters. We use the same architecture hyperparameters as in Gregor et al. (2015).

After training we evaluate test set log likelihood, which are lower bounds on the true value. See Table 1 which reports both approximations and lower bounds of $\log p(\mathbf{x})$ for various methods. The VGP achieves the highest known results on log-likelihood using DRAW, reporting a value of **-80.11** compared to the original highest of -80.97.² The VGP also achieves the highest known results among the class of non-structure exploiting models using the DLGM, with a value of -81.90 compared to the previous best of -82.90 reported by Burda et al. (2015).

5 Conclusion

We present the variational Gaussian process (VGP), a variational model which adapts its shape to match complex posterior distributions. The VGP draws samples from a tractable distribution, and posits a Bayesian nonparametric prior over transformations from the tractable distribution to mean-field parameters. The VGP adaptively learns the transformations from the space of all continuous mappings—it is a universal approximator and achieves powerful flexibility in practice.

²While superceding DRAW’s previous best results, we have not fully explored initialization in the VGP for learning the model. We hypothesize the VGP may achieve better results given more tuning and computation.

References

- James Bergstra, Olivier Breuleux, Frédéric Bastien, Pascal Lamblin, Razvan Pascanu, Guillaume Desjardins, Joseph Turian, David Warde-Farley, and Yoshua Bengio. Theano: a CPU and GPU math expression compiler. In *Proceedings of the Python for Scientific Computing Conference (SciPy)*, 2010.
- David M Blei and John D Lafferty. Dynamic topic models. In *International Conference on Machine Learning*, 2006.
- Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.
- Bob Carpenter, Matthew D. Hoffman, Marcus Brubaker, Daniel Lee, Peter Li, and Michael Betancourt. The Stan Math Library: Reverse-mode automatic differentiation in C++. *arXiv preprint arXiv:1509.07164*, 2015.
- John P Cunningham, Krishna V Shenoy, and Maneesh Sahani. Fast Gaussian process methods for point process intensity estimation. In *International Conference on Machine Learning*. ACM, 2008.
- Karol Gregor, Ivo Danihelka, Andriy Mnih, Charles Blundell, and Daan Wierstra. Deep autoregressive networks. In *International Conference on Machine Learning*, 2014.
- Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. DRAW: A recurrent neural network for image generation. In *International Conference on Machine Learning*, 2015.
- G. Hinton and D. Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Computational Learning Theory*, pages 5–13. ACM, 1993.
- Geoffrey E Hinton and Richard S Zemel. Autoencoders, minimum description length, and helmholtz free energy. *Advances in neural information processing systems*, pages 3–3, 1994.
- Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14:1303–1347, 2013.
- Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- Diederik P Kingma and Max Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.
- Neil Lawrence. *Variational Inference in Probabilistic Models*. PhD thesis, 2000.
- Neil Lawrence. Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *The Journal of Machine Learning Research*, 6:1783–1816, 2005.
- Iain Murray and Ruslan R Salakhutdinov. Evaluating probabilities under high-dimensional latent variable models. In *Advances in neural information processing systems*, pages 1137–1144, 2009.
- Radford M Neal. Learning stochastic feedforward networks. *Department of Computer Science, University of Toronto*, 1990.
- Michael Osborne. *Bayesian Gaussian processes for sequential prediction, optimisation and quadrature*. PhD thesis, Oxford University New College, 2010.
- Joaquin Quiñero-Candela and Carl Edward Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research*, 6:1939–1959, 2005.
- Tapani Raiko, Yao Li, Kyunghyun Cho, and Yoshua Bengio. Iterative neural autoregressive distribution estimator nade-k. In *Advances in Neural Information Processing Systems*, pages 325–333, 2014.
- Rajesh Ranganath, Sean Gerrish, and David M Blei. Black box variational inference. In *Artificial Intelligence and Statistics*, 2014.

- Rajesh Ranganath, Dustin Tran, and David M. Blei. Hierarchical variational models. *arXiv preprint arXiv:1511.02386*, 2015.
- Carl Edward Rasmussen and Christopher K I Williams. *Gaussian processes for machine learning*. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2006.
- Danilo Jimenez Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, 2015.
- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 2014.
- Herbert Robbins and S Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, 1951.
- Ruslan Salakhutdinov and Hugo Larochelle. Efficient learning of deep Boltzmann machines. In *International Conference on Artificial Intelligence and Statistics*, pages 693–700, 2010.
- Ruslan Salakhutdinov and Iain Murray. On the quantitative analysis of deep belief networks. In *International Conference on Machine Learning*, 2008.
- Tim Salimans, Diederik P Kingma, and Max Welling. Markov chain Monte carlo and variational inference: Bridging the gap. In *International Conference on Machine Learning*, 2015.
- Aad Van Der Vaart and Harry Van Zanten. Information rates of nonparametric gaussian process methods. *The Journal of Machine Learning Research*, 12:2095–2119, 2011.
- S. Waterhouse, D. MacKay, and T. Robinson. Bayesian methods for mixtures of experts. *Neural Information Processing Systems*, pages 351–357, 1996.

A Proof of Theorem 1

Theorem 1. Let $q(\mathbf{z}; \boldsymbol{\theta}, \mathcal{D})$ denote the *variational Gaussian process*. For any posterior distribution $p(\mathbf{z} | \mathbf{x})$ with a finite number of latent variables and continuous quantile function (inverse CDF), there exist a set of parameters $(\boldsymbol{\theta}, \mathcal{D})$ such that

$$\text{KL}(q(\mathbf{z}; \boldsymbol{\theta}, \mathcal{D}) \| p(\mathbf{z} | \mathbf{x})) = 0.$$

Proof. Let the mean-field distribution be given by degenerate delta distributions

$$q(\mathbf{z}_i | f_i) = \delta_{f_i}(\mathbf{z}_i).$$

Let the size of the latent input be equivalent to the number of latent variables $c = d$ and fix $\sigma_{\text{ARD}}^2 = 1$ and $\omega_j = 1$. Furthermore for simplicity, we assume that $\boldsymbol{\xi}$ is drawn uniformly on the d -dimensional hypercube. Then if we let P^{-1} denote the inverse posterior cumulative distribution function, the optimal f denoted f^* such that

$$\text{KL}(q(\mathbf{z}; \boldsymbol{\theta}) \| p(\mathbf{z} | \mathbf{x})) = 0$$

is

$$f^*(\boldsymbol{\xi}) = P^{-1}(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_d).$$

Define O_m to be the set of points $k/2^m$ for $k = 0$ to 2^m , and define S_m to be the d -dimensional product of O_m . Let \mathcal{D}_m be the set containing the pairs $(s_i, f^*(s_i))$, for each element s_i in S_m . Denote f^m as the GP mapping conditioned on the dataset \mathcal{D}_m , this random mapping satisfies $f^m(s_i) = f^*(s_i)$ for all $s_i \in S_m$ by the noise free prediction property of Gaussian processes (Rasmussen and Williams, 2006). Then by continuity, as $m \rightarrow \infty$, f^m converges to f^* . \square

A broad condition under which the quantile function of a distribution is continuous is if that distribution has positive density with respect to the Lebesgue measure.

The rate of convergence for finite sizes of the variational data can be studied via posterior contraction rates for GPs under random covariates (Van Der Vaart and Van Zanten, 2011). Only an additional assumption using stronger continuity conditions for the posterior quantile and the use of Matern covariance functions is required for the theory to be applicable in the variational setting.

B Variational objective

We derive the tractable lower bound to the model evidence $\log p(\mathbf{x})$ presented. We first penalize the **ELBO** with an expected KL term:

$$\begin{aligned}\log p(\mathbf{x}) &\geq \mathcal{L} = \mathbb{E}_{q_{\text{VGP}}}[\log p(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\text{VGP}}(\mathbf{z}) \| p(\mathbf{z})) \\ &\geq \mathbb{E}_{q_{\text{VGP}}}[\log p(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\text{VGP}}(\mathbf{z}) \| p(\mathbf{z})) - \mathbb{E}_{q_{\text{VGP}}}\left[\text{KL}(q(\boldsymbol{\xi}, f | \mathbf{z}) \| r(\boldsymbol{\xi}, f | \mathbf{z}))\right].\end{aligned}$$

We can combine all terms into the expectations as follows:

$$\begin{aligned}\tilde{\mathcal{L}} &= \mathbb{E}_{q(\mathbf{z}, \boldsymbol{\xi}, f)}\left[\log p(\mathbf{x} | \mathbf{z}) - \log q(\mathbf{z}) + \log p(\mathbf{z}) - \log q(\boldsymbol{\xi}, f | \mathbf{z}) + \log r(\boldsymbol{\xi}, f | \mathbf{z})\right] \\ &= \mathbb{E}_{q(\mathbf{z}, \boldsymbol{\xi}, f)}\left[\log p(\mathbf{x} | \mathbf{z}) - \log q(\mathbf{z} | f(\boldsymbol{\xi})) + \log p(\mathbf{z}) - \log q(\boldsymbol{\xi}, f) + \log r(\boldsymbol{\xi}, f | \mathbf{z})\right],\end{aligned}$$

where we apply the product rule $q(\mathbf{z})q(\boldsymbol{\xi}, f | \mathbf{z}) = q(\mathbf{z} | f(\boldsymbol{\xi}))q(\boldsymbol{\xi}, f)$. Recombining terms as KL divergences, and written with parameters $(\boldsymbol{\theta}, \boldsymbol{\phi})$, this recovers the auto-encoded variational objective in [Section 3](#):

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\theta})}[\log p(\mathbf{x} | \mathbf{z})] - \mathbb{E}_{q(\mathbf{z}, \boldsymbol{\xi}, f; \boldsymbol{\theta})}\left[\text{KL}(q(\mathbf{z} | f(\boldsymbol{\xi})) \| p(\mathbf{z})) + \text{KL}(q(\boldsymbol{\xi}, f; \boldsymbol{\theta}) \| r(\boldsymbol{\xi}, f | \mathbf{z}; \boldsymbol{\phi}))\right].$$

We now analyze the KL terms. Recall that we specify the auxiliary model $r(\boldsymbol{\xi}, f | \mathbf{z}) = \mathcal{N}((\boldsymbol{\xi}, f(\boldsymbol{\xi}))^\top | \mathbf{z}; \mathbf{m}, \mathbf{S})$, where $\mathbf{m} \in \mathbb{R}^{c+d}$, $\mathbf{S} \in \mathbb{R}^{c+d}$. The variational prior for the **VGP** is also jointly Gaussian:

$$q(\boldsymbol{\xi}, f) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\xi} \\ f(\boldsymbol{\xi}) \end{bmatrix}; \begin{bmatrix} \mathbf{0} \\ \tilde{\mathbf{m}} \end{bmatrix}, \begin{bmatrix} \mathbf{I} & \mathbf{C} \\ \mathbf{C}^\top & \tilde{\mathbf{K}} \end{bmatrix}\right),$$

where the vector $\tilde{\mathbf{m}}$ has elements $\tilde{\mathbf{m}}_i = \mathbf{K}_{\xi_s} \mathbf{K}_{s_s}^{-1} \mathbf{t}_i$, and the matrix $\tilde{\mathbf{K}}$ is diagonal with elements $\tilde{\mathbf{K}}_{ii} = \mathbf{K}_{\xi\xi} - \mathbf{K}_{\xi_s} \mathbf{K}_{s_s}^{-1} \mathbf{K}_{s_s}^\top$. Applying the law of total expectations, the cross-correlation is

$$\mathbf{C} = \mathbb{E}[\boldsymbol{\xi} f(\boldsymbol{\xi})^\top] - \mathbb{E}[\boldsymbol{\xi}] \mathbb{E}[f(\boldsymbol{\xi})]^\top = \mathbb{E}[\boldsymbol{\xi} \mathbb{E}_{f|\boldsymbol{\xi}}[f(\boldsymbol{\xi}) | \boldsymbol{\xi}]^\top] = \mathbb{E}\left[\boldsymbol{\xi} \left[\prod_{i=1}^d \mathbf{K}_{\xi_s} \mathbf{K}_{s_s}^{-1} \mathbf{t}_i\right]^\top\right].$$

This expectation is analytically intractable as $\boldsymbol{\xi}$ appears non-linearly in the kernel. In the special case when the kernel is linear, the expression is analytic and corresponds to probabilistic principal components analysis ([Lawrence, 2005](#)). In general we apply a Monte Carlo estimate using standard normal samples of $\boldsymbol{\xi}$; for black box gradients we obtain these samples for free as they already follow the process of sampling \mathbf{z} (producing intermediate samples $\boldsymbol{\xi}$).

Let $\mathbf{m}_q, \boldsymbol{\Sigma}_q$ denote the concatenated mean vector and covariance matrix for $q(\boldsymbol{\xi}, f)$. Because $q(\boldsymbol{\xi}, f)$ and $r(\boldsymbol{\xi}, f | \mathbf{z})$ are both Gaussian, the KL has an analytic form:

$$\begin{aligned}\text{KL}(q(\boldsymbol{\xi}, f; \mathbf{m}_q, \boldsymbol{\Sigma}_q) \| r(\boldsymbol{\xi}, f | \mathbf{z}; \mathbf{m}, \mathbf{S})) &= \\ &= \frac{1}{2} \left((\mathbf{m} - \mathbf{m}_q)^\top \mathbf{S}^{-1} (\mathbf{m} - \mathbf{m}_q) + \text{tr}(\mathbf{S}^{-1} \boldsymbol{\Sigma}_q + \log \mathbf{S} - \log \boldsymbol{\Sigma}_q) - (c + d) \right).\end{aligned}$$

Since \mathbf{S} is a diagonal matrix, inversion is trivial—the whole expression is simple to compute and backpropagate gradients.

We now consider the first KL term. The KL divergence between the mean-field $q(\mathbf{z} | f(\boldsymbol{\xi}))$ and the model prior $p(\mathbf{z})$ is analytically tractable for certain popular models. For example, in the deep latent Gaussian model ([Rezende et al., 2014](#)) and DRAW ([Gregor et al., 2015](#)), both the mean-field distribution and model prior are Gaussian, leading to an analytic KL term similar to the above. In general, when the KL is intractable, we combine the term with the reconstruction term, and maximize the variational objective

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \mathbb{E}_{q(\mathbf{z}, \boldsymbol{\xi}, f; \boldsymbol{\theta})}[\log p(\mathbf{x}, \mathbf{z}) - q(\mathbf{z} | f(\boldsymbol{\xi}))] - \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\theta})}[\text{KL}(q(\boldsymbol{\xi}, f; \boldsymbol{\theta}) \| r(\boldsymbol{\xi}, f | \mathbf{z}; \boldsymbol{\phi}))], \quad (3)$$

which leads to slightly higher variance in the stochastic gradients during optimization.

C Gradients of the variational objective

We derive gradients for the variational objective (Eq.2). This follows trivially by backpropagation:

$$\begin{aligned}\nabla_{\theta} \tilde{\mathcal{L}}(\theta, \phi) &= \mathbb{E}_{\mathcal{N}(\xi)}[\mathbb{E}_{w(\epsilon)}[\nabla_{\theta} \mathbf{f}(\xi) \nabla_{\mathbf{f}} z(\epsilon) \nabla_z \log p(\mathbf{x} | z)]] \\ &\quad - \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\nabla_{\theta} \text{KL}(q(z | \mathbf{f}(\xi; \theta)) \| p(z)) \right] \right] \\ &\quad - \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\nabla_{\theta} \text{KL}(q(\xi, f; \theta) \| r(\xi, f | z(\epsilon; \mathbf{f}(\xi; \theta)); \phi)) \right] \right] \\ \nabla_{\phi} \tilde{\mathcal{L}}(\theta, \phi) &= -\mathbb{E}_{\mathcal{N}(\xi)}[\mathbb{E}_{w(\epsilon)}[\nabla_{\phi} \text{KL}(q(\xi, f; \theta) \| r(\xi, f | z(\epsilon; \mathbf{f}(\xi; \theta)); \phi))]],\end{aligned}$$

where we assume the KL terms are analytically written from Appendix B and gradients are propagated similarly through their computational graph.

We also derive gradients for the general variational bound of Eq.3—it assumes that the first KL term, measuring the divergence between q and the prior for p , is not necessarily tractable. Following the reparameterizations described in Section 3.2, this variational objective can be rewritten as

$$\begin{aligned}\tilde{\mathcal{L}}(\theta, \phi) &= \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\log p(\mathbf{x}, z(\epsilon; \mathbf{f}(\xi; \theta))) - \log q(z(\epsilon; \mathbf{f}(\xi; \theta)) | \mathbf{f}(\xi; \theta)) \right] \right] \\ &\quad - \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\text{KL}(q(\xi, f; \theta) \| r(\xi, f | z(\epsilon; \mathbf{f}(\xi; \theta)); \phi)) \right] \right].\end{aligned}$$

We calculate gradients by backpropagating over the nested reparameterizations:

$$\begin{aligned}\nabla_{\theta} \tilde{\mathcal{L}}(\theta, \phi) &= \mathbb{E}_{\mathcal{N}(\xi)}[\mathbb{E}_{w(\epsilon)}[\nabla_{\theta} \mathbf{f}(\xi) \nabla_{\mathbf{f}} z(\epsilon) [\nabla_z \log p(\mathbf{x}, z) - \nabla_z \log q(z | \mathbf{f})]]] \\ &\quad - \mathbb{E}_{\mathcal{N}(\xi)} \left[\mathbb{E}_{w(\epsilon)} \left[\nabla_{\theta} \text{KL}(q(\xi, f; \theta) \| r(\xi, f | z(\epsilon; \phi))) \right] \right] \\ \nabla_{\phi} \tilde{\mathcal{L}}(\theta, \phi) &= -\mathbb{E}_{\mathcal{N}(\xi)}[\mathbb{E}_{w(\epsilon)}[\nabla_{\phi} \text{KL}(q(\xi, f; \theta) \| r(\xi, f | z(\epsilon; \mathbf{f}(\xi; \theta)); \phi))]].\end{aligned}$$

D Computational and storage complexity

The algorithm has $\mathcal{O}(d + m^3 + LN^2)$ complexity, where d is the number of latent variables, m is the size of the variational data, and L is the number of layers of the neural networks with N the average hidden layer size. In particular, the algorithm is linear in the number of latent variables, which is competitive with other variational methods. The number of variational and auxiliary parameters has $\mathcal{O}(c + LN)$ complexity—storing the kernel hyperparameters and the neural network parameters. Note that unlike most GP literature, we require no low rank constraints such as the use of inducing variables (Quiñero-Candela and Rasmussen, 2005) for scalable computation.

If massive sizes of variational data are required, e.g., when its cubic complexity due to inversion of a $m \times m$ matrix becomes the bottleneck during computation, we can scale it further. Consider fixing the variational inputs to lie on a grid. For stationary kernels, this allows us to exploit Toeplitz structure for fast $m \times m$ matrix inversion. In particular, one can embed the Toeplitz matrix into a circulant matrix and apply conjugate gradient combined with fast Fourier transforms in order to compute inverse-matrix vector products in $\mathcal{O}(m \log m)$ computation and $\mathcal{O}(m)$ storage (Cunningham et al., 2008). For product kernels, we can further exploit Kronecker structure to allow fast $m \times m$ matrix inversion in $\mathcal{O}(Pm^{1+1/P})$ operations and $\mathcal{O}(Pm^{2/P})$ storage, where $P > 1$ is the number of kernel products (Osborne, 2010). The automatic relevance determination (ARD) kernel specifically leads to $\mathcal{O}(cm^{1+1/c})$ complexity, which is linear in m .