



First Optimize, Then Discretize for SciML

Marius Zeinhofer

Workshop on Artificial Intelligence for Theoretical Sciences
at the Kavli Institute of Theoretical Sciences, Beijing, China

ETH Zürich, Seminar for Applied Mathematics,
Switzerland

November 14, 2024

Today's Talk

Part I – Function Space Optimization for SciML

Discretize function space algorithms in the tangent space of a neural network ansatz¹

- Develop optimizers that provably mimic function space dynamics
- Design new algorithms and understand existing ones within this framework

Part II – Function Space Optimization at Scale

Scale the function space algorithms to millions of parameters²

- Leverage the Kronecker structure of linear layer Jacobians to approximate curvature matrices

¹J. Müller and M. Zeinhofer (2024). "Position: Optimization in SciML Should Employ the Function Space Geometry". In: *Forty-first International Conference on Machine Learning*.

²F. Dangel, J. Müller, and M. Zeinhofer (2024). "Kronecker-Factored Approximate Curvature for Physics-Informed Neural Networks". In: *Advances in Neural Information Processing Systems*.

Table of Contents

1. Function Space Optimization for SciML

2. Scalability

3. Numerical Experiments

4. Summary

Scientific Machine Learning

When referring to Scientific Machine Learning (SciML) we mean:

- Operator Learning³
- Variational Monte Carlo for Schrödinger's equation with NN ansatz⁴
- PINNs, Deep Ritz, etc.⁵
- Anything involving PDEs and neural networks

Observation

Introducing PDEs to objective functions complicates neural network training drastically.

³Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar (2021). "Fourier Neural Operator for Parametric Partial Differential Equations". In: *International Conference on Learning Representations*.

⁴D. Pfau, J. S. Spencer, A. G. Matthews, and W. M. C. Foulkes (2020). "Ab Initio Solution of the Many-Electron Schrödinger Equation with Deep Neural Networks". In: *Physical Review Research* 2.3, p. 033429.

⁵M. Raissi, P. Perdikaris, and G. E. Karniadakis (2019). "Physics-Informed Neural Networks: A Deep Learning Framework for Solving Forward and Inverse Problems Involving Nonlinear Partial Differential Equations". In: *Journal of Computational Physics* 378, pp. 686–707.

Optimization in SciML is a Challenge

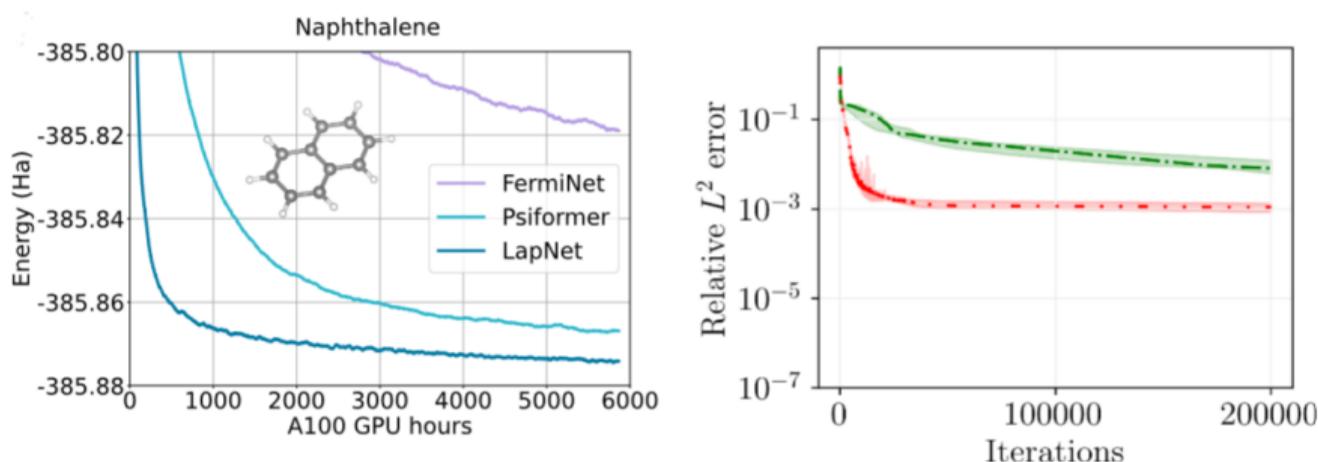


Figure: Illustration of neural network wavefunction optimization⁶ and PINN optimization⁷.

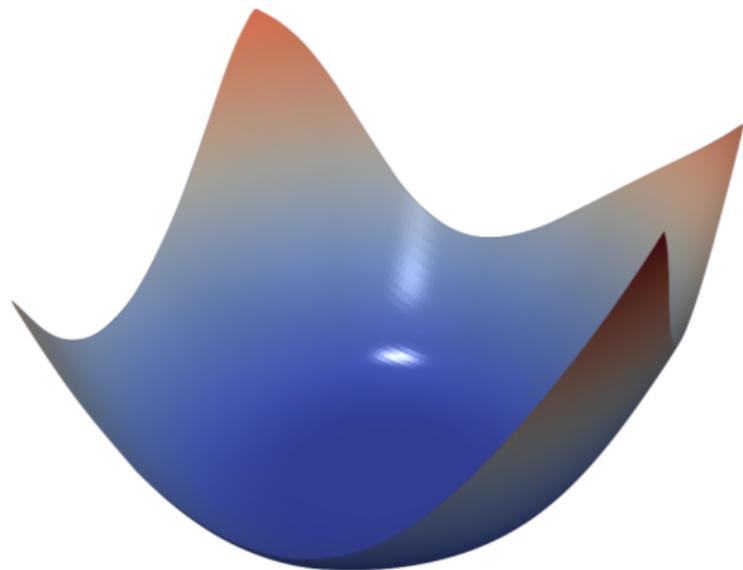
Proposition

Abandon first-order optimization. Exploit geometry in function space.

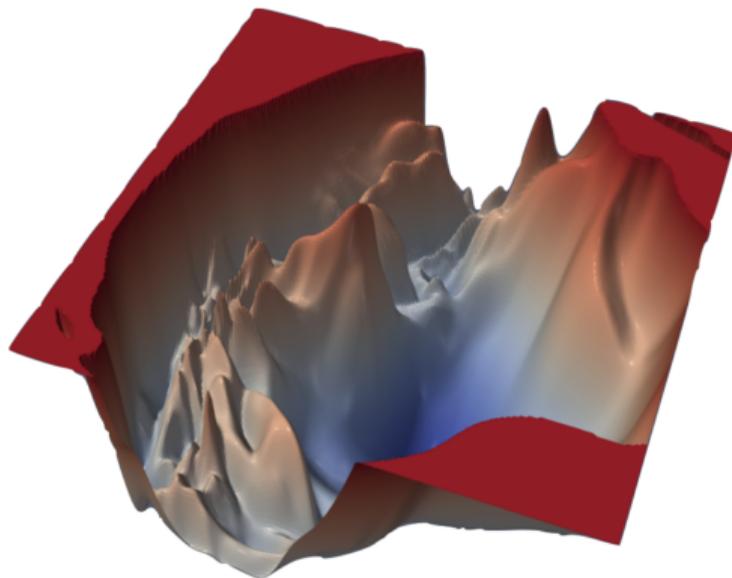
⁶R. Li et al. (2024). "A computational framework for neural network-based variational Monte Carlo with Forward Laplacian". In: *Nature Machine Intelligence*, pp. 1–11.

⁷J. Müller and M. Zeinhofer (2023). "Achieving High Accuracy with PINNs via Energy Natural Gradients". In: *International Conference on Machine Learning*.

Geometry in Function Space vs Parameter Space⁸



$u \mapsto E(u)$
“Function Space”



$\theta \mapsto L(\theta) = E(u_\theta)$
“Parameter Space”

⁸H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein (2018). “Visualizing the loss landscape of neural nets”. In: *Advances in neural information processing systems* 31.

Abstract Function Space Optimization⁹

Informal Roadmap: “First Optimize, Then Discretize”

Given: “a SciML problem” and a neural network ansatz. We propose:

- (i) **Continuous Formulation:** Formulate minimization/saddle point/... problem in a Hilbert space \mathcal{H}

$$E : \mathcal{H} \rightarrow \mathbb{R}.$$

- (ii) **Optimize:** Decide for an appropriate iterative algorithm *in function space* \mathcal{H}

$$u_{k+1} = u_k + d_k, \quad k = 0, 1, 2, \dots$$

- (iii) **Discretize:** Project update directions on the tangent space of neural network ansatz.

⁹J. Müller and M. Zeinhofer (2024). “Position: Optimization in SciML Should Employ the Function Space Geometry”. In: *Forty-first International Conference on Machine Learning*.

Recap

Newton's Method for Minimization in Function Spaces

For (general) objective $E : \mathcal{H} \rightarrow \mathbb{R}$ do second-order Taylor expansion (in function space)

$$E(u + d) \approx E(u) + DE(u)d + \frac{1}{2}D^2E(u)(d, d).$$

Explicitly minimize in the variable d :

$$d = u - D^2E(u)^{-1}[DE(u)]$$

- Think of the expression $D^2E(u)^{-1}[DE(u)]$ as solving a PDE¹⁰.

¹⁰The Euler-Lagrange equations of the Taylor expansion of E .

Newton: Continuous Formulation

- Suppose we aim to solve Poisson's equation

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= g && \text{on } \partial\Omega \end{aligned}$$

- Reformulate as a minimization problem¹¹

$$\min_{u \in H^2(\Omega)} E(u) = \frac{1}{2} \|\Delta u + f\|_{L^2(\Omega)}^2 + \frac{1}{2} \|u - g\|_{L^2(\partial\Omega)}^2$$

- Ignore neural network ansatz classes.

¹¹A formulation as a root finding problem or a variational formulation make sense, too.

Newton: Optimize

- For the quadratic minimization problem

$$u^* = \operatorname{argmin}_{u \in H^2(\Omega)} E(u) = \frac{1}{2} \|\Delta u + f\|_{L^2(\Omega)}^2 + \frac{1}{2} \|u - g\|_{L^2(\partial\Omega)}^2$$

- Newton's method in function space has one step convergence

$$\begin{aligned} u^* &= u - D^2 E(u)^{-1} [DE(u)] \\ &= u + \underbrace{(u^* - u)}_{=d}. \end{aligned}$$

Goal

Discretize such that we follow the optimal step $d = u^* - u$?

Newton: Discretize I – Bilinear Forms

Question

How to discretize the Newton step $D^2E(u)^{-1}[DE(u)]$?

- $D^2E(u) : H^2(\Omega) \times H^2(\Omega) \rightarrow \mathbb{R}$ is a *bilinear form*:

$$D^2E(u)(v, w) = \int_{\Omega} \Delta v \Delta w \, dx + \int_{\partial\Omega} v w \, ds$$

- $DE(u) : H^2(\Omega) \rightarrow \mathbb{R}$ is a *linear functional*:

$$DE(u)(v) = \int_{\Omega} (\Delta u + f) \Delta v \, dx + \int_{\partial\Omega} (u - g) v \, ds$$

- The expression $D^2E(u)^{-1}(DE(u))$ means: Find v such that

$$D^2E(u)(v, w) = DE(u)(w), \quad \text{for all } w \in H^2(\Omega)$$

Newton: Discretize II – Galerkin Methods

Question

How to discretize an equation of the form: $D^2E(u)(v, w) = DE(u)(w)$, for all w ?

- Choose basis functions $v_1, \dots, v_n \in H^2(\Omega)$ and transfer to a matrix equation:

$$G_{ij} = D^2E(u)(v_j, v_i), \quad b_i = DE(u)(v_i)$$

- Solve the matrix equation

$$G\alpha = b$$

- Under suitable conditions, $v = \sum_{i=1}^n \alpha_i v_i$ is an approximate solution.

Naming convention

Discretizing a bilinear form with a finite dimensional vector space is a *Galerkin Method*.

Newton: Discretize III – Galerkin in Tangent Space

- Given a neural network ansatz $\{u_\theta \mid \theta \in \Theta\}$, as basis function choose

$$\partial_{\theta_1} u_\theta, \dots, \partial_{\theta_p} u_\theta$$

- Discretize with a Galerkin ansatz using the basis functions $\partial_{\theta_1} u_\theta, \dots, \partial_{\theta_p} u_\theta$:

$$G(\theta)_{ij} = D^2 E(u_\theta)(\partial_{\theta_j} u_\theta, \partial_{\theta_i} u_\theta), \quad \nabla L(\theta)_i = DE(u_\theta)(\partial_{\theta_i} u_\theta)$$

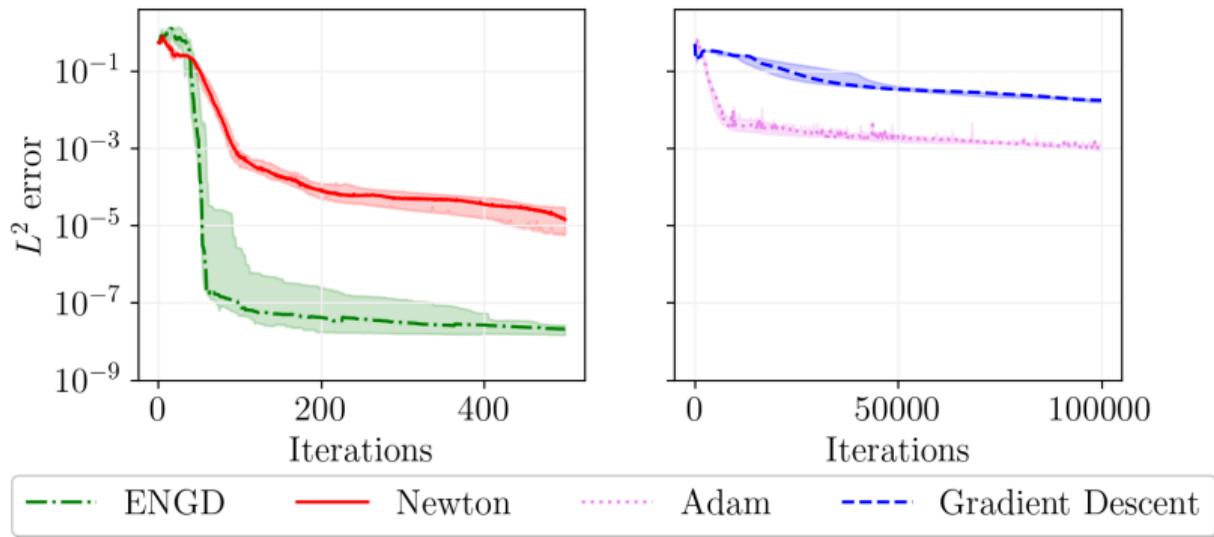
- This yields an optimization algorithm:

$$\theta_{k+1} = \theta_k - \eta_k G(\theta_k)^\dagger \nabla L(\theta_k).$$

- If $E(u) = \frac{1}{2} \|\Delta u + f\|_{L^2(\Omega)}^2 + \frac{1}{2} \|u - g\|_{L^2(\partial\Omega)}^2$ we get:

$$G(\theta)_{ij} = \int_{\Omega} \Delta \partial_{\theta_i} u_\theta \Delta \partial_{\theta_j} u_\theta \, dx + \int_{\partial\Omega} \partial_{\theta_i} u_\theta \partial_{\theta_j} u_\theta \, ds$$

This Yields a Highly Efficient Method¹²



$$\theta_{k+1} = \theta_k - \eta_k G(\theta_k)^\dagger \nabla L(\theta_k), \quad k = 0, 1, 2, \dots$$
$$G(\theta)_{ij} = \int_{\Omega} \Delta \partial_{\theta_i} u_{\theta} \Delta \partial_{\theta_j} u_{\theta} dx + \int_{\partial\Omega} \partial_{\theta_i} u_{\theta} \partial_{\theta_j} u_{\theta} ds$$

¹²J. Müller and M. Zeinhofer (2023). "Achieving High Accuracy with PINNs via Energy Natural Gradients". In: *International Conference on Machine Learning*.

L^2 Gradient Flow for Rayleigh quotient: Cont. Form

- Suppose we want to minimize the Rayleigh quotient

$$R(\psi) = \frac{\langle H\psi, \psi \rangle}{\langle \psi, \psi \rangle},$$

where H is the Hamiltonian (or any symmetric PDE operator).

- Rewrite the problem using $\tilde{\psi} = \psi / \|\psi\|_{L^2}$ as

$$\tilde{R}(\tilde{\psi}) = \langle H\tilde{\psi}, \tilde{\psi} \rangle$$

- The transformation $\psi / \|\psi\|_{L^2}$ is for mathematical convenience.

L^2 Gradient Flow for Rayleigh quotient: Optimize

- Do L^2 gradient flow with explicit Euler discretization on \tilde{R}

$$\tilde{\psi}_{k+1} = \tilde{\psi}_K - \eta_k T^{-1}(D\tilde{R}(\tilde{\psi}_k)),$$

- $T : L^2 \rightarrow (L^2)^*$ is the Riesz isometry¹³ of L^2

$$T(\psi)(\phi) = \langle \psi, \phi \rangle$$

- $D\tilde{R}(\tilde{\psi})\phi = \langle H\tilde{\psi}, \phi \rangle$ is the Fréchet (functional) derivative of \tilde{R} .

¹³Can be viewed & discretized as bilinear form

L^2 Gradient Flow for Rayleigh quotient: Discretize

- The neural network ansatz and tangent space/Galerkin space are

$$\{\psi_\theta / \|\psi_\theta\|_{L^2} \mid \theta \in \Theta\}, \quad \text{span}\left\{\partial_{\theta_1} \left(\frac{\psi_\theta}{\|\psi_\theta\|}\right), \dots, \partial_{\theta_P} \left(\frac{\psi_\theta}{\|\psi_\theta\|}\right)\right\}$$

- Discretizing the Riesz map T yields

$$F_{ij} = T(\partial_{\theta_i}(\psi_\theta / \|\psi_\theta\|))(\partial_{\theta_j}(\psi_\theta / \|\psi_\theta\|)) = \langle \partial_{\theta_i}(\psi_\theta / \|\psi_\theta\|), \partial_{\theta_j}(\psi_\theta / \|\psi_\theta\|) \rangle$$

- Discretizing the functional derivative $D\tilde{R}$ yields the gradient of the loss
 $L(\theta) = R(\psi_\theta)$
- The final algorithm is $\theta_{k+1} = \theta_k - \eta_k F^\dagger \nabla L(\theta_k)$

Fact

One can compute that this is *stochastic reconfiguration*.

Translation: Function Space to Parameter Space

We can carry out the previous reasoning for different function space optimization methods¹⁴

Function Space	Parameter Space	Name in Literature
Newton	Generalized Gauss-Newton	ENGD ¹⁵
Gradient Descent	Natural Gradient Descent ¹⁶	—
Lagrange-Newton	Competitive Gradient Descent ¹⁷	CPINNs ¹⁸
Gauss-Newton	Gauss-Newton	GNNGD ¹⁹

¹⁴J. Müller and M. Zeinhofer (2024). “Position: Optimization in SciML Should Employ the Function Space Geometry”. In: *Forty-first International Conference on Machine Learning*.

¹⁵J. Müller and M. Zeinhofer (2023). “Achieving High Accuracy with PINNs via Energy Natural Gradients”. In: *International Conference on Machine Learning*.

¹⁶S.-I. Amari (1998). “Natural Gradient Works Efficiently in Learning”. In: *Neural Computation* 10.2, pp. 251–276.

¹⁷F. Schäfer and A. Anandkumar (2019). “Competitive gradient descent”. In: *Advances in Neural Information Processing Systems* 32.

¹⁸Q. Zeng, S. H. Bryngelson, and F. T. Schaefer (2023). “Competitive Physics Informed Networks”. In: *International Conference on Learning Representations*.

¹⁹A. Jnini, F. Vella, and M. Zeinhofer (2024). “Gauss-Newton Natural Gradient Descent for Physics-Informed Computational Fluid Dynamics”. In: *arXiv preprint arXiv:2402.10680*.

A General Class of Minimization Algorithms

- Many optimization methods are of the form:

$$u_{k+1} = u_k + T_{u_k}^{-1}[DE(u_k)]$$

where T_k is a linear, continuous, bijective map $T_{u_k} : \mathcal{H} \rightarrow \mathcal{H}^*$.

- Gradient Descent, Newton, Gauss-Newton, Lagrange-Newton, ...
- Assumption: T_{u_k} is coercive, i.e.,

$$\langle T_{u_k} v, v \rangle \geq \alpha \|v\|_{\mathcal{H}}^2$$

for some $\alpha > 0$.

Projection Theorem

Theorem

For the discretized algorithm $\theta_{k+1} = \theta_k - \eta_k G(\theta_k)^\dagger \nabla L(\theta_k)$ it holds

$$u_{\theta_{k+1}} = u_{\theta_k} - \eta_k \Pi_{u_{\theta_k}} \underbrace{[T_{u_{\theta_k}}^{-1}(DE(u_{\theta_k}))]}_{=d_k} + \epsilon_k$$

where Π_u denotes the Galerkin projection onto the tangent space. The term ϵ_k corresponds to an error vanishing quadratically in the step and step size length

$$\epsilon_k = O(\eta_k^2 \|G(\theta_k)^\dagger \nabla L(\theta_k)\|^2).$$

Proof.

Taylor expansion & Céa's Lemma. □

Visualization of the Projection Theorem

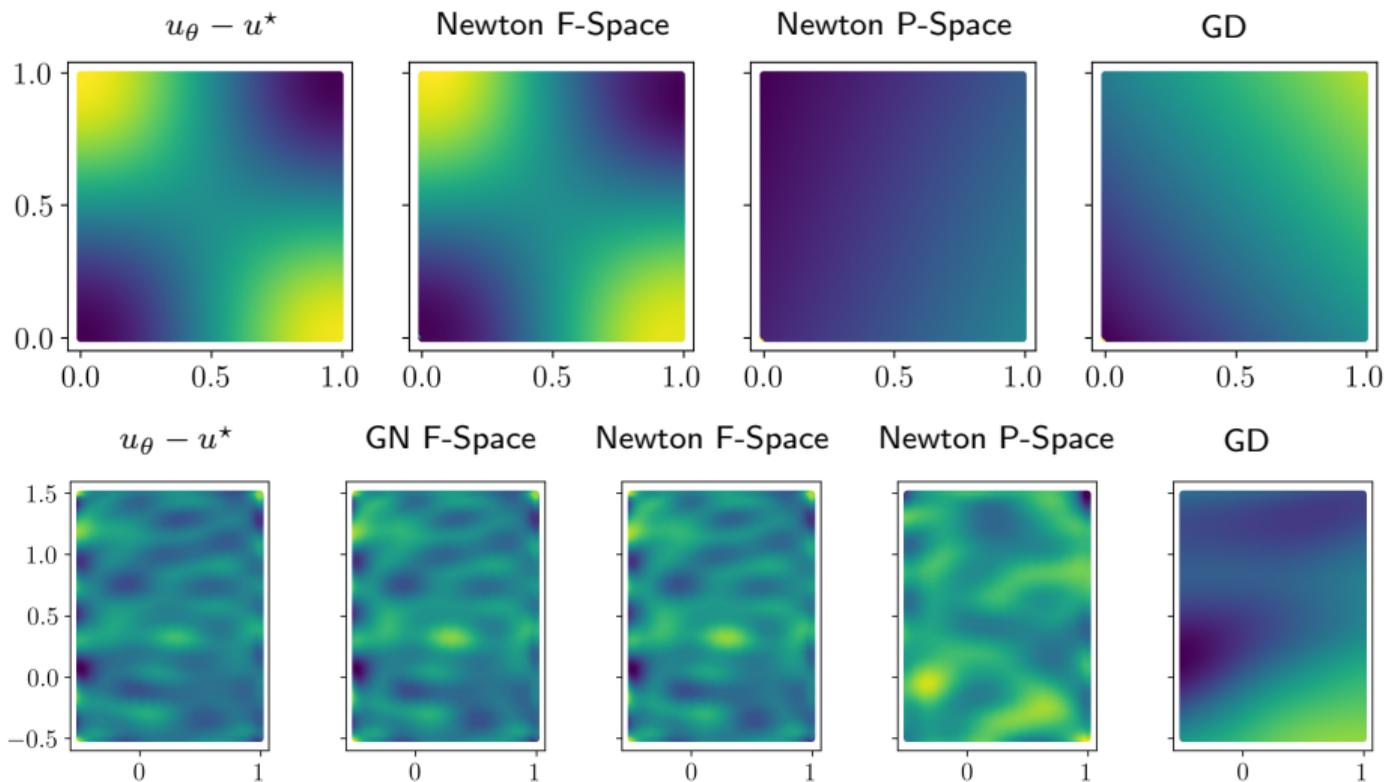


Table of Contents

1. Function Space Optimization for SciML

2. Scalability

3. Numerical Experiments

4. Summary

The Price of Non-local Ansatz Functions

Recall the Structure of the ENGD Matrix

In every step of optimization the following matrix needs to be assembled and inverted

$$G(\theta)_{ij} = D^2 E(u_\theta)(\partial_{\theta_i} u_\theta, \partial_{\theta_j} u_\theta), \quad i, j = 1, \dots, p.$$

For neural network ansatz functions the matrix is dense, rank deficient and ill-conditioned.

- A direct solve works well for small networks.
- Naive application is doomed to fail, already for networks of modest size, say $p > 10^4$ trainable parameters.

What are our Options?

- Matrix-free solution methods for the system $G(\theta)d = \nabla L(\theta)$.
 - The matrix vector-product $G(\theta)v$ can be computed at comparable cost to the gradient $\nabla L(\theta)$ ^{20, 21}, hence use CG.
 - Drawback: CG iterates long due to ill-conditioning of $G(\theta)$.
- Derive cheap-to-invert approximations.
 - Derivatives of linear network layers yield Kronecker structure²². Leverage to build an approximation of $G(\theta)$!
 - The details fill the remaining talk.
- Combine both ideas, i.e., use preconditioned CG.

²⁰N. N. Schraudolph (2002). “Fast curvature matrix-vector products for second-order gradient descent”. In: *Neural computation* 14.7, pp. 1723–1738.

²¹A. Jnini, F. Vella, and M. Zeinhofer (2024). “Gauss-Newton Natural Gradient Descent for Physics-Informed Computational Fluid Dynamics”. In: *arXiv preprint arXiv:2402.10680*.

²²J. Martens (2020). “New Insights and Perspectives on the Natural Gradient Method”. In: *The Journal of Machine Learning Research* 21.1, pp. 5776–5851.

Kronecker-Factor Approximation

- We discard cross-interactions between layers and use a Kronecker Approximation²³:

$$\begin{aligned}G(\theta) &\approx \text{diag}(G^{(1)}, \dots, G^{(L)}) \\ &\approx \text{diag}(A^{(1)} \otimes B^{(1)}, \dots, A^{(L)} \otimes B^{(L)}).\end{aligned}$$

- Use the property of Kronecker product

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$$

- Ignoring bias terms, for layer l mapping from $\mathbb{R}^n \rightarrow \mathbb{R}^m$ the matrix $G^{(l)}(\theta)$ is of size (nm, nm) but the Kronecker factors are of size

$$A^{(l)} \in \mathbb{R}^{n \times n}, \quad B^{(l)} \in \mathbb{R}^{m \times m}$$

²³We omit the θ dependency for brevity only.

Where does the Kronecker Structure Come From?

Structure of ENGD Matrix

Ignoring boundary terms, one can show that the ENGD matrix for Poisson is

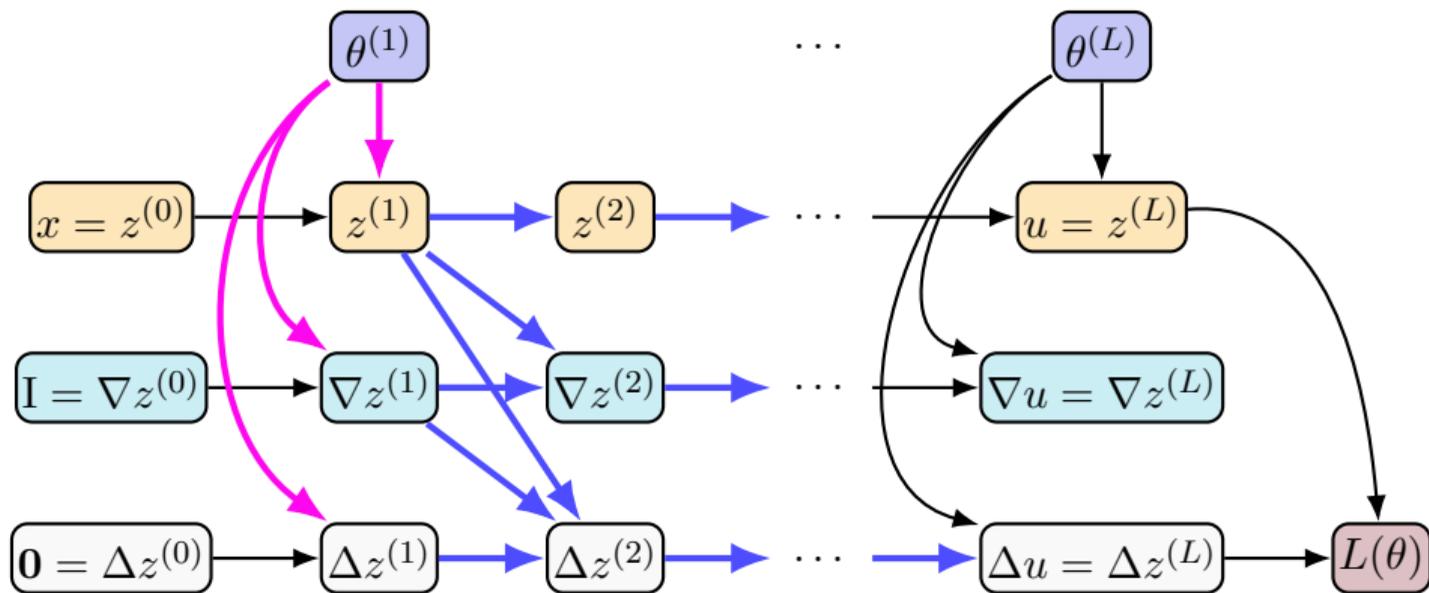
$$G^{(l)}(\theta) = \frac{1}{N_\Omega} \sum_{n=1}^{N_\Omega} \mathbf{J}_{\theta^{(l)}} \Delta u_\theta(x_n)^\top \mathbf{J}_{\theta^{(l)}} \Delta u_\theta(x_n)$$

- Key for Kronecker factors: Jacobian of $\theta \mapsto \theta x$, where θ is weight matrix

$$\mathbf{J}_\theta(\theta x) = x^\top \otimes \text{Id}$$

- Need: children of $\theta^{(l)}$ in computational graph of $\theta^{(l)} \mapsto \Delta u_\theta(x_n)$.

Taylor Mode alias Forward Laplacian²⁴ Graph



$$\theta^{(l)} \mapsto \Delta u_{\theta}(x) = [(z^{(l)}, \nabla z^{(l)}, \Delta z^{(l)}) \mapsto \Delta u_{\theta}(x)] \circ [\theta^{(l)} \mapsto (z^{(l)}, \nabla z^{(l)}, \Delta z^{(l)})]$$

²⁴R. Li et al. (2024). "A computational framework for neural network-based variational Monte Carlo with Forward Laplacian". In: *Nature Machine Intelligence*, pp. 1–11.

The Full K-FAC Approximation

Exact block of the ENGD matrix of a Laplace operator

$$G_{\Omega}^{(l)}(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{s=1}^S \sum_{s'=1}^S Z_{n,s}^{(l-1)} Z_{n,s'}^{(l-1)\top} \otimes g_{n,s}^{(l)} g_{n,s'}^{(l)\top}.$$

K-FAC for the ENGD matrix of a Laplace operator

$$G_{\Omega}^{(l)}(\theta) \approx \left(\frac{1}{NS} \sum_{n,s=1}^{N,S} Z_{n,s}^{(l-1)} Z_{n,s}^{(l-1)\top} \right) \otimes \left(\frac{1}{N} \sum_{n,s=1}^{N,S} g_{n,s}^{(l)} g_{n,s}^{(l)\top} \right) =: A_{\Omega}^{(l)} \otimes B_{\Omega}^{(l)}$$

$$Z_{n,1}^{(l)} := z_n^{(l)}, \quad Z_{n,2}^{(l)} := \partial_{x_1} z_n^{(l)}, \dots, Z_{n,1+d}^{(l)} := \partial_{x_d} z_n^{(l)}, \quad Z_{n,2+d}^{(l)} := \Delta z_n^{(l)}$$
$$g_{n,s}^{(l)} := J_{Z_{n,s}^{(l)}} \Delta u_n$$

Further Details & Generalizations

- The K-FAC approximation works for more than just the Poisson equation:
 - general nonlinear PDEs using PINN formulation
 - Variational formulations (deep Ritz)
 - Neural Operators? Requires more work.
- Typically combined with
 - Trust-region methods for step-size choice
 - Exponential moving average for K-FAC-matrices
 - Infrequent updates of the K-FAC-matrices
 - Momentum on the natural gradient
- Find all details in the pre-print²⁵.

²⁵F. Dangel, J. Müller, and M. Zeinhofer (2024). "Kronecker-Factored Approximate Curvature for Physics-Informed Neural Networks". In: *Advances in Neural Information Processing Systems*.

Table of Contents

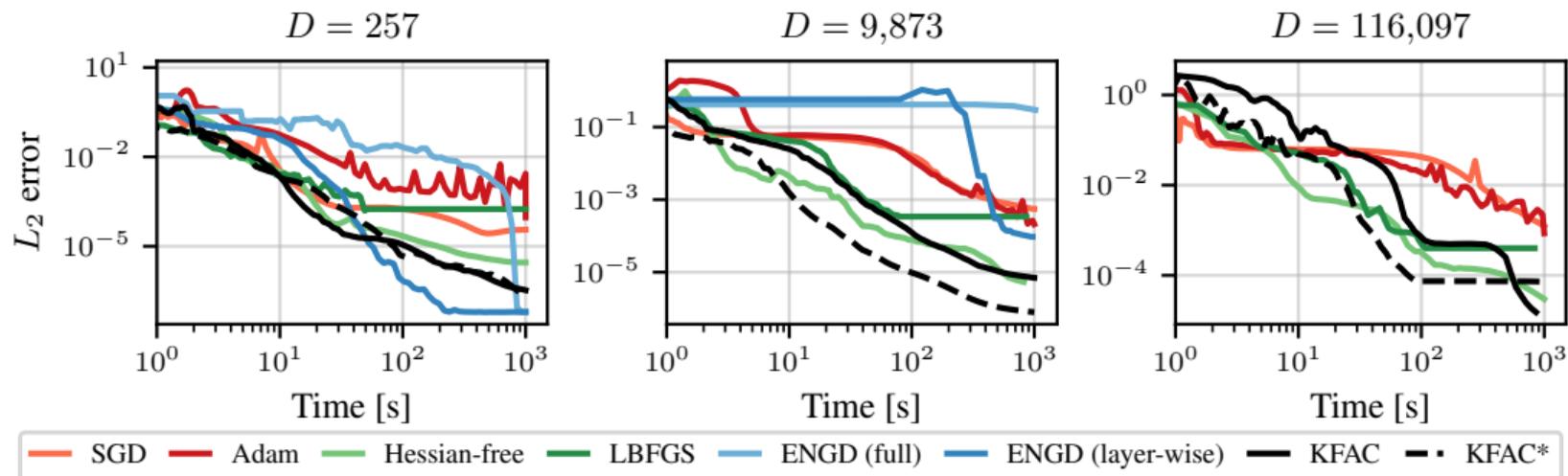
1. Function Space Optimization for SciML

2. Scalability

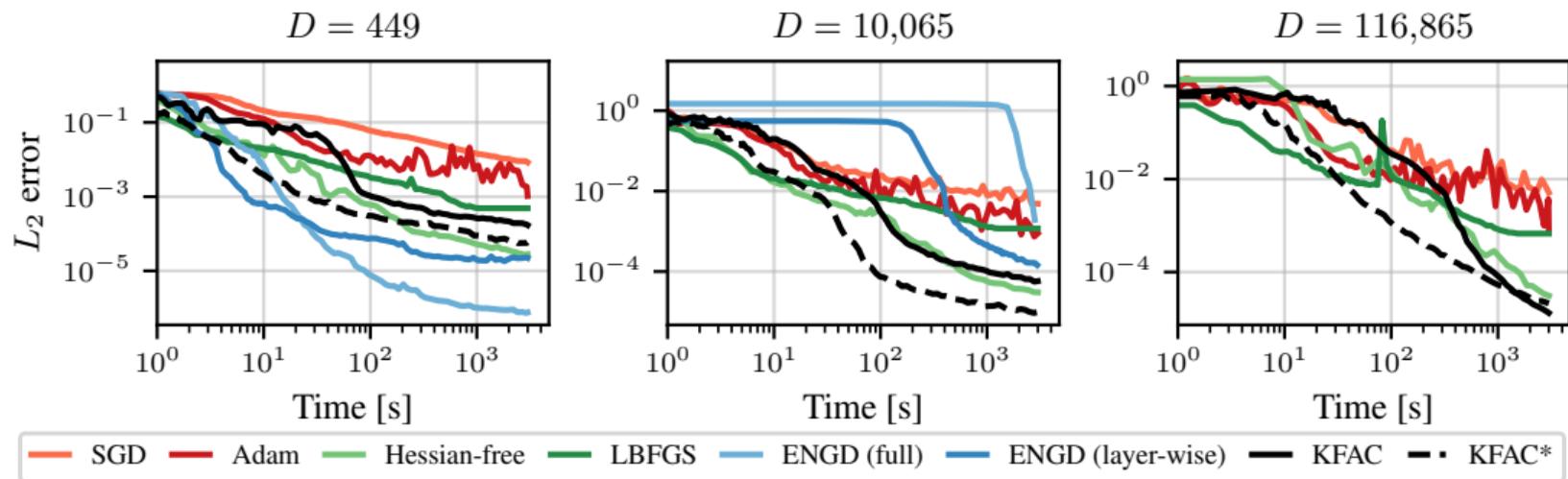
3. Numerical Experiments

4. Summary

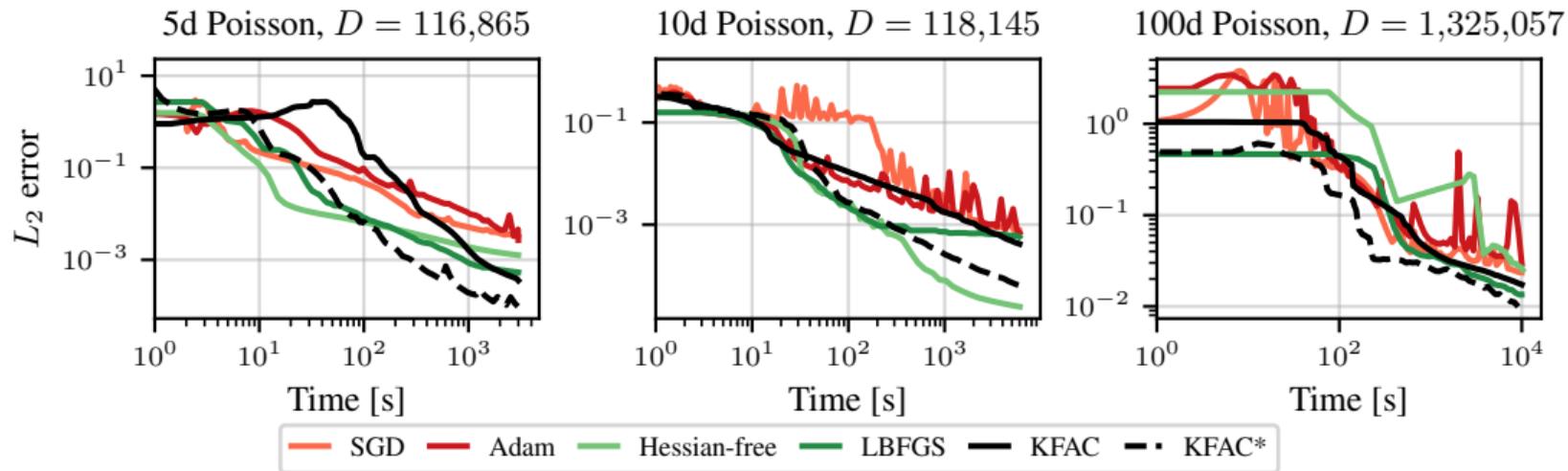
Poisson Equation Two Dimensions



Heat Equation Equation (4+1) Dimensions



High Dimensional Poisson Equation



Discussion

- + The K-FAC approximation allows to scale natural gradient methods to millions of parameters.
- + Even preliminary numerical experiments compare favorable to established optimizers.
- + The approximation generalizes to other PDEs (see pre-print²⁶ for details).
 - The numerical experiments are preliminary, more equations need to follow.
 - The K-FAC approximation is architecture dependent and complicated to handle.

²⁶F. Dangel, J. Müller, and M. Zeinhofer (2024). "Kronecker-Factored Approximate Curvature for Physics-Informed Neural Networks". In: *Advances in Neural Information Processing Systems*.

Table of Contents

1. Function Space Optimization for SciML
2. Scalability
3. Numerical Experiments
- 4. Summary**

When is the Functional Viewpoint Useful?

- We discussed difficulty of NN optimization for loss functions including PDE operators.
- Infinite dimensional algorithms can exploit geometric structure.
- Obtain an algorithm by Galerkin discretization in tangent space.
- Scalability can be non-trivial.

When can this be useful?

The problem possesses a continuous formulation & contains PDE terms.

When does it not help?

First order methods are successful e.g., supervised deep learning applications.

References I

-  Amari, S.-I. (1998). “Natural Gradient Works Efficiently in Learning”. In: *Neural Computation* 10.2, pp. 251–276.
-  Dangel, F., J. Müller, and M. Zeinhofer (2024). “Kronecker-Factored Approximate Curvature for Physics-Informed Neural Networks”. In: *Advances in Neural Information Processing Systems*.
-  Eschenhagen, R. et al. (2024). “Kronecker-factored approximate curvature for modern neural network architectures”. In: *Advances in Neural Information Processing Systems* 36.
-  Hu, T., B. Jin, and Z. Zhou (2023). “Solving Poisson Problems in Polygonal Domains with Singularity Enriched Physics Informed Neural Networks”. In: *arXiv preprint arXiv:2308.16429*.

References II

-  Jnini, A., F. Vella, and M. Zeinhofer (2024). “Gauss-Newton Natural Gradient Descent for Physics-Informed Computational Fluid Dynamics”. In: *arXiv preprint arXiv:2402.10680*.
-  Li, H. et al. (2018). “Visualizing the loss landscape of neural nets”. In: *Advances in neural information processing systems* 31.
-  Li, R. et al. (2024). “A computational framework for neural network-based variational Monte Carlo with Forward Laplacian”. In: *Nature Machine Intelligence*, pp. 1–11.
-  Li, Z. et al. (2021). “Fourier Neural Operator for Parametric Partial Differential Equations”. In: *International Conference on Learning Representations*.
-  Martens, J. (2020). “New Insights and Perspectives on the Natural Gradient Method”. In: *The Journal of Machine Learning Research* 21.1, pp. 5776–5851.
-  Müller, J. and M. Zeinhofer (2023). “Achieving High Accuracy with PINNs via Energy Natural Gradients”. In: *International Conference on Machine Learning*.

References III

-  Müller, J. and M. Zeinhofer (2024). “Position: Optimization in SciML Should Employ the Function Space Geometry”. In: *Forty-first International Conference on Machine Learning*.
-  Pfau, D. et al. (2020). “Ab Initio Solution of the Many-Electron Schrödinger Equation with Deep Neural Networks”. In: *Physical Review Research* 2.3, p. 033429.
-  Raissi, M., P. Perdikaris, and G. E. Karniadakis (2019). “Physics-Informed Neural Networks: A Deep Learning Framework for Solving Forward and Inverse Problems Involving Nonlinear Partial Differential Equations”. In: *Journal of Computational Physics* 378, pp. 686–707.
-  Schäfer, F. and A. Anandkumar (2019). “Competitive gradient descent”. In: *Advances in Neural Information Processing Systems* 32.
-  Schraudolph, N. N. (2002). “Fast curvature matrix-vector products for second-order gradient descent”. In: *Neural computation* 14.7, pp. 1723–1738.

References IV

-  Zeng, Q., S. H. Bryngelson, and F. T. Schaefer (2023). “Competitive Physics Informed Networks”. In: *International Conference on Learning Representations*.