



AdaGeo: Adaptive Geometric Learning for Optimization and Sampling



Gabriele Abbati¹, Alessandra Tosi², Michael A Osborne¹, Seth Flaxman³

¹University of Oxford, ²Mind Foundry Ltd, ³Imperial College London

Abstract

In high-dimensional optimization and sampling, well-known issues such as slow-mixing, non-convexity and correlations can arise.

We propose **AdaGeo**, a preconditioning framework for **adaptively** learning the **geometry** of parameter space.

We use the **Gaussian Process latent variable model** (GP-LVM) to represent a lower-dimensional embedding of the parameters, identifying the underlying **Riemannian manifold** on which the optimization or sampling is taking place. Samples or optimization steps are then proposed based on the geometry of the manifold.

We apply our framework to **stochastic gradient descent** (SGD) and **stochastic gradient Langevin dynamics** (SGLD) and show performance improvements for both optimization and sampling.

GP-Latent Variable Models

Latent variable models relate a set of observed variables $\Theta \subset \mathbb{R}^D$ to a set of lower-dimensional unobserved (or **latent**) variables $\Omega \subset \mathbb{R}^Q$ through the mapping f :

$$\theta = f(\omega) + \eta \quad (6)$$

where $\theta \in \Theta$, $\omega \in \Omega$ and η is a noise term. Gaussian Process Latent Variable models assume a **Gaussian Process prior** for f .

For differentiable kernels the **Jacobian** J of f is normally distributed (if its columns are assumed to be independent):

$$p(\mathbf{J} | \Theta, \Omega, \beta) = \prod_{i=1}^D \mathcal{N}(\mathbf{J}_{i,:} | \mu_{\mathbf{J}_{i,:}}, \Sigma_{\mathbf{J}}), \quad (7)$$

where $\mu_{\mathbf{J}_{i,:}}$ and $\Sigma_{\mathbf{J}}$ are respectively the posterior mean and covariance of the Jacobian given θ and the mapping f (derivation in [1]).

Gradients of a function $g(\theta) : \mathbb{R}^D \rightarrow \mathbb{R}$ can then be transformed as:

$$\nabla_{\omega} g(\mathbf{f}(\omega)) = \mu_{\mathbf{J}} \nabla_{\theta} g(\theta). \quad (8)$$

AdaGeo Optimization

The **minimization** of an objective function $g(\theta)$:

$$\theta^* = \arg \min_{\theta} g(\theta) \quad (1)$$

can be solved with **gradient-based** schemes of the form $\theta_{t+1} = \theta_t - \Delta \theta_t$ where $\Delta \theta_t = \Delta \theta_t(\nabla g(\theta))$. After t steps the GP-LVM is trained on the set $\Theta = \{\theta_1, \dots, \theta_t\}$ and the updates are computed into the **latent space**.

E.g., a vanilla stochastic gradient descent:

$$\theta_{t+1} = \theta_t - \frac{\epsilon_t}{N_b} \sum_{i=1}^{N_b} \nabla_{\theta} g(\theta_{ti}) \quad (2)$$

becomes, using eq. 8 for the gradients:

$$\begin{aligned} \omega_{t+1} &= \omega_t - \frac{\epsilon_t}{N_b} \sum_{i=1}^{N_b} \nabla_{\omega} g(\omega_{ti}) \\ \theta_{t+1} &= f(\omega_{t+1}). \end{aligned} \quad (3)$$

Latent and classic updates have to be **alternated** so that to acquire and exploit geometrical information.

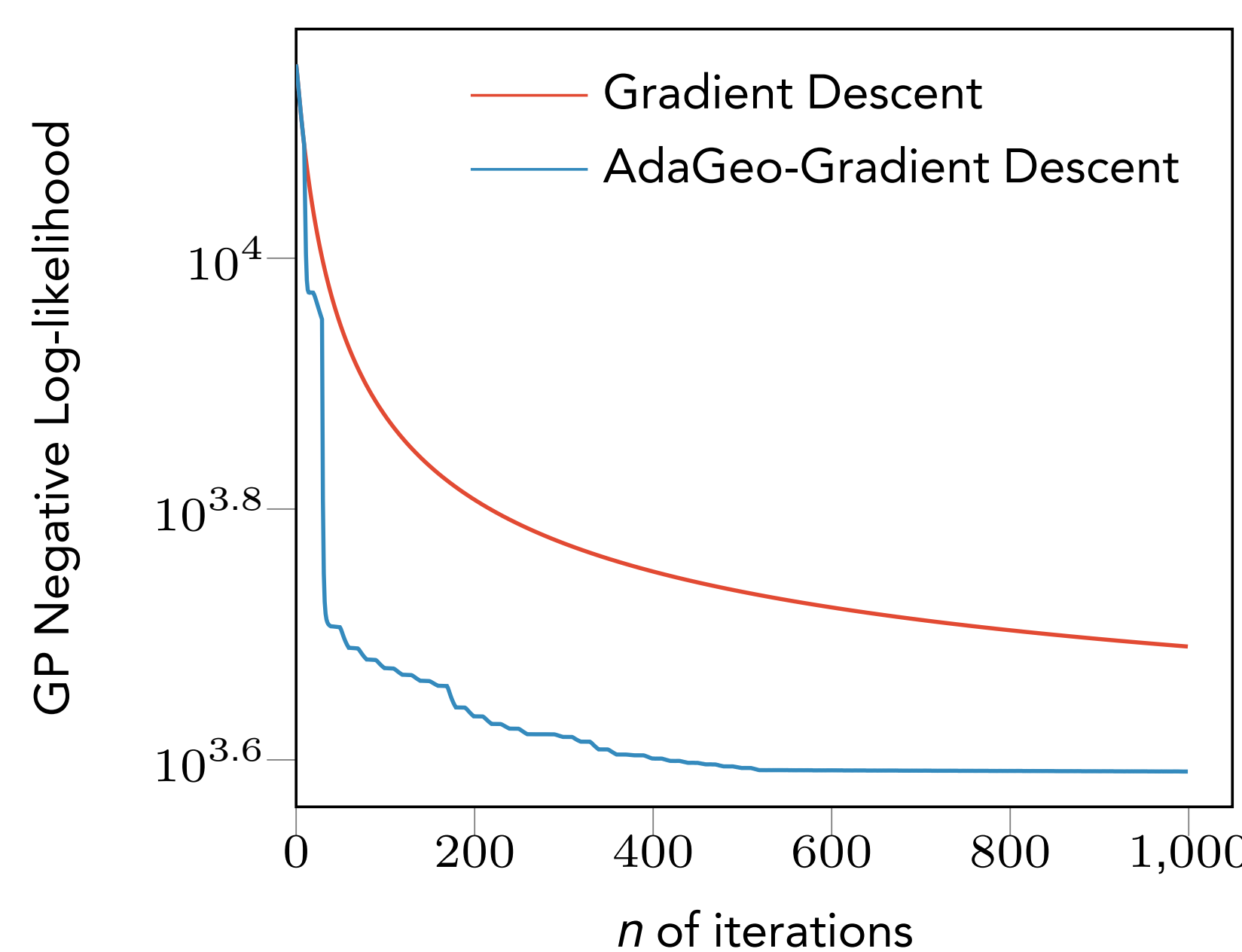
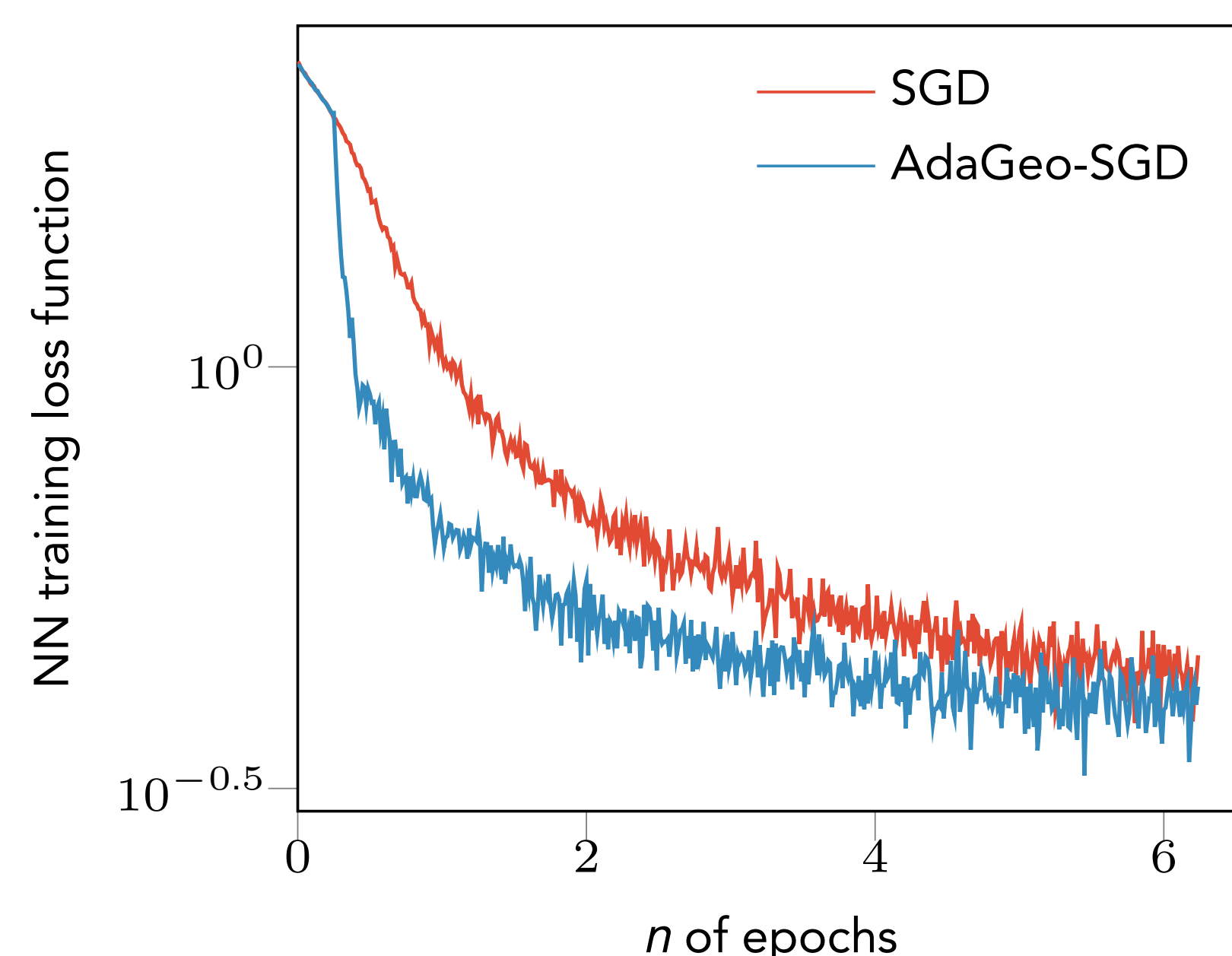


Fig. 1: Above: neural network loss function during training. Below: Negative log-likelihood during training of a Gaussian process.

AdaGeo-SGLD Sampling

Stochastic gradient Langevin dynamics (SGLD)[2] implements Bayesian **posterior sampling** of a distribution $p(x, \theta)$ by building an iterative method of the form:

$$\Delta \theta_t = \frac{\epsilon_t}{2} \left(\nabla_{\theta} \log p(\theta_t) + \frac{N}{n} \sum_{i=1}^n \nabla_{\theta} \log p(\mathbf{x}_i | \theta_t) \right) + \eta_t, \quad (4)$$

where $\eta_t \sim \mathcal{N}(\mathbf{0}, \epsilon_t \mathbf{I})$. Analogously as the optimization case, the update can be brought onto the **lower-dimensional latent space**:

$$\Delta \omega_t = \frac{\epsilon_t}{2} \left(\nabla_{\omega} \log p(\mathbf{f}(\omega_t)) + \frac{N}{n} \sum_{i=1}^n \nabla_{\omega} \log p(\mathbf{x}_i | \mathbf{f}(\omega_t)) \right) + \eta_t, \quad (5)$$

where $\eta_t \sim \mathcal{N}(\mathbf{0}, \epsilon_t \mathbf{I})$.

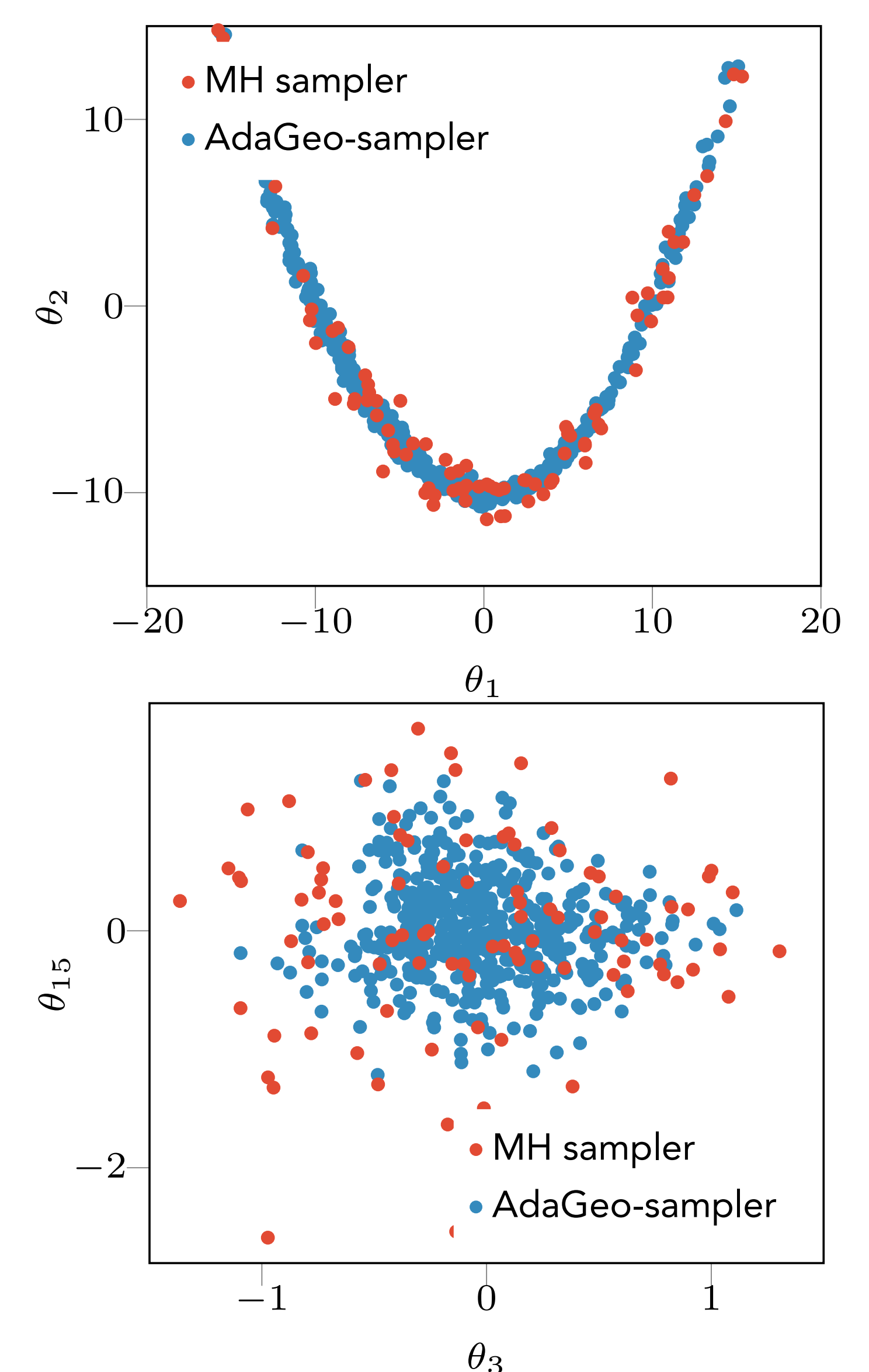


Fig. 2: Sampling with AdaGeo-SGLD from a 50-dimensional banana distribution:

$$p(\theta) \propto \exp \left(-\frac{\theta_1^2}{200} - \frac{(\theta_2 - b\theta_1^2 + 100b)^2}{2} - \sum_{j=3}^D \theta_j^2 \right).$$

References

- [1] Alessandra Tosi et al., *Metrics for probabilistic geometries*, UAI, page 800, 2014
- [2] Max Welling et al., *Bayesian learning via stochastic gradient Langevin dynamics*, ICML, pages 681–688, 2011

Results

Sampling We employed AdaGeo-SGLD to sample from a 50-dimensional **banana** distribution: **Metropolis-Hastings** yields the first 100 samples, after 10^4 burn-in iterations and with a thinning factor of 250. **AdaGeo-SGLD**, using a 5-dimensional latent space, is then employed for the next 250 samples, with 10^3 burn-in steps and a thinning factor of 100. As seen in figure 2 the main features of the distribution are preserved. A relevant speed up and better autocorrelation plots are obtained as well.

Optimization First, the training of a one-layer **neural network** implementing **logistic regression** on the **MNIST** dataset was speeded up: the 7850 weights were modeled through a 9-dimensional latent space and **AdaGeo** was used to accelerate an **SGD** scheme, where 20 classic updates were alternated with 30 latent ones until convergence. Second, a **Gaussian Process** with a total of 9 hyper-parameters is trained using **gradient descent** with Nesterov momentum: here we alternate 15 classic updates with 15 latent ones, using a 3-dimensional latent space.



Contact

Gabriele Abbati
University of Oxford
Email: gabb@robots.ox.ac.uk
Website: <http://www.robots.ox.ac.uk/~gabb/>