Perturbative Black Box Variational Inference

Cheng Zhang^{*1}, Robert Bamler^{*1}, Manfred Opper², Stephan Mandt^{*1} 1. Disney Research, 2. Technical University of Berlin * Equal contribution

Introduction

- We establish a unified view on black box variational inference with generalized divergences as a form of *biased importance sampling*.
- We use these insights to construct a new variational bound with favorable properties with respect to a variance-bias trade-off.
- In our experiments, the resulting posterior covariances are closer to the true posterior, and likelihoods on heldout data are higher than with traditional black box variational inference.

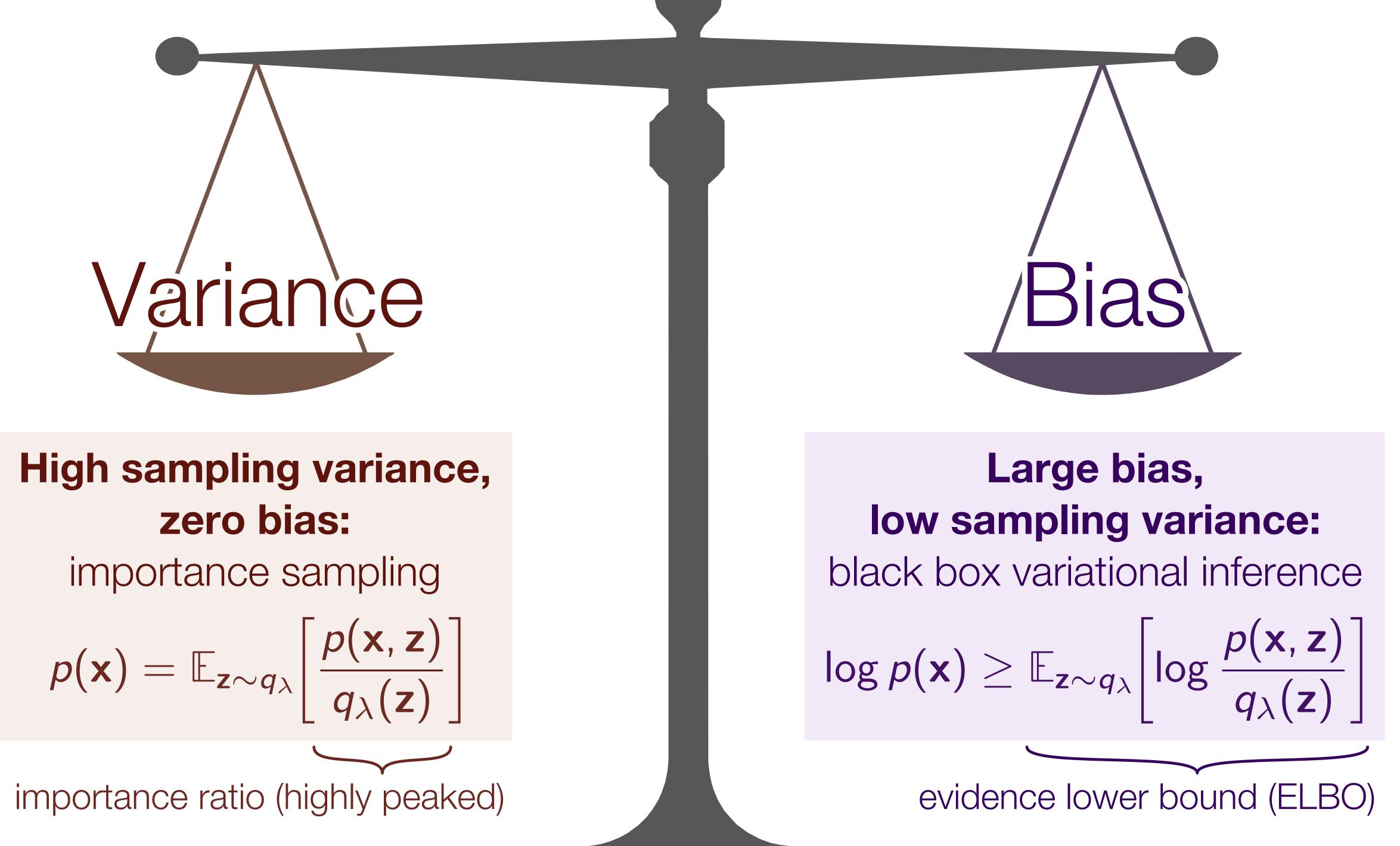
Variational Inference as Biased Importance Sampling

Model with joint p.d.f. p(x, z)observed laten \rightarrow seek posterior $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x},\mathbf{z})}{p(\mathbf{x})}$

Problem: intractable denominator $p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$

Black box variational inference (BBVI) estimates a lower bound $\mathcal{L}(\lambda)$ on log $p(\mathbf{x})$ based on Monte Carlo samples from a variational distribution $q_{\lambda}(\mathbf{z})$. Taking $\mathcal{L}(\lambda)$ as a proxy for log $p(\mathbf{x})$ results in a bias and a sampling variance.

Variance-bias trade-off:



A Unified View

For any concave function f with $f(\xi) \leq \xi$.

- \rightarrow importance sampling: f = id
- \rightarrow (traditional) Kullback-Leibler BBVI: $f = \log + \text{const.}$
- \rightarrow BBVI with alpha-divergence: $f^{(\alpha)}(\xi) \propto \xi^{1-\alpha}$

Observation: $\xi \equiv \frac{p(\mathbf{x}, \mathbf{z})}{q_{\lambda}(\mathbf{z})}$ is highly peaked in **z**-space.

- \Rightarrow If $f(\xi)$ depends algebraically on ξ , as in the alpha bound, then the sampling variance is high and reparemeterization gradients are noisy.
- \Rightarrow If $f(\xi)$ depends only on $\log \xi$, as in the KL bound, then the sampling variance is lower and reparemeterization gradients are less noisy. However, the logarithm introduces a bias, i.e., the KL-bound is less tight.

Perturbative Black Box Variational Inference (PVI)

Aim: Lower bound with small bias and small gradient noise.

Taylor expansion of $\frac{p(x,z)}{q_{\lambda}(z)} = \exp[\log p(x,z) - \log q_{\lambda}(z)]$ around a reference value of e^{-V_0} :

$$\mathcal{L}_{f_{V_0}^{(n)}}(\lambda) \;\equiv\; e^{-V_0}\sum_{k=0}^n rac{1}{k!} \mathbb{E}_{\mathsf{z}\sim q_\lambda} \Big[ig(V_0 + \log p(\mathsf{x}, \mathsf{z}) - \log q_\lambda(\mathsf{z})ig)^k \Big]$$

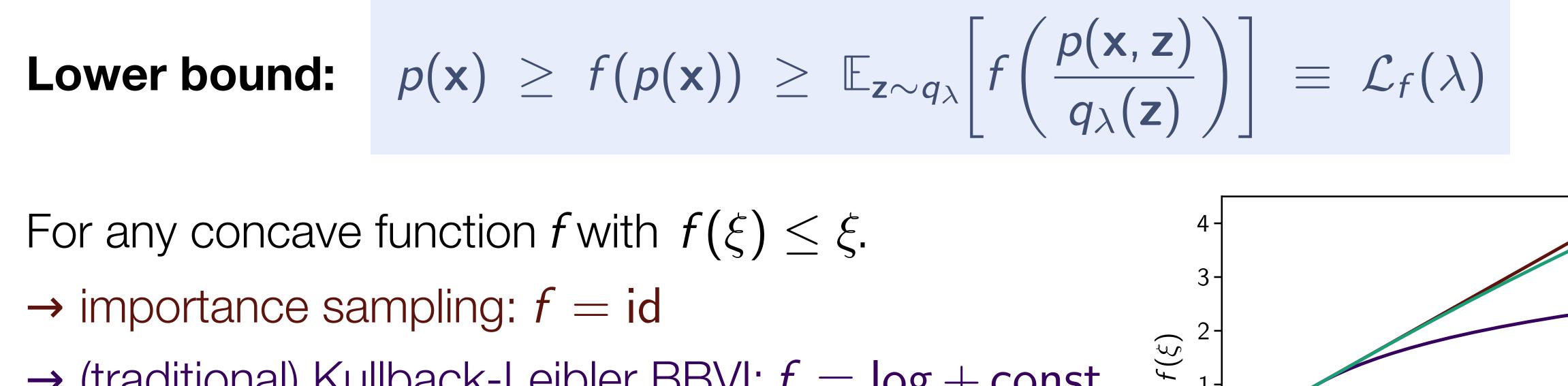
We show that $p(\mathbf{x}) \geq \mathcal{L}_{f^{(n)}}(\lambda)$ for all odd n.

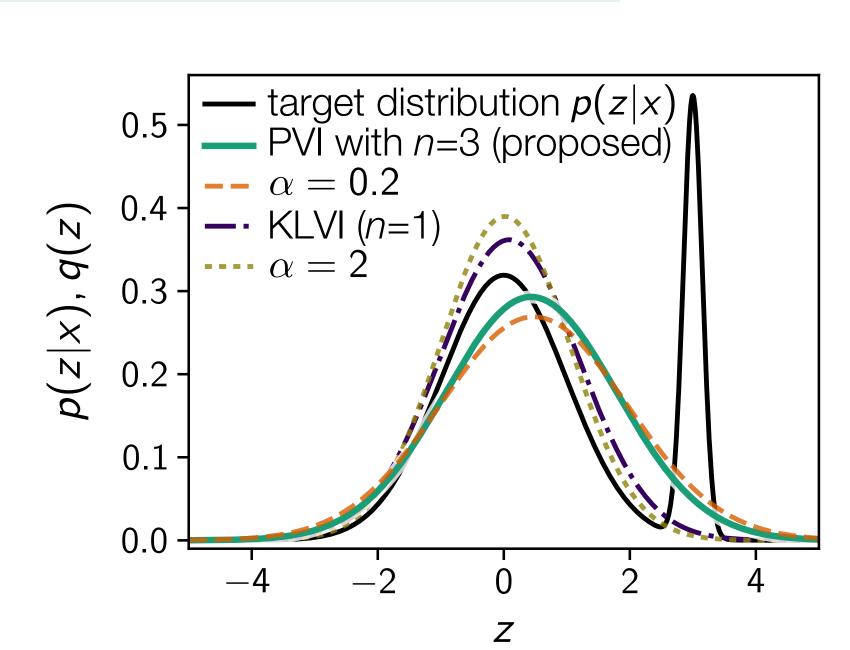
 $\rightarrow n \rightarrow \infty$: importance sampling

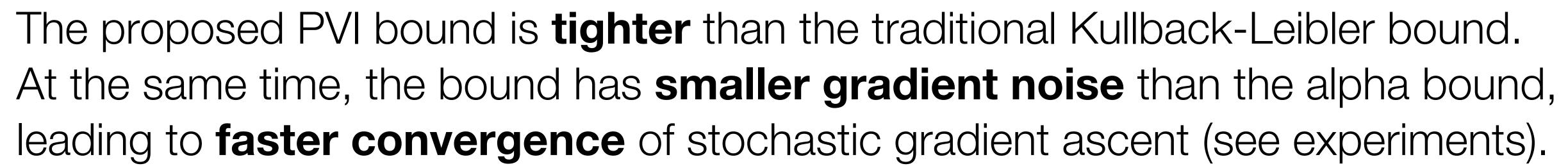
 $\rightarrow n = 1$: (traditional) Kullback-Leibler VI

 $\rightarrow n = 3$: proposed

The proposed PVI bound is **tighter** than the traditional Kullback-Leibler bound.







Experiments KLVI

With increasing dimensionality of the latent space, the gradient noise grows exponentially for the alpha-VI bound, and only algebraically for the proposed PVI bound (green line).

Gaussian process classification with the UCI data sets

We consistently ob **lower error rates** on held out test set

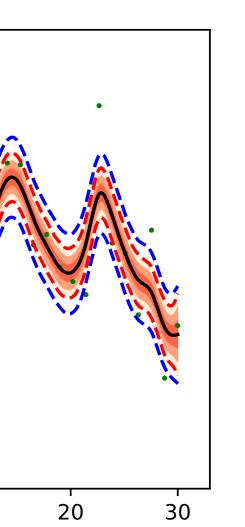
gradient noise, see above.

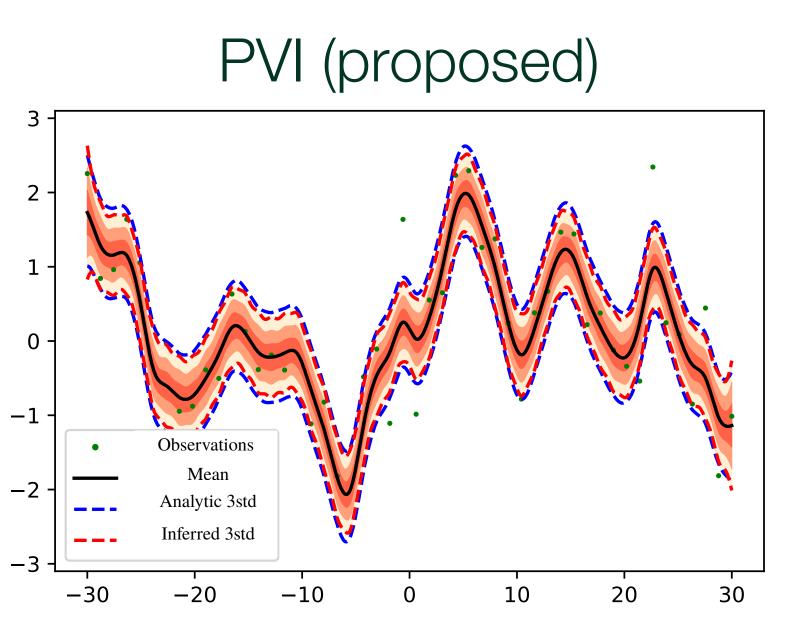
Variational autoencoder (VAE)

DISNEØ Research

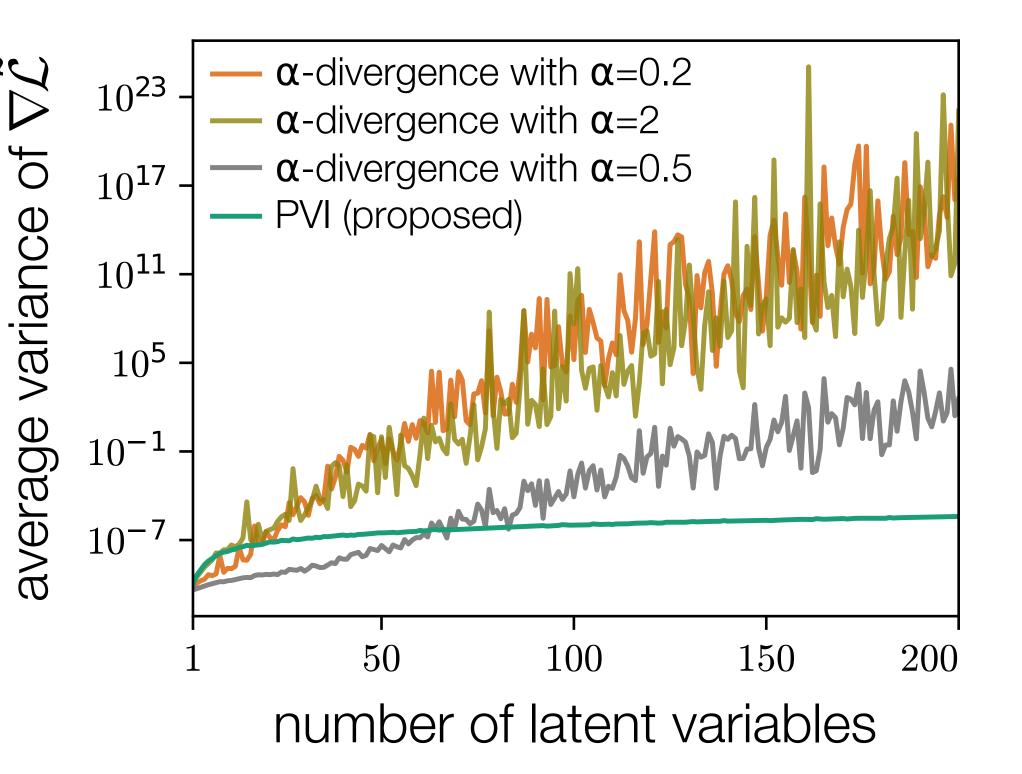


Gaussian process regression with synthetic data





	average posterior variance
Analytic	0.0415
KLVI	0.0176
PVI (proposed)	0.0355



btain	Data set	Crab	Pima	Heart	Sonar
S	KLVI	0.22	0.245	0.148	0.212
ets.	PVI (proposed)	0.11	0.240	0.133	0.173

