



AdaGeo: Adaptive Geometric Learning for Optimization and Sampling

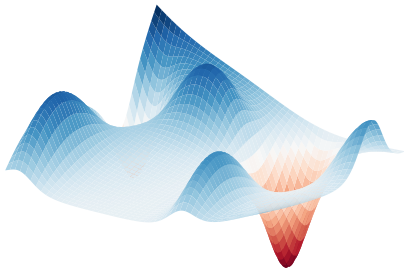
Gabriele Abbati¹, Alessandra Tosi², Seth Flaxman³, Michael A Osborne¹

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

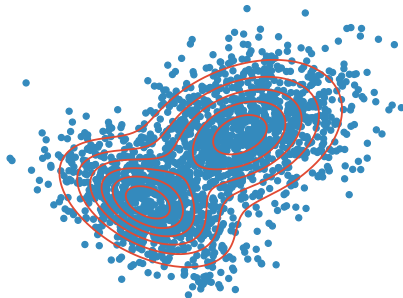
Afternoon Meeting on Bayesian Computation 2018
University of Reading

High-dimensional Problems

- **Gradient-based optimization**



- **MCMC Sampling**



Issues arising from high dimensionality:

- non-convexity
- strong correlations
- multimodality



Gradient-based optimization

- AdaGrad
- AdaDelta
- Adam
- RMSProp

MCMC Sampling

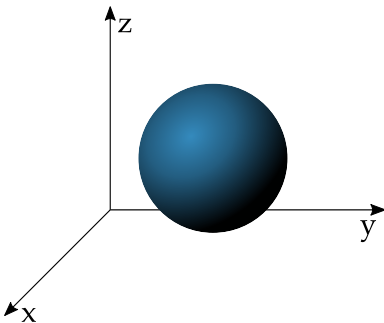
- Hamiltonian Monte Carlo
- Particle Monte Carlo
- Stochastic gradient Langevin dynamics

All of these methods focus on computing clever updates for optimization algorithms or for Markov chains.

Novelty: to the best of our knowledge, no dimensionality reduction approaches were applied in this direction before.

After t steps of optimization or sampling, we assume the obtained points in the parameter space to be lying on a **manifold**.

We then feed them to a dimensionality reduction method to find a **lower-dimensional representation**.



3D example: if the sampler/optimizer algorithm keeps on returning proposals on a sphere surface, that information might be used to our advantage

Can we perform better if the algorithm acts with **knowledge** of the manifold?



Latent Variable Models

Latent Variable Models describe a set Θ through a lower-dimensional latent set Ω

Latent Variable Models

$$\Theta = \{\theta_1, \dots, \theta_N \in \mathbb{R}^D\} \xleftarrow[\text{with } Q < D]{\text{map } f} \Omega = \{\omega_1, \dots, \omega_N \in \mathbb{R}^Q\}$$

where:

- θ : observed variables/parameters
- ω : latent variables
- f : mapping
- D, Q : dimensionalities of Θ and Ω respectively



Latent Variable Models

Latent Variable Models describe a set Θ through a lower-dimensional latent set Ω

Latent Variable Models

$$\Theta = \{\theta_1, \dots, \theta_N \in \mathbb{R}^D\} \xleftarrow[\text{with } Q < D]{\text{map } f} \Omega = \{\omega_1, \dots, \omega_N \in \mathbb{R}^Q\}$$

$$\text{mapping: } \theta = f(\omega) + \eta \quad \text{with} \quad \eta \sim \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I})$$

Dimensionality reduction \longrightarrow **Manifold identification**

The lower-dimensional manifold on which the samples lie is characterized through the latent set



Gaussian Process Latent Variable Model

The choice of the dimensionality reduction method fell on the **Gaussian Process Latent Variable Model**^[1].

GPLVM: Gaussian Process prior over mapping f in

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

Motivation:

- Analytically sound mathematical tool
- Full distribution over the mapping f
- Full distribution over the derivatives of the mapping f

^[1]Lawrence, N., Probabilistic non-linear principal component analysis with Gaussian process latent variable models. *Journal of machine learning research* (2005)

Gaussian Process

Gaussian Process^[2]: a collection of random variables, any finite number of which have a joint Gaussian distribution.

If a real-valued stochastic process f is a GP, it will be denoted as

$$f(\cdot) \sim \mathcal{GP}(m(\cdot), k(\cdot, \cdot))$$

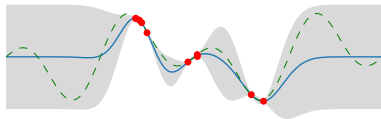
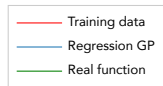
A Gaussian Process is fully specified by

- a mean function $m(\cdot)$
- a covariance function $k(\cdot, \cdot)$

where

$$m(\omega) = \mathbb{E}[f(\omega)],$$

$$k(\omega, \omega') = \mathbb{E}[(f(\omega) - m(\omega))(f(\omega') - m(\omega'))]$$



^[2]Rasmussen, C. E., Williams, C. K. I., Gaussian Processes for Machine Learning, the MIT Press (2006)

Gaussian Process Latent Variable Model

GPLVM: Gaussian Process prior over mapping \mathbf{f} in

$$\boldsymbol{\theta} = \mathbf{f}(\boldsymbol{\omega}) + \boldsymbol{\eta}$$

The likelihood of the data $\boldsymbol{\Theta}$ given the latent $\boldsymbol{\Omega}$ is given by

- ① marginalizing the mapping
- ② optimizing the latent variables

Resulting likelihood:

$$\begin{aligned} p(\boldsymbol{\Theta} \mid \boldsymbol{\Omega}, \beta) &= \prod_{j=1}^D \mathcal{N}(\boldsymbol{\theta}_{:,j} \mid \mathbf{0}, \mathbf{K} + \beta^{-1} \mathbf{I}) \\ &= \prod_{j=1}^D \mathcal{N}(\boldsymbol{\theta}_{:,j} \mid \mathbf{0}, \tilde{\mathbf{K}}) \end{aligned}$$

With the resulting noise model being:

$$\theta_{i,j} = \tilde{\mathbf{K}}_{(\boldsymbol{\omega}_i, \boldsymbol{\Omega})} \tilde{\mathbf{K}}^{-1} \boldsymbol{\Theta}_{:,j} + \eta_j$$

Gaussian Process Latent Variable Model

GPLVM: Gaussian Process prior over mapping \mathbf{f} in

$$\boldsymbol{\theta} = \mathbf{f}(\boldsymbol{\omega}) + \boldsymbol{\eta}$$

For differentiable kernels $k(\cdot, \cdot)$, the Jacobian \mathbf{J} of the mapping \mathbf{f} can be computed analytically:

$$J_{ij} = \frac{\partial f_i}{\partial \omega_j}$$

But as previously said, GPLVM can yield the full (Gaussian) distribution over the Jacobian. If the rows of \mathbf{J} are assumed to be independent:

$$p(\mathbf{J} \mid \boldsymbol{\Omega}, \beta) = \prod_{i=1}^D \mathcal{N}(\mathbf{J}_{i,:} \mid \boldsymbol{\mu}_{\mathbf{J}_{i,:}}, \boldsymbol{\Sigma}_{\mathbf{J}}),$$

Recap

- 1 After t iterations the optimization or sampling algorithm has yielded a set of observed points $\Theta = \{\theta_1, \dots, \theta_N \in \mathbb{R}^D\}$ in the parameter space
- 2 A GPLVM is trained on Θ in order to build a latent space Ω that describes the lower-dimensional manifold on which the optimization/sampling is allegedly taking place. We can:

- move from the latent space Ω to the observed space Θ :

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

$$\Theta \leftarrow \Omega$$

but not viceversa (f is not invertible)

- taking the gradients of a generic function $g: \Theta \rightarrow \mathbb{R}$ from the observed space Θ to the latent space Ω :

$$\nabla_{\omega} g(\omega) = \nabla_{\theta} g(\theta)$$

$$\Omega \leftarrow \Theta$$

- In this case a principal component of Θ is given by the vector of its distribution



Recap

- 1 After t iterations the optimization or sampling algorithm has yielded a set of observed points $\Theta = \{\theta_1, \dots, \theta_N \in \mathbb{R}^D\}$ in the parameter space
- 2 A GPLVM is trained on Θ in order to build a latent space Ω that describes the lower-dimensional manifold on which the optimization/sampling is allegedly taking place. We can:

- **move** from the latent space Ω to the observed space Θ :

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

$$\Theta \leftarrow \Omega$$

but not viceversa (f is not invertible)

- bring the **gradients** of a generic function $g : \Theta \rightarrow \mathbb{R}$ from the observed space Θ to the latent space Ω :

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

$$\Omega \leftarrow \Theta$$

In this case a punctual estimate of \mathbf{J} is given by the mean of its distribution.

Recap

- 1 After t iterations the optimization or sampling algorithm has yielded a set of observed points $\Theta = \{\theta_1, \dots, \theta_N \in \mathbb{R}^D\}$ in the parameter space
- 2 A GPLVM is trained on Θ in order to build a latent space Ω that describes the lower-dimensional manifold on which the optimization/sampling is allegedly taking place. We can:

- **move** from the latent space Ω to the observed space Θ :

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

$$\Theta \leftarrow \Omega$$

but not viceversa (f is not invertible)

- bring the **gradients** of a generic function $g : \Theta \rightarrow \mathbb{R}$ from the observed space Θ to the latent space Ω :

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

$$\Omega \leftarrow \Theta$$

In this case a punctual estimate of \mathbf{J} is given by the mean of its distribution.

Recap

- 1 After t iterations the optimization or sampling algorithm has yielded a set of observed points $\Theta = \{\theta_1, \dots, \theta_N \in \mathbb{R}^D\}$ in the parameter space
- 2 A GPLVM is trained on Θ in order to build a latent space Ω that describes the lower-dimensional manifold on which the optimization/sampling is allegedly taking place. We can:

- **move** from the latent space Ω to the observed space Θ :

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

$$\Theta \leftarrow \Omega$$

but not viceversa (f is not invertible)

- bring the **gradients** of a generic function $g : \Theta \rightarrow \mathbb{R}$ from the observed space Θ to the latent space Ω :

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

$$\Omega \leftarrow \Theta$$

In this case a punctual estimate of \mathbf{J} is given by the mean of its distribution.

AdaGeo Gradient-based Optimization

Minimization problem:

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

Iterative scheme solution (e.g. (stochastic) gradient descent):

$$\theta_{t+1} = \theta_t - \Delta \theta_t (\nabla_{\theta} g)$$

We propose, after having learned a latent representation with GPLVM, to move the problem onto the latent space Ω

Minimization problem:

$$\omega^* = \arg \min_{\omega} g(f(\omega))$$

Iterative scheme solution (e.g. (stochastic) gradient descent):

$$\omega_{t+1} = \omega_t - \Delta \omega_t (\nabla_{\omega} g)$$

AdaGeo Gradient-based Optimization

Minimization problem:

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

Iterative scheme solution (e.g. (stochastic) gradient descent):

$$\theta_{t+1} = \theta_t - \Delta \theta_t (\nabla_{\theta} g)$$

We propose, after having learned a latent representation with GPLVM, to move the problem onto the latent space Ω

Minimization problem:

$$\omega^* = \arg \min_{\omega} g(f(\omega))$$

Iterative scheme solution (e.g. (stochastic) gradient descent):

$$\omega_{t+1} = \omega_t - \Delta \omega_t (\nabla_{\omega} g)$$

Algorithm 1 AdaGeo gradient-based optimization (minimization)

- 1: **while** convergence is not reached **do**
- 2: Perform T_θ iterations with classic updates on the parameter space Θ :

$$\Delta\theta_t = \Delta\theta_t(\nabla_\theta g(\theta))$$

$$\theta_{t+1} = \theta_t - \Delta\theta_t$$

- 3: Train the GP-LVM model on the samples $\Theta = \{\theta_1, \dots, \theta_{T_\theta}\}$
- 4: Continue performing T_ω using the AdaGeo optimizer:

$$\Delta\omega_t = \Delta\omega_t(\nabla_\omega g(f(\omega)))$$

$$\omega_{t+1} = \omega_t - \Delta\omega_t$$

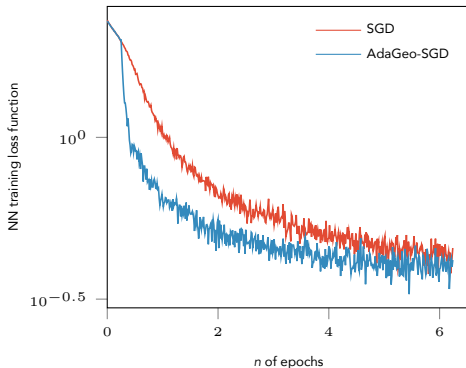
and moving back to the parameter space with

$$\theta_{t+1} = f(\omega_{t+1}).$$

- 5: **end while.**
-

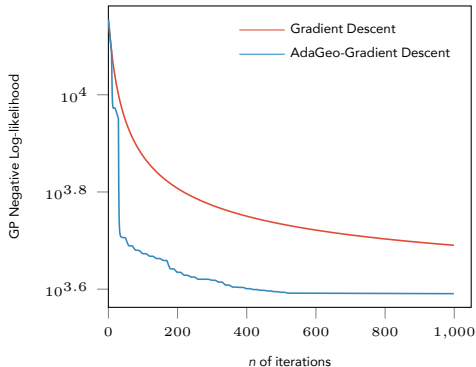
Experiment: Logistic Regression on MNIST

- Neural Network with a single hidden layer implementing logistic regression on MNIST
- Dimension of parameter space:
 $D = 7850$
- Dimension of latent space:
 $Q = 9$
- Iterations: $T_\theta = 20$ and $T_\omega = 30$



Experiment: Gaussian Process Training

- Concrete compressive strength dataset^[3]: regression task with 8 real variables
- Composite kernel (RBF, Matérn, linear and bias)
- Dimension of parameter space:
 $D = 9$
- Dimension of latent space:
 $Q = 3$
- Iterations: $T_{\theta} = 15$ and $T_{\omega} = 15$



^[3]Lichman, M., UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science (2013)

AdaGeo Bayesian Sampling

Bayesian sampling framework:

- a **dataset** $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is given
- \mathbf{X} is modeled with a **generative model** whose likelihood is

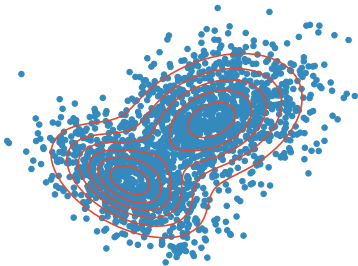
$$p(\mathbf{X}, \theta) = \prod_{i=1}^N p(\mathbf{x}_i, \theta),$$

parameterized by the vector $\theta \in \mathbb{R}^D$, with prior $p(\theta)$.

Performing statistical inference means getting insights on the **posterior distribution**

$$p(\theta | \mathbf{X}) = \frac{p(\mathbf{X} | \theta)p(\theta)}{p(\mathbf{X})}$$

analytically or **approximately** through **sampling**.



AdaGeo Bayesian Sampling

Bayesian sampling framework:

- a **dataset** $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is given
- \mathbf{X} is modeled with a **generative model** whose likelihood is

$$p(\mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{x}_i, \boldsymbol{\theta}),$$

parameterized by the vector $\boldsymbol{\theta} \in \mathbb{R}^D$, with prior $p(\boldsymbol{\theta})$.

Unfortunately the denominator is often intractable. One possible approach is to approximate the integral

$$p(\mathbf{X}) = \int_{\boldsymbol{\Theta}} p(\mathbf{X}, \boldsymbol{\theta}) d\boldsymbol{\theta} = \int_{\boldsymbol{\Theta}} p(\mathbf{X}|\boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

through **Markov Chain Monte Carlo** or similar methods.

Stochastic Gradient Langevin Dynamics

Stochastic gradient Langevin dynamics^[4] combines **stochastic optimization** and the physical concept of **Langevin dynamics** to build a posterior sampler

At each time t a mini-batch is extracted and the parameters are updated as:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \Delta \boldsymbol{\theta}_t,$$

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

$$\boldsymbol{\eta}_t \sim \mathcal{N}(\mathbf{0}, \epsilon_t \mathbf{I})$$

with the learning rate ϵ_t satisfying:

$$\sum_{t=1}^{\infty} \epsilon_t = \infty, \quad \sum_{t=1}^{\infty} \epsilon_t^2 < \infty$$

^[4]Welling, M., Teh, Y. W., Bayesian learning via stochastic gradient Langevin dynamics. *Proceedings of the 28th International Conference on Machine Learning (ICML)* (2011)

AdaGeo - Stochastic Gradient Langevin Dynamics

Analogously as before, we propose to:

- 1 Pick your favourite sampler and produce the first t samples to build the set $\Theta = \{\theta_1, \dots, \theta_N\}$
- 2 Train a GPLVM on Θ to learn the latent space Ω
- 3 Move the updates onto the latent space with **AdaGeo - Stochastic Gradient Langevin Dynamics**:

$$\omega_{t+1} = \omega_t + \Delta\omega_t,$$

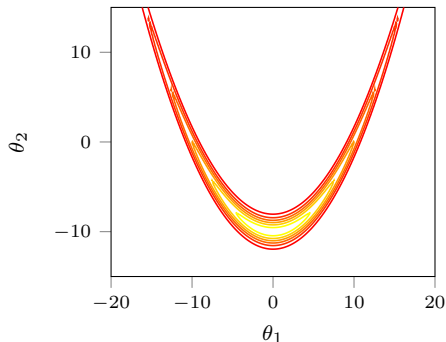
$$\Delta\omega_t = \frac{\epsilon_t}{2} \left(\nabla_{\omega} \log p(f(\omega_t)) + \frac{N}{n} \sum_{i=1}^n \nabla_{\omega} \log p(\mathbf{x}_{ti} | f(\omega_t)) \right) + \eta_t,$$

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

Experiment: Sampling from the Banana Distribution

The “banana” distribution has this formula:

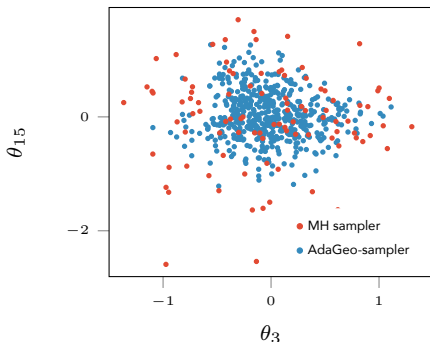
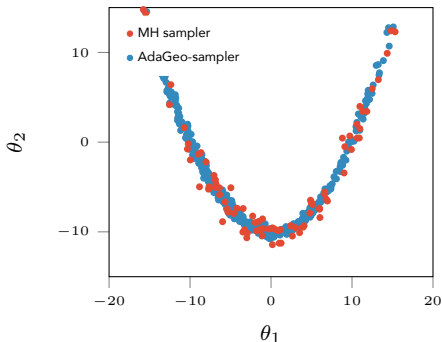
$$p(\boldsymbol{\theta}) \propto \exp \left(-\frac{\theta_1^2}{200} - \frac{(\theta_2 - b\theta_1^2 + 100b)^2}{2} - \sum_{j=3}^D \theta_j^2 \right)$$



θ_1 and θ_2 present the interaction shown on the left, while the other variables produce Gaussian noise

Experiment: Sampling from the Banana Distribution

- A Metropolis-Hastings returns the first 100 samples drawn from a 50-dimensional banana distribution
- AdaGeo-SGLD is then employed to sample from a 5-dimensional latent space



Bonus Round: Riemannian Extensions (theory only)

If a **covariance function** of a Gaussian Process is **differentiable**, then it is straightforward to show that the **mapping** f is also **differentiable**.

Under this assumption we can compute the latent **metric tensor** \mathbf{G} , which will give further information about the **geometry** of the latent space (distances, geodetic lines etc.)

If \mathbf{J} is the Jacobian of the mapping f , then

$$\mathbf{G} = \mathbf{J}^\top \mathbf{J}$$

This yields a distribution over the metric tensor^[5]:

$$\mathbf{G} \sim \mathcal{W}_Q(D, \Sigma_{\mathbf{J}}, \mathbb{E}[\mathbf{J}^\top] \mathbb{E}[\mathbf{J}])$$

and a punctual estimate can be obtained with

$$\mathbb{E}[\mathbf{J}^\top \mathbf{J}] = \mathbb{E}[\mathbf{J}^\top] \mathbb{E}[\mathbf{J}] + D \Sigma_{\mathbf{J}}.$$

^[5]Tosi, A., Hauberg, S., Vellido, A., Lawrence, N. D., Metrics for probabilistic geometries. *Uncertainty in Artificial Intelligence* (2014)

Stochastic gradient Riemannian Langevin dynamics^[6] puts together the advantages of exploiting a known **Riemannian geometry** with the scalability of the **stochastic optimization** approaches:

$$\begin{aligned}
 \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t + \Delta \boldsymbol{\theta}_t, \\
 \Delta \boldsymbol{\theta}_t &= \frac{\epsilon_t}{2} \boldsymbol{\mu}(\boldsymbol{\theta}_t) + \mathbf{G}^{-\frac{1}{2}}(\boldsymbol{\theta}_t) \boldsymbol{\eta}_t \\
 \boldsymbol{\eta}_t &\sim \mathcal{N}(\mathbf{0}, \epsilon_t \mathbf{I}),
 \end{aligned}$$

where

$$\begin{aligned}
 \boldsymbol{\mu}(\boldsymbol{\theta})_j &= \left(\mathbf{G}^{-1}(\boldsymbol{\theta}) \left(\nabla_{\boldsymbol{\theta}} \log p(\boldsymbol{\theta}) + \frac{N}{n} \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}_{ti} | \boldsymbol{\theta}) \right) \right)_j \\
 &\quad - 2 \sum_{k=1}^D \left(\mathbf{G}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{G}(\boldsymbol{\theta})}{\partial \theta_k} \mathbf{G}^{-1}(\boldsymbol{\theta}) \right)_{jk} + \sum_{k=1}^D (\mathbf{G}^{-1}(\boldsymbol{\theta}))_{jk} \text{Tr} \left(\mathbf{G}^{-1}(\boldsymbol{\theta}) \frac{\partial \mathbf{G}(\boldsymbol{\theta})}{\partial \theta_k} \right)
 \end{aligned}$$

^[6]Patterson, S., Teh, Y. W., Stochastic gradient Riemannian Langevin dynamics on the probability simplex. *Advances in Neural Information Processing Systems* (2013)

Bonus Round: AdaGeo - Stochastic Gradient Riemannian Langevin Dynamics

Analogously as the SGLD case, we can move now the update in the latent space with **AdaGeo - Stochastic gradient Riemannian Langevin dynamics**:

$$\omega_{t+1} = \omega_t + \Delta\omega_t$$

$$\Delta\omega_t = \frac{\epsilon_t}{2} \mu(\omega_t) + \mathbf{G}_\omega^{-\frac{1}{2}}(\omega_t) \eta_t$$

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

where

$$\begin{aligned} \mu(\omega)_j = & \left(\mathbf{G}_\omega^{-1}(\omega) \left(\nabla_\omega \log p(\mathbf{f}(\omega)) + \frac{N}{n} \sum_{i=1}^n \nabla_\omega \log p(\mathbf{x}_{ti} | \mathbf{f}(\omega)) \right) \right)_j \\ & - 2 \sum_{k=1}^Q \left(\mathbf{G}_\omega^{-1}(\omega) \frac{\partial \mathbf{G}_\omega(\omega)}{\partial \omega_k} \mathbf{G}_\omega^{-1}(\omega) \right)_{jk} + \sum_{k=1}^Q (\mathbf{G}_\omega^{-1}(\omega))_{jk} \text{Tr} \left(\mathbf{G}_\omega^{-1}(\omega) \frac{\partial \mathbf{G}_\omega(\omega)}{\partial \omega_k} \right) \end{aligned}$$



Conclusions

- We develop a **generic framework** for combining **dimensionality reduction** techniques with **sampling and optimization** methods
- We contribute to **gradient-based optimization** methods by coupling them with appropriate dimensionality reduction techniques. In particular, we improve the performances of gradient descent and stochastic gradient descent, when training respectively a Gaussian Process and a neural network
- We contribute to **Markov Chain Monte Carlo** by developing a AdaGeo version of stochastic gradient Langevin dynamics; the information gathered through the latent space are employed to compute the steps of the Markov chain
- We extend the approach to stochastic gradient **Riemannian** Langevin dynamics, thanks to the geometric tensor naturally recovered by the GP-LVM model



Thank you

Reference: Abbati, G., Tosi, A., Flaxman, S., Osborne, M. A. (2018). AdaGeo: Adaptive Geometric Learning for Optimization and Sampling. In *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics (AISTATS)*, to appear.