

Natural gradient descent with momentum

Nonlinear Approximation for High-Dimensional Problems 2025
Workshop in honor of Albert Cohen

Agustín Somacal

In collaboration with **Anthony Nouy**

École Centrale de Nantes

Laboratoire de Mathématiques Jean Leray



Albert Cohen



Welcome to the home page of Albert Cohen.

Address: [Laboratoire Jacques-Louis Lions](#),
Sorbonne Université
4, Place Jussieu
75005 Paris, France

Phone: (33)(1)44277195
email : albert.cohen@sorbonne-universite.fr

[Research interests](#)

Abridged [curriculum vitae](#)

List of [publications](#)

Cours [de M2 2023/2024](#)

Activity group [SMAT-SIGMA](#)

The conference [Curves and Surfaces 2022](#)

The conference [Foundations of Computational Mathematics 2023](#)

The [FoCM](#) society

Things I like to do ...

... [below see level](#)

... [above see level](#)

Albert Cohen



Welcome to the home page of Albert Cohen.

Address: [Laboratoire Jacques-Louis Lions](#),
Sorbonne Université
4, Place Jussieu
75005 Paris, France

Phone: (33)(1)44277195
email : albert.cohen@sorbonne-universite.fr

[Research interests](#)

Abridged [curriculum vitae](#)

List of [publications](#)

Cours [de M2 2023/2024](#)

Activity group [SMAT-SIGMA](#)

The conference [Curves and Surfaces 2022](#)

The conference [Foundations of Computational Mathematics 2023](#)

The [FoCM](#) society

Things I like to do ...

... [below see level](#)

... [above see level](#)

Outline

1. What is natural gradient and why we may need it?

- Problem setting
- From gradient descent to Newton's method.
- From Newton's method to natural gradient.
- Toy examples to gain intuition.

2. What is momentum and when we may need it?

3. How to combine momentum and natural gradient.

- Two toy examples
- Two less toy examples



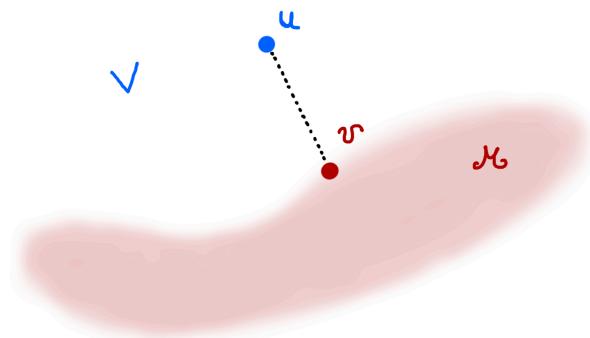
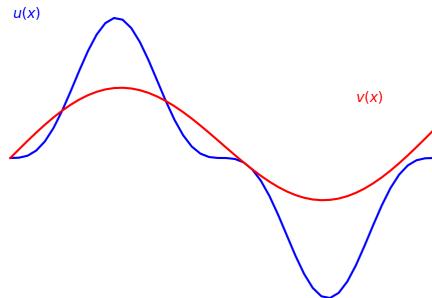
Problem formulation

Objective: approximate target function $u \in V$ by $v \in \mathcal{M} \subset V$

Target function

$$u : \mathbb{R}^d \rightarrow \mathbb{R} \in V$$

Hilbert space
 $L^2(\Omega), H^1(\Omega), \dots$



Problem formulation

Objective: approximate target function $u \in V$ by $v \in \mathcal{M} \subset V$

Target function

$$\textcolor{blue}{u} : \mathbb{R}^d \rightarrow \mathbb{R} \in V$$

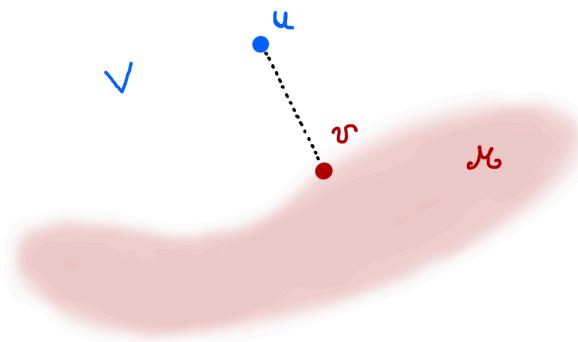
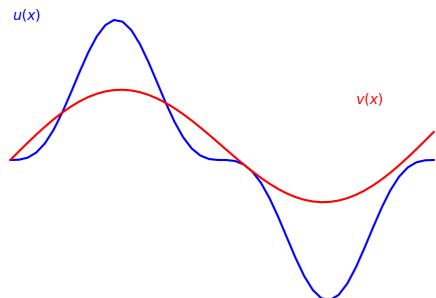
Hilbert space
 $L^2(\Omega), H^1(\Omega), \dots$

Approximation
manifold

$$\mathcal{M} := \{v_\theta(x) = \textcolor{red}{A}(\theta)(x); \theta \in \mathbb{R}^p\}$$

Linear model

$$\begin{aligned}\textcolor{red}{A}(\theta)(x) &= \theta_1 \phi_1(x) + \dots + \theta_p \phi_p(x) \\ &= \theta^T \Phi(x)\end{aligned}$$



Problem formulation

Objective: approximate target function $u \in V$ by $v \in \mathcal{M} \subset V$

Target function

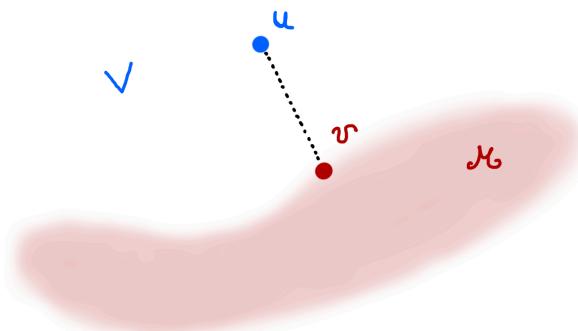
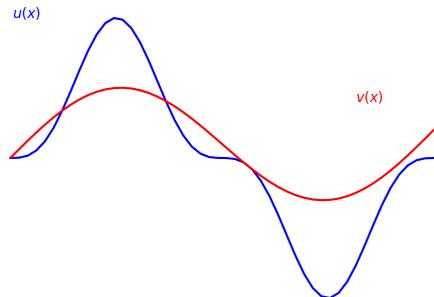
$$\textcolor{blue}{u} : \mathbb{R}^d \rightarrow \mathbb{R} \in V$$

Hilbert space
 $L^2(\Omega), H^1(\Omega), \dots$

Approximation
manifold

$$\mathcal{M} := \{v_\theta(x) = \textcolor{red}{A}(\theta)(x); \theta \in \mathbb{R}^p\}$$

Linear model,
Neural network,
...



Minimization problem

Objective: approximate $u \in V$ by $v \in \mathcal{M}$.

Continuous problem

$$\begin{aligned}\mathcal{L}_u(\textcolor{red}{v}) &= \frac{1}{2} \|\textcolor{blue}{u} - \textcolor{red}{v}\|_V^2 \\ &= \frac{1}{2} \langle \textcolor{blue}{u} - \textcolor{red}{v}, \textcolor{blue}{u} - \textcolor{red}{v} \rangle_V \\ &= \frac{1}{2} \int (\textcolor{blue}{u}(x) - \textcolor{red}{v}(x))^2 d\mu(x)\end{aligned}$$

Minimization problem

Objective: approximate $u \in V$ by $v \in \mathcal{M}$.

Continuous problem

$$\begin{aligned}\mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_V^2 \\ &= \frac{1}{2} \langle \mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v} \rangle_V \\ &= \frac{1}{2} \int (\mathbf{u}(x) - \mathbf{v}(x))^2 d\mu(x)\end{aligned}$$

Discrete problem

$$\begin{aligned}\mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_m^2 \\ &= \frac{1}{2} \langle \mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v} \rangle_m \\ &= \frac{1}{2m} \sum_{i=1}^m (\mathbf{u}(x_i) - \mathbf{v}(x_i))^2\end{aligned}$$

Minimization problem

Objective: approximate $u \in V$ by $v \in \mathcal{M}$.

Continuous problem

$$\begin{aligned}\mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_V^2 \\ &= \frac{1}{2} \langle \mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v} \rangle_V \\ &= \frac{1}{2} \int (\mathbf{u}(x) - \mathbf{v}(x))^2 d\mu(x)\end{aligned}$$

Discrete problem

$$\begin{aligned}\mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_m^2 \\ &= \frac{1}{2} \langle \mathbf{u} - \mathbf{v}, \mathbf{u} - \mathbf{v} \rangle_m \\ &= \frac{1}{2m} \sum_{i=1}^m (\mathbf{u}(x_i) - \mathbf{v}(x_i))^2\end{aligned}$$

Functional perspective

$$\mathbf{v}^* = \arg \min_{\mathbf{v} \in \mathcal{M}} \mathcal{L}_u(\mathbf{v})$$

Parameter perspective

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \mathcal{L}_u(A(\theta))$$

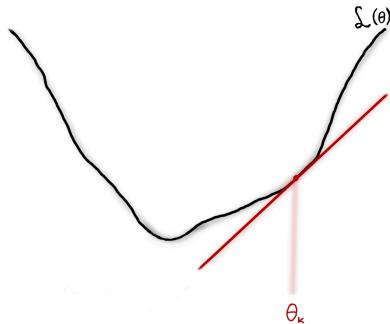
$$\mathcal{L}_u(A(\theta)) = (\mathcal{L}_u \circ A)(\theta) =: \mathcal{L}(\theta)$$

Gradient descent

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k .

$$\mathcal{L}(\theta) \approx \mathcal{L}(\theta_k) + \langle \nabla_{\theta} \mathcal{L}(\theta_k), \theta - \theta_k \rangle_{\mathbb{R}^p}$$

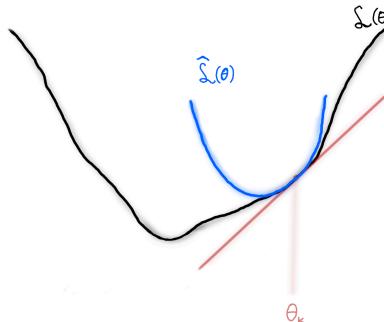


Gradient descent

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k plus **penalization on the distance** traveled on each step.

$$\mathcal{L}(\theta) \approx \mathcal{L}(\theta_k) + \langle \nabla_{\theta} \mathcal{L}(\theta_k), \theta - \theta_k \rangle_{\mathbb{R}^p} + \frac{1}{2s} \rho(\theta, \theta_k)$$

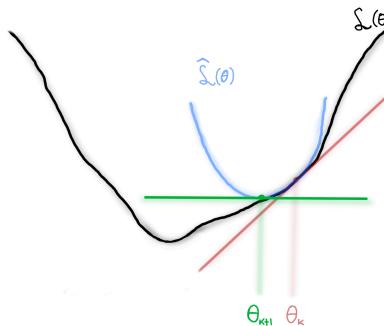


Gradient descent

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k plus **penalization on the distance** traveled on each step.

$$0 = \nabla_{\theta} \left[\mathcal{L}(\theta_k) + \langle \nabla_{\theta} \mathcal{L}(\theta_k), \theta - \theta_k \rangle_{\mathbb{R}^p} + \frac{1}{2s} \rho(\theta, \theta_k) \right]$$



Gradient descent

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k plus **penalization on the distance** traveled on each step.

$$\frac{1}{2} \nabla_{\theta} \rho(\theta, \theta_k) = -s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Gradient descent

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k plus **penalization on the distance** traveled on each step.

$$\frac{1}{2} \nabla_{\theta} \rho(\theta, \theta_k) = -s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Gradient descent

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_{\mathbb{R}^p}^2$$

$$\nabla_{\theta} \rho(\theta, \theta_k) = 2(\theta - \theta_k)$$

$$\theta = \theta_k - s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Preconditioned gradient

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k plus **penalization on the distance** traveled on each step.

$$\frac{1}{2} \nabla_{\theta} \rho(\theta, \theta_k) = -s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Gradient descent

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_{\mathbb{R}^p}^2$$

$$\nabla_{\theta} \rho(\theta, \theta_k) = 2(\theta - \theta_k)$$

$$\theta = \theta_k - s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Preconditioned gradient

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_M^2$$

$$\nabla_{\theta} \rho(\theta, \theta_k) = 2M(\theta - \theta_k)$$

$$\theta = \theta_k - s M^{-1} \nabla_{\theta} \mathcal{L}(\theta_k)$$

Newton's method

Iteratively improve approximation by minimizing $\mathcal{L}(\theta_k)$.

Taylor expansion around current iterate θ_k plus **penalization on the distance** traveled on each step.

$$\frac{1}{2} \nabla_{\theta} \rho(\theta, \theta_k) = -s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Gradient descent

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_{\mathbb{R}^p}^2$$

$$\nabla_{\theta} \rho(\theta, \theta_k) = 2(\theta - \theta_k)$$

$$\theta = \theta_k - s \nabla_{\theta} \mathcal{L}(\theta_k)$$

Newton's method

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_H^2$$

$$\nabla_{\theta} \rho(\theta, \theta_k) = 2H(\theta - \theta_k)$$

$$\theta = \theta_k - s H^{-1} \nabla_{\theta} \mathcal{L}(\theta_k)$$

Functional gradient and geometric intuition

Some properties [Gruhlke, Robert, Anthony Nouy, and Philipp Trunschke. 2024].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \end{aligned}$$

Functional gradient and geometric intuition

Some properties [Gruhlke, Robert, Anthony Nouy, and Philipp Trunschke. 2024].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \end{aligned}$$

Functional gradient and geometric intuition

Some properties [Gruhlke, Robert, Anthony Nouy, and Philipp Trunschke. 2024].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \end{aligned}$$

$$\mathcal{L}_{\textcolor{blue}{u}}(\textcolor{red}{v}) = \frac{1}{2} \|\textcolor{blue}{u} - \textcolor{red}{v}\|_{L^2(\Omega)}^2 \quad \longrightarrow \quad V \ni \nabla \mathcal{L}_u = \textcolor{blue}{u} - \textcolor{red}{v}$$

Natural gradient from Newton's method

[Amari, Shun-ichi. 1998] [Martens, James 2020].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \\ &= \left\langle H_V \mathcal{L} \frac{\partial A}{\partial \theta_i}, \frac{\partial A}{\partial \theta_j} \right\rangle_V + \left\langle \nabla \mathcal{L}_u, \frac{\partial^2 A}{\partial \theta_i \partial \theta_j} \right\rangle_V \\ &= G + \langle \nabla \mathcal{L}_u, H_A \rangle_V \end{aligned}$$

Natural gradient from Newton's method

[Amari, Shun-ichi. 1998] [Martens, James 2020].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \\ &= \left\langle H_V \mathcal{L} \frac{\partial A}{\partial \theta_i}, \frac{\partial A}{\partial \theta_j} \right\rangle_V + \left\langle \nabla \mathcal{L}_u, \frac{\partial^2 A}{\partial \theta_i \partial \theta_j} \right\rangle_V \\ &= G + \langle \nabla \mathcal{L}_u, H_A \rangle_V \end{aligned}$$

Natural gradient from Newton's method

[Amari, Shun-ichi. 1998] [Martens, James 2020].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \\ &= \left\langle H_V \mathcal{L} \frac{\partial A}{\partial \theta_i}, \frac{\partial A}{\partial \theta_j} \right\rangle_V + \left\langle \nabla \mathcal{L}_u, \frac{\partial^2 A}{\partial \theta_i \partial \theta_j} \right\rangle_V \\ &= G + \langle \nabla \mathcal{L}_u, H_A \rangle_V \end{aligned}$$

$$G_{ij} = \int \frac{\partial A}{\partial \theta_i}(x) [H_V \mathcal{L}(\textcolor{red}{v})](x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x).$$

Natural gradient from Newton's method

[Amari, Shun-ichi. 1998] [Martens, James 2020].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \\ &= \left\langle H_V \mathcal{L} \frac{\partial A}{\partial \theta_i}, \frac{\partial A}{\partial \theta_j} \right\rangle_V + \left\langle \nabla \mathcal{L}_u, \frac{\partial^2 A}{\partial \theta_i \partial \theta_j} \right\rangle_V \\ &= G + \langle \nabla \mathcal{L}_u, H_A \rangle_V \end{aligned}$$

$$G_{ij} = \int \frac{\partial A}{\partial \theta_i}(x) [H_V \mathcal{L}(\textcolor{red}{v})](x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x).$$

$$\mathcal{L}_{\textcolor{blue}{u}}(\textcolor{red}{v}) = \frac{1}{2} \| \textcolor{blue}{u} - \textcolor{red}{v} \|_{L^2(\Omega)}^2 \quad \longrightarrow \quad H_V \mathcal{L}(\textcolor{red}{v})(x) = 1$$

Natural gradient from Newton's method

[Amari, Shun-ichi. 1998] [Martens, James 2020].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \\ &= \left\langle H_V \mathcal{L} \frac{\partial A}{\partial \theta_i}, \frac{\partial A}{\partial \theta_j} \right\rangle_V + \left\langle \nabla \mathcal{L}_u, \frac{\partial^2 A}{\partial \theta_i \partial \theta_j} \right\rangle_V \\ &= G + \langle \nabla \mathcal{L}_u, H_A \rangle_V \end{aligned}$$

$$G_{ij} = \int \frac{\partial A}{\partial \theta_i}(x) [H_V \mathcal{L}(\textcolor{red}{v})](x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x).$$

$$\mathcal{L}_{\textcolor{blue}{u}}(\textcolor{red}{v}) = \frac{1}{2} \| \textcolor{blue}{u} - \textcolor{red}{v} \|_{L^2(\Omega)}^2 \quad \longrightarrow \quad H_V \mathcal{L}(\textcolor{red}{v})(x) = 1$$

$$G_{ij} = \int \left[\frac{\partial A}{\partial \theta_i} \frac{\partial A}{\partial \theta_j} \right] (x) dx$$

Natural gradient from Newton's method

[Amari, Shun-ichi. 1998] [Martens, James 2020].

$$\begin{aligned} H_{ij} &= \frac{\partial^2 \mathcal{L}}{\partial \theta_i \partial \theta_j} = \frac{\partial}{\partial \theta_i} \left[\frac{\partial \mathcal{L}}{\partial \theta_j} \right] = \frac{\partial}{\partial \theta_i} \left[\frac{\partial}{\partial \theta_j} (\mathcal{L}_{\textcolor{blue}{u}} \circ A) \right] = \frac{\partial}{\partial \theta_i} \left[\int_{\Omega} \nabla \mathcal{L}_u(\textcolor{red}{v}_{\theta})(x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x) \right] \\ &= \frac{\partial}{\partial \theta_i} \left[\left\langle \nabla \mathcal{L}_u, \frac{\partial A}{\partial \theta_j} \right\rangle_V \right] \\ &= \left\langle H_V \mathcal{L} \frac{\partial A}{\partial \theta_i}, \frac{\partial A}{\partial \theta_j} \right\rangle_V + \left\langle \nabla \mathcal{L}_u, \frac{\partial^2 A}{\partial \theta_i \partial \theta_j} \right\rangle_V \\ &= G + \langle \nabla \mathcal{L}_u, H_A \rangle_V \quad \text{Model linearization} \end{aligned}$$

$$G_{ij} = \int \frac{\partial A}{\partial \theta_i}(x) [H_V \mathcal{L}(\textcolor{red}{v})](x) \frac{\partial A}{\partial \theta_j}(x) d\mu(x).$$

$$\mathcal{L}_{\textcolor{blue}{u}}(\textcolor{red}{v}) = \frac{1}{2} \| \textcolor{blue}{u} - \textcolor{red}{v} \|_{L^2(\Omega)}^2 \quad \longrightarrow \quad H_V \mathcal{L}(\textcolor{red}{v})(x) = 1$$

$$G_{ij} = \int \left[\frac{\partial A}{\partial \theta_i} \frac{\partial A}{\partial \theta_j} \right](x) dx$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1} \nabla_{\theta} \mathcal{L}$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1} \nabla_{\theta} \mathcal{L}$$

$$\begin{aligned} \frac{\partial \mathbf{v}_{\theta}}{\partial s} &= \frac{\partial \mathbf{v}_{\theta}}{\partial \theta} \frac{\partial \theta}{\partial s} = \frac{\partial A}{\partial \theta}^T M^{-1} \nabla_{\theta} \mathcal{L} = -\frac{\partial A}{\partial \theta}^T M^{-1} \left\langle \nabla \mathcal{L}, \frac{\partial A}{\partial \theta} \right\rangle_V \\ &= -\left(G^{-\frac{1}{2}} \frac{\partial A}{\partial \theta} \right)^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \left\langle \nabla \mathcal{L}, G^{-\frac{1}{2}} \frac{\partial A}{\partial \theta} \right\rangle_V \\ &= -E^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, E \rangle_V \end{aligned}$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1} \nabla_{\theta} \mathcal{L}$$

$$\begin{aligned}\frac{\partial \mathbf{v}_{\theta}}{\partial s} &= \frac{\partial \mathbf{v}_{\theta}}{\partial \theta} \frac{\partial \theta}{\partial s} = \frac{\partial A}{\partial \theta}^T M^{-1} \nabla_{\theta} \mathcal{L} = -\frac{\partial A}{\partial \theta}^T M^{-1} \left\langle \nabla \mathcal{L}, \frac{\partial A}{\partial \theta} \right\rangle_V \\ &= -\left(G^{-\frac{1}{2}} \frac{\partial A}{\partial \theta} \right)^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \left\langle \nabla \mathcal{L}, G^{-\frac{1}{2}} \frac{\partial A}{\partial \theta} \right\rangle_V \\ &= -E^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, E \rangle_V\end{aligned}$$

$$\langle E, E \rangle_V = \int_{\Omega} \left(G^{-\frac{1}{2}} \frac{\partial A}{\partial \theta} \right)^T \left(G^{-\frac{1}{2}} \frac{\partial A}{\partial \theta} \right) d\mu(x) = \int_{\Omega} \frac{\partial A}{\partial \theta}^T G^{-1} \frac{\partial A}{\partial \theta} d\mu(x) = I$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1} \nabla_{\theta} \mathcal{L}$$

Preconditioned gradient flow
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1} \nabla_{\theta} \mathcal{L}$$

Preconditioned gradient flow
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Gradient descent $M = I$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} I^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - s G^{-1} \nabla \mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1} \nabla_{\theta} \mathcal{L}$$

Preconditioned gradient flow
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Gradient descent $M = I$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - sG^{-1}\nabla\mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1}\nabla_{\theta}\mathcal{L}$$

Preconditioned gradient flow
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Gradient descent $M = I$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient $M = G$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} G^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - sG^{-1}\nabla\mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1}\nabla_{\theta}\mathcal{L}$$

Preconditioned gradient flow
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Gradient descent $M = I$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient $M = G$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient flow dynamics

$$\rho(\theta, \theta_k) = \|\theta - \theta_k\|_G^2 \quad \longrightarrow \quad \theta = \theta_k - sG^{-1}\nabla\mathcal{L}(\theta_k)$$

Preconditioned gradient flow
in **parameter space**

$$\frac{\partial \theta}{\partial s} = -M^{-1}\nabla_{\theta}\mathcal{L}$$

Preconditioned gradient flow
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G^{\frac{1}{2}} M^{-1} G^{\frac{1}{2}} \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Gradient descent $M = I$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -\textcolor{green}{E}^T G \langle \nabla \mathcal{L}, \textcolor{green}{E} \rangle_V$$

Natural gradient $M = G$
in **functional space**

$$\frac{\partial \textcolor{red}{v}_{\theta}}{\partial s} = -P_{\mathcal{T}_k} \nabla \mathcal{L}$$



Toy example

Gradient descent trajectory.

$$\mathbf{u} \in L^2([0, 1])$$

$$\mathcal{L}_u(\mathbf{v}) = \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|^2$$

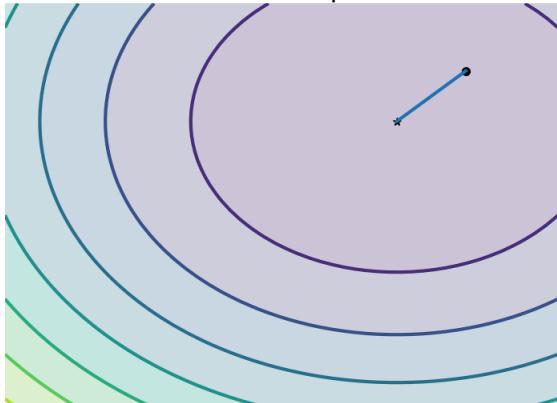
$$\mathbf{v}_\theta(x) = \mathbf{A}(\theta)(x) = \theta^T \Phi(x)$$

$$= \theta^T \begin{bmatrix} 1 \\ \sqrt{2} \sin(2\pi x) \end{bmatrix}$$

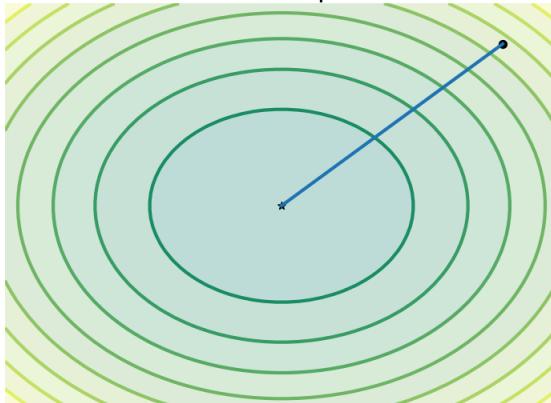
$$\frac{\partial A}{\partial \theta_i}(\theta)(x) = \Phi_i(x)$$

$$G_{ij} = \int \left[\frac{\partial A}{\partial \theta_i} \frac{\partial A}{\partial \theta_j} \right] (x) dx = \delta_{ij}$$

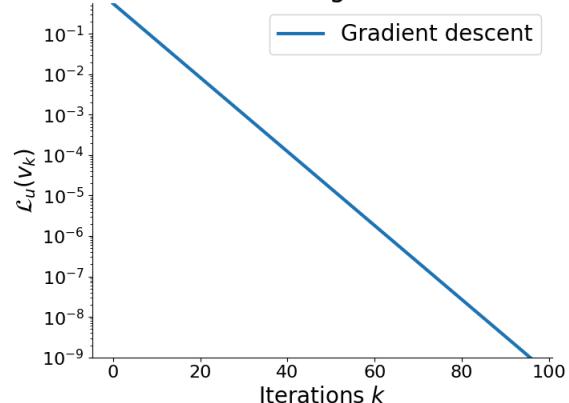
Parameter space



Function space



Convergence

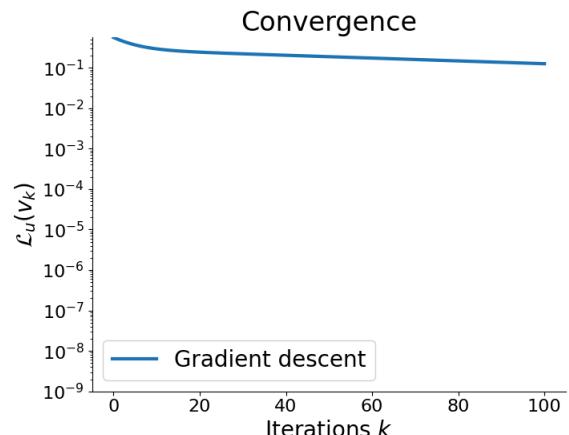
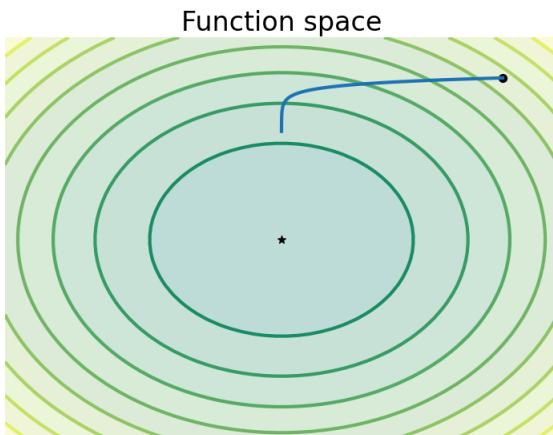
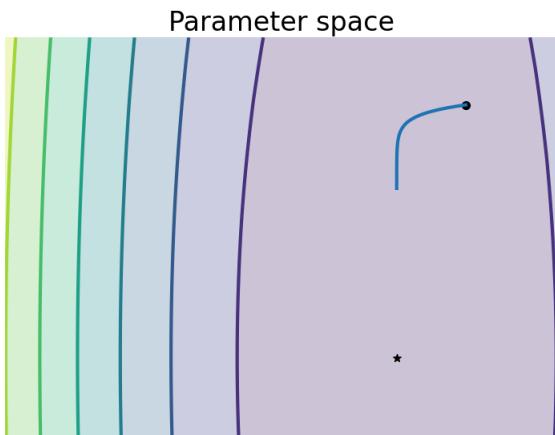


Toy example

Gradient descent is biased in functional space.

$$\begin{aligned}\mathbf{u} &\in L^2([0, 1]) \\ \mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|^2 \\ \mathbf{v}_\theta(x) &= \mathbf{A}(\theta)(x) = \theta^T \mathbf{B} \Phi(x) \\ &= \theta^T \mathbf{B} \begin{bmatrix} 1 \\ \sqrt{2} \sin(2\pi x) \end{bmatrix}\end{aligned}$$

$$\begin{aligned}\mathbb{R}^{p \times p} &\ni \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0.2 \end{bmatrix} \\ \frac{\partial A}{\partial \theta_i}(\theta)(x) &= \mathbf{B}_i \Phi_i(x) \\ G_{ij} &= \int \left[\frac{\partial A}{\partial \theta_i} [\mathbf{B}^T \mathbf{B}]_{ij} \frac{\partial A}{\partial \theta_j} \right] (x) dx = [\mathbf{B}^T \mathbf{B}]_{ij}\end{aligned}$$

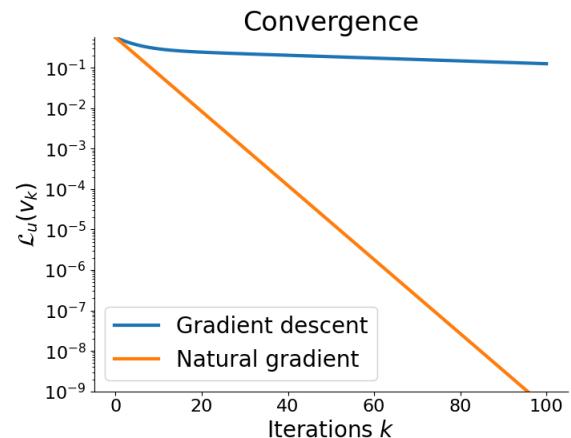
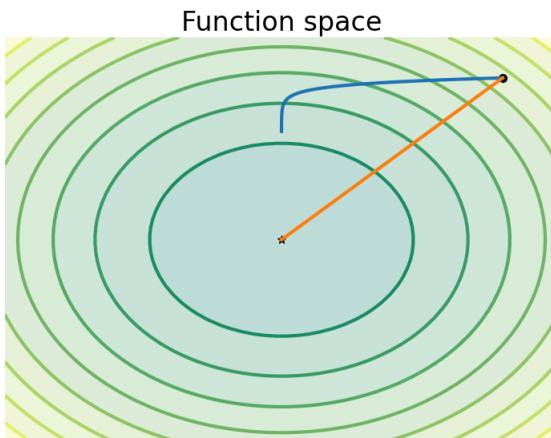
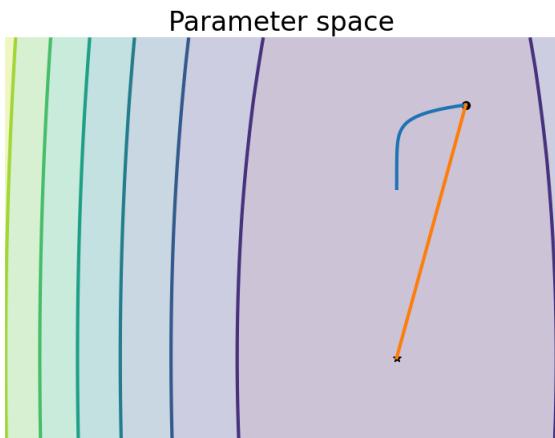


Toy example

Natural gradient descent.

$$\begin{aligned}\mathbf{u} &\in L^2([0, 1]) \\ \mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|^2 \\ \mathbf{v}_\theta(x) &= \mathbf{A}(\theta)(x) = \theta^T \mathbf{B} \Phi(x) \\ &= \theta^T \mathbf{B} \begin{bmatrix} 1 \\ \sqrt{2} \sin(2\pi x) \end{bmatrix}\end{aligned}$$

$$\begin{aligned}\mathbb{R}^{p \times p} &\ni \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0.2 \end{bmatrix} \\ \frac{\partial \mathbf{A}}{\partial \theta_i}(\theta)(x) &= \mathbf{B}_i \Phi_i(x) \\ G_{ij} &= \int \left[\frac{\partial \mathbf{A}}{\partial \theta_i} [\mathbf{B}^T \mathbf{B}]_{ij} \frac{\partial \mathbf{A}}{\partial \theta_j} \right] (x) dx = [\mathbf{B}^T \mathbf{B}]_{ij}\end{aligned}$$



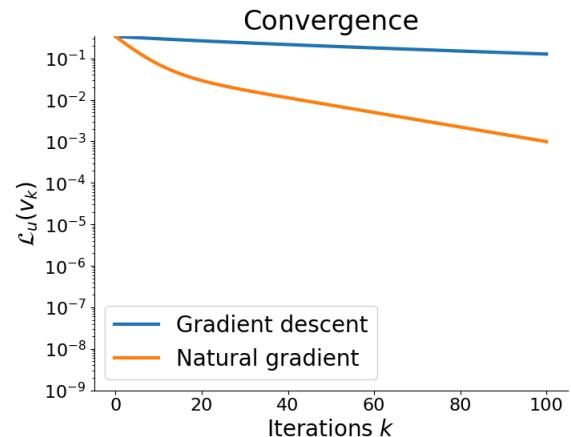
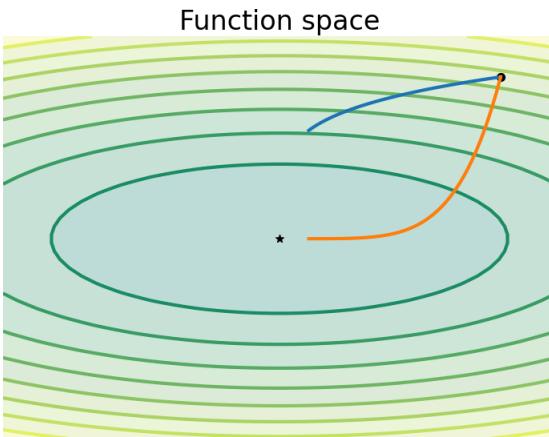
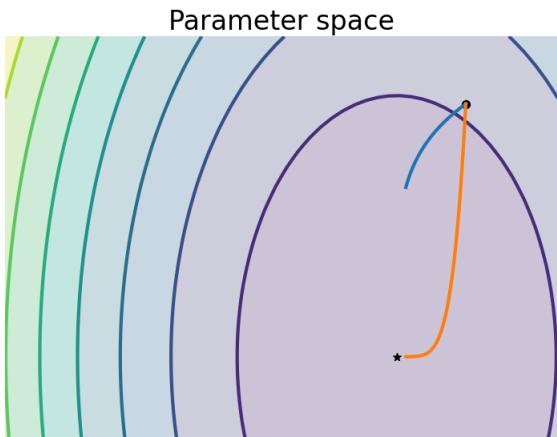
Toy example

Non isotropic loss.

$$\begin{aligned}\mathbf{u} &\in L^2([0, 1]) \\ \mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_{\mathbf{K}}^2 \\ \mathbf{v}_{\theta}(x) &= \mathbf{A}(\theta)(x) = \theta^T \mathbf{B} \Phi(x)\end{aligned}$$

$$= \theta^T \mathbf{B} \left[\frac{1}{\sqrt{2} \sin(2\pi x)} \right]$$

$$\begin{aligned}\mathbb{R}^{p \times p} &\ni \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0.2 \end{bmatrix} \\ \frac{\partial A}{\partial \theta_i}(\theta)(x) &= \mathbf{B}_i \Phi_i(x) \\ G_{ij} &= \int \left[\frac{\partial A}{\partial \theta_i} [\mathbf{B}^T \mathbf{K} \mathbf{B}]_{ij} \frac{\partial A}{\partial \theta_j} \right] (x) dx = [\mathbf{B}^T \mathbf{K} \mathbf{B}]_{ij}\end{aligned}$$



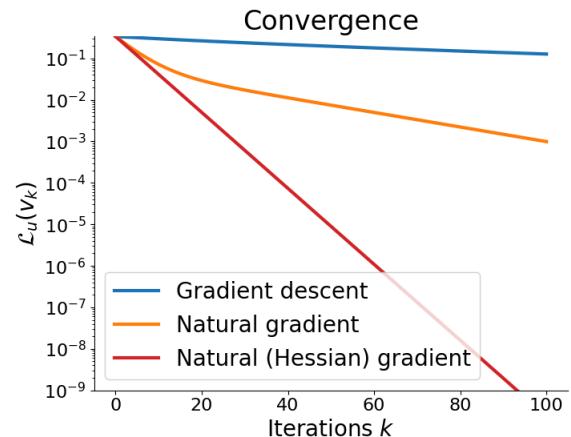
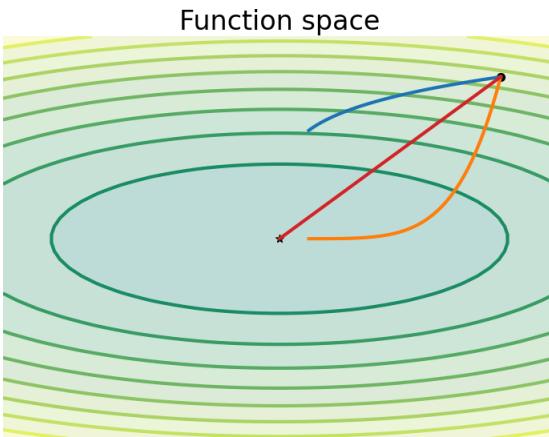
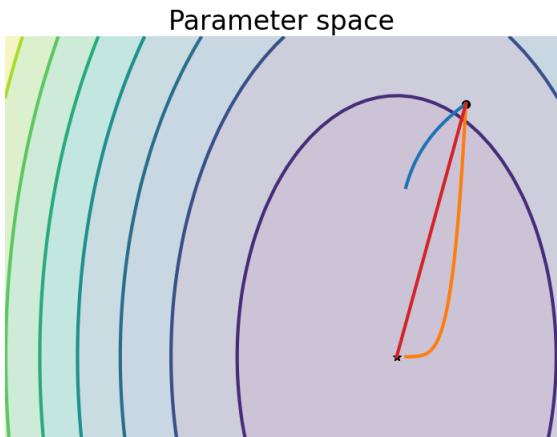
Toy example

Natural gradient descent with loss hessian.

$$\begin{aligned} \mathbf{u} &\in L^2([0, 1]) \\ \mathcal{L}_u(\mathbf{v}) &= \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_{\mathbf{K}}^2 \\ \mathbf{v}_{\theta}(x) &= \mathbf{A}(\theta)(x) = \theta^T \mathbf{B} \Phi(x) \end{aligned}$$

$$= \theta^T \mathbf{B} \left[\frac{1}{\sqrt{2} \sin(2\pi x)} \right]$$

$$\begin{aligned} \mathbb{R}^{p \times p} &\ni \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0.2 \end{bmatrix} \\ \frac{\partial \mathbf{A}}{\partial \theta_i}(\theta)(x) &= \mathbf{B}_i \Phi_i(x) \\ G_{ij} &= \int \left[\frac{\partial \mathbf{A}}{\partial \theta_i} [\mathbf{B}^T \mathbf{K} \mathbf{B}]_{ij} \frac{\partial \mathbf{A}}{\partial \theta_j} \right] (x) dx = [\mathbf{B}^T \mathbf{K} \mathbf{B}]_{ij} \end{aligned}$$



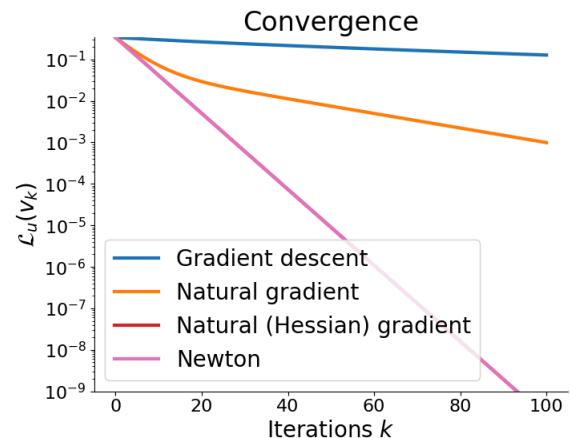
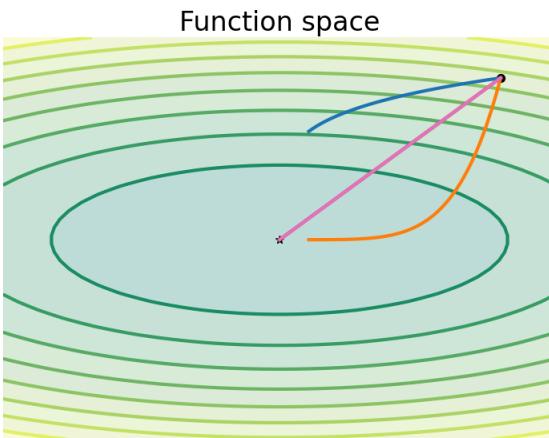
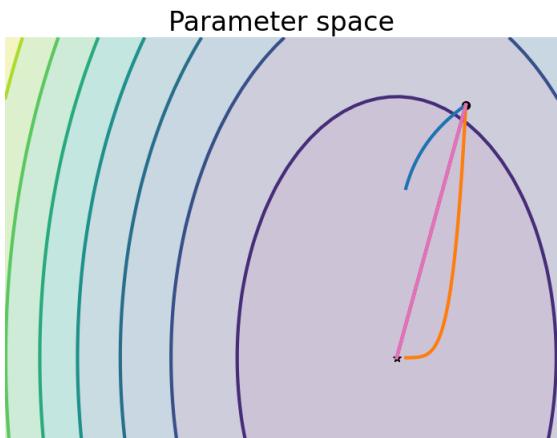
Toy example

Natural gradient and Newton method are equivalent for linear models.

$$\begin{aligned}\textcolor{blue}{u} &\in L^2([0, 1]) \\ \mathcal{L}_u(\textcolor{red}{v}) &= \frac{1}{2} \|\textcolor{blue}{u} - \textcolor{red}{v}\|_{\textcolor{blue}{K}}^2 \\ \textcolor{red}{v}_\theta(x) &= A(\theta)(x) = \theta^T \textcolor{red}{B} \Phi(x)\end{aligned}$$

$$= \theta^T \textcolor{red}{B} \left[\frac{1}{\sqrt{2} \sin(2\pi x)} \right]$$

$$\begin{aligned}\mathbb{R}^{p \times p} &\ni \textcolor{red}{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0.2 \end{bmatrix} \\ \frac{\partial A}{\partial \theta_i}(\theta)(x) &= \textcolor{red}{B}_i \Phi_i(x) \\ G_{ij} &= \int \left[\frac{\partial A}{\partial \theta_i} [\textcolor{red}{B}^T K \textcolor{blue}{B}]_{ij} \frac{\partial A}{\partial \theta_j} \right] (x) dx = [\textcolor{red}{B}^T K \textcolor{blue}{B}]_{ij}\end{aligned}$$

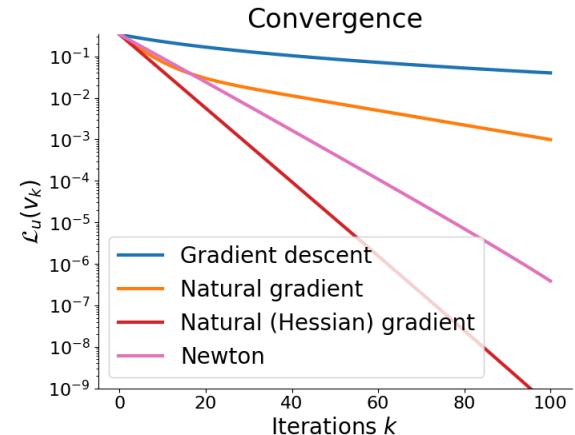
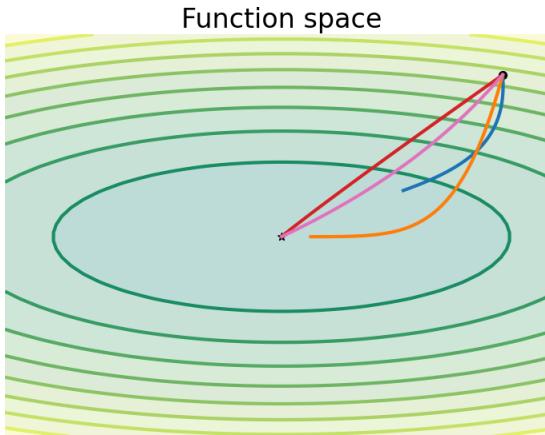
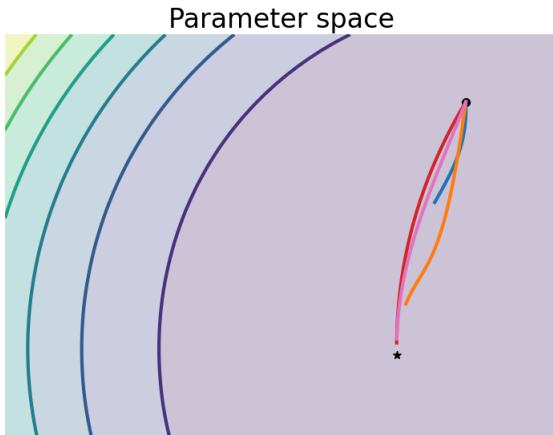


Toy example

Nonlinear manifold.

$$\begin{aligned}\textcolor{blue}{u} &\in L^2([0, 1]) \\ \mathcal{L}_u(\textcolor{red}{v}) &= \frac{1}{2} \|\textcolor{blue}{u} - \textcolor{red}{v}\|_{\textcolor{blue}{K}}^2 \\ v_\theta(x) &= A(\theta)(x) = \theta^T (\textcolor{red}{B} + \textcolor{violet}{Q}\theta)\Phi(x)\end{aligned}$$

$$\begin{aligned}\mathbb{R}^{p \times p} &\ni \textcolor{red}{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0.2 \end{bmatrix}, \textcolor{violet}{Q} \in \mathbb{R}^{p \times p \times p} \\ \frac{\partial A}{\partial \theta_i}(\theta)(x) &= (\textcolor{red}{B}_i + \textcolor{violet}{Q}_i\theta)\Phi_i(x) \\ G_{ij}(\theta) &= [(\textcolor{red}{B} + \textcolor{violet}{Q}_i\theta)^T \textcolor{blue}{K} (\textcolor{red}{B} + \textcolor{violet}{Q}_i\theta)]_{ij}\end{aligned}$$



Why we need momentum

Using momentum to **scape local minima** and **correct the biased step** of NGD.

Beyond L^2 loss

Classification
problem: KL-div

$$\mathcal{L}_u(\mathbf{v}) = \int \mathbf{v}(x) \log \frac{\mathbf{v}(x)}{\mathbf{u}(x)} dx$$

Stochastic setting
empirical loss

$$\begin{aligned} \mathcal{L}_u(\mathbf{v}_k) &= \|\mathbf{u} - \mathbf{v}_k\|_m^2 \\ &\frac{1}{2m} \sum_{i=1}^m (\mathbf{u}(x_{I_i^k}) - \mathbf{v}(x_{I_i^k}))^2 \end{aligned}$$

PDE residual

$$\begin{aligned} \mathcal{L}(\mathbf{v}) &= \|R(\mathbf{v})\|^2 \\ &= \|-\epsilon \partial_{xx} \mathbf{v} + \partial_x \mathbf{v} - 1\|^2 \end{aligned}$$

Escape local minima

$$\frac{\partial \mathbf{v}_\theta}{\partial s} = -P_{\mathcal{T}_k} \nabla \mathcal{L}$$

Non-linear model

$$\theta_{k+1} = \theta_k - s G_k^{-1} \nabla_\theta \mathcal{L}(\theta_k)$$

$$\theta_s = \theta_k - \int_0^s G_k^{-1}(\theta_\ell) \nabla_\theta \mathcal{L}(\theta_\ell) d\ell$$



Momentum dynamics

From gradient flow to momentum [Polyak, B.T. 1964] [Nesterov, Yurii. 1983].

$$\frac{d\theta}{ds} = -\nabla_{\theta}\mathcal{L}$$

Heavy-ball

$$\theta_{k+1} = \theta_k + p_k$$

$$p_k = \beta p_{k-1} - \alpha \nabla_{\theta}\mathcal{L}_u(\theta_k)$$

$$\frac{d^2\theta}{ds^2} = -\gamma \frac{d\theta}{ds} - \nabla_{\theta}\mathcal{L}$$

Nestorov

$$y_k = \theta_k + \beta(\theta_k - \theta_{k-1})$$

$$\theta_{k+1} = y_k - \alpha \nabla_{\theta}\mathcal{L}_u(y_k)$$

$$\theta_{k+1} = \theta_k + \beta p_{k-1} - \alpha \nabla_{\theta}\mathcal{L}_u(\theta_k)$$

$$\theta_{k+1} = \theta_k + \beta(\theta_k - \theta_{k-1}) - \alpha \nabla_{\theta}\mathcal{L}_u(y_k)$$

Momentum dynamics in functional space

From momentum in parameter space to functional space.

Heavy-ball

$$\theta_{k+1} - \theta_k = \beta p_{k-1} - \alpha \nabla_\theta \mathcal{L}_u(\theta_k)$$

Nestorov

$$\theta_{k+1} - \theta_k = \beta(\theta_k - \theta_{k-1}) - \alpha \nabla_\theta \mathcal{L}_u(y_k)$$

$$\mathcal{T}_k \ni \mathbf{d}v_k = P_{\mathcal{T}_k}[\beta p_{k-1} - \alpha \nabla \mathcal{L}(v_k)]$$

$$\mathcal{T}_k \ni \mathbf{d}v_k = P_{\mathcal{T}_k}[\beta(v_k - v_{k-1}) - \alpha \nabla \mathcal{L}(v_k)]$$

$$P_{\mathcal{T}_k} \nabla \mathcal{L}(v_k)$$

$$P_{\mathcal{T}_k} p_{k-1}$$

$$P_{\mathcal{T}_k} (v_k - v_{k-1})$$

Momentum dynamics in functional space

From momentum in parameter space to functional space.

Heavy-ball

Nestorov

$$\theta_{k+1} - \theta_k = \beta p_{k-1} - \alpha \nabla_{\theta} \mathcal{L}_u(\theta_k)$$

$$\theta_{k+1} - \theta_k = \beta(\theta_k - \theta_{k-1}) - \alpha \nabla_{\theta} \mathcal{L}_u(y_k)$$

$$\mathcal{T}_k \ni \mathbf{d}v_k^{HB} = P_{\mathcal{T}_k}[\beta p_{k-1} - \alpha \nabla \mathcal{L}(\mathbf{v}_k)]$$

$$\mathcal{T}_k \ni \mathbf{d}v_k^N = P_{\mathcal{T}_k}[\beta(\mathbf{v}_k - \mathbf{v}_{k-1}) - \alpha \nabla \mathcal{L}(\mathbf{v}_k)]$$

$$P_{\mathcal{T}_k} \nabla \mathcal{L}(\mathbf{v}_k)$$

$$P_{\mathcal{T}_k} p_{k-1}$$

$$P_{\mathcal{T}_k} (\mathbf{v}_k - \mathbf{v}_{k-1})$$

$$G_k^{-1} \nabla_{\theta} \mathcal{L}(\theta_k)$$

$$G_k^{-1} G_{k,k-1} p_{k-1}$$

$$G_k^{-1} \int \left[\frac{\partial A}{\partial \theta} \Bigg|_{\theta_k} (\mathbf{v}_k - \mathbf{v}_{k-1}) \right] (x) dx$$

$$G_{k,k-1} = \int \left[\frac{\partial A}{\partial \theta} \Bigg|_{\theta_k} \frac{\partial A}{\partial \theta} \Bigg|_{\theta_{k-1}} \right] (x) dx$$

Momentum dynamics in functional space

From momentum in parameter space to functional space.

Heavy-ball

Nestorov

$$\theta_{k+1} - \theta_k = \beta p_{k-1} - \alpha \nabla_{\theta} \mathcal{L}_u(\theta_k)$$

$$\theta_{k+1} - \theta_k = \beta(\theta_k - \theta_{k-1}) - \alpha \nabla_{\theta} \mathcal{L}_u(y_k)$$

$$\mathcal{T}_k \ni \mathbf{d}v_k^{HB} = \beta P_{\mathcal{T}_k} \mathbf{p}_{k-1} - \alpha P_{\mathcal{T}_k}^{H_{\mathcal{L}}} \nabla \mathcal{L}(\mathbf{v}_k) \quad \mathcal{T}_k \ni \mathbf{d}v_k^N = \beta P_{\mathcal{T}_k} (\mathbf{v}_k - \mathbf{v}_{k-1}) - \alpha P_{\mathcal{T}_k}^{H_{\mathcal{L}}} \nabla \mathcal{L}(\mathbf{v}_k)$$

$$P_{\mathcal{T}_k}^{H_{\mathcal{L}}} \nabla \mathcal{L}(\mathbf{v}_k)$$

$$P_{\mathcal{T}_k} \mathbf{p}_{k-1}$$

$$P_{\mathcal{T}_k} (\mathbf{v}_k - \mathbf{v}_{k-1})$$

$$G_k^{-1} \nabla_{\theta} \mathcal{L}(\theta_k) \quad G_k^{-1} G_{k,k-1} p_{k-1} \quad G_k^{-1} \int \left[\frac{\partial A}{\partial \theta} \Bigg|_{\theta_k} (\mathbf{v}_k - \mathbf{v}_{k-1}) \right] (x) dx$$
$$G_k = \int \left[\frac{\partial A}{\partial \theta} H_{\mathcal{L}} \frac{\partial A}{\partial \theta} \right] (x) dx \quad G_{k,k-1} = \int \left[\frac{\partial A}{\partial \theta} \Bigg|_{\theta_k} \frac{\partial A}{\partial \theta} \Bigg|_{\theta_{k-1}} \right] (x) dx$$

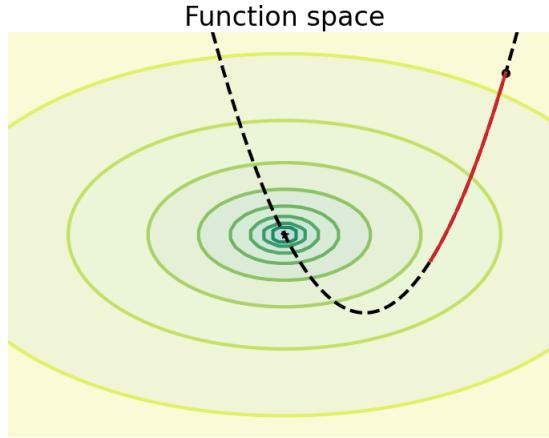
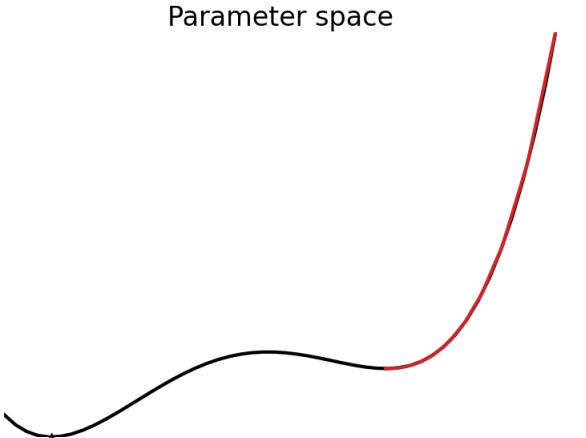
Toy example

Escaping local minima.

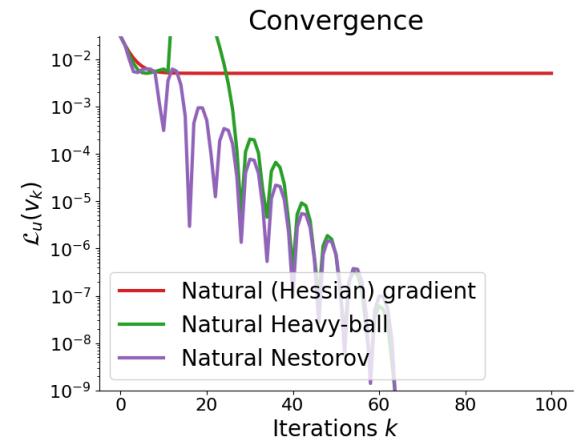
$$\textcolor{blue}{u} \in L^2([0, 1])$$

$$\mathcal{L}_u(\textcolor{red}{v}) = \frac{1}{2} \|\textcolor{blue}{u} - \textcolor{red}{v}\|_{\textcolor{blue}{K}}^2$$

$$v_{\theta}(x) = \theta_1 \textcolor{red}{b}^T \Phi(x) + \theta_1^2 b^{\perp T} \Phi(x)$$



$$\begin{aligned} \mathbf{d}v_k^{HB} &= P_{\mathcal{T}_k}[\beta \textcolor{violet}{p}_{k-1} - \alpha \nabla \mathcal{L}(\textcolor{red}{v}_k)] \\ \mathbf{d}v_k^N &= P_{\mathcal{T}_k}[\beta(\textcolor{red}{v}_k - \textcolor{violet}{v}_{k-1}) - \alpha \nabla \mathcal{L}(\textcolor{red}{v}_k)] \end{aligned}$$



Toy example

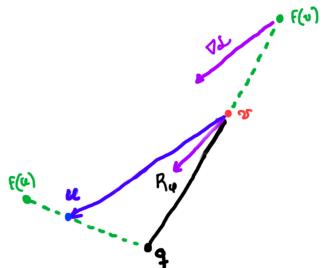
Not L^2 loss.

$$\mathbf{u} \in L^2([0, 1])$$

$$\mathcal{L}_u(\mathbf{v}) = \frac{1}{2} \|f(\mathbf{u}) - f(\mathbf{v})\|_K^2$$

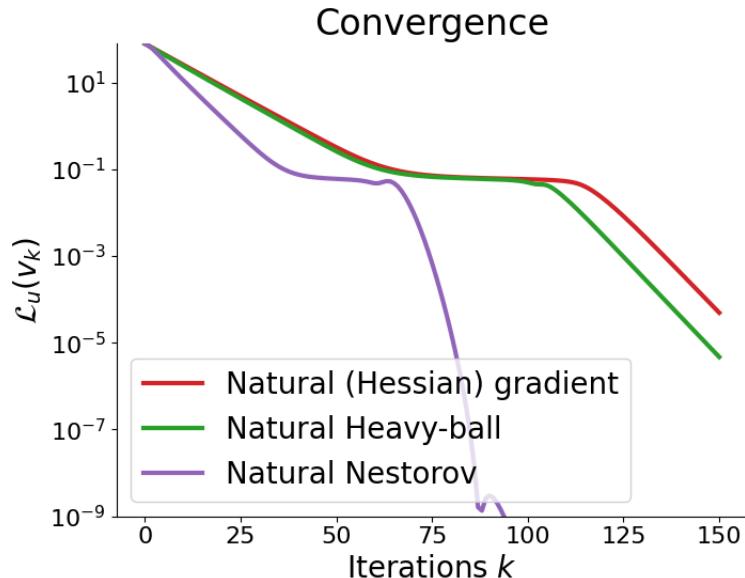
$$f(\mathbf{v}) = (1 + \omega \|\mathbf{v} - \mathbf{q}\|^2)(\mathbf{v} - \mathbf{q}) + \mathbf{q}$$

$$\mathbf{q} = R_\varphi(\mathbf{u} - \mathbf{v}) + \mathbf{v}$$



$$\mathbf{d}v_k^{HB} = P_{\mathcal{T}_k}[\beta \mathbf{p}_{k-1} - \alpha \nabla \mathcal{L}(\mathbf{v}_k)]$$

$$\mathbf{d}v_k^N = P_{\mathcal{T}_k}[\beta(\mathbf{v}_k - \mathbf{v}_{k-1}) - \alpha \nabla \mathcal{L}(\mathbf{v}_k)]$$



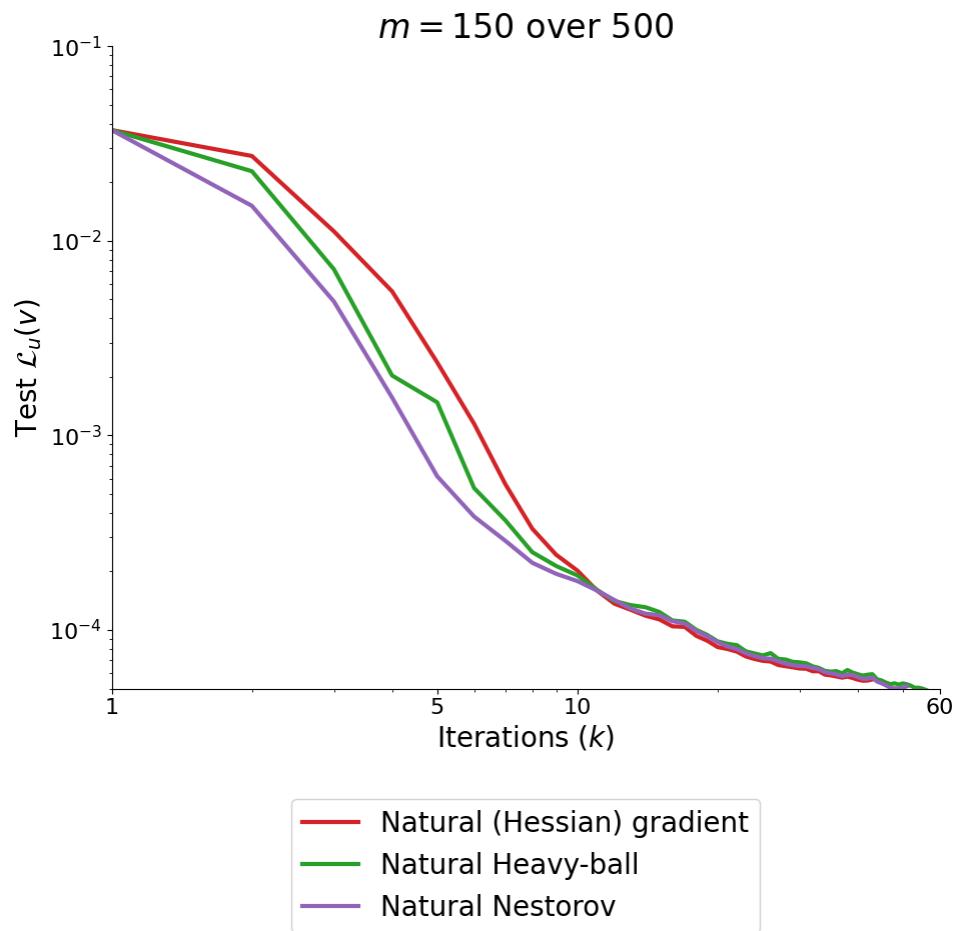
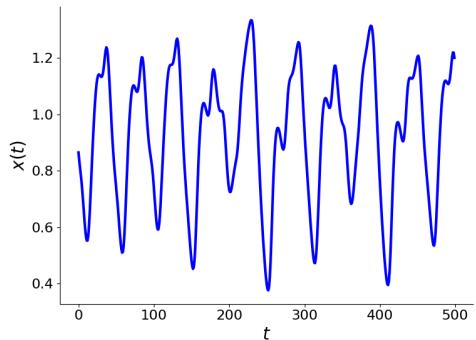
Mackey Glass

A less toy example [Park, H, S.-I Amari, and K Fukun (2000)].

Mackey Glass caotic time series:

- $z(t + 1) = (1 - b)z(t) + a \frac{z(t-\tau)}{1+z(t-\tau)^{10}}$
- Input: $z(t), z(t - 6), z(t - 12), z(t - 18)$
- Output: $z(t + 6)$

Model (shallow NN 10 neurons): $v_\theta : \mathbb{R}^4 \rightarrow \mathbb{R}$



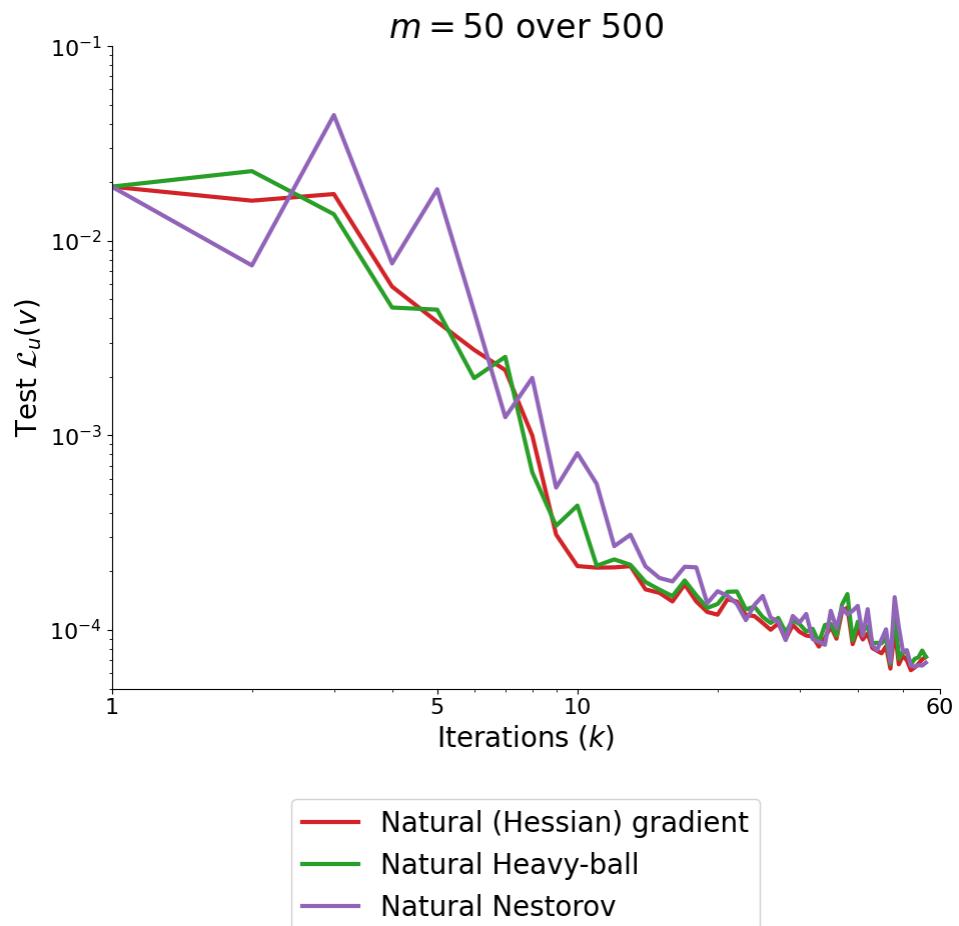
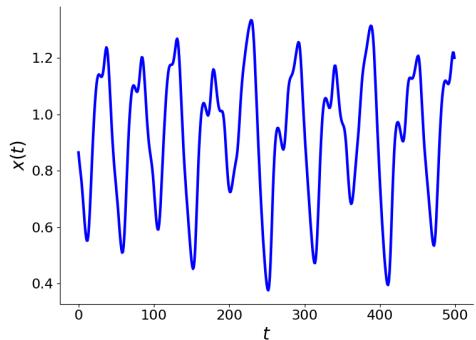
Mackey Glass

A less toy example [Park, H, S.-I Amari, and K Fukun (2000)].

Mackey Glass caotic time series:

- $z(t + 1) = (1 - b)z(t) + a \frac{z(t-\tau)}{1+z(t-\tau)^{10}}$
- Input: $z(t), z(t - 6), z(t - 12), z(t - 18)$
- Output: $z(t + 6)$

Model (shallow NN 10 neurons): $v_\theta : \mathbb{R}^4 \rightarrow \mathbb{R}$



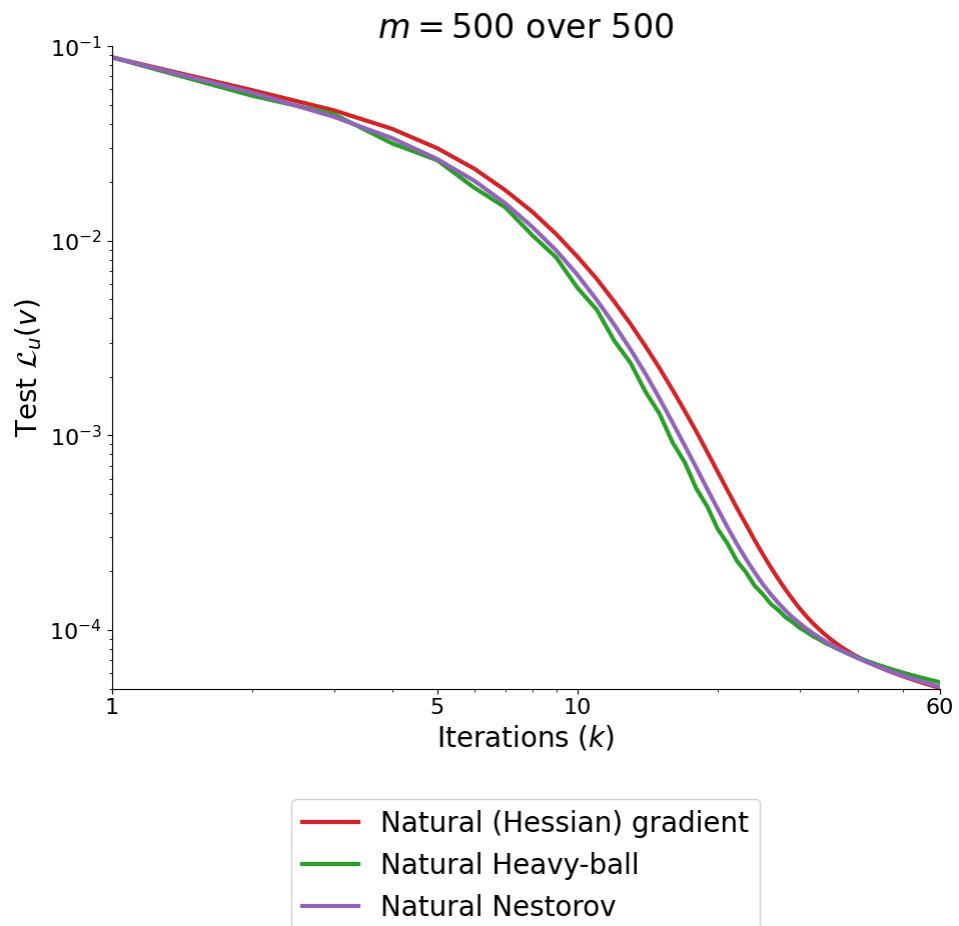
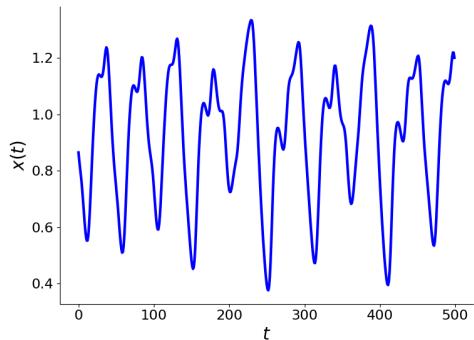
Mackey Glass

A less toy example [Park, H, S.-I Amari, and K Fukun (2000)].

Mackey Glass caotic time series:

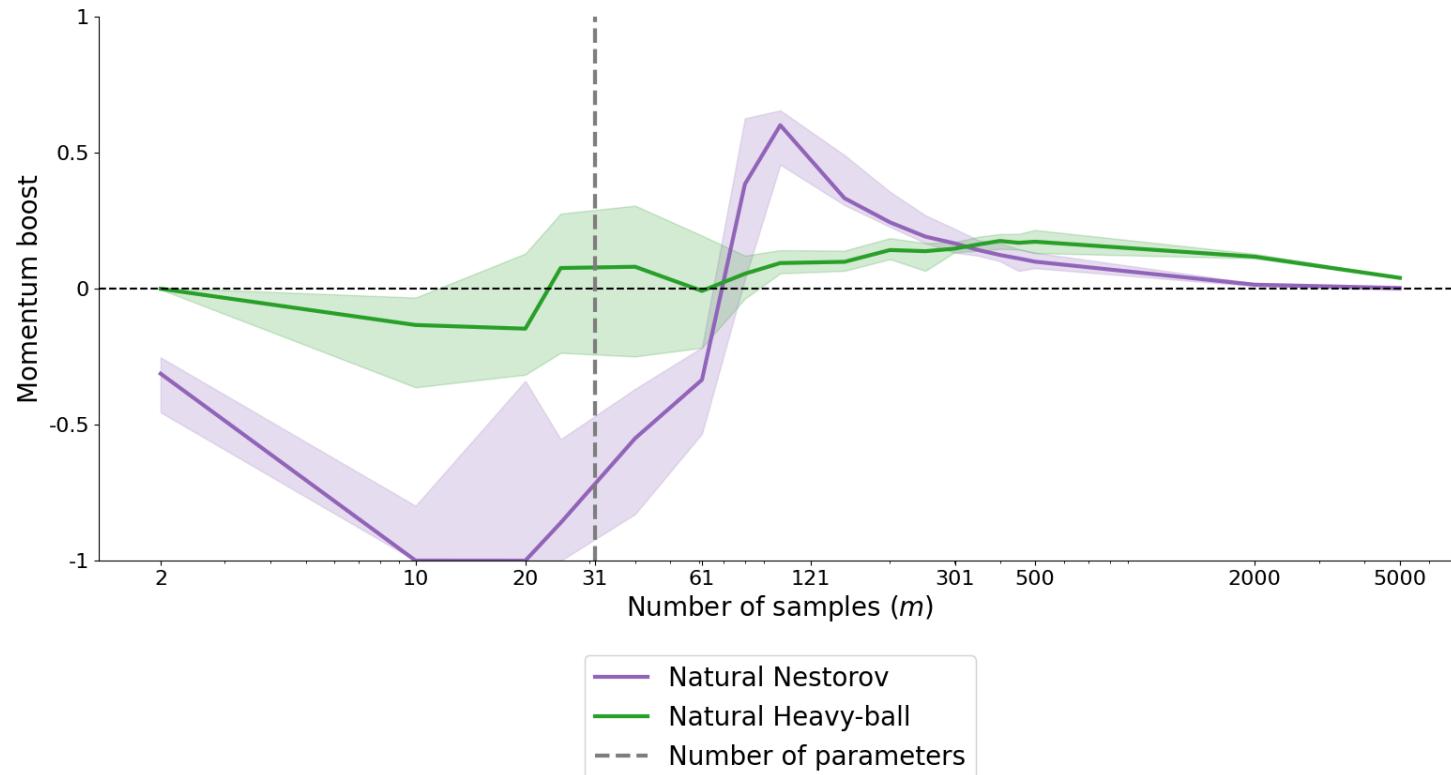
- $z(t + 1) = (1 - b)z(t) + a \frac{z(t-\tau)}{1+z(t-\tau)^{10}}$
- Input: $z(t), z(t - 6), z(t - 12), z(t - 18)$
- Output: $z(t + 6)$

Model (shallow NN 10 neurons): $v_\theta : \mathbb{R}^4 \rightarrow \mathbb{R}$



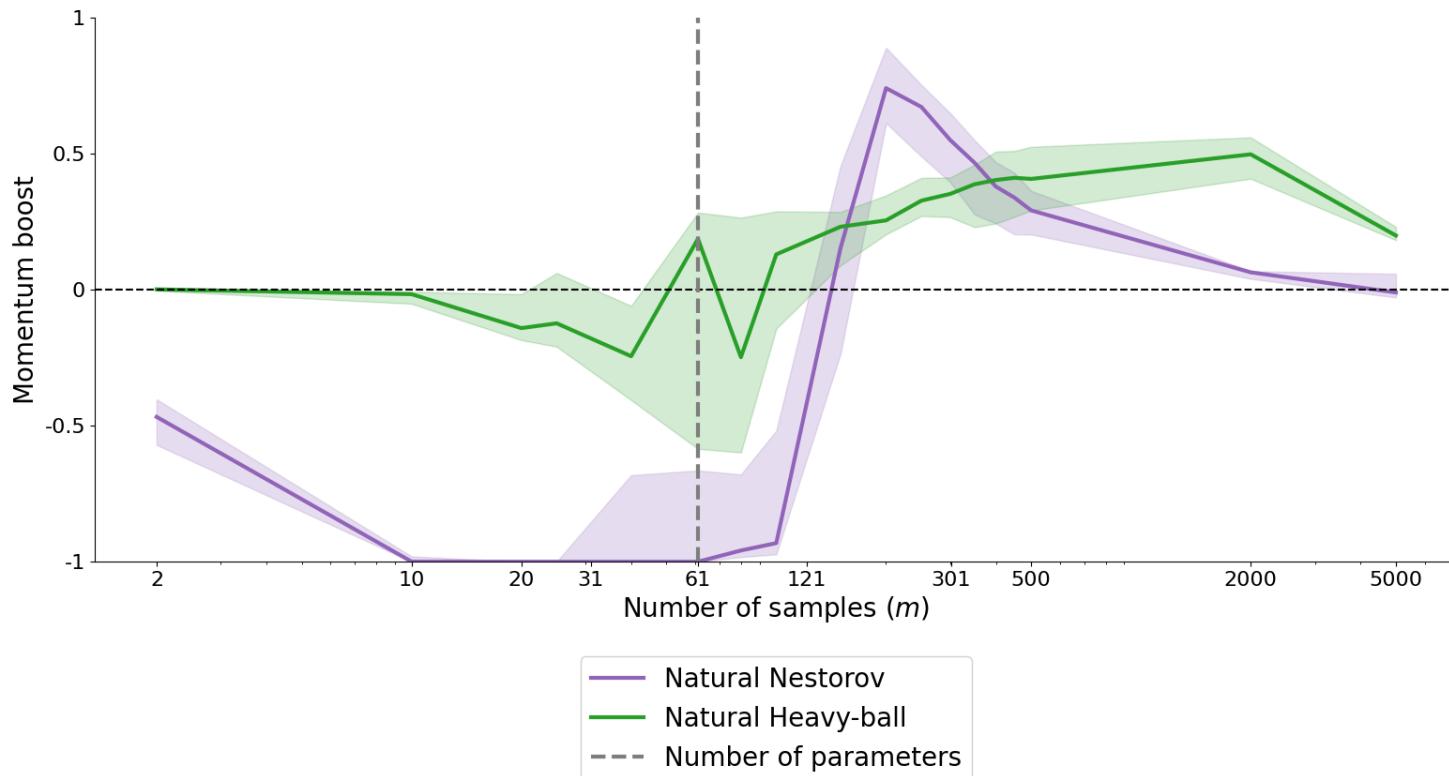
Mackey Glass: sampling

Effect of natural gradient direction estimation with limited number of samples.



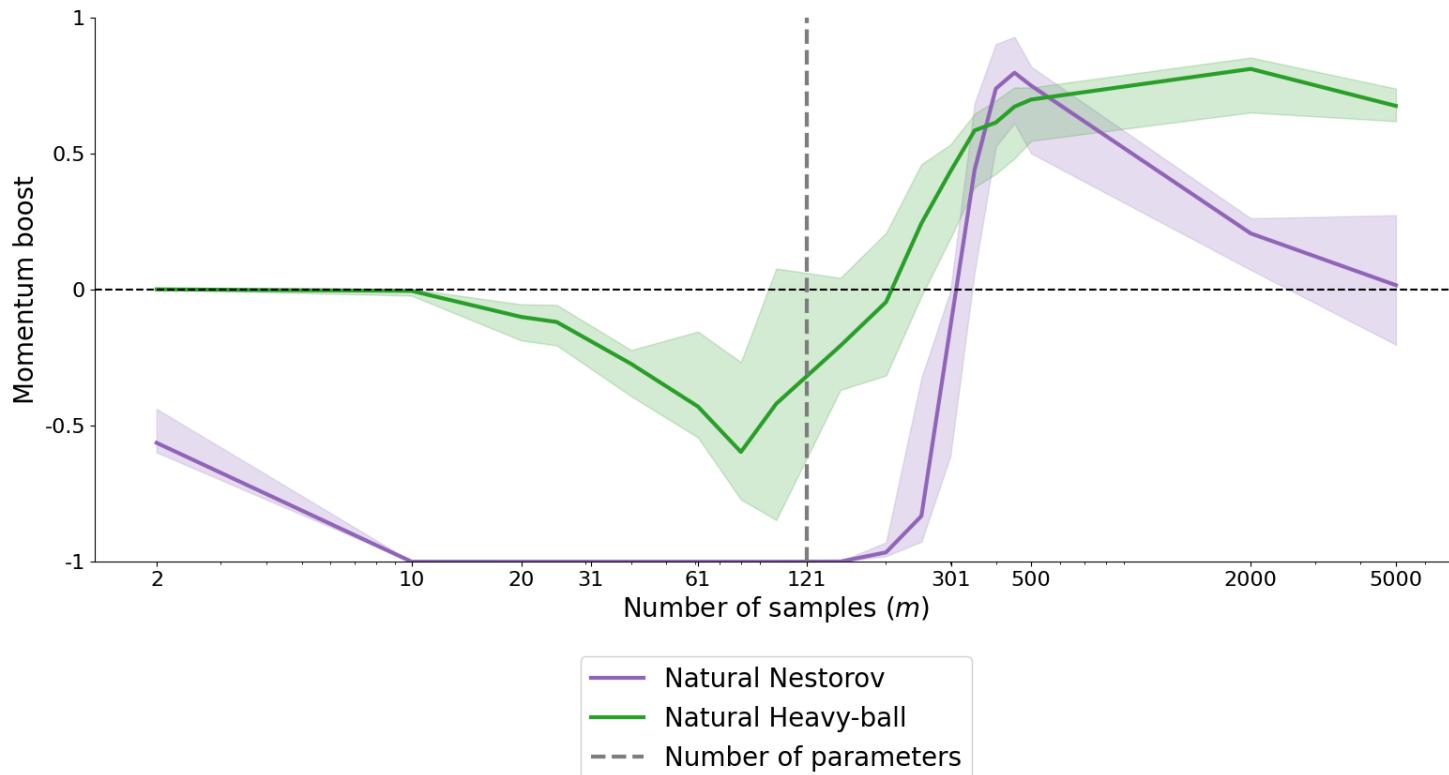
Mackey Glass: sampling

Effect of natural gradient direction estimation with limited number of samples.



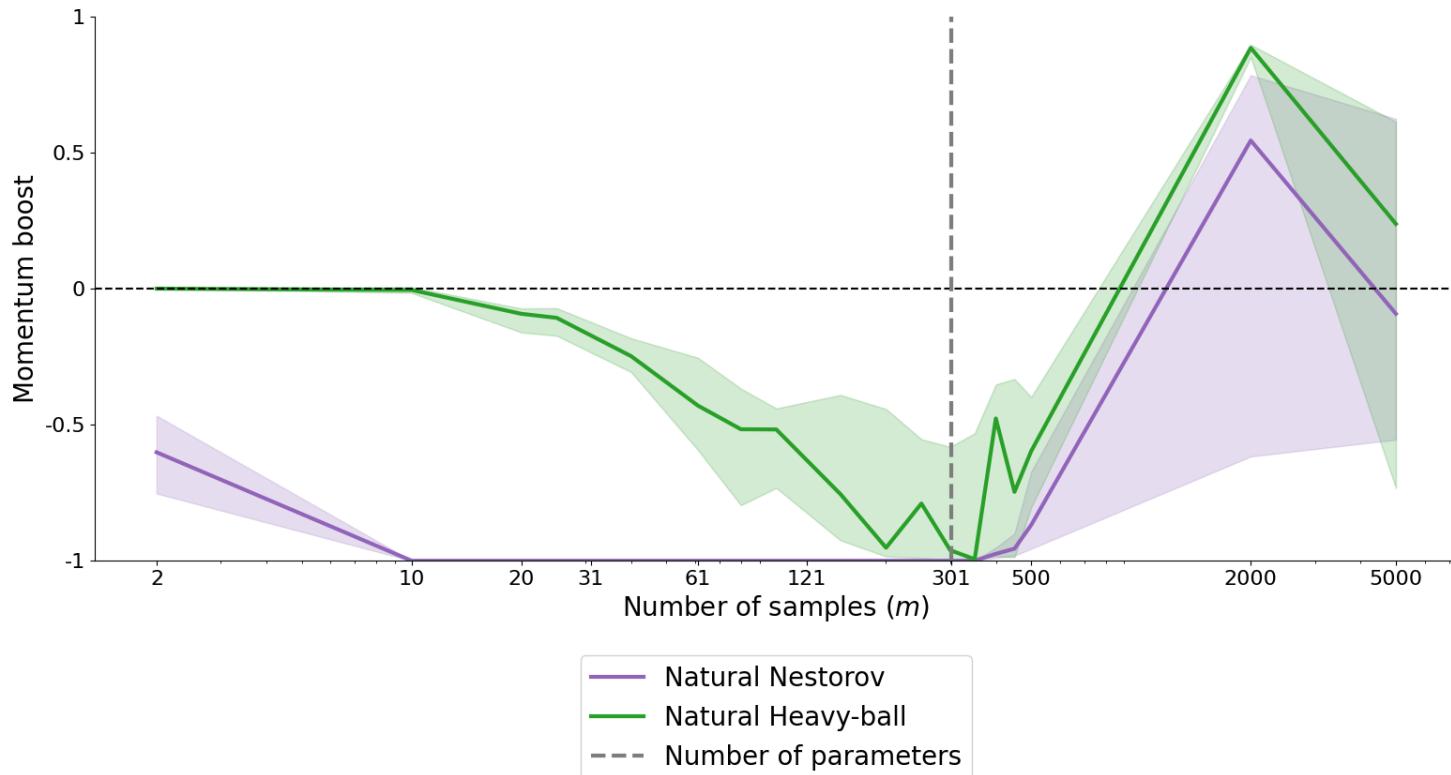
Mackey Glass: sampling

Effect of natural gradient direction estimation with limited number of samples.



Mackey Glass: sampling

Effect of natural gradient direction estimation with limited number of samples.



Classification task

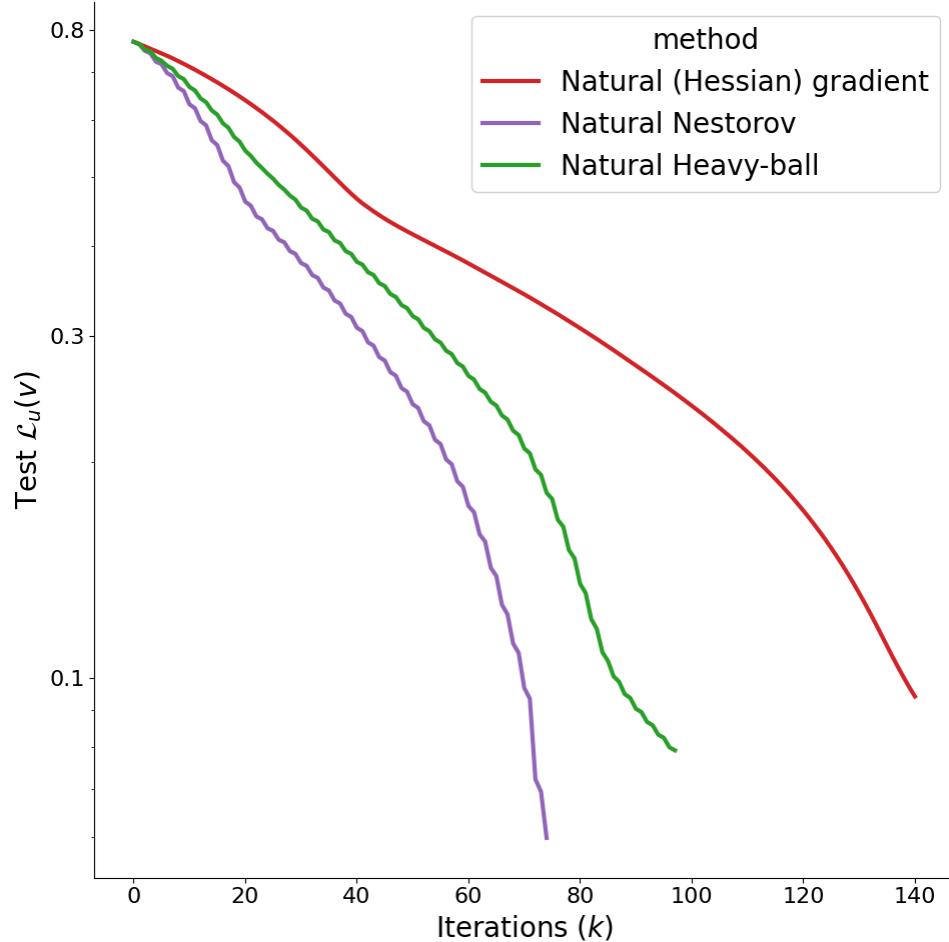
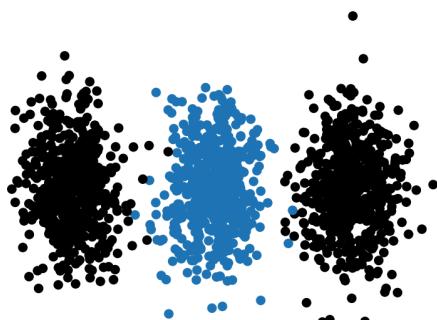
Another less toy example from [Park, H, S.-I Amari, S. Ito, K Fukumizu (2000)].

Mackey Glass caotic time series:

- Input: $z \in \mathbb{R}^2$
- Output: $c \in \{0, 1\}$

$$\mathcal{L}_u(\mathbf{v}) = \int -u \log(v) - (1-u) \log(1-v)$$

Model (Shallow NN): $\mathbf{v}_\theta : \mathbb{R}^2 \rightarrow [0, 1]$



Classification task

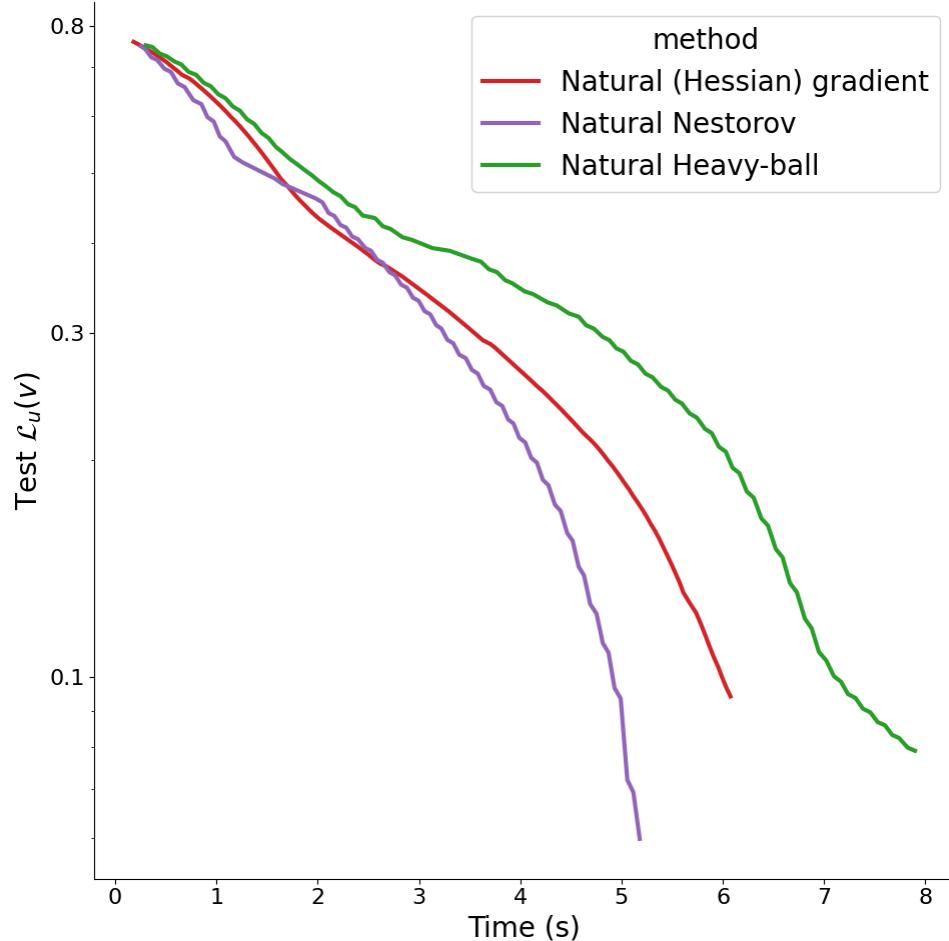
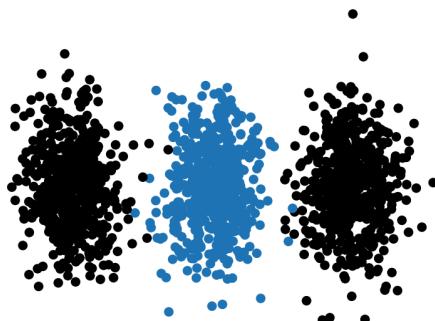
Another less toy example from [Park, H, S.-I Amari, and K Fukumizu (2000)].

Mackey Glass caotic time series:

- Input: $z \in \mathbb{R}^2$
- Output: $c \in \{0, 1\}$

$$\mathcal{L}_u(v) = -u \log(v) - (1-u) \log(1-v)$$

Model (Shallow NN): $v_\theta : \mathbb{R}^2 \rightarrow [0, 1]$





Thank you!

(2024) R. Gruhlke, A. Nouy, P. Trunschke.	NGD and optimal sampling.
(2020) J. Martens.	NGD Review.
(2000) H. Park, S. Amari, K. Fukumizu.	NGD Experiments.
(1998) S. Amari.	NGD Initial paper.
(1983) Y. Nesterov.	Nestorov acceleration.
(1964) B.T. Polyak.	Heavy Ball acceleration.