

ANaGRAM: A Natural Gradient Relative to Adapted Model for efficient PINNs learning

ICLR 2025

Nilo Schwencke, Cyril Furtlehner

TAU Team, INRIA Saclay—A&O, LISN, Paris-Saclay University—CNRS

March 31, 2025



Physics Informed Neural Networks (PINNs)

Problem statement

We aim to solve:

$$\begin{cases} D(u) = f \in L^2(\Omega \rightarrow \mathbb{R}, \mu) & \text{in } \Omega \\ B(u) = g \in L^2(\partial\Omega \rightarrow \mathbb{R}, \sigma) & \text{on } \partial\Omega \end{cases}.$$

PINNs key idea

Optimize a neural network $u|_{\theta}$ on the loss:

$$\begin{aligned} \hat{\ell}_{D,B}(\theta) := & \frac{1}{2S_D} \sum_{i=1}^{S_D} \left(D[u|_{\theta}](x_i^D) - f(x_i^D) \right)^2 \\ & + \frac{1}{2S_B} \sum_{i=1}^{S_B} \left(B[u|_{\theta}](x_i^B) - g(x_i^B) \right)^2, \end{aligned}$$

using autodiff to compute D, B (Raissi et al., 2019).

Problem

This leads to low accuracy with SGD.

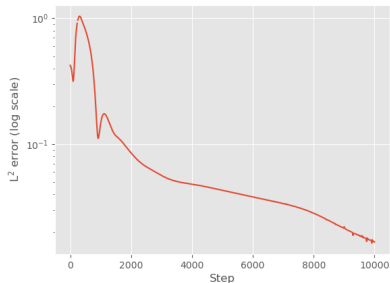


Figure: PINN solution under standard Adam optimization, to Laplace equation in 2 D.

Natural Gradient

Reinterpreting quadratic loss

Consider the loss of a classical quadratic regression problem, with batch (x_i) :

$$\hat{\ell}(\boldsymbol{\theta}) := \frac{1}{2S} \sum_{i=1}^S (u_{|\boldsymbol{\theta}}(x_i) - f(x_i))^2.$$

In the population limit:

$$\hat{\ell}(\theta) \xrightarrow{S \rightarrow \infty} \mathcal{L}(u_{|\theta}); \quad \mathcal{L}(u) := \frac{1}{2} \|u - f\|_{L^2(\Omega)}^2$$

This yields the Fréchet derivative:

$$d\mathcal{L}|_u(h) = \langle \underbrace{u - f}_{\nabla \mathcal{L}|_u}, h \rangle_{L^2(\Omega)},$$

and thus the gradient flow:

$$\begin{cases} u_0 \in L^2(\Omega) \\ \dot{u}_t = -\nabla \mathcal{L}|_{u_t} = f - u_t \end{cases}.$$

Solution: $u_t = f - e^{-t}(u_0 - f)$.

A functional geometry perspective

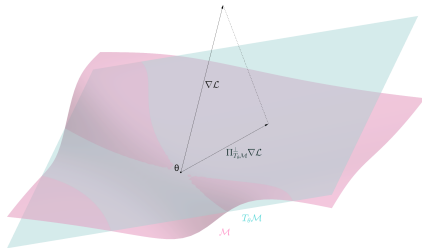
Natural gradient in functional space

The functional space is constrained to:

- $\mathcal{M} := \text{Im } u = \{u_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \mathbb{R}^P\}$
- $T_{\boldsymbol{\theta}}\mathcal{M} := \text{Im } du_{|\boldsymbol{\theta}} = \text{Span}(\partial_p u_{\boldsymbol{\theta}})$

The Natural Gradient is then defined as:

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \eta du_{|\boldsymbol{\theta}_t}^\dagger \left(\Pi_{T_{\boldsymbol{\theta}_t}^{\perp} \mathcal{M}} \nabla \mathcal{L}|_{u_{|\boldsymbol{\theta}_t}} \right),$$



Definition-Proposition (Schwencke and Furtlehner (2025))

The **Natural Neural Tangent Kernel (NNTK)** is the kernel of the projection $\Pi_{T_{\theta}\mathcal{M}} : L^2(\Omega) \rightarrow L^2(\Omega)$ onto $T_{\theta}\mathcal{M}$. It is given by the formula:

$$NNTK_{\theta}(x, y) := \sum_{1 \leq p, q \leq P} (\partial_p u|_{\theta}(x)) \ G_{\theta pq}^{\dagger} (\partial_q u|_{\theta}(y))^t; \quad G_{\theta p, q} := \langle \partial_p u|_{\theta}, \partial_q u|_{\theta} \rangle_{L^2(\Omega)}.$$

Corollary

The Natural Gradient update rewrites: $\theta_{t+1} \leftarrow \theta_t - \eta G_{\theta_t}^{\dagger} \nabla \ell(\theta_t)$; $\ell(\theta) := \mathcal{L}(u|_{\theta})$.

Shortcomings

- Computation of the Gram matrix G_{θ_t} is quadratic in the number of parameters.
- Inversion of G_{θ_t} is cubic

We introduce a the empirical Natural Gradient that scales linearly with the number of parameters.

empirical Natural Gradient

(N)NTK in a nutshell

The functional dynamic of (N)GD on the empirical loss $\hat{\ell}$ is described by (Jacot et al., 2018; Rudner et al., 2019):

$$\frac{du_{\theta_t}}{dt}(x) = - \sum_{i=1}^S (N)NTK_{\theta_t}(x, x_i)(u_{|\theta_t}(x_i) - y_i),$$

Key Idea

The empirical dynamics takes place in:

$$\hat{T}_{\theta}\mathcal{M} := \text{Span} \left(NNTK_{\theta}(x_i, \cdot) : (x_i)_{1 \leq i \leq N} \right).$$

We can define the empirical Natural Gradient:

$$\theta_{t+1} = \theta_t - \eta du_{|\theta_t}^{\dagger} \left(\Pi_{\hat{T}_{\theta_t}\mathcal{M}}^{\perp} \nabla \mathcal{L}_{|u_{|\theta_t}} \right).$$

Theorem (ANaGRAM)

Under mild assumptions:

$$du_{|\theta_t}^{\dagger} \left(\Pi_{\hat{T}_{\theta_t}\mathcal{M}}^{\perp} \nabla \mathcal{L}_{|u_{|\theta_t}} \right) = \hat{\phi}_{\theta_t}^{\dagger} \widehat{\nabla \mathcal{L}}_{\theta_t},$$

with: for all $1 \leq p \leq P, 1 \leq i \leq S$

- $\hat{\phi}_{\theta_{t i}, p} := \partial_p u_{|\theta_t}(x_i)$
- $\widehat{\nabla \mathcal{L}}_{\theta_{t i}} := \nabla \mathcal{L}_{|u_{|\theta_t}}(x_i)$

Key fact

$\hat{\phi}_{\theta_t}^{\dagger}$ can be computed with a SVD, with complexity $O(\min(P^2 S, S^2 P))$.

Corollary

There exist P points (\hat{x}_i) such that:

$$\Pi_{\hat{T}_{\theta}\mathcal{M}}^{\perp} \nabla \mathcal{L}_{|u_{|\theta}} = \Pi_{T_{\theta}\mathcal{M}}^{\perp} \nabla \mathcal{L}_{|u_{|\theta}}.$$

Application to PINNs

PINNs are a quadratic regression problem

Natural Gradient of PINNs

Figure: Illustration of PINNs Natural Gradient

PINNs are a quadratic regression problem with model: $(D, B) \circ u :$

$$\begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H} & \rightarrow L^2(\Omega \rightarrow \mathbb{R}, \mu) \times \\ & & L^2(\partial\Omega \rightarrow \mathbb{R}, \sigma) \\ \boldsymbol{\theta} & \mapsto u|_{\boldsymbol{\theta}} & \mapsto (D[u|_{\boldsymbol{\theta}}], B[u|_{\boldsymbol{\theta}}]) \end{cases}$$

Experiments

2 D Laplace equation

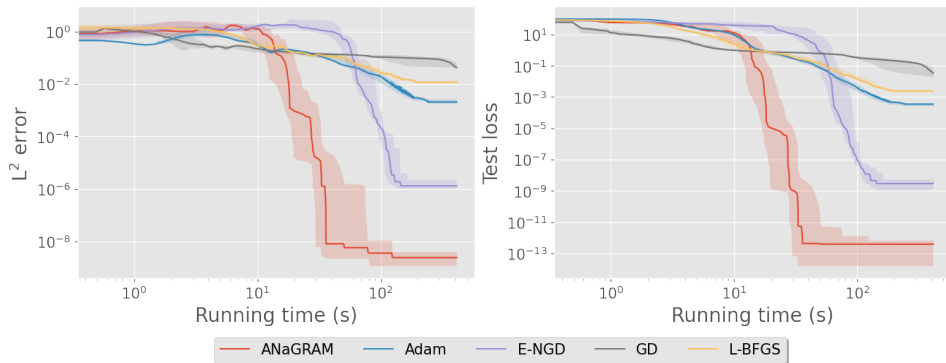


Figure: Performance comparison w.r.t running time for Laplace equation in 2 D:

$$\begin{cases} \Delta u = -2\pi^2 \sin(\pi x_1) \sin(\pi x_2) & \text{in } [0, 1]^2 \\ u = 0 & \text{on } \partial[0, 1]^2 \end{cases}$$

1+1 D Heat equation

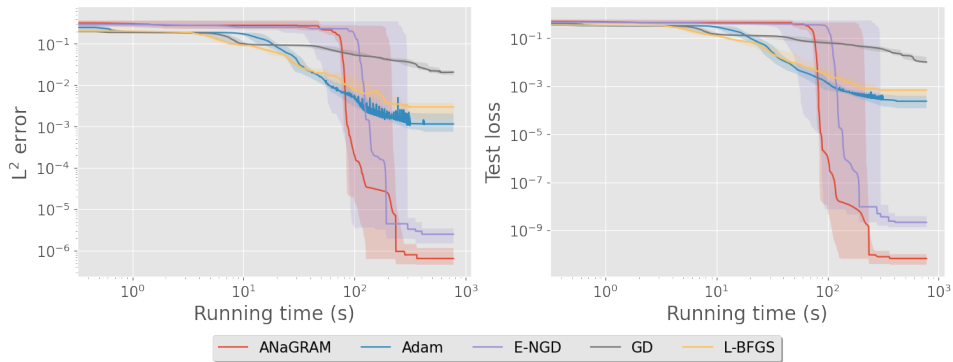


Figure: Performance comparison w.r.t running time for Heat equation in 1+1 D:

$$\begin{cases} \partial_t u - \frac{1}{4} \partial_{xx} u = 0 & \text{in } [0, 1]^2 \\ u = 0 & \text{on } [0, 1] \times \{0, 1\} \\ u = \sin(\pi x) & \text{on } \{0\} \times [0, 1] \end{cases}$$

5 D Laplace equation

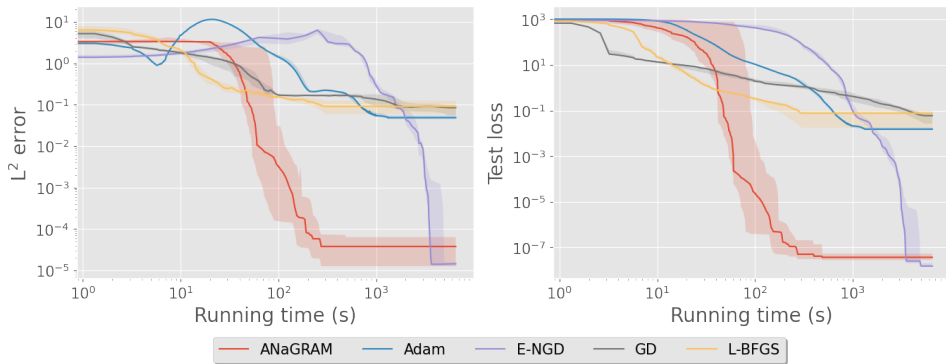


Figure: Performance comparison w.r.t running time for Laplace equation in 5 D:

$$\begin{cases} \Delta u = \pi^2 \sum_{k=1}^5 \sin(\pi x_k) & \text{in } \Omega = [0, 1]^5 \\ u = \sum_{k=1}^5 \sin(\pi x_k) & \text{on } \partial\Omega \end{cases}$$

1+1 D Allen-Cahn equation

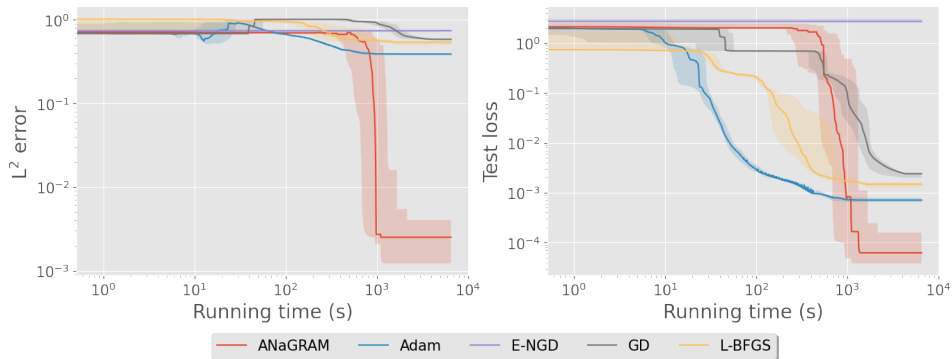


Figure: Performance comparison w.r.t running time for Allen-Cahn equation in 1+1 D:

$$\begin{cases} \partial_t u - 10^{-3} \partial_{xx} u - 5(u - u^3) = 0 & \text{in } \Omega = [0, 1] \times [-1, 1] \\ u = -1 & \text{on } \partial\Omega_{\text{border}} = [0, 1] \times \{-1, 1\} \\ u(0, x) = x^2 \cos(\pi x) & \text{on } \partial\Omega_0 = \{0\} \times [-1, 1] \end{cases}$$

Conclusion and Perspectives

Conclusions

- Anagram gives a theoretically founded simplification to any natural-gradient algorithm lowering the complexity from $O(P^3)$ to $O(\min(PN^2, P^2N))$, which is above stochastic gradient descent only by a factor $\min(P, N)$.
- In the case of PINNs, we prove that natural gradient correspond to an optimal linear update following the Green's function.
- Empirical results are improved by several orders of magnitude.
- The SVD cut-off factor appears to be a pivotal hyper-parameter of the algorithm.

Perspectives

- Design of an optimal collocation points procedure, coupled with SVD cut-off factor adaptation strategy.
- Establish theoretical connections with classical algorithms, such as FEMs, FDMs, *etc.*
- Include data assimilation in this theoretical setting, and understand its regularizing effect.
- Include common optimization techniques (e.g. Momentum)
- Extend to order 2 methods
- Extend it to Operator learning
- Application to HJB equation

Thank you for your attention !

- JACOT, A., F. GABRIEL, AND C. HONGLER (2018): “Neural Tangent Kernel: Convergence and Generalization in Neural Networks,” *Advances in neural information processing systems*, 31.
- RAISSI, M., P. PERDIKARIS, AND G. KARNIADAKIS (2019): “Physics-Informed Neural Networks: A Deep Learning Framework for Solving Forward and Inverse Problems Involving Nonlinear Partial Differential Equations,” *Journal of Computational Physics*, 378, 686–707.
- RUDNER, T. G., F. WENZEL, Y. W. TEH, AND Y. GAL (2019): “The Natural Neural Tangent Kernel: Neural Network Training Dynamics under Natural Gradient Descent,” in *4th Workshop on Bayesian Deep Learning (NeurIPS 2019)*.
- SCHWENCKE, N. AND C. FURTLEHNER (2025): “ANaGRAM: A Natural Gradient Relative to Adapted Model for Efficient PINNs Learning,” in *The Thirteenth International Conference on Learning Representations*.