

---

# AMSTRAMGRAM: ADAPTIVE MULTI-CUTOFF STRATEGY MODIFICATION FOR ANaGRAM

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Recent works have shown that natural gradient methods can significantly outperform standard optimizers when training physics-informed neural networks (PINNs). In this paper, we analyze the training dynamics of PINNs optimized with ANaGRAM, a natural-gradient-inspired approach employing singular value decomposition with cutoff regularization. Building on this analysis, we propose a multi-cutoff adaptation strategy that further enhances ANaGRAM’s performance. Experiments on benchmark PDEs validate the effectiveness of our method, which allows to reach machine precision on some experiments. To provide theoretical grounding, we develop a framework based on spectral theory that explains the necessity of regularization and extend previous shown connections with Green’s functions theory.

## 1 INTRODUCTION

Physics-informed neural networks (PINNs) have recently emerged as a promising alternative for the numerical solution of partial differential equations (PDEs) (Raissi et al., 2019). By leveraging neural networks as universal function approximators (Leshno et al., 1993), PINNs replace traditional mesh-based discretizations with sampling-based collocation methods, enabling a straightforward extension to high-dimensional domains. This mesh-free formulation not only circumvents the “curse of dimensionality” inherent in grid-based approaches, but also allows continuous evaluation of the solution throughout the domain without explicit mesh generation (Cuomo et al., 2022).

Despite these advantages, achieving low training error with PINNs remains a major challenge (Wang et al., 2023; Urbán et al., 2025; Kiyani et al., 2025; De Ryck et al., 2024). Open questions include how to select and distribute collocation points, how to balance the PDE residual against boundary-condition penalties, and which optimization strategies most effectively minimize the composite loss (Krishnapriyan et al., 2021; Wang et al., 2021; McClenny & Braga-Neto, 2022). **Recent work has pursued machine-precision accuracy through hierarchical or multilevel PINN architectures (??), demonstrating that PINNs can successfully tackle ill-posed problems that challenge conventional finite element methods. These approaches improve accuracy primarily through architectural design. In contrast, our work focuses on accuracy gains driven purely by the optimizer, offering a complementary perspective.**

A different line of research has recently reexamined PINNs from the perspective of functional geometry (Müller & Zeinhofer, 2023; 2024; Jnini et al., 2024), providing a mathematically principled view of the training dynamics. In this vein, the ANaGRAM algorithm (Schwencke & Furtlehner, 2025) applies a natural-gradient update (Amari, 1998; Ollivier, 2015), based on a reinterpretation and generalization of the neural tangent kernel (NTK; Jacot et al. (2018)) as the kernel of the projection onto the neural network’s tangent space. This leads to a notion of the empirical natural gradient that projects the true functional gradient onto the empirical tangent space, yielding significantly faster convergence and lower errors compared to standard optimizers on PDE benchmarks.

Nevertheless, while ANaGRAM improves over standard optimizers, it still falls short of the accuracy attained by classical mesh-based methods, such as the finite element method (Grossmann et al., 2024). Moreover, its final performance is highly affected by the way the pseudo-inverse of the feature matrix is computed. In particular, ANaGRAM sets a fixed level of *cutoff*: a value below which the singular values of the feature matrix are ignored, *i.e.* it controls how much loss signal is incorporated into an

054 update. ANaGRAM’s cutoff is currently chosen manually, as no automatic selection procedure has  
 055 been proposed.

056 In this paper, we study the performance and training dynamics of ANaGRAM, with a particular  
 057 focus on the role of the chosen cutoff. Typically, the training loss of ANaGRAM exhibits the slow  
 058 convergence at the early iterations followed by a sudden drop at the end of the training – similar  
 059 behavior is shown by the eNGD method (Müller & Zeinhofer, 2023). We discover that it is closely  
 060 connected to what we further refer as the *flattening phenomenon*, which we define and characterize  
 061 using the *reconstruction error*: a novel metric that measures how much of the loss signal is lost  
 062 by different choices of cutoffs. Relying on the adaptive multi-cutoff strategy, our new algorithm  
 063 AMStramGRAM manages to capitalize on this phenomenon, resulting in a significant improvement  
 064 (of several orders of magnitude) on various PDE benchmarks. To complement our empirical findings,  
 065 we also present a functional-analytic view linking cutoff (and ridge regularization) to (generalized)  
 066 Green operator theory, clarifying why cutoff regularization is essential and not just a mere fix to  
 067 stabilize training.

## 068 2 PROBLEM STATEMENT

### 069 2.1 DIFFERENTIAL OPERATORS AND PHYSICS-INFORMED NEURAL NETWORKS (PINNS)

070 Let  $\Omega \subset \mathbb{R}^d$  be a domain. We introduce two operators,  $D$  and  $B$ , defined on a Hilbert space  $\mathcal{H}$  of  
 071 real-valued functions, acting respectively on  $\Omega$  and on its boundary  $\partial\Omega$ :

$$072 \quad D : \begin{cases} \mathcal{H} & \rightarrow \mathbf{L}^2(\Omega \rightarrow \mathbb{R}, \mu) \\ u & \mapsto D[u] \end{cases}, \quad B : \begin{cases} \mathcal{H} & \rightarrow \mathbf{L}^2(\partial\Omega \rightarrow \mathbb{R}, \sigma) \\ u & \mapsto B[u] \end{cases}. \quad (1)$$

073 Here,  $D$  denotes a differential operator, while  $B$  represents a boundary operator. A function  $u \in \mathcal{H}$  is  
 074 said to be a *classical solution* to the *Partial Differential Equation* (PDE) associated with  $D$  and  $B$  if  
 075 it satisfies

$$076 \quad \begin{cases} D(u) = f \in \mathbf{L}^2(\Omega \rightarrow \mathbb{R}, \mu), & \text{in } \Omega, \\ B(u) = g \in \mathbf{L}^2(\partial\Omega \rightarrow \mathbb{R}, \sigma), & \text{on } \partial\Omega, \end{cases} \quad (2)$$

077 A *physics-informed neural network* (PINN) approximates the solution  $u$  by a parametric model  $u_\theta$ ,  
 078 where  $u_\theta$  is a neural network with parameters  $\theta \in \mathbb{R}^P$ . The learning objective is to minimize the  
 079 empirical loss

$$080 \quad \ell_{D,B}(\theta) := \frac{1}{2S_D} \sum_{i=1}^{S_D} (D[u_\theta](x_i^D) - f(x_i^D))^2 + \frac{1}{2S_B} \sum_{i=1}^{S_B} (B[u_\theta](x_i^B) - g(x_i^B))^2. \quad (3)$$

### 081 2.2 PINNS OPTIMIZERS

082 Training PINNs is notoriously challenging. Issues such as spectral bias, where networks struggle  
 083 to learn high-frequency components, and the difficulty of balancing residual and boundary loss  
 084 terms—often with vastly different magnitudes— result in unsatisfactory performance of standard  
 085 deep learning optimizers (Wang et al., 2021; De Ryck et al., 2024; Krishnapriyan et al., 2021; Liu  
 086 et al., 2024).

087 To mitigate these challenges, researchers have proposed various strategies. These include adaptive  
 088 sampling approaches that focus on regions with high error (Krishnapriyan et al., 2021), dynamic loss  
 089 weighting schemes (McClenny & Braga-Neto, 2022), and architectural modifications (Wang et al.,  
 090 2024). Another promising line of research has focused on modifying the optimizers. In particular,  
 091 two main branches of optimization approaches for PINNs have emerged:

- 092 (i) **Second-Order Methods.** These methods, based on Quasi-Newton techniques, particularly  
 093 the BFGS algorithm (Nocedal & Wright, 1999, Chapter 6) and its memory-efficient ap-  
 094 proximation L-BFGS (Liu & Nocedal, 1989), address some of the training difficulties by  
 095 considering the curvature of the loss landscape. This curvature arises from the non-linearities  
 096 of both the neural network and the differential operators (Rathore et al., 2024). Recently,  
 097 Urbán et al. (2025) extended this approach by modifying the self-scaled BFGS (SSBFGS;

Al-Baali, 1998) and self-scaled Broyden (SSBroyden; Al-Baali & Khalfan, 2005), along with other computational enhancements such as point resampling (Wu et al., 2023) and boundary condition enforcement (Wang et al., 2023), achieving state-of-the-art results (Kiyani et al., 2025).

- (ii) **Natural Gradient Methods.** In contrast to second-order methods, natural gradient methods are **first-order** techniques<sup>1</sup> that provide a principled way to incorporate the geometry and metric structure of the problem space. Initially introduced in the context of information geometry by Amari (1998) and later extended by Ollivier (2015), these methods were introduced for PINNs by Müller & Zeinhofer (2023). In subsequent work, Schwencke & Furtlehner (2025) connected these methods to kernel methods, yielding an efficient implementation they linked to Green’s function theory (Duffy, 2015).

### 2.3 NATURAL GRADIENT METHODS FOR PINNS

As a preliminary observation highlighted in Schwencke & Furtlehner (2025, Section 4.1), PINNs can be interpreted as a quadratic regression problem. This viewpoint arises naturally once the parametric model  $u_\theta$  is replaced with the following compound model:

$$(D, B) \circ u : \begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H} & \rightarrow \mathbf{L}^2(\Omega, \mu) \times \mathbf{L}^2(\partial\Omega, \sigma) \\ \boldsymbol{\theta} & \mapsto u_\theta & \mapsto (D[u_\theta], B[u_\theta]) \end{cases}. \quad (4)$$

For ease of exposition, and without loss of generality, we restrict attention to regression in  $\mathbf{L}^2(\Omega, \mu)$ . Given  $f \in \mathbf{L}^2(\Omega, \mu)$ , we define the associated empirical loss

$$\ell(\boldsymbol{\theta}) := \frac{1}{2S} \sum_{i=1}^S (u_\theta(x_i) - f(x_i))^2, \quad (5)$$

which can be seen as a discretization of the functional loss

$$\mathcal{L}(u) := \frac{1}{2} \|u - f\|_{\mathbf{L}^2(\Omega, \mu)}^2, \quad u \in \mathbf{L}^2(\Omega, \mu). \quad (6)$$

The natural gradient approach seeks to compute the optimal update direction in function space and then pull it back to parameter space. A single Fréchet derivative of the functional loss Equation (6) yields  $\nabla \mathcal{L}|_u = u - f$ . The key insight is that admissible updates are constrained to the tangent space of the parametric model,

$$T_{\boldsymbol{\theta}} \mathcal{M} := \text{Im}(du_{\boldsymbol{\theta}}) = \text{Span}(\partial_p u_{\boldsymbol{\theta}} : 1 \leq p \leq P) \subset \mathcal{H}, \quad (7)$$

where  $\mathcal{M} := \text{Im}(u) = \{u_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \mathbb{R}^P\} \subset \mathcal{H}$  is the manifold of functions parametrized by  $\boldsymbol{\theta}$ . Thus, the optimal update in function space is the projection of  $\nabla \mathcal{L}|_u$  onto the tangent space (cf. Figure 5),

$$u_{\boldsymbol{\theta}_{t+1}} \leftarrow u_{\boldsymbol{\theta}_t} - \eta \Pi_{T_{\boldsymbol{\theta}_t} \mathcal{M}}(\nabla \mathcal{L}_{u_{\boldsymbol{\theta}_t}}); \quad \boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \eta du_{\boldsymbol{\theta}_t}^\dagger(\Pi_{T_{\boldsymbol{\theta}_t} \mathcal{M}}(\nabla \mathcal{L}_{u_{\boldsymbol{\theta}_t}})), \quad (8)$$

where the second equation is simply the pullback of the functional update to parameter space. We prove in Section H.1 that this update is equivalent to the Gram–matrix formulation:

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \eta G_{\boldsymbol{\theta}_t}^\dagger \nabla \ell(\boldsymbol{\theta}_t); \quad G_{\boldsymbol{\theta}_{t,p,q}} := \langle \partial_p u_{\boldsymbol{\theta}_t}, \partial_q u_{\boldsymbol{\theta}_t} \rangle_{\mathbf{L}^2(\Omega, \mu)}. \quad (9)$$

### 2.4 ANAGRAM: EMPIRICAL NATURAL GRADIENT

The  $O(P^3)$  complexity of matrix inversion in Equation (9) renders a direct implementation prohibitively expensive. ANaGRAM (Schwencke & Furtlehner, 2025) circumvents this by exploiting a motivated approximation. The key observation is that the update can be expressed in terms of the empirical feature matrix  $\widehat{\phi} \in \mathbb{R}^{S \times P}$  and the empirical functional residuals  $\widehat{\mathcal{L}}_\theta \in \mathbb{R}^S$ :

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \eta \widehat{\phi}^\dagger \widehat{\nabla \mathcal{L}}_\theta; \quad \widehat{\phi}_{i,p} := \partial_p u_\theta(x_i); \quad (\widehat{\nabla \mathcal{L}}_\theta)_i := u_\theta(x_i) - f(x_i). \quad (10)$$

<sup>1</sup>contrary to a widespread misconception, which arises from their analogy in the context of information theory

Here, the pseudo-inverse is computed via singular value decomposition (SVD):  $\hat{\phi}^\dagger = \hat{U} \hat{\Delta}^\dagger \hat{V}^T$  with  $\hat{\phi} = \hat{V} \hat{\Delta} \hat{U}^T$ , where  $\hat{U} \in \mathbb{R}^{P \times r_{\text{svd}}}$ ,  $\hat{\Delta} \in \mathbb{R}^{r_{\text{svd}} \times r_{\text{svd}}}$ ,  $\hat{V} \in \mathbb{R}^{S \times r_{\text{svd}}}$ , and  $r_{\text{svd}} = \min(P, S)$ . This reduces computational cost to  $O(\min(P, S)^2)$ , which is tractable in practice. A comparable complexity was later obtained by Guzmán-Cordero et al. (2025) using a Cholesky factorization approach.

For further details on the derivation of the empirical natural gradient, we refer to Schwencke & Furtlehner (2025). In what follows, we adopt a slight abuse of notation by omitting the explicit dependence on  $\theta$  whenever it is clear from context. When iteration indices matter, we explicitly write  $t$  to emphasize the connection to  $\theta_t$ .

## 2.5 REGULARIZATION

As discussed in Section G.1, the type of problem we consider is ill-conditioned, which necessitates the use of regularization. We distinguish between two main regularization schemes: (i) *ridge regression*, which consists in adding a factor  $\alpha^2 I_d$  (or, according to conventions,  $\alpha^{-2} I_d$ ) to the Gram matrix  $G_\theta$  in Equation (9) (or its approximation  $\hat{G}_\theta$ ), thereby making it invertible or (ii) *cutoff regularization*, a scheme that applies a binary threshold (used in ANaGRAM):

$$\hat{\Delta}_{t,i}^\dagger = \begin{cases} \hat{\Delta}_{t,i}^{-1}, & \text{if } \hat{\Delta}_{t,i} \geq \alpha, \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

Here  $\alpha$  denotes the cutoff threshold. This regularization is the focus of our analysis in Section 3. For completeness, we provide a geometric interpretation of each scheme in Section G. We further show that cutoff regularization extends previously established connections between natural gradient methods and Green’s function theory (Schwencke & Furtlehner, 2025). In particular, we obtain:

**Theorem 1.** *The generalized Green’s function of the operator  $D$  in the regularized space  $\mathcal{H}_{D, \mathcal{H}_0}^\alpha$  is given, for all  $x, y \in \Omega$ , by*

$$g_D(x, y) := D[k_D(x, \cdot)](y), \quad (12)$$

where  $\mathcal{H}_{D, \mathcal{H}_0}^\alpha$  is a regularized space with reproducing kernel  $k_D$ , defined in Section G.4. [As a consequence, we show that AMStrAMGRAM converges in the NTK-regime as explained in Section G.5.](#)

## 3 INSIGHTS ON ANaGRAM’S TRAINING DYNAMICS

In this section, we will look at relevant quantities of interest to understand this empirical phenomenon.

### 3.1 RECONSTRUCTION ERROR OF FUNCTIONAL GRADIENT

Let  $\theta \in \mathbb{R}^P$ , the empirical feature matrix  $\hat{\phi} \in \mathbb{R}^{S \times P}$ , and the empirical functional gradient  $\widehat{\nabla} \mathcal{L} \in \mathbb{R}^S$  as defined in Equation (10). Let us consider various empirical tangent spaces formed by taking different ranges of right singular vectors of  $\hat{\phi} = \hat{U} \hat{\Delta} \hat{V}^T$ , i.e.  $\widehat{T}_N^M \mathcal{M} = \text{Span}(\hat{V}_{t,i} : M \leq i \leq N)$ . For  $1 \leq N \leq r_{\text{svd}}$ , reconstruction error measures how much information from the functional gradient signal is lost when considering only first  $N$  components in SVD (the error caused by the projection onto the empirical tangent space  $\widehat{T}_N^0 \mathcal{M}$ ) is defined as follows

$$\text{RCE}_N^S := \frac{1}{\sqrt{S}} \left\| \widehat{V} \Pi_N^0 \widehat{V}^T \widehat{\nabla} \mathcal{L} - \widehat{\nabla} \mathcal{L} \right\|_{\mathbb{R}^S} = \frac{1}{\sqrt{S}} \left\| \Pi_{\widehat{T}_N^0 \mathcal{M}}^\perp \widehat{\nabla} \mathcal{L} - \widehat{\nabla} \mathcal{L} \right\|, \quad (13)$$

where we define  $\Pi_N^M \in \mathbb{R}^{r_{\text{svd}} \times r_{\text{svd}}}$  as a projection operator onto  $\widehat{T}_N^M \mathcal{M}$ :

$$\Pi_N^M = \sum_{p=M+1}^N \mathbf{e}^{(p)} \mathbf{e}^{(p)T}, \quad (14)$$

with  $(\mathbf{e}^{(p)})_{1 \leq p \leq r_{\text{svd}}}$  being the canonical basis of  $\mathbb{R}^{r_{\text{svd}}}$ .

**Proposition 1.**  *$\text{RCE}_N^S$  is a non-increasing function of  $N$ , i.e. for all  $1 \leq M, N \leq r_{\text{svd}}$ :*

$$M \leq N \implies \text{RCE}_M^S \geq \text{RCE}_N^S. \quad (15)$$

Furthermore, assuming that  $(x_i)_{i=1}^S$  are i.i.d sampled from  $\mu$ , we have  $\mu$ -almost surely

$$\lim_{S \rightarrow \infty} RCE_N^S = \left\| \nabla \mathcal{L}_{u_\theta} - \Pi_{T_N^0 \mathcal{M}}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega, \mu)} = \left\| \Pi_{[T_N^0 \mathcal{M}]^\perp}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega, \mu)}, \quad (16)$$

where  $T_N^M \mathcal{M} = \text{Span}(V_{t,i} : M \leq i \leq N)$ , while  $(V_{t,i})_{1 \leq i \leq r_{\text{svd}}}$  are the right singular-vectors of the differential  $du_\theta$  ordered in a decreasing order according to their associated singular values.

*Remark 1.* Note that  $\hat{V}_{t,i} \in \mathbb{R}^S$  for  $i \in 1, \dots, N$ , the right singular vectors of  $\hat{\phi}$ , can be seen as discretized versions of  $V_{t,i}$  from Proposition 1. Indeed, a weak convergence holds, i.e.  $\forall h \in \mathcal{H}$ ,  $\frac{1}{S} \sum_{j=1}^S \hat{V}_{t,i,j} h_j = \frac{1}{S} \sum_{j=1}^S V_{t,i}(x_j) h(x_j) \xrightarrow{S \rightarrow \infty} \langle V_{t,i}, h \rangle_{L^2}$ .

Proof of Proposition 1 can be found in Section H.3. From Proposition 1 RCE is related to the concept of *expressivity bottleneck* illustrated in Verbockhaven et al. (2024), and measures what part of the learning signal is not captured by truncating at  $N$  components for natural gradient computation. Therefore, this metric allows us to explicitly estimate and compare different cutoff choices. Note that this metric incurs no additional computational cost since ANaGRAM already computes the required SVD.

### 3.2 EMPIRICAL OBSERVATIONS: FLATTENING

Here we illustrate the evolution of training loss and reconstruction error, where Figure 1 schematically outlines key stages of ANaGRAM’s training dynamics. The plot of a real experiment is provided in Section E.

Let  $\alpha$  is a cutoff level (also referred to as precision) and  $r_{\text{cutoff}}$  denote the number of components retained by the cutoff, i.e.,  $r_{\text{cutoff}}(t) = \max\{j : \hat{\Delta}_{t,j} \geq \alpha\}$ . In Figure 1, we observe different stages of the training. First, the reconstruction error is above the wanted precision (Figure 6a). As the training progresses, the training loss drops and the reconstruction error drops until reaching the cutoff precision (Figure 6b). Eventually, the reconstruction error drops below the cutoff threshold (Figure 6c). During this phase, the training loss (corresponding to the RCE for 0 component (green line in the figure)) is not decreasing a lot.

Then, a phenomenon that we call "flattening" occurs: once the reconstruction error is small compared to the cutoff precision value, reconstruction error *flattens* over the interval  $[N_{\text{flat}}, r_{\text{cutoff}}]$ , where  $N_{\text{flat}}$  is the smallest number such as

$$RCE_{N_{\text{flat}}}^S - RCE_{r_{\text{cutoff}}}^S \approx 0. \quad (17)$$

Eventually, the phenomenon propagates toward low numbers of retained components (Figure 6e) and  $N_{\text{flat}} = 0$ . Reconstruction error is now constant for all retained components and the training ends with training loss at cutoff precision. We refer a reader to Section H.3 to have a more theoretical insight on what is happening during the flattening.

*Remark 2.* This phenomenon sheds light on the sharp drop in training loss observed near the end of optimization, as reported in Schwencke & Furtlehner (2025). By combining Equations (5), (10) and (13) and using that  $\Pi_0^0 = 0$ , we obtain

$$RCE_0^{S^2} \stackrel{13}{=} \frac{1}{S} \left\| \hat{V} \Pi_0^0 \hat{V}^T \widehat{\nabla \mathcal{L}} - \widehat{\nabla \mathcal{L}} \right\|_{\mathbb{R}^S}^2 = \frac{1}{S} \left\| \widehat{\nabla \mathcal{L}} \right\|_{\mathbb{R}^S}^2 \stackrel{10}{=} \frac{1}{S} \sum_{i=1}^S (u_\theta(x_i) - f(x_i))^2 \stackrel{5}{=} \ell(\theta). \quad (18)$$

Thus, the last iteration of flattening is **directly responsible for the sudden drop of train loss** at the end of the training.

*Remark 3.* We see that for higher precision than the cutoff value ( $N > r_{\text{cutoff}}$ ), the RCE is still decreasing as we increase the number of components kept. This indicates that there is still information to capture in the functional eigenspace composed of components associated to lower eigenvalues, see also Section H.3.

The final interesting observation is that

$$RCE_0^S - RCE_{r_{\text{cutoff}}}^S \approx 0 \quad \Leftrightarrow \quad \Pi_{r_{\text{cutoff}}}^0 \hat{V}^T \widehat{\nabla \mathcal{L}} \approx 0. \quad (19)$$

Thus, the flattening phenomenon means that the projection of the signal onto the first  $r_{\text{cutoff}}$  components retained by the cutoff is negligible. In other words, the optimization has extracted all the *usable* signal from these components at this cutoff level.

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

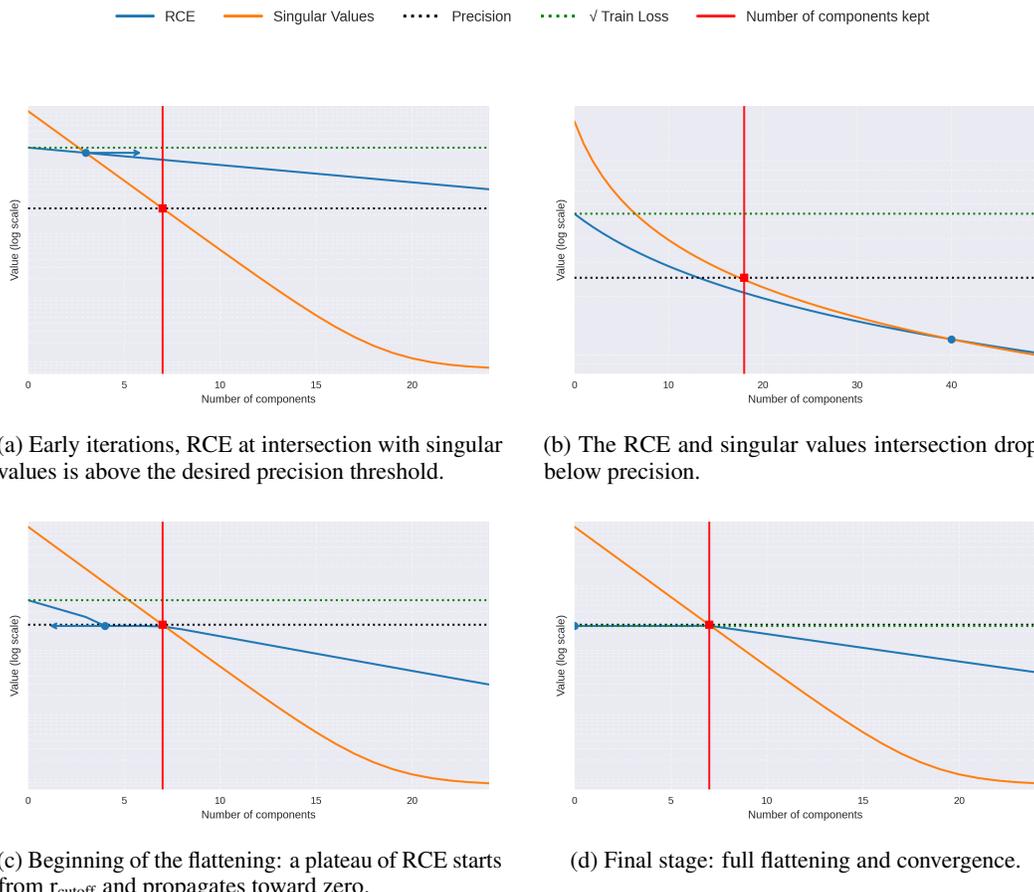


Figure 1: **ANaGRAM training dynamics.** Legend (top) and four key phases: (a) initial evolution, (b) reconstruction-singular value intersection passes target precision, (c) emergence of the flattening regime, (d) complete flattening yielding final loss level. Despite changing scale, target precision is constant and fixed across all plots. The number of ANaGRAM’s retained components  $r_{\text{cutoff}}$  is at intersection of precision line with singular values curve.

### 3.3 INCOMPLETE FLATTENING AND ADAPTIVE STRATEGIES

In practice, for some experiments we observe that the flattening may remain incomplete with  $\lim_{t \rightarrow \infty} N_{\text{flat}} = N_{\text{flat}}^{\infty} > 0$ : the system remains in a state similar to that shown in Figure 1c and never (at least not within a reasonable number of iterations) reaches the configuration illustrated in Figure 1d. A natural question arises: *what happens if we adjust the cutoff to retain exactly  $N_{\text{flat}}^{\infty}$  components?*

If we try this trick in practice (see Figure 7), then a single natural gradient step with an adjusted cutoff can be enough to get immediate and complete flattening ( $N_{\text{flat}} = 0$ ) and eventually dramatically reduce training loss. This abrupt flattening when restricting cutoff to low number of feature is typically accompanied by a learning rate found by the line search to be very close to one. A possible explanation is that this may represent an iteration in the *lazy training* regime (NTK and the feature matrix are nearly constant), where we regress linearly (and thus fast) based on learned features. [Link between flattening and lazy training is described in Appendix J.](#)

This empirical insight motivates the use of an adaptive algorithm: by dynamically adjusting cutoffs, we can hope to accelerate convergence and achieve higher precision.

---

## 4 ALGORITHMIC DESIGN: EXPLOITING FLATTENING

Building upon the empirical analysis presented in Section 3, we develop a principled algorithm that controls and exploits the flattening phenomenon identified in ANaGRAM’s training dynamics. Our approach is based on tracking the relationship between reconstruction error and singular values to automatically determine well-adapted cutoff in order to reach the target precision (error)  $\epsilon$  at the end of the training. This well-adapted cutoff should vary from one iteration to another to adjust to the currently learned weights and training dynamics in such a way to avoid early flattening (if flattening happens too early, the training stagnates at higher values of losses) and when intersection between RCE and singular values goes below the target precision  $\epsilon$ , we enforce the flattening, so that the final training loss also drops to  $\epsilon$ .

### 4.1 ADAPTIVE CUTOFF STRATEGY

In what follows, we suggest an adaptive cutoff rank  $r_{\text{cutoff}}$  that indicates how much components of  $\hat{\Delta}$  are retained for the next update of ANaGRAM. Our algorithm operates by dynamically selecting cutoff ranks based on the relationship between reconstruction error and singular values:

$$r_{\text{cutoff}}(t) = \begin{cases} r_{\text{int}}(t) := \max \{j : \text{RCE}_j^S(t) \leq \hat{\Delta}_{t,j}\} & \text{if } \text{RCE}_{r_{\text{int}}(t)}^S(t) > \epsilon \text{ (intersection rank),} \\ r_{\epsilon}(t) := \max \{j : \text{RCE}_j^S(t) \geq \epsilon\} & \text{if } \text{RCE}_{r_{\text{int}}(t)}^S(t) \leq \epsilon \text{ (precision rank).} \end{cases} \quad (20)$$

The algorithm terminates when  $r_{\epsilon}(t) = 0$ , indicating that the reconstruction error  $\text{RCE}_0^S$  that is equal to the training error is indeed below the predefined precision threshold.

**Target accuracy** The parameter  $\epsilon$  specifies the target residual precision that the algorithm aims to achieve. By design, if the algorithm converges (i.e., if the final rank reaches 0), the resulting training residual is guaranteed to lie at or below this prescribed level. As shown in Figure 2, once the smallest retained singular value reaches  $\epsilon$ , the algorithm enters the flattening regime, ensuring that the residual stabilizes beneath this threshold.

In practice, when an appropriate network architecture and a sufficiently informative set of collocation points are used (as illustrated in our Heat and Laplace experiments), choosing  $\epsilon$  close to machine precision typically yields the best results. Conversely, if convergence is not observed,  $\epsilon$  can be increased by a few orders of magnitude, accepting a lower target precision but ensuring convergence.

For ease of presentation, we provide only the core elements of AMStraMGRAM in Algorithm 1 consisting in adaptively choosing, which  $r_{\text{cutoff}}$  to apply for  $\hat{\Delta}$  at each update of ANaGRAM. The final algorithm is explained in Section C. Final Algorithm ?? addresses some irregularities observed in evolution of RCE and singular values that we explain in more details in Section C.4.

### 4.2 GEOMETRICAL INTERPRETATION OF THE ADAPTIVE STRATEGY

The algorithm exploits the geometric relationship between the empirical tangent space and the functional gradient. By tracking the intersection, we maximize the projection of the functional gradient onto the empirical tangent space while staying out of flattening. Once the intersection reach the precision level, we exploit the flattening phenomenon to achieve prescribed precision.

According to Proposition 1, the reconstruction error  $\text{RCE}_N^S$  measures how much of the functional gradient signal remains to be captured by the first  $N$  components. The intersection point thus represents the good balance between signal capture and phase transition.

## 5 EXPERIMENTAL RESULTS

We first compare in Table 1 our method implemented in JAX<sup>2</sup> with the ANaGRAM method (Schwencke & Furtlehner, 2025) on the benchmark problems presented in their paper, with modified datasets. As we see, for every equation, we perform better.

---

<sup>2</sup><https://anonymous.4open.science/r/AMStraMGRAM-8D1B/>

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

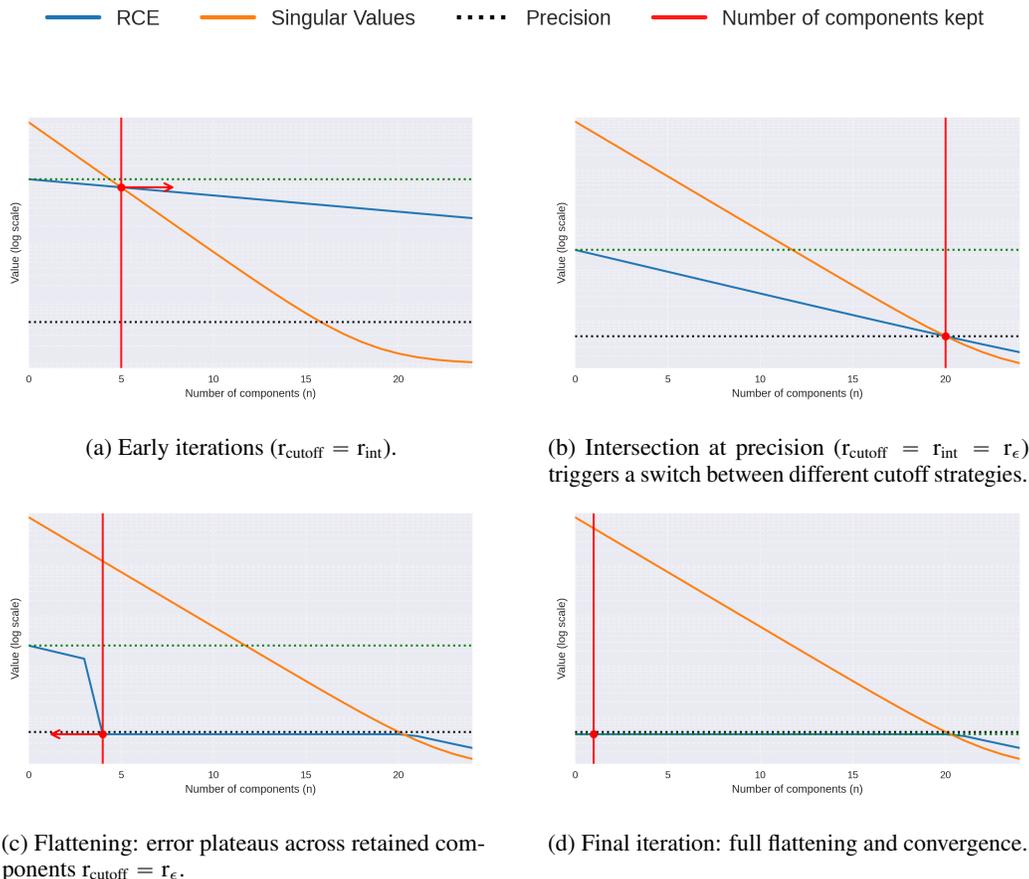


Figure 2: **Dynamics of the adaptive multi-cutoff strategy in AMStramGRAM.** Progression from (a) initial exploration, (b) intersection reaches precision, (c) flattening onset, to (d) converged state. Red arrows (when present) indicate the retained rank dynamics (pointing right – increasing, pointing left – decreasing). Legends are shown below.

Table 1: Performance comparison between AMStramGRAM (our method) and ANaGRAM Schwencke & Furtlehner (2025). The adaptive strategy demonstrates significant improvements across all benchmark problems, with  $L_2$  error improvements of up to 8 orders of magnitude.

Experiment	Mean Squared Error (MSE)		$L_2$ Error	
	Ours	ANaGRAM	Ours	ANaGRAM
Heat Equation	<b>6.29e-29</b> ± <b>6.78e-30</b>	8.56e-11 ± 7.05e-11	<b>2.32e-14</b> ± <b>1.14e-14</b>	1.28e-06 ± 1.75e-06
Laplace 2D	<b>1.46e-28</b> ± <b>1.87e-29</b>	4.27e-13 ± 4.66e-13	<b>2.24e-15</b> ± <b>2.52e-16</b>	3.49e-09 ± 3.58e-09
Laplace 5D	<b>2.04e-08</b> ± <b>1.16e-08</b>	6.37e-08 ± 7.01e-08	<b>2.12e-05</b> ± <b>8.15e-06</b>	4.00e-05 ± 2.93e-05
Allen–Cahn	<b>3.19e-11</b> ± <b>2.37e-11</b>	2.19e-04 ± 4.16e-04	<b>5.87e-05</b> ± <b>6.25e-06</b>	4.32e-03 ± 5.93e-03

We then compare our method with the baseline methods from Urbán et al. (2025) on the benchmark problems presented in their paper. Note that in our case we do not need to enforce boundary constraints. The methodology of sampling is also slightly different, as we sample the data from a fixed grid, following the methodology of Schwencke & Furtlehner (2025), while in Urbán et al. (2025) they perform batching of randomly sampled points.

Table 2: Performance comparison between AMStraMGRAM (our method) and baseline Urbán et al. (2025) methods. Our method demonstrates improvements across benchmark problems, without requiring enforcement of boundary constraints.

Experiment	Mean Squared Error (MSE)		$L_2$ Error	
	Ours	SSBroyden*	Ours	SSBroyden*
One-dimensional Burgers (1DB)	<b>2.99e-12</b> $\pm$ <b>9.26e-13</b>	2.92e-10 $\pm$ 1.45e-10	<b>1.5e-06</b> $\pm$ <b>9.43e-7</b>	1.59e-06 $\pm$ 1.02e-6
Non-Linear Poisson (k=1)	<b>8.51e-24</b> $\pm$ <b>2.24e-24</b>	3.03e-16 $\pm$ 3.82e-16	6.81e-10 $\pm$ 1.41e-09	<b>9.29e-12</b> $\pm$ <b>5.85e-12</b>
Allen–Cahn (AC)	3.19e-11 $\pm$ 2.37e-11	<b>6.42e-12</b> $\pm$ <b>5.52e-12</b>	5.87e-05 $\pm$ 6.25e-06	<b>3.94e-06</b> $\pm$ <b>1.72e-06</b>

\* refer to method from Urbán et al. (2025) with adaptive sampling and hard constraint enforcement on boundary conditions.

## 6 LIMITATIONS

Despite its effectiveness, AMStraMGRAM can exhibit overfitting, particularly in problems with sharp features like the Allen–Cahn equation. The algorithm drives the training error to machine precision on the sampled points, but the learned function may develop high-frequency oscillations between them, especially in regions of high curvature where the approximation is the most challenging. These artifacts, visible as “overfitting lines” in Figure 3, are an imprint of the sampling lattice (see regions around  $x = \pm 0.5$ ). They arise because the SVD cutoff effectively projects the update onto a low-rank subspace of the tangent space. This subspace is often aligned with the grid axes, leading to anisotropic smoothing that perfectly fits the data on the grid lines but interpolates poorly in the under-sampled regions between them. Once the flattening phase begins, the training enters a quasi-linear regime that can “lock in” these geometric artifacts.

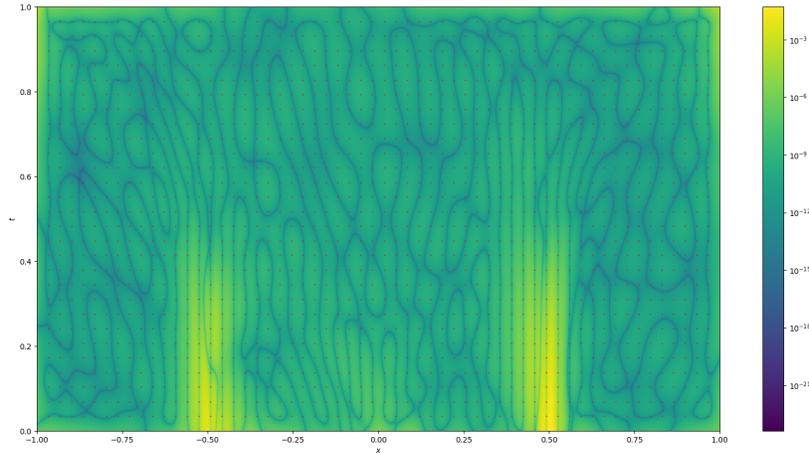


Figure 3: Allen–Cahn overfitting: residual lines align with sampling lines. Low-rank (post-cutoff) tangent projections fit exactly on sampled fibers while interpolation between them inherits weakly constrained oscillations in regions of steep interface curvature.

This phenomenon highlights that while our method significantly improves on ANaGRAM, the quality of the final solution remains fundamentally limited by the sampling strategy. Mitigating such overfitting requires co-designing the sampler and the optimizer. Potential remedies include adaptive sampling, where new collocation points are added in regions of high reconstruction error, or curriculum-based approaches that progressively refine the sampling grid.

**Disentangling solution accuracy from sampling error.** Our method attains machine precision on the residual and not on the error. This last metric depends crucially on the choice of points. Now that limitations of the optimizer are not the main bottleneck anymore, the investigation of a proper sampling/quadrature strategy should be the next step to further reduce  $L^2$  error. We illustrate the relationship between collocation density (expressed as number of points per dimension) and error in Figure 4.

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

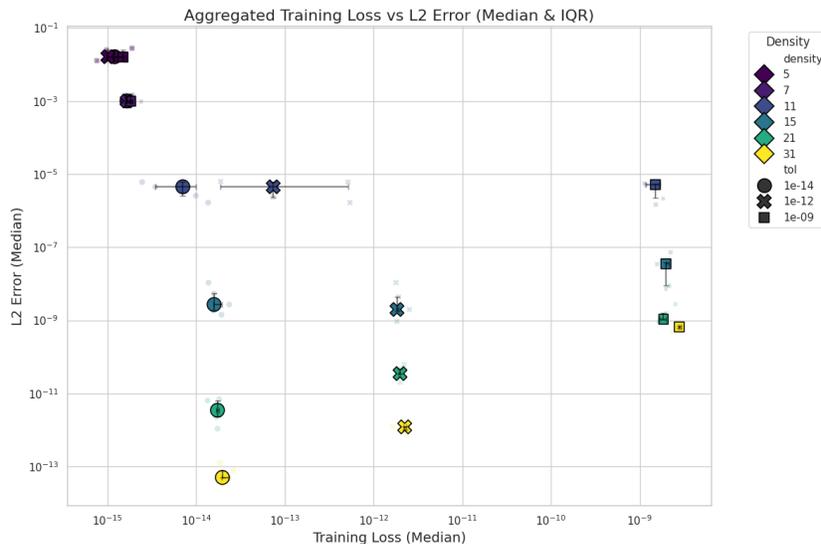


Figure 4: Relationship between collocation density and error for the Heat equation with fixed model size. The mesh density corresponds to the number of points per dimension in the interior of the domain. The ‘tol’ values correspond to the desired precision (cutoff). We observe that increasing the number of points improves the  $L_2$  error without significant change in the residual, suggesting that the optimizer effectively minimizes the loss on the available points, and the remaining error is dominated by the discretization (sampling) error.

**Stability Under Resampling** While AMStraMGRAM employs a fixed grid to ensure deterministic updates and stable cutoff tracking, standard PINN training often relies on minibatch sampling. In a batchwise setting, the truncated natural-gradient update becomes stochastic not only because of noise in gradient estimation, but also due to fluctuations in the empirical feature matrix. In Appendix I, we provide additional experiments illustrating that different stochastic mini-batches induce different empirical tangent spaces. Our results indicate that when the cutoff rank is either very low or very high, the choice of samples has only a minor effect on the resulting empirical tangent space. In contrast, for intermediate cutoff values, different batches can produce noticeably different update directions, and consequently lead to distinct training dynamics.

These observations suggest that developing principled stochastic batching schemes or adaptive sampling strategies for our algorithm constitutes an interesting direction for future work.

## 7 CONCLUSION

In this work, we have introduced AMStraMGRAM, an adaptive multi-cutoff strategy that enhances the ANaGRAM natural gradient method for training PINNs. Our work provides an analytical framework to explain ANaGRAM’s convergence behavior, uncovering a *flattening* phenomenon that clarifies its training dynamics. The proposed algorithm automatically adjusts cutoff regularization. Notably, AMStraMGRAM exhibits “overfitting” as demonstrated in Allen-Cahn experiments. These results underscore the potential of natural gradient optimization for PINNs while highlighting the critical role of sampling strategies in realizing their full accuracy.

Future research will focus on integrating residual-based methods to further stabilize training, establishing rigorous dynamics analysis for the feature-development phase, and extending the approach to higher-dimensional PDEs and complex geometries. Exploring the interplay between network architecture and optimization—as well as further developing sampling techniques—will be essential to address the fundamental challenge of balancing optimization power with data representation. Ultimately, our findings suggest that with careful algorithmic design, PINNs can achieve the precision required for practical scientific computing, paving the way for mesh-free methods in computational science.

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593

---

## REFERENCES

- Ben Adcock and Daan Huybrechs. Frames and Numerical Approximation. *SIAM Review*, 61(3): 443–473, January 2019. ISSN 0036-1445, 1095-7200. doi: 10.1137/17M1114697.
- Ben Adcock and Daan Huybrechs. Frames and numerical approximation II: Generalized sampling, July 2020.
- M. Al-Baali. Numerical Experience with a Class of Self-Scaling Quasi-Newton Algorithms. *Journal of Optimization Theory and Applications*, 96(3):533–553, March 1998. ISSN 0022-3239, 1573-2878. doi: 10.1023/A:1022608410710.
- Mehiddin Al-Baali and Humaid Khalfan. Wide interval for efficient self-scaling quasi-Newton algorithms. *Optimization Methods and Software*, 20(6):679–691, December 2005. ISSN 1055-6788, 1029-4937. doi: 10.1080/10556780410001709448.
- Shun-Ichi Amari. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Yurij M. Berezansky, Zinovij G. Sheftel, and Georgij F. Us. *Functional Analysis. Vol. II*, volume 86 of *Operator Theory Advances and Applications*. Birkhäuser, 1996.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- Salvatore Cuomo, Vincenzo Schiano Di Cola, Fabio Giampaolo, Gianluigi Rozza, Maziar Raissi, and Francesco Piccialli. Scientific Machine Learning Through Physics-Informed Neural Networks: Where we are and What’s Next. *Journal of Scientific Computing*, 92(3):88, July 2022. ISSN 1573-7691. doi: 10.1007/s10915-022-01939-z.
- Tim De Ryck, Florent Bonnet, Siddhartha Mishra, and Emmanuel de Bézenac. An operator preconditioning perspective on training in physics-informed machine learning, May 2024.
- Dean G. Duffy. *Green’s Functions with Applications*. Chapman and Hall/CRC, 2015.
- Tamara G Grossmann, Urszula Julia Komorowska, Jonas Latz, and Carola-Bibiane Schönlieb. Can physics-informed neural networks beat the finite element method? *IMA Journal of Applied Mathematics*, 89(1):143–174, January 2024. ISSN 0272-4960. doi: 10.1093/imamat/hxae011.
- Andrés Guzmán-Cordero, Felix Dangel, Gil Goldshlager, and Marius Zeinhofer. Improving Energy Natural Gradient Descent through Woodbury, Momentum, and Randomization, May 2025.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
- Anas Jnini, Flavio Vella, and Marius Zeinhofer. Gauss-Newton Natural Gradient Descent for Physics-Informed Computational Fluid Dynamics, February 2024.
- Andreas Kirsch. *An Introduction to the Mathematical Theory of Inverse Problems*, volume 120 of *Applied Mathematical Sciences*. Springer International Publishing, Cham, 2021. ISBN 978-3-030-63342-4 978-3-030-63343-1. doi: 10.1007/978-3-030-63343-1.
- Elham Kiyani, Khemraj Shukla, Jorge F. Urbán, Jérôme Darbon, and George Em Karniadakis. Which Optimizer Works Best for Physics-Informed Neural Networks and Kolmogorov-Arnold Networks?, April 2025.
- Rainer Kress. *Linear Integral Equations*, volume 82 of *Applied Mathematical Sciences*. Springer, New York, NY, 2014. ISBN 978-1-4614-9592-5 978-1-4614-9593-2. doi: 10.1007/978-1-4614-9593-2.
- Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, and Michael W Mahoney. Characterizing possible failure modes in physics-informed neural networks. In *Advances in Neural Information Processing Systems*, volume 34, pp. 26548–26560. Curran Associates, Inc., 2021.

---

594 Moshe Leshno, Vladimir Ya Lin, Allan Pinkus, and Shimon Schocken. Multilayer feedforward  
595 networks with a nonpolynomial activation function can approximate any function. *Neural networks*,  
596 6(6):861–867, 1993.

597  
598 Dong C. Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization.  
599 *Mathematical Programming*, 45(1):503–528, August 1989. ISSN 1436-4646. doi: 10.1007/  
600 BF01589116.

601 Songming Liu, Chang Su, Jiachen Yao, Zhongkai Hao, Hang Su, Youjia Wu, and Jun Zhu. Precondi-  
602 tioning for Physics-Informed Neural Networks, February 2024.

603  
604 Levi McClenny and Ulisses Braga-Neto. Self-Adaptive Physics-Informed Neural Networks using a  
605 Soft Attention Mechanism, April 2022.

606  
607 Johannes Müller and Marius Zeinhofer. Achieving high accuracy with PINNs via energy natural  
608 gradient descent. In *International Conference on Machine Learning*, pp. 25471–25485. PMLR,  
609 2023.

610  
611 Johannes Müller and Marius Zeinhofer. Position: Optimization in SciML Should Employ the Function  
612 Space Geometry. In *Forty-First International Conference on Machine Learning*, February 2024.

613  
614 Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 1999.

615  
616 Yann Ollivier. Riemannian metrics for neural networks I: Feedforward networks. *Information and*  
617 *Inference: A Journal of the IMA*, 4(2):108–153, 2015.

618  
619 Vern I. Paulsen and Mrinal Raghupathi. *An Introduction to the Theory of Reproducing Kernel Hilbert*  
620 *Spaces*, volume 152. Cambridge university press, 2016.

621  
622 M. Raissi, P. Perdikaris, and G.E. Karniadakis. Physics-informed neural networks: A deep learning  
623 framework for solving forward and inverse problems involving nonlinear partial differential  
624 equations. *Journal of Computational Physics*, 378:686–707, February 2019. ISSN 00219991. doi:  
625 10.1016/j.jcp.2018.10.045.

626  
627 Pratik Rathore, Weimu Lei, Zachary Frangella, Lu Lu, and Madeleine Udell. Challenges in Training  
628 PINNs: A Loss Landscape Perspective, February 2024.

629  
630 Nilo Schwencke and Cyril Furtlehner. ANaGRAM: A natural gradient relative to adapted model for  
631 efficient PINNs learning. In *The Thirteenth International Conference on Learning Representations*,  
632 2025.

633  
634 Jorge F. Urbán, Petros Stefanou, and José A. Pons. Unveiling the optimization process of physics  
635 informed neural networks: How accurate and competitive can PINNs be? *Journal of Computational*  
636 *Physics*, 523:113656, February 2025. ISSN 0021-9991. doi: 10.1016/j.jcp.2024.113656.

637  
638 Manon Verbockhaven, Sylvain Chevallier, and Guillaume Charpiat. Growing tiny networks: Spotting  
639 expressivity bottlenecks and fixing them optimally. *arXiv preprint arXiv:2405.19816*, 2024.

640  
641 Sifan Wang, Yujun Teng, and Paris Perdikaris. Understanding and mitigating gradient flow pathologies  
642 in physics-informed neural networks. *SIAM Journal on Scientific Computing*, 43(5):A3055–A3081,  
643 2021.

644  
645 Sifan Wang, Shyam Sankaran, Hanwen Wang, and Paris Perdikaris. An Expert’s Guide to Training  
646 Physics-informed Neural Networks, August 2023.

647  
648 Sifan Wang, Bowen Li, Yuhan Chen, and Paris Perdikaris. PirateNets: Physics-informed Deep  
649 Learning with Residual Adaptive Networks. *Journal of Machine Learning Research*, 25(402):  
650 1–51, 2024. ISSN 1533-7928.

651  
652 Chenxi Wu, Min Zhu, Qinyang Tan, Yadhu Kartha, and Lu Lu. A comprehensive study of non-  
653 adaptive and residual-based adaptive sampling for physics-informed neural networks. *Computer*  
654 *Methods in Applied Mechanics and Engineering*, 403:115671, 2023.

## A ILLUSTRATION OF NATURAL GRADIENT

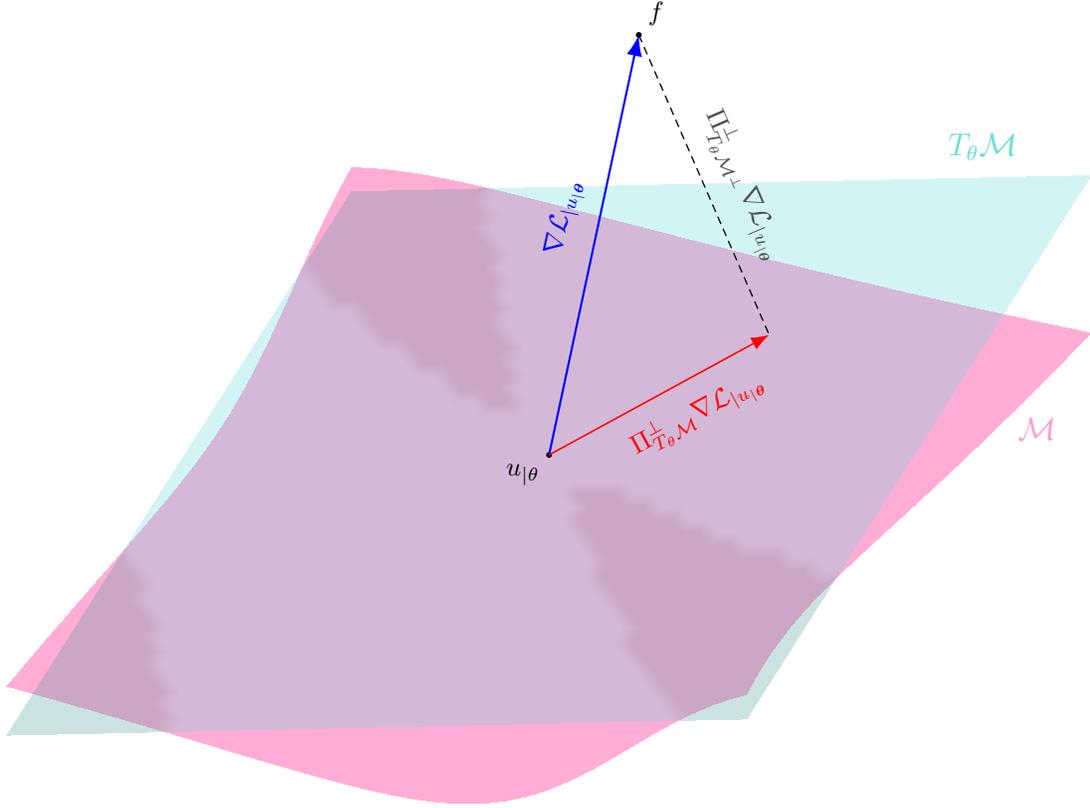


Figure 5: Illustration of the orthogonal projection of the functional gradient onto the tangent space. While the ideal update direction would be the functional gradient  $\nabla \mathcal{L}|_{u_\theta}$  (shown in blue), our model constrains us to follow directions within the tangent space  $T_\theta \mathcal{M}$  (shown as a green plane). The optimal feasible direction is thus the orthogonal projection  $\Pi_{T_\theta \mathcal{M}}^\perp (\nabla \mathcal{L}|_{u_\theta})$  (shown in red).

## B OUR VOCABULARY

- **Domain** ( $\Omega$ ).
- **Boundary** ( $\partial\Omega$ ).
- **Differential operators** ( $D, B$ ).
- **Cutoff** ( $\alpha_t$ ). A threshold below which the components of the matrix  $\hat{\Delta}$  are truncated, *i.e.*

$$\hat{\Delta} \leftarrow \begin{cases} \hat{\Delta} & \text{if } \hat{\Delta} \geq \alpha_t, \\ 0 & \text{else.} \end{cases}$$
- **Full rank** ( $\mathbf{r}_{\text{svd}}$ ). A full rank of feature matrix  $\hat{\phi}$  that we assume, without loss of generality, to be equal to  $\min(P, S)$ .
- **Rank** ( $\mathbf{r}_{\text{cutoff}}$ ). A number of  $\hat{\Delta}$  components that are retained when computing a pseudo-inverse of  $\hat{\Delta}$  in ANaGRAM. Depending on a current regime of the training and a desired effect, it can be set at  $r_{\text{int}}$  or  $r_\epsilon$ .
- **Flattening**. The phenomenon described in Section 3.2, when reconstruction error starts to stabilize for a range of possible ranks.
- **Flat cutoff** ( $N_{\text{flat}}$ ). A number of components that corresponds to the beginning of flattening in reconstruction error curve.

- 702 • **Feature matrix** ( $\hat{\phi} \in \mathbb{R}^{P \times S}$ ). It is defined by a jacobian  $\partial_p u_{\theta}(x_i)$ , which is used in an
- 703 ANaGRAM's update to "project" a functional gradient onto parameter space of  $\theta$ .
- 704 • **Precision** ( $\epsilon$ ). A hyperparameter of AMStraMGRAM that prescribes a target error level that
- 705 the algorithm should achieve.
- 706 • **Intersection rank** ( $r_{\text{int}}$ ). Defined in Equation (20), roughly speaking it corresponds to
- 707 a number of components at which reconstruction error and singular values curves are
- 708 intersecting.
- 709 • **Precision rank** ( $r_{\epsilon}$ ). Defined in Equation (20), it corresponds to a number of components at
- 710 which reconstruction error curve and precision level are intersecting.
- 711 • **Functional gradient** ( $\nabla \mathcal{L}$ ). A Frechet derivative of squared  $L^2$  loss  $\mathcal{L}$ , its negative gives
- 712 the "ideal" update direction in non-parametric case.
- 713 • **Empirical functional gradient** ( $\widehat{\nabla \mathcal{L}} \in \mathbb{R}^S$ ). A vector obtain by evaluating  $\nabla \mathcal{L}$  on some
- 714 finite number of samples  $x_i \in \Omega$ , for  $i \in 1, \dots, S$ .
- 715 • **Parametric model** ( $u_{\theta}$ ). A function parametrized with  $\theta$  that serves to approximate a
- 716 solution to a problem (regression or PDE). Typically, it is a neural network, where  $\theta$  are its
- 717 full set of weights.
- 718 • **Differential of the model** ( $du_{\theta}$ ). Defined as  $du_{\theta}(h) = \sum_{p=1}^P h_p \frac{\partial u}{\partial \theta_p} = \lim_{\epsilon \rightarrow 0} \frac{u_{|\theta+\epsilon h} - u_{\theta}}{\epsilon}$ . It
- 719 measures how much  $u_{\theta}$  changes in a given direction  $h$ .
- 720 • **Tangent space** ( $T_{\theta} \mathcal{M}$ ). Image of a differential of the model, giving a space of possible
- 721 updates for a model  $u_{\theta}$ .
- 722 • **SVD components of  $\hat{\phi}$**  ( $\widehat{U}, \widehat{\Delta}, \widehat{V}$ ). In particular,  $\hat{\phi} = \widehat{U} \widehat{\Delta} \widehat{V}^T$ , where  $\widehat{U} \in \mathbb{R}^{P \times S}$  is a left
- 723 singular vector matrix,  $\widehat{\Delta} \in \mathbb{R}^{\text{rsvd} \times \text{rsvd}}$  is a diagonal matrix with singular values on a diagonal
- 724 ordered in a decreasing order and  $\widehat{V}$  is a right singular vector matrix.
- 725 • **Functional singular vectors** ( $V_{t,i}$ ). Right singular vectors of the differential  $du_{\theta}$ .
- 726 • **Empirical tangent space** ( $T_N^M \mathcal{M}$ ). A subspace of tangent space  $T_{\theta} \mathcal{M}$ , restricted to a span
- 727 of the right functional singular vectors  $V_{t,i}$  corresponding to a range of components from  $M$
- 728 to  $N$ , i.e.  $\text{Span}(V_{t,i} : 1 \leq M \leq N \leq N)$ .
- 729 • **Discretized empirical tangent space** ( $\widehat{T_N^M \mathcal{M}}$ ). A version of  $T_N^M \mathcal{M}$  discretized on a set of
- 730 samples  $\{x_i\}_{i=1}^S$  coming from  $\Omega$ .
- 731 • **Reconstruction error** ( $\text{RCE}_N^S$ ). A measure identifying the portion of the functional gradient
- 732 signal that is lost when restricting  $\widehat{\nabla \mathcal{L}}$  to  $\widehat{T_N^Q \mathcal{M}}$ .
- 733 • **Feature development phase**. The early phase in the training, during which high volatility
- 734 is observed in both quantities of interest with high sensitivity to the choice of  $r_{\text{cutoff}}$ .
- 735 • **Flattening phase**. The later phase in the training, during which reconstruction error starts to
- 736 flatten for some values of  $N$ , at the same time singular values dominate over reconstruction
- 737 error for all retained components, resulting in a drop of training loss.

## 743 C PRACTICAL IMPLEMENTATION CONSIDERATIONS

744 While the principled algorithm discussed in the main paper and summarized in Algorithm 1 provides  
 745 a sound framework, empirical observations reveal that additional mechanisms are necessary for robust  
 746 performance across diverse PDE problems. This section describes additional modifications to make  
 747 the algorithm more practical.

### 750 C.1 THE DUAL CUTOFF STRATEGY: ADDRESSING EMPIRICAL CHALLENGES

751 Our experiments reveal that the single cutoff approach, while theoretically elegant, suffers from  
 752 numerical instabilities and incomplete convergence in practice. We observed three critical issues:

- 753 1. **Ignition failure:** The intersection between reconstruction error and singular values some-
- 754 times fails to evolve, preventing the algorithm from reaching lower error values.

- 
- 756           2. **Retreating dynamics:** The intersection rank may decrease during training, disrupting  
757           convergence.  
758  
759           3. **Incomplete flattening:** Without additional stabilization, the flattening phenomenon may  
760           not complete, leading to suboptimal final accuracy.  
761

762  
763 To address these challenges, we introduce a dual cutoff strategy inspired by the staged design of  
764 rocket launches:  
765

## 766 C.2 THREE-PHASE TRAINING DYNAMICS

### 767 C.2.1 IGNITION PHASE

768 We initialize two cutoffs:  
769

- 770 • **Minimum cutoff ( $r_{\min}$ ):** Set at the intersection point  $r_{\text{int}}(t)$   
771
- 772 • **Maximum cutoff ( $r_{\max}$ ):** Set at the "elbow" of the singular value curve (see algorithm 4)  
773

774  
775 The algorithm performs two natural gradient steps per iteration, one with each cutoff. If the intersec-  
776 tion position remains static after both updates, we increment  $r_{\max}$  by one to promote exploration of  
777 additional gradient components.  
778

779 This phase ends when  $r_{\min}$  reaches  $r_{\max}$ —an event we term **liftoff**.  
780

### 781 C.2.2 ASCENT PHASE

782 During ascent, both cutoffs track the moving intersection, but with a stability mechanism:  
783

$$784 \quad r_{\max}(t) = \max(r_{\max}(t-1), r_{\text{int}}(t)). \quad (21)$$

785 This monotonicity constraint prevents the intersection rank from falling to zero, which would disrupt  
786 training dynamics.  
787

### 788 C.2.3 STAGE SEPARATION AND PRECISION LOCKING

789 When  $\text{RCE}_{r_{\text{int}}(t)}^S(t) \leq \epsilon$ , we trigger **stage separation**:  
790

- 791 •  $r_{\min}$  is fixed at the precision level:  $r_{\min} = r_{\epsilon}(t)$   
792
- 793 •  $r_{\max}$  continues tracking the intersection to maintain stability  
794

795 The algorithm continues until  $r_{\min} = 0$  (**booster return**), indicating complete convergence. The final  
796 algorithm that combines all three stages is mentioned in Algorithm ??.  
797

---

810 **Algorithm 1:** Sketch of the Adaptative MultiCutoff Strategy for ANaGRAM (AMStraMGRAM)

---

811 **Input:**  $u_{\theta} : \mathbb{R}^P \rightarrow L^2(\Omega, \mu)$ ,  $\theta_0 \in \mathbb{R}^P$ ,  $f \in L^2(\Omega, \mu)$ ,  $(x_i) \in \Omega^S$ ,  $\epsilon > 0$ ,  $T_{\max} \in \mathbb{N}$

812 // Initialization

813 1  $t \leftarrow 0$

814 2  $\hat{\phi}_0 \leftarrow (\partial_p u_{\theta_0}(x_i))_{i,p}$  for  $i \in 1, \dots, S$  and  $p \in 1, \dots, P$

815 3  $\hat{U}_0, \hat{\Delta}_0, \hat{V}_0^T \leftarrow \text{SVD}(\hat{\phi}_0)$

816 4  $\widehat{\nabla} \mathcal{L}_0 \leftarrow (u_{\theta_0}(x_i) - f(x_i))_i$  for  $i \in 1, \dots, S$

817 5 Compute  $(\text{RCE}_j^S)$  for all  $j \in 1, \dots, r_{\text{svd}}$  following Equation (13)

818 6 **repeat**

819     // Compute adaptive ranks

820     7 Compute  $r_{\text{int}}$  and  $r_{\epsilon}$  using expressions from Equation (20)

821     // Determine a final cutoff rank

822     8 **if**  $\text{RCE}_{r_{\text{int}}}^S > \epsilon$  **then**

823         9      $r_{\text{cutoff}} \leftarrow r_{\text{int}}$                              // Track intersection

824     10 **else**

825         11      $r_{\text{cutoff}} \leftarrow r_{\epsilon}$                              // Lock on precision

826     // Natural gradient step

827     12 Set  $\hat{\Delta}_t \leftarrow \begin{cases} \hat{\Delta}_{t,i} & \text{if } i \leq r_{\text{cutoff}}, \\ 0 & \text{else;} \end{cases}$

828     13 Get new  $\theta_{t+1}$  after one ANaGRAM step with Equation (10)

829     // Update for next iteration

830     14  $\hat{\phi}_{t+1} \leftarrow (\partial_p u_{\theta_{t+1}}(x_i))_{i,p}$

831     15  $\hat{U}_{t+1}, \hat{\Delta}_{t+1}, \hat{V}_{t+1}^T \leftarrow \text{SVD}(\hat{\phi}_{t+1})$

832     16  $\widehat{\nabla} \mathcal{L}_{t+1} \leftarrow (u_{\theta_{t+1}}(x_i) - f(x_i))_i$

833     17 Recompute  $\text{RCE}_j^S$  for all  $j \in 1, \dots, r_{\text{svd}}$  following Equation (13)

834     18  $t \leftarrow t + 1$

835 19 **until**  $r_{\epsilon} = 0$  or  $t \geq T_{\max}$

836 **Output:**  $\theta_t$

---

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

---

### C.3 COMPLETE PRACTICAL ALGORITHM

---

**Algorithm 2:** AMStrAMGRAM : Adaptive Multicutoff Strategy Modification for ANaGRAM

---

**Input:** •  $u : \mathbb{R}^P \rightarrow \mathbb{L}^2(\Omega, \mu)$  // neural network architecture

- $\theta_0 \in \mathbb{R}^P$  // initialization of the neural network
- $f \in \mathbb{L}^2(\Omega, \mu)$  // target function of the quadratic regression
- $(x_i) \in \Omega^S$  // a batch in  $\Omega$
- $\epsilon > 0$  // precision level of the optimization

```

874 1 begin Initialization
875 2    $\lambda \leftarrow False$  // Liftoff indicator
876 3    $\hat{\phi}_{\theta_0} \leftarrow (\partial_p u_{\theta_0}(x_i))_{1 \leq i \leq S, 1 \leq p \leq P}$  // Computed via auto-differentiation
877 4    $\hat{U}_{\theta_0}, \hat{\Delta}_{\theta_0}, \hat{V}_{\theta_0}^t \leftarrow \text{SVD}(\hat{\phi}_{\theta_0})$ 
878 5    $\widehat{\nabla} \mathcal{L}_{\theta_0} \leftarrow (u_{\theta_0}(x_i) - f(x_i))_{1 \leq i \leq S}$ 
879 6    $\text{RCE}_0^S \leftarrow \text{ReconstructionErrors}(\hat{V}_{\theta_0}^t, \widehat{\nabla} \mathcal{L}_{\theta_0})$ 
880 7    $r_{\max 0} \leftarrow \text{FindElbow}((1, \dots, r_{\text{svd}}), \hat{\Delta}_{\theta_0})$ 
881
882 8 repeat
883 9    $r_{1t} \leftarrow \# \left\{ \text{RCE}_{0_j}^S \leq \hat{\Delta}_{\theta_{t_j}} : 1 \leq j \leq r_{\text{svd}} \right\}$ 
884 10   $r_{2t} \leftarrow \# \left\{ \text{RCE}_{0_j}^S \geq \epsilon : 1 \leq j \leq r_{\text{svd}} \right\}$ 
885 11  /* with # standing for the cardinal */
886 12   $r_{\min t} \leftarrow \min(r_{1t}, r_{2t})$ 
887 13   $r_{\max t} \leftarrow \max(r_{1t}, r_{\max t-1})$ 
888 14  if not  $\lambda_t$  then
889 15  | if  $r_{\min t} \geq r_{\max t}$  then
890 16  | |  $\lambda_t \leftarrow True$ 
891 17  | else if  $r_{\min t-1} = r_{\min t}$  then
892 18  | |  $r_{\max t} \leftarrow r_{\max t} + 1$ 
893 19  | foreach  $r_{\text{cutoff}} \in \{r_{\max t}, r_{\min t}\}$  do
894 20  | |  $\hat{\Delta}_{\theta_t} \leftarrow (\hat{\Delta}_{\theta_{t,p}} \text{ if } p \geq r_{\text{cutoff}} \text{ else } 0)_{1 \leq p \leq P}$ 
895 21  | |  $\widehat{\nabla} \mathcal{L}_{\theta_t} \leftarrow (u_{\theta_t}(x_i) - f(x_i))_{1 \leq i \leq S}$ 
896 22  | |  $d_{\theta_t} \leftarrow \hat{V}_{\theta_t} \hat{\Delta}_{\theta_t}^{\dagger} \hat{U}_{\theta_t}^t \widehat{\nabla} \mathcal{L}_{\theta_t}$ 
897 23  | |  $\eta_t \leftarrow \arg \min_{\eta \in \mathbb{R}^+} \sum_{1 \leq i \leq S} (f(x_i) - u_{\theta_t - \eta d_{\theta_t}}(x_i))^2$  // via line search
898 24  | |  $\theta_{t+1} \leftarrow \theta_t - \eta_t d_{\theta_t}$ 
899 25  | |  $\hat{\phi}_{\theta_{t+1}} \leftarrow (\partial_p u_{\theta_{t+1}}(x_i))_{1 \leq i \leq S, 1 \leq p \leq P}$  // Computed
900 26  | | via auto-differentiation
901 27  | |  $\hat{U}_{\theta_{t+1}}, \hat{\Delta}_{\theta_{t+1}}, \hat{V}_{\theta_{t+1}}^t \leftarrow \text{SVD}(\hat{\phi}_{\theta_{t+1}})$ 
902
903 28 until  $r_{1t} = 0$  or  $t \geq T_{\max}$ 

```

---

### C.4 EMPIRICAL JUSTIFICATION FOR DESIGN CHOICES

The dual cutoff strategy addresses specific empirical challenges we observed:

**Dual gradient steps:** Without the second cutoff, training dynamics sometimes stagnate. The dual approach provides both stability (via  $r_{\min}$ ) and exploration (via  $r_{\max}$ ).

**Elbow initialization:** The elbow point marks where singular values cease contributing meaningful signal, providing a natural upper bound for exploration.

**Monotonic  $r_{\max}$ :** Prevents catastrophic retreat of the intersection point, which we observed in complex equations like Allen-Cahn.

918 **Stage separation timing:** Triggered precisely when the intersection error drops below target precision,  
 919 ensuring optimal utilization of the flattening phenomenon.  
 920

921 We see in the next section how this practical algorithm successfully improve empirical robustness.  
 922  
 923

## 924 D ALGORITHMIC DETAILS

---

### 928 Algorithm 3: Find elbow

---

```

929 1 Function FindElbow
930   Input:-  $(x_i) \in \mathbb{R}^m$  // an increasing sequence of  $m \in \mathbb{N}$  points in  $\mathbb{R}$ 
931           -  $\hat{f} \in \mathbb{R}^m$  // a decreasing function evaluated at points  $(x_i)$ 
932
933   /* Clockwise normal vector to  $(x_m - x_1, \hat{f}_m - \hat{f}_1)$  */
934   2  $\vec{n} \leftarrow (\hat{f}_m - \hat{f}_1, x_1 - x_m) \in \mathbb{R}^2$ 
935   3  $(s_j)_{1 \leq j \leq m} \leftarrow \left( \left\langle \vec{n}, (x_j - x_1, \hat{f}_j - \hat{f}_1) \right\rangle_{\mathbb{R}^2} \right)_{1 \leq j \leq m}$ 
936   Output:  $\arg \max_{1 \leq j \leq m} s_j$ 
937
938 4 end

```

---

### 942 Algorithm 4: Reconstruction Errors

---

```

943 1 Function ReconstructionErrors
944   Input:-  $\hat{V}^t \in \mathbb{R}^{r_{\text{svd}}, S}$  // right singular vectors of the Jacobian  $\hat{\phi}$ 
945           -  $\widehat{\nabla \mathcal{L}} \in \mathbb{R}^S$  // Evaluated functional gradient
946
947   2 begin Initialization
948     3  $\hat{\Sigma} \leftarrow 0 \in \mathbb{R}^S$  // cumulative approximation of  $\widehat{\nabla \mathcal{L}}$ 
949     4  $\text{RCE}^S \leftarrow 0 \in \mathbb{R}^{r_{\text{svd}}}$  // cumulated reconstruction errors
950     5  $\hat{c} \leftarrow \hat{V}^t \widehat{\nabla \mathcal{L}} \in \mathbb{R}^{r_{\text{svd}}}$ 
951   6 end
952   7 foreach  $j \in (1, \dots, r_{\text{svd}})$  do
953     8  $\hat{\Sigma} \leftarrow \hat{\Sigma} + \hat{c}_j$ 
954     9  $\text{RCE}_j^S \leftarrow \|\hat{\Sigma} - \hat{c}\|_2$ 
955   10 end
956   Output:  $\text{RCE}^S$ 
957
958 11 end

```

---

## 962 E EMPIRICAL EXAMPLE OF ANAGRAM TRAINING DYNAMICS

963  
 964  
 965 In Figure 6, we analyze ANaGRAM’s training on the heat equation with a fixed cutoff threshold  
 966  $\alpha = 10^{-3}$  and line search for the learning rate. The training loss coincides with  $\|\widehat{\nabla \mathcal{L}}\|^2$ . We can  
 967 see the flattening phenomenon to occur on Iteration 120 and completed at 150. As discussed in the  
 968 main paper, sometimes the flattening can be incomplete, and for many iterations remain without any  
 969 further progress ( $N_{\text{flat}}$  never reaching zero). In this case, changing a cutoff threshold results in an  
 970 immediate and complete flattening for all first components up to  $r_{\text{cutoff}}$ , which is demonstrated in  
 971 Figure 7 for Iteration 120 of Figure 6.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

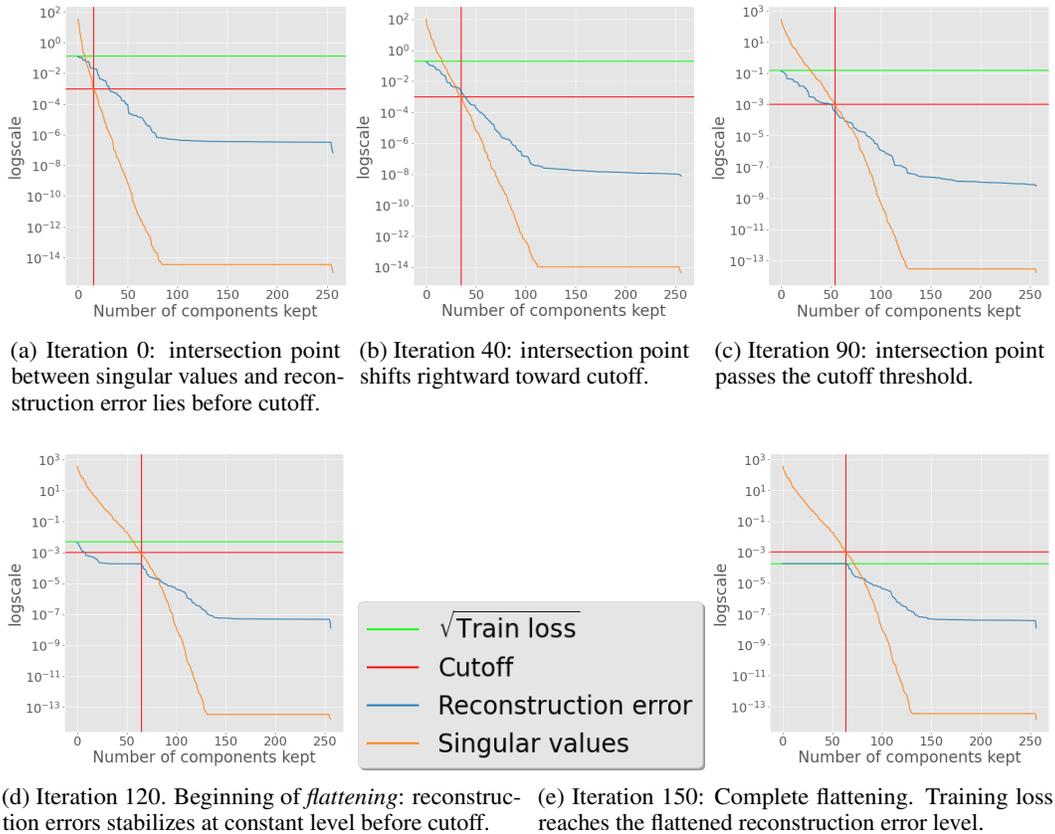


Figure 6: Evolution of quantities of interest during ANaGRAM training on heat equation. The dynamics reveal two distinct phases culminating in reconstruction error flattening.

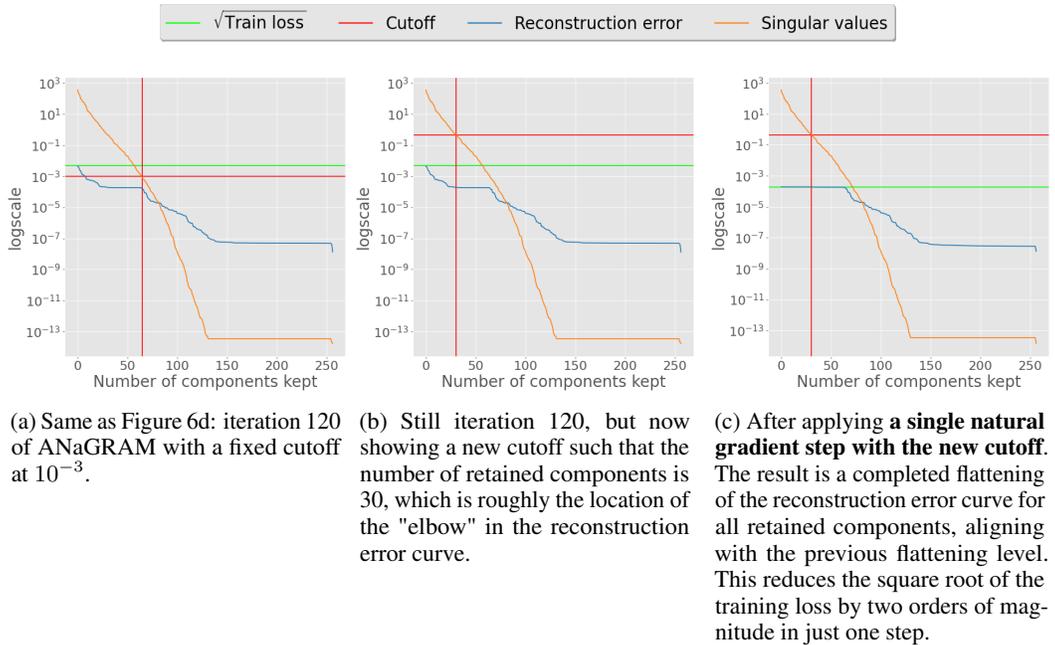


Figure 7: Illustration of “instant flattening” through adaptive cutoff adjustment. A single step with adjusted cutoff completes the flattening process.

## F DEEP DIVE ON SELECTED EXPERIMENTS

In this section we look at curves of training and estimations obtained with AMStramGRAM on benchmark of PDEs.

### F.1 ONE DIMENSIONAL BURGERS EQUATION

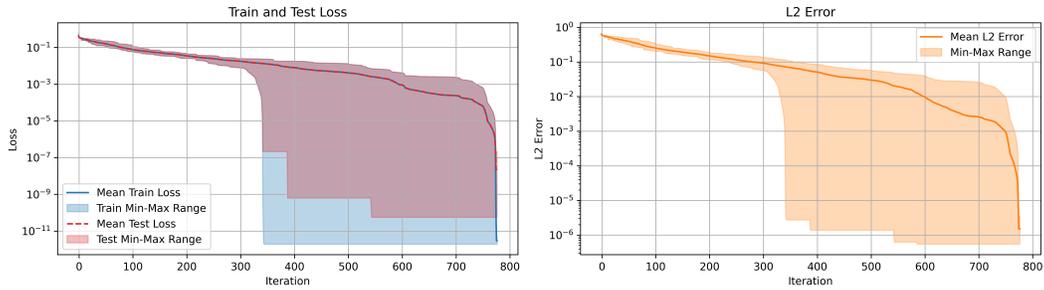
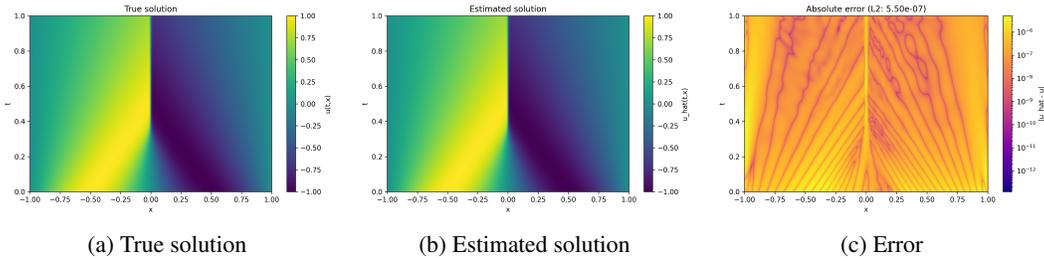


Figure 8: Training metrics for the One-Dimensional Burgers equation, showing convergence behavior with our adaptive multi-cutoff strategy.



(a) True solution

(b) Estimated solution

(c) Error

Figure 9: Results for One Dimensional Burgers Equation with cutoff  $10^{-6}$ .

### F.2 HEAT EQUATION

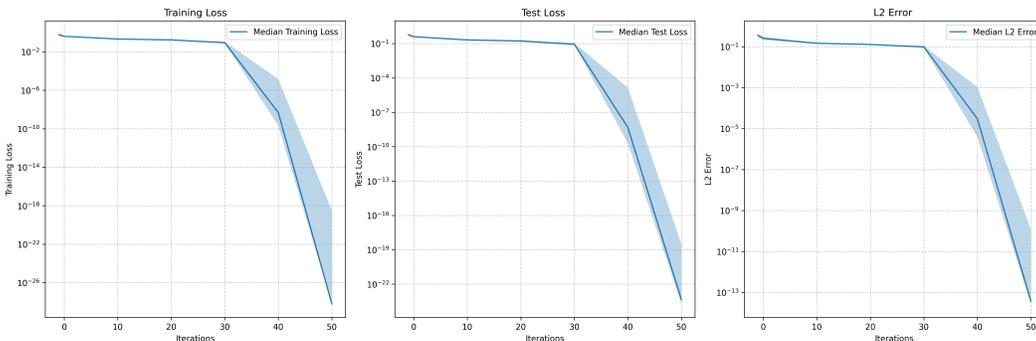


Figure 10: Convergence results for the Heat equation showing the  $L_2$  error over iterations. Our method (AMStramGRAM) converges faster and reaches a lower final error than ANaGRAM and baselines. Variability across runs is due to differing feature development speed from the random initialization.

1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133

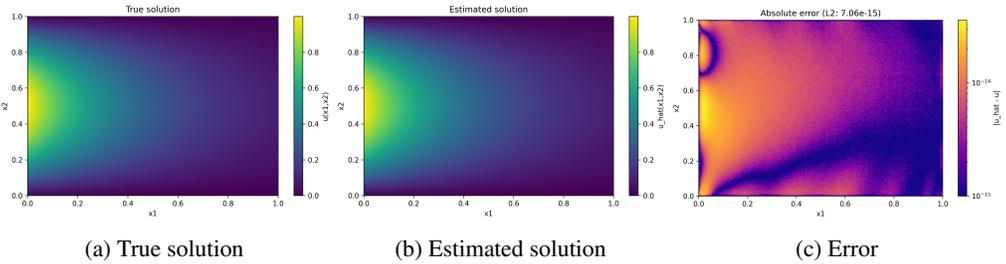


Figure 11: Results for the Heat equation (solution cutoff  $10^{-14}$ ). The error remains uniformly low over the domain, illustrating the effectiveness of the adaptive multi-cutoff strategy.

### F.3 LAPLACE EQUATIONS (L2D AND L5D)

For the Laplace equation in 2D, our method also demonstrates remarkable performance improvements over the baselines. The convergence is both faster and reaches a significantly lower error plateau.

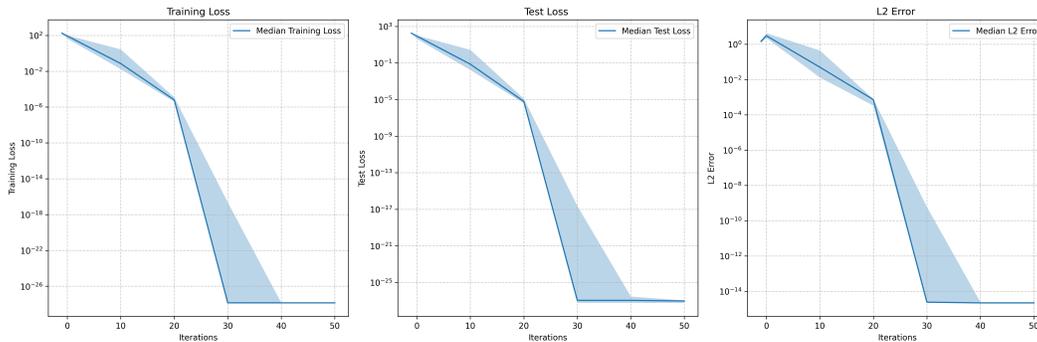


Figure 12: Convergence results for the Laplace 2D problem, showing the  $L_2$  error over iterations. Our method (AMStrAMGRAM) achieves both faster convergence and lower final error compared to ANaGRAM and other baseline methods. The observed variance between runs can be explained by different speed of convergence depending on the initialization.

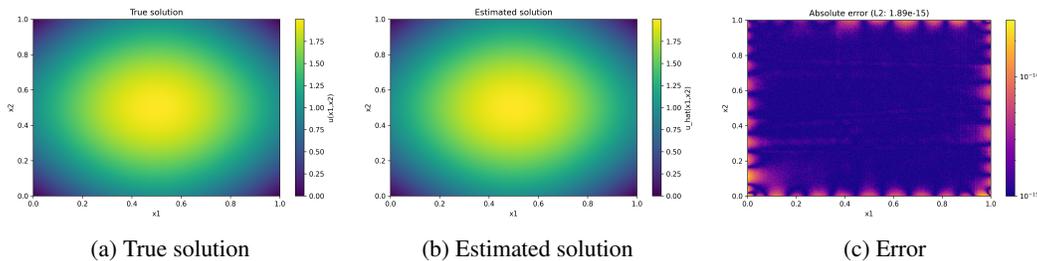
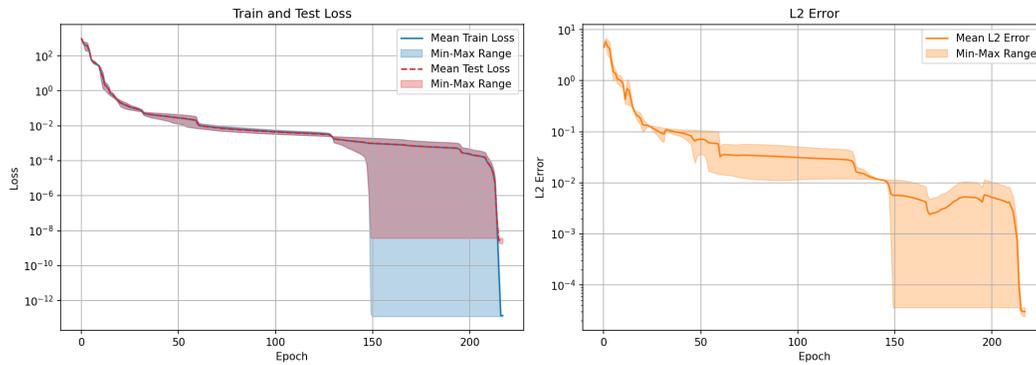


Figure 13: Results for Laplace 2D Equation with cutoff  $10^{-6}$ .

1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145

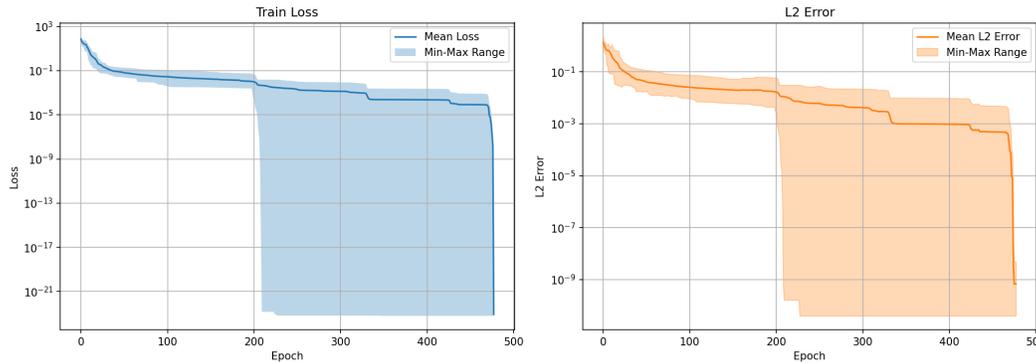


1146  
 1147 Figure 14: Convergence results for the Laplace 5D problem, showing the  $L_2$  error over iterations. Our method (AMStraMGRAM) achieves faster convergence but not lower final error compared to ANaGRAM and other baseline methods. We see that seeds change the speed of convergence of the algorithm

1153 F.4 NON LINEAR POISSON EQUATION

1154  
 1155 To compare ourselves with Urbán et al. (2025), we select (K=1).

1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170



1171 Figure 15: Convergence results for the Non Linear Poisson equation, showing the  $L_2$  error over iterations. Our method (AMStraMGRAM) achieves both faster convergence and lower final error compared to ANaGRAM and other baseline methods.

1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187

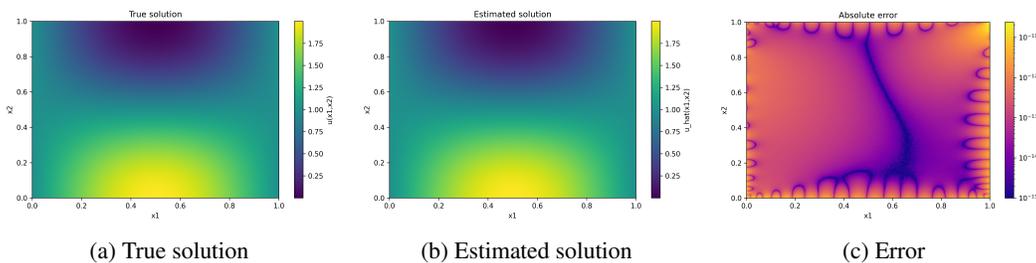


Figure 16: Results for the Nonlinear Poisson equation (cutoff  $10^{-4}$ ).

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

## F.5 ALLEN-CAHN EQUATION

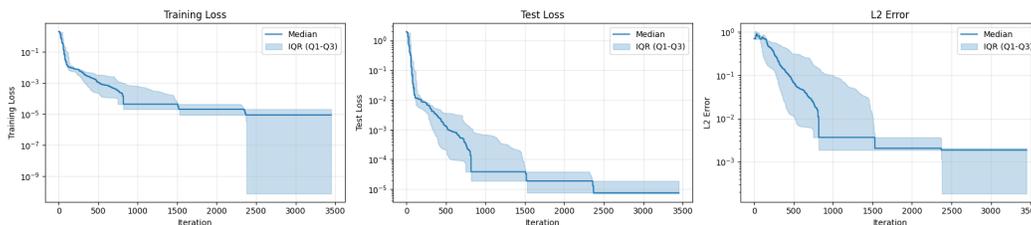


Figure 17: Training curves for the Allen-Cahn equation, showing the evolution of loss and error over iterations.

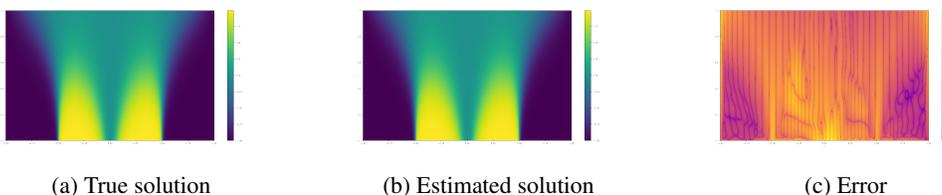


Figure 18: Results on the Allen-Cahn equation, showing the error distribution (left), model prediction (middle), and true solution (right). The error is mostly present in regions with the "sharpest" transitions, which exemplifies the challenge of accurately capturing sharp interfaces still remains even for our advanced optimization approach.

## G GEOMETRICAL INTERPRETATION OF REGULARIZATIONS

### G.1 WHY REGULARIZATION IS NECESSARY

We recall that our goal is to solve the operator equation  $D[u] = f$  by minimizing the squared residual

$$\|D[u] - f\|_{L^2(\Omega, \mu)}^2. \quad (22)$$

For simplicity, assume  $D$  is linear. Then the mapping

$$u \in C^\infty(\Omega) \mapsto \|D[u]\|_{L^2(\Omega, \mu)} \quad (23)$$

defines a semi-norm on  $C^\infty(\Omega)$ . We can "upgrade" this semi-norm into a true norm by introducing the following generalized Sobolev norm:

$$\|\cdot\|_{\tilde{\mathcal{H}}_D} : \begin{cases} C^\infty(\Omega \rightarrow \mathbb{R}) & \rightarrow \mathbb{R}^+ \\ u & \mapsto \sqrt{\|u\|_{L^2(\Omega \rightarrow \mathbb{R}, \mu)}^2 + \|D[u]\|_{L^2(\Omega \rightarrow \mathbb{R}, \mu)}^2} \end{cases} \quad (24)$$

Clearly, for any  $u$ ,

$$\|u\|_{L^2(\Omega, \mu)} \leq \|u\|_{\tilde{\mathcal{H}}_D}, \quad (25)$$

which guarantees that  $\|\cdot\|_{\tilde{\mathcal{H}}_D}$  is *definite*, i.e.  $\|u\|_{\tilde{\mathcal{H}}_D} = 0 \iff u = 0$ .

Completing  $C^\infty(\Omega)$  with respect to  $\|\cdot\|_{\tilde{\mathcal{H}}_D}$  yields a generalized Sobolev space  $(\mathcal{H}_D, \|\cdot\|_{\tilde{\mathcal{H}}_D})$ . This Hilbert space is the largest subspace of  $L^2(\Omega, \mu)$  on which  $D$  is continuous. Indeed, for every  $u \in \mathcal{H}_D$ ,

$$\|D[u]\|_{L^2(\Omega, \mu)} \leq \|u\|_{\mathcal{H}_D}. \quad (26)$$

Since our goal is to solve  $D[u] = f$ , we need  $D$  to be continuously invertible. That is, we need the reverse inequality of Equation (26) to hold (up to a constant  $\alpha > 0$ ). Formally, if  $D$  were algebraically

1242 invertible (bijective as a mapping), this condition would read:  
 1243

$$\begin{aligned}
 & \left( \exists \alpha > 0, \forall u \in \mathcal{H}_D, \|u\|_{\mathcal{H}_D} \leq \alpha \|D[u]\|_{L^2(\Omega \rightarrow \mathbb{R}, \mu)} \right) \\
 & \iff \left( \exists \alpha > 0, \forall u \in \mathcal{H}_D, \|D^{-1}[D[u]]\|_{\mathcal{H}_D} \leq \alpha \|D[u]\|_{L^2(\Omega \rightarrow \mathbb{R}, \mu)} \right) \quad (27) \\
 & \iff \left( \exists \alpha > 0, \forall f \in L^2(\Omega \rightarrow \mathbb{R}, \mu), \|D^{-1}[f]\|_{\mathcal{H}_D} \leq \alpha \|f\|_{L^2(\Omega \rightarrow \mathbb{R}, \mu)} \right)
 \end{aligned}$$

1251 **Operator ill-conditioning.** Even if  $D$  is bijective, Equation (27) may fail to hold, i.e.  $D$  can be  
 1252 ill-conditioned. Suppose there exists a subspace  $\mathcal{H}_K \subset \mathcal{H}_D$  such that  $D$  acts compactly on  $\mathcal{H}_K$  with  
 1253 infinite rank. Then  $D$  admits a singular value decomposition (Kress, 2014, Theorem 15.16): for  
 1254  $u \in \mathcal{H}_K$ ,

$$D[u] = \sum_{n \in \mathbb{N}} e_n \lambda_n \langle v_n, u \rangle_{\mathcal{H}_D}, \quad (28)$$

1259 with  $(v_n)$  orthonormal in  $\mathcal{H}_D$ ,  $(e_n)$  orthonormal in  $L^2(\Omega, \mu)$ , and  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ .

1261 For Equation (27) to hold, we would need  $\inf_n \lambda_n > 0$ , contradicting  $\lambda_n \rightarrow 0$ . This is exactly the  
 1262 classical inverse problem setting:  $D$  is bijective but ill-conditioned, and regularization is unavoidable.  
 1263 Among the many schemes developed, Tikhonov regularization is the canonical example (Kirsch,  
 1264 2021).

1266 **Non-bijectivity.** If  $D$  is not bijective, two additional issues may occur.

1267 **NON-SURJECTIVITY.** If  $\text{Im } D$  is a closed subspace, we can still obtain a solution by replacing the  
 1270 target  $f$  with its projection  $\Pi_{\text{Im } D} f$ . Note that minimizing  $\|D[u] - f\|_{L^2(\Omega, \mu)}^2$  yields precisely this  
 1271 least-squares solution.

1272 **NON-INJECTIVITY.** The lack of injectivity is a much more subtle issue. Since  $D$  is linear and  
 1273 continuous, its null space  $\text{Ker } D$  is a closed subspace of  $\mathcal{H}_D$ . In principle, one could restrict the  
 1274 domain of  $D$  to  $\text{Ker } D^\perp$  to make it injective. The problem, however, is that identifying  $\text{Ker } D$  is  
 1275 typically just as hard as solving the original problem itself, since it amounts to characterizing all  
 1276  $u \in \mathcal{H}_D$  such that  $D[u] = 0$ . Therefore, unless one can rely on theoretical results that explicitly  
 1277 describe  $\text{Ker } D$ , or construct a subspace  $\mathcal{H}_0 \subset \mathcal{H}_D$  for which  $\text{Ker } D \cap \mathcal{H}_0$  is explicitly known (so  
 1278 that  $D$  can be restricted to  $\mathcal{H}_0$ ), it is generally impossible to “get rid of”  $\text{Ker } D$  in practice.

1282 On the other hand, if we do not filter out  $\text{Ker } D$ , this has the unwanted consequence of introducing  
 1283 “spurious” low-energy signals. To be concrete, suppose we approximate our solution in a space  
 1284  $\mathcal{H}_K$  with orthonormal basis  $(u_n)_{n \in \mathbb{N}}$ . Assume there exists a subsequence  $(u_n^S) \notin \text{Ker } D$  converging  
 1285 towards  $\text{Ker } D$ . Since  $\text{Ker } D$  is closed (by continuity of  $D$ ), this means

$$\lim_{n \rightarrow \infty} \|\Pi_{\text{Ker } D} u_n^S - u_n^S\|_{\mathcal{H}_D}^2 = 0. \quad (29)$$

1288 Equivalently, after extraction, this can be rewritten for all  $n \in \mathbb{N}$  as

$$\frac{\|\Pi_{\text{Ker } D} u_n^S\|_{\mathcal{H}_D}^2}{\|u_n^S\|_{\mathcal{H}_D}^2} \geq 1 - 2^{-n}. \quad (30)$$

Now consider normalized vectors  $u_n^S / \|u_n^S\|_{\mathcal{H}_D}$ . We have

$$\begin{aligned}
0 < \left\| D \left[ \frac{u_n^S}{\|u_n^S\|_{\mathcal{H}_D}} \right] \right\|_{\mathcal{H}_D}^2 &= \left\| D \left[ \frac{\Pi_{\text{Ker } D^\perp} u_n^S + \Pi_{\text{Ker } D} u_n^S}{\|u_n^S\|_{\mathcal{H}_D}} \right] \right\|_{\mathcal{H}_D}^2 \\
&= \left\| D \left[ \frac{\Pi_{\text{Ker } D^\perp} u_n^S}{\|u_n^S\|_{\mathcal{H}_D}} \right] + \underbrace{D \left[ \frac{\Pi_{\text{Ker } D} u_n^S}{\|u_n^S\|_{\mathcal{H}_D}} \right]}_{=0} \right\|_{\mathcal{H}_D}^2 \\
&= \left\| D \left[ \frac{\Pi_{\text{Ker } D^\perp} u_n^S}{\|u_n^S\|_{\mathcal{H}_D}} \right] \right\|_{\mathcal{H}_D}^2 \\
&\stackrel{(26)}{\leq} \left\| \frac{\Pi_{\text{Ker } D^\perp} u_n^S}{\|u_n^S\|_{\mathcal{H}_D}} \right\|_{\mathcal{H}_D}^2 \\
&= 1 - \frac{\|\Pi_{\text{Ker } D} u_n^S\|_{\mathcal{H}_D}^2}{\|u_n^S\|_{\mathcal{H}_D}^2} \\
&\stackrel{(30)}{\leq} 2^{-n}.
\end{aligned} \tag{31}$$

In particular, if (for simplicity) the normalized  $(u_n / \|u_n\|_{\mathcal{H}_D})$  are right singular vectors of  $D$ , then the vectors  $(u_n^S / \|u_n^S\|_{\mathcal{H}_D})$  will correspond to singular values vanishing at least as fast as  $(2^{-n})$ . Crucially, however, these vanishing singular values do not reflect an intrinsic ill-conditioning of  $D$ , but rather an *artificial* ill-conditioning induced by the choice of approximation space  $\mathcal{H}_K$ . In other words, the spurious instability arises from how we approximate the operator, not from the operator itself. For more details on this approximation-induced phenomenon, see Adcock & Huybrechs (2019; 2020).

These remarks highlight the *inevitable need for regularization* in practice. In the next section, we will provide a geometric interpretation of the two regularization schemes introduced in Section 2.5, emphasizing how fundamentally different they are in nature.

*Remark 4.* The above discussion becomes even more critical when we restrict ourselves to a finite-dimensional approximation space  $\mathcal{H}_{\text{app}} \subset \mathcal{H}_D$ . In this case, the restriction  $D_{\text{app}}$  is automatically compact, since it is of finite rank. As a consequence, both types of ill-conditioning described above may occur simultaneously. This highlights once again why regularization is not merely convenient but *unavoidable* in numerical practice.

## G.2 RIDGE-REGRESSION

Returning to the definition given in Section 2.5, recall that Ridge regression amounts to adding  $\alpha^2 I_d$  (for some  $\alpha > 0$ ) to the Gram matrix  $G_\theta$  introduced in Equation (9):

$$\theta_{t+1} \leftarrow \theta_t - \eta G_{\theta_t}^\dagger \nabla \ell(\theta_t); \quad G_{\theta_t p, q} := \langle \partial_p u_{\theta_t}, \partial_q u_{\theta_t} \rangle_{L^2(\Omega, \mu)}. \tag{9}$$

We can reformulate this observation in the following way: given our model

$$u : \mathbb{R}^P \rightarrow L^2(\Omega \rightarrow \mathbb{R}, \mu), \tag{32}$$

consider the *regularized model*

$$u^\alpha : \begin{cases} \mathbb{R}^P & \rightarrow L^2(\Omega, \mu) \times \mathbb{R}^P \\ \theta & \mapsto (u_\theta, \alpha \theta). \end{cases} \tag{33}$$

The Gram matrix of this regularized model is exactly  $G_\theta + \alpha^2 I_d$ . Suppose further that regression is performed with respect to some function  $f \in L^2(\Omega, \mu)$ . Then we must adapt the objective to the regularized model, replacing  $f$  with the pair

$$(f, \alpha \theta) \in L^2(\Omega, \mu) \times \mathbb{R}^P. \tag{34}$$

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

A straightforward computation shows that, for all  $1 \leq p \leq \min(P, S)$ ,

$$\begin{aligned} \langle \partial_p u_{\theta}^{\alpha}, (f, \alpha \theta) - u_{\theta}^{\alpha} \rangle_{L^2(\Omega, \mu) \times \mathbb{R}^P} &= \langle \partial_p u_{\theta}, f - u_{\theta} \rangle_{L^2(\Omega, \mu)} + \underbrace{\alpha \langle e^{(p)}, \theta - \theta \rangle_{\mathbb{R}^P}}_{=0} \\ &= \langle \partial_p u_{\theta}, f - u_{\theta} \rangle_{L^2(\Omega, \mu)}. \end{aligned} \quad (35)$$

Thus, regression of  $(f, \alpha \theta)$  with the regularized model is exactly equivalent to Ridge regression. Equivalently, Ridge regression corresponds to replacing the original model  $u$  by the regularized model  $u^{\alpha}$ , and replacing the objective  $f$  by  $(f, \alpha \theta)$ . From this point of view, the choice of  $\alpha \theta$  as the secondary target may be interpreted as a *default assumption* in the absence of prior information on the parameters: one simply uses the current parameters as a reference target.

We can now extract several fundamental facts:

1. As  $\alpha \rightarrow 0$ , the regularized model  $u^{\alpha}$  tends in operator norm to the unregularized model  $(u, 0)$  (i.e.  $u$  by abuse of notation). Indeed,

$$\sup_{\|\theta\|_{\mathbb{R}^P}=1} \|du_{\theta}^{\alpha} - (du_{\theta}, 0)\|_{L^2(\Omega, \mu) \times \mathbb{R}^P} = \alpha \sup_{\|\theta\|_{\mathbb{R}^P}=1} \|\theta\|_{\mathbb{R}^P} = \alpha. \quad (36)$$

2. The model  $u^{\alpha}$  is injective and continuous. Since  $du_{\theta}$  is continuous (as  $\mathbb{R}^P$  is finite-dimensional), the only possible source of non-injectivity is  $\text{Ker } du_{\theta}^{\alpha}$ . But

$$\text{Ker } du_{\theta}^{\alpha} = \text{Ker } du_{\theta} \cap \text{Ker}(\alpha I_{\mathbb{R}^P}) \subset \text{Ker}(\alpha I_{\mathbb{R}^P}) = \{0\}, \quad (37)$$

hence injectivity. Restricting  $u^{\alpha}$  to its image makes it algebraically bijective, and the inverse is continuous since

$$\alpha \|\theta\|_{\mathbb{R}^P} \leq \|du_{\theta}^{\alpha}\|_{L^2(\Omega, \mu) \times \mathbb{R}^P}. \quad (38)$$

By the equivalence stated in Equation (27), this implies that  $(du_{\theta}^{\alpha})^{-1}$  is continuous. Consequently,  $\text{Im } du_{\theta}^{\alpha}$  is closed in  $L^2(\Omega, \mu) \times \mathbb{R}^P$ , since it is the inverse image of a closed set under  $(du_{\theta}^{\alpha})^{-1}$ . Therefore least-squares solution is well-defined.

3. The least-squares solution of  $u^{\alpha} = (f, 0)$  is influenced by  $\alpha$  as follows:  $(f, 0)$  is projected onto

$$\text{Im } du_{\theta}^{\alpha} = \text{Span}((\partial_p u_{\theta}, \alpha e^{(p)}) : 1 \leq p \leq P). \quad (39)$$

In particular, even if  $f \in \text{Im } du_{\theta}$  and  $f \neq 0$ , we still have  $(f, 0) \notin \text{Im } du_{\theta}^{\alpha}$  (since  $du_{\theta}(0) = 0$ ). Consequently,

$$\left( \Pi_{\text{Im } du_{\theta}^{\alpha}}^{\perp}(f, 0) \right)_1 \neq f, \quad (40)$$

where the subscript 1 denotes projection onto the first component in  $L^2(\Omega, \mu) \times \mathbb{R}^P$ .

We illustrate these phenomena in Figure 19a.

Building on the above analysis, we now show that Ridge regression can be extended to the functional setting. To this end, let us reconsider the operator  $D : \mathcal{H}_D \rightarrow L^2(\Omega, \mu)$  introduced in Section G.1. Analogously to what we did for the parametric model  $u$ , we define the *regularized operator* at level  $\alpha > 0$  as

$$D^{\alpha} : \begin{cases} \mathcal{H}_D & \rightarrow L^2(\Omega, \mu) \times \mathcal{H}_D \\ u & \mapsto (D[u], \alpha u) \end{cases}. \quad (41)$$

The corresponding target becomes the *regularized objective*  $(f, \alpha u)$ .

At this level of generality, the equivalence with Gram-matrix regularization no longer holds, since we are dealing with infinite-dimensional operators for which no direct Gram-matrix representation exists. Nevertheless, the fundamental properties remain valid, namely:

1. When  $\alpha \rightarrow 0$ , the regularized operator  $D^{\alpha}$  converges to  $(D, 0)$  in the operator-norm sense, i.e. to  $D$  by a mild abuse of notation. Indeed, we have

$$\sup_{\|u\|_{\mathcal{H}_D}=1} \|D^{\alpha}[u] - (D, 0)\|_{L^2(\Omega, \mu) \times \mathcal{H}_D} = \alpha \sup_{\|u\|_{\mathcal{H}_D}=1} \|u\|_{\mathcal{H}_D} = \alpha. \quad (42)$$

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

2. The operator  $D^\alpha$  is injective and continuous. Indeed,  $D$  is continuous by the very construction of  $\mathcal{H}_D$  (see Section G.1), and injectivity follows since

$$\text{Ker } D^\alpha = \text{Ker } D \cap \text{Ker}(\alpha I_{\mathcal{H}_D}) \subseteq \text{Ker}(\alpha I_{\mathcal{H}_D}) = \{0\}. \quad (43)$$

Restricting  $D^\alpha$  to its image makes it algebraically bijective, and the inverse is continuous: we have  $\alpha \|u\|_{\mathcal{H}_D} \leq \|D^\alpha[u]\|_{L^2(\Omega, \mu) \times \mathcal{H}_D}$ , which by the equivalence in Equation (27) implies that  $(D^\alpha)^{-1}$  is continuous. Consequently,  $\text{Im } D^\alpha$  is closed in  $L^2(\Omega, \mu) \times \mathcal{H}_D$ , since it is the inverse image of a closed set under  $(D^\alpha)^{-1}$ . Therefore least-squares solution is well-defined.

3. Least-squares solutions of the regularized problem  $D^\alpha[u] = (f, 0)$  are impacted by  $\alpha$  in the following way: we are projecting  $(f, 0)$  onto

$$\text{Im } D^\alpha = \text{Span} \left( (D[h], \alpha h) : h \in \mathcal{H}_D \right). \quad (44)$$

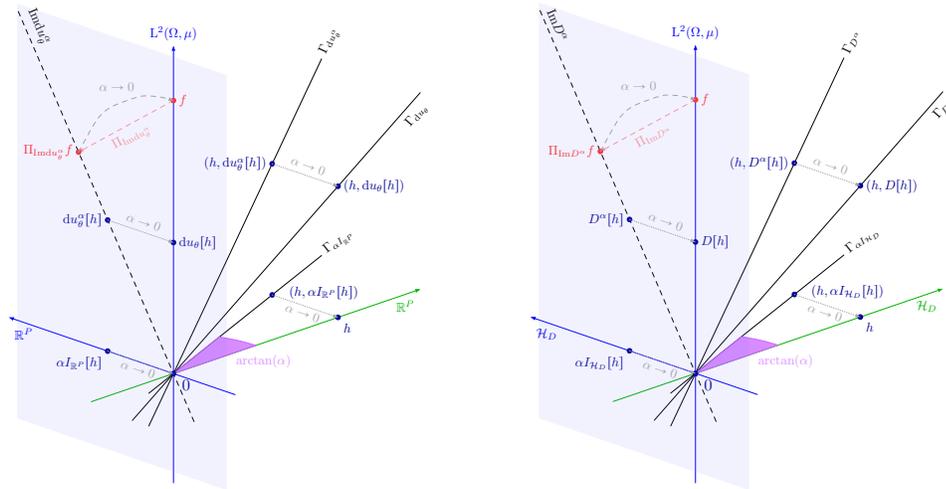
In particular, even if  $f \in \text{Im } D$  with  $f \neq 0$ , we have  $(f, 0) \notin \text{Im } D^\alpha$  (since  $D[0] = 0$ ), and hence

$$\left( \Pi_{\text{Im } D^\alpha}^\perp (f, 0) \right)_1 \neq f, \quad (45)$$

where the subscript 1 denotes the first coordinate in  $L^2(\Omega, \mu) \times \mathcal{H}_D$ .

We illustrate these phenomena in Figure 19b.

In summary, Ridge regression can be interpreted as a modification of the operator  $D$ , rendering it injective and continuously invertible on its image. However, this comes at a price: the regularized solutions are *never* exact solutions of the original equation  $D[u] = f$ , even when  $\alpha$  is arbitrarily small, since we are in fact solving a different operator equation. This marks a fundamental distinction from cutoff regularization, which instead acts directly on the approximation space, as we shall see in the next section.



(a) **Illustration of parametric Ridge regression.**

The green region represents the solution space, while the blue regions denote the target spaces. As  $\alpha \rightarrow 0$ , the regularized graph  $\Gamma_{Du^\alpha}$  of  $du_\theta^\alpha$  approaches the graph  $\Gamma_{Du}$  of  $du_\theta$ , with the angle between them vanishing at rate  $\arctan(\alpha)$ . The key consequence is that the projection of the objective  $f$  onto  $\text{Im } du_\theta^\alpha$  follows a non-linear path as  $\alpha \rightarrow 0$ , coinciding with  $\Pi_{\text{Im } du_\theta} f$  only asymptotically.

(b) **Illustration of functional Ridge regression.**

The green region represents the solution space, while the blue regions denote the target spaces. As  $\alpha \rightarrow 0$ , the regularized graph  $\Gamma_{D^\alpha}$  of  $D^\alpha$  approaches the graph  $\Gamma_{\mathcal{H}_D}$  of  $D$ , with the angle between them vanishing at rate  $\arctan(\alpha)$ . The key consequence is that the projection of the objective  $f$  onto  $\text{Im } D^\alpha$  follows a non-linear path as  $\alpha \rightarrow 0$ , coinciding with  $\Pi_{\text{Im } D} f$  only asymptotically.

Figure 19: Illustrations of Ridge regression.

---

### 1458 G.3 CUTOFF REGRESSION

1459  
1460 As in Section G.2, let us return to the setting of Section 2.5. In Equation (11), we introduced cutoff  
1461 regularization from the SVD perspective: given the differential  $du_\theta$  of the model  $u$ , at the point  $\theta$ ,  
1462 and its singular value decomposition  $du_\theta = V_\theta \Delta_\theta U_\theta^T$ , the cutoff-regularized pseudo-inverse  $du_\theta^{\dagger\alpha}$  at  
1463 level  $\alpha > 0$  is defined as

$$1464 \quad du_\theta^{\dagger\alpha} := U_\theta \Delta_\theta^{\dagger\alpha} V_\theta^T; \quad \Delta_{\theta,p}^{\dagger\alpha} := \begin{cases} \Delta_{\theta,p}^{-1} & \text{if } \Delta_{\theta,p} \geq \alpha \\ 0 & \text{otherwise} \end{cases}, \quad 1 \leq p \leq P. \quad (46)$$

1465  
1466  
1467 Let us reinterpret this construction. Denote by  $N_\alpha \in \mathbb{N}$  the number of singular values larger than  $\alpha$ .  
1468 Equivalently, assuming  $(\Delta_{\theta,p})_{1 \leq p \leq P}$  is non-increasing,

$$1469 \quad N_\alpha := \arg \max_{p \in \mathbb{N}} \{ \Delta_{\theta,p} \geq \alpha \}. \quad (47)$$

1470 Define

$$1471 \quad \Theta_\alpha := \text{Span}\{U_{\theta,p} : 1 \leq p \leq N_\alpha\}, \quad T_{N_\alpha}^0 \mathcal{M} := \text{Span}\{V_{i,p} : 1 \leq p \leq N_\alpha\}, \quad (48)$$

1472 so that  $T_{N_\alpha}^0 \mathcal{M} = du_\theta(\Theta_\alpha)$ . We then have

$$1473 \quad \left( du_{\theta|_{\Theta_\alpha}}^{T_{N_\alpha}^0 \mathcal{M}} \right)^{-1} = du_\theta^{\dagger\alpha}, \quad (49)$$

1474 meaning that the restriction  $du_\theta^\alpha := du_{\theta|_{\Theta_\alpha}}$  of  $du_\theta$  to the domain  $\Theta_\alpha$  becomes invertible once its  
1475 codomain is restricted to its image  $T_{N_\alpha}^0 \mathcal{M}$ , with inverse given precisely by the cutoff pseudo-inverse  
1476  $du_\theta^{\dagger\alpha}$ . Moreover, for any  $h \in \Theta_\alpha$ ,

$$1477 \quad \|du_\theta(h)\|_{L^2(\Omega, \mu)} = \|V_\theta \Delta_\theta U_\theta^T h\|_{L^2(\Omega, \mu)} \stackrel{V_\theta \text{ unitary}}{=} \|\Delta_\theta U_\theta^T h\|_{\mathbb{R}^P} \stackrel{h \in \Theta_\alpha}{\geq} \alpha \|U_\theta^T h\|_{\mathbb{R}^P} \stackrel{U_\theta \text{ unitary}}{=} \alpha \|h\|_{\mathbb{R}^P}. \quad (50)$$

1478 In other words, Equation (27) is satisfied by  $du_\theta^\alpha$ .

1479 Thus, while ridge regularization modifies the model itself, cutoff regularization instead restricts the  
1480 domain of the model so that, on this restricted domain, Equation (27) holds and the model becomes  
1481 invertible. We summarize the fundamental properties:

1482 1. We have

$$1483 \quad \bigcap_{\alpha > 0} (\mathbb{R}^P \setminus \Theta_\alpha) = \text{Ker } du_\theta, \quad (51)$$

1484 that is,  $\lim_{\alpha \rightarrow 0} \mathbb{R}^P \setminus \Theta_\alpha = \text{Ker } du_\theta$ , since for all  $\alpha > \beta$  we have  $\Theta_\alpha \subset \Theta_\beta$  and then  
1485  $\mathbb{R}^P \setminus \Theta_\beta \subset \mathbb{R}^P \setminus \Theta_\alpha$ . Similarly,  $\lim_{\alpha \rightarrow 0} T_{N_\alpha}^0 \mathcal{M} = \text{Im } du_\theta$ . Moreover, for each  $\alpha > 0$ , the  
1486 restriction  $du_\theta^\alpha$  coincides with  $du_\theta$  on  $\Theta_\alpha$ .

1487 2. By Equation (50),  $du_\theta^\alpha$  is injective and continuous. Restricting it to its image  $T_{N_\alpha}^0 \mathcal{M}$   
1488 makes it bijective and bicontinuous, with inverse exactly the cutoff pseudo-inverse  $du_\theta^{\dagger\alpha}$ . In  
1489 particular  $du(\Theta_\alpha)$  is closed in  $L^2(\Omega, \mu)$ , since it is the inverse image of a closed set under  
1490  $du_\theta^{\dagger\alpha}$ . Therefore least-squares solution is well-defined.

1491 3. Solving the least-squares problem  $du_\theta^\alpha = f$  is now altered in the following way: the target  
1492  $f$  is first projected onto  $T_{N_\alpha}^0 \mathcal{M} = \text{Im } du_\theta^\alpha$ . In particular, if for some  $\alpha > 0$  we already have  
1493  $f \in \text{Im } du_\theta^\alpha$ , then the regularized least-squares formulation recovers an *exact solution* to the  
1494 problem. This stands in sharp contrast with Ridge regression, where such exact recovery  
1495 can only occur *asymptotically* in the limit  $\alpha \rightarrow 0$ .

1506  
1507 As in Section G.2, we now need to reinterpret the cutoff regularization in order to extend it to the  
1508 functional setting. Let us return once more to the operator  $D : \mathcal{H}_D \rightarrow L^2(\Omega, \mu)$  introduced in  
1509 Section G.1. In general, one cannot define an SVD for such an operator (except when it is compact).  
1510 We must therefore appeal to the spectral theorem for bounded self-adjoint operators, which relies on  
1511 the notion of a *projection-valued measure* (also called a resolution of the identity). For our purposes,  
it will be sufficient to simply state the definition.

1512 **Definition 1** (Projection-valued measure). Let  $(X, \mathcal{A})$  be a measurable space, where  $\mathcal{A}$  denotes its  
 1513  $\sigma$ -algebra, and let  $\mathcal{H}$  be a Hilbert space. A *projection-valued measure* (PVM) is a map  
 1514

$$1515 \pi : \mathcal{A} \rightarrow \mathcal{L}_b(\mathcal{H} \rightarrow \mathcal{H}),$$

1517 where  $\mathcal{L}_b(\mathcal{H} \rightarrow \mathcal{H})$  denotes the set of bounded operators on  $\mathcal{H}$ , such that for every  $A \in \mathcal{A}$ ,  $\pi(A)$  is  
 1518 an orthogonal projection on  $\mathcal{H}$ , and the following properties hold:  
 1519

- 1520 1.  $\pi(\emptyset) = 0$  and  $\pi(X) = I_{\mathcal{H}}$ , where  $I_{\mathcal{H}}$  is the identity operator on  $\mathcal{H}$ ;
- 1521 2.  $\pi(A \cap B) = \pi(A)\pi(B)$  for all  $A, B \in \mathcal{A}$ ;
- 1522 3. For every countable family  $(A_i)_{i=1}^{\infty}$  of disjoint sets in  $\mathcal{A}$ ,

$$1523 \pi\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \pi(A_i),$$

1524 where the series converges in the strong operator topology.  
 1525

1526 Since projection-valued measures are measures, one can define integrals with respect to them. We  
 1527 refer to (Berezansky et al., 1996, Chapter 13) for details. We may now state the spectral theorem.  
 1528

1529 **Theorem 2.** *Let  $\mathcal{H}$  be a Hilbert space and let  $A : \mathcal{H} \rightarrow \mathcal{H}$  be a self-adjoint operator. Then there  
 1530 exists a projection-valued measure  $\pi$  on the Borel  $\sigma$ -algebra of  $\mathbb{R}$  such that*  
 1531

$$1532 A = \int_{\mathbb{R}} \lambda \pi(d\lambda) = \int_{\sigma(A)} \lambda \pi(d\lambda), \quad (52)$$

1533 where  $\sigma(A)$  denotes the spectrum of  $A$ .  
 1534

1535 A proof can be found in (Berezansky et al., 1996, Theorem 4.1, Section 4.1, Chapter 13). In particular,  
 1536 since  $\pi$  is a projection-valued measure, we have by Definition 1:  
 1537

$$1538 I_{\mathcal{H}} = \int_{\mathbb{R}} \pi(d\lambda). \quad (53)$$

1539 Since  $D$  is continuous, we can define its adjoint  $D^* : L^2(\Omega, \mu) \rightarrow \mathcal{H}_D$ , and hence the self-adjoint  
 1540 operator  $D^*D : \mathcal{H}_D \rightarrow \mathcal{H}_D$ . Applying Theorem 2, we obtain a projection-valued measure  $\pi_D$  on  $\mathbb{R}$   
 1541 endowed with its Borel  $\sigma$ -algebra, such that  
 1542

$$1543 D^*D = \int_{\mathbb{R}_+} \lambda \pi_D(d\lambda), \quad I_{\mathcal{H}_D} = \int_{\mathbb{R}_+} \pi_D(d\lambda), \quad (54)$$

1544 where the integration is restricted to  $\mathbb{R}_+$  since  $D^*D$  is a positive operator. We can then define  
 1545

$$1546 \Pi_D^\alpha := \int_{\alpha^2}^{+\infty} \pi_D(d\lambda), \quad (55)$$

1547 which is an orthogonal projection in  $\mathcal{H}_D$  since  $\pi_D$  is a projection-valued measure. We then define  
 1548 the regularized space  $\mathcal{H}_D^\alpha$  at level  $\alpha > 0$  by  
 1549

$$1550 \mathcal{H}_D^\alpha := \text{Im } \Pi_D^\alpha \subset \mathcal{H}_D. \quad (56)$$

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

For any  $u \in \mathcal{H}_D^\alpha$ , we compute

$$\begin{aligned}
\|D[u]\|_{L^2(\Omega, \mu)}^2 &= \langle D[u], D[u] \rangle_{L^2(\Omega, \mu)} = \langle u, D^* D[u] \rangle_{\mathcal{H}_D} \\
&= \left\langle u, \int_{\mathbb{R}_+} \lambda \pi_D(d\lambda) u \right\rangle_{\mathcal{H}_D} \\
&\stackrel{u \in \mathcal{H}_D^\alpha}{=} \left\langle u, \int_{\mathbb{R}_+} \lambda_1 \pi_D(d\lambda_1) \int_{\alpha^2}^{+\infty} \pi_D(d\lambda_2) u \right\rangle_{\mathcal{H}_D} \\
&\stackrel{\pi_D \text{ PVM}}{=} \left\langle u, \int_{\alpha^2}^{+\infty} \lambda \pi_D(d\lambda) u \right\rangle_{\mathcal{H}_D} \\
&= \int_{\alpha^2}^{+\infty} \lambda \langle u, \pi_D(d\lambda) u \rangle_{\mathcal{H}_D} \\
&\geq \alpha^2 \int_{\alpha^2}^{+\infty} \langle u, \pi_D(d\lambda) u \rangle_{\mathcal{H}_D} \stackrel{u \in \mathcal{H}_D^\alpha}{=} \alpha^2 \langle u, u \rangle_{\mathcal{H}_D} = \alpha^2 \|u\|_{\mathcal{H}_D}^2.
\end{aligned} \tag{57}$$

That is,

$$\|D[u]\|_{L^2(\Omega, \mu)} \geq \alpha \|u\|_{\mathcal{H}_D}, \tag{58}$$

so that Equation (27) is verified on  $\mathcal{H}_D^\alpha$ . We denote

$$D^\alpha := D|_{\mathcal{H}_D^\alpha}, \tag{59}$$

the restriction of  $D$  to the domain  $\mathcal{H}_D^\alpha$ . We can now list the fundamental properties:

1. We have

$$\bigcap_{\alpha > 0} (\mathcal{H}_D \setminus \mathcal{H}_D^\alpha) = \text{Ker } D \tag{60}$$

that is,  $\lim_{\alpha \rightarrow 0} \mathcal{H}_D \setminus \mathcal{H}_D^\alpha = \text{Ker } D$ , since for all  $\alpha > \beta$ ,  $\mathcal{H}_D^\alpha \subset \mathcal{H}_D^\beta$  by Property 3 of Definition 1. Moreover, by continuity of  $D$ , we also have  $\lim_{\alpha \rightarrow 0} D[\mathcal{H}_D^\alpha] = \text{Im } D$ . Finally, for each  $\alpha > 0$ ,  $D^\alpha$  coincides with  $D$  on  $\mathcal{H}_D^\alpha$  by construction.

2. As established by Equation (27),  $D^\alpha$  is injective and continuous. When restricted to its image, it is therefore bijective and bicontinuous, hence invertible. In particular  $D[\mathcal{H}_D^\alpha]$  is closed in  $L^2(\Omega, \mu)$ , since it is the inverse image of a closed set under  $(D^\alpha)^{-1}$ . Therefore least-squares solution is well-defined.
3. The least-squares solution of  $D^\alpha = f$  is now modified as follows: one projects  $f$  onto  $\text{Im } D^\alpha$ . In particular, if for some  $\alpha > 0$  we already have  $f \in \text{Im } D^\alpha$ , then the regularized least-squares formulation recovers an *exact solution* to the problem  $D[u] = f$ . This stands in sharp contrast with Ridge regression, where such exact recovery can only occur *asymptotically* in the limit  $\alpha \rightarrow 0$ .

#### G.4 CONNECTION TO GREEN'S FUNCTION

To further highlight the difference between the two regularization schemes, we now reinterpret them through the lens of Green's functions of the operator  $\bar{D}$ . Schwencke & Furtlehner (2025, Theorem 2) established in the finite-dimensional case a connection between the natural gradient for PINNs and Green's functions. Their proof relies on Schwencke & Furtlehner (2025, Proposition 3), which will be our starting point. We restate the relevant definitions and results for completeness.

**Definition 2** (Schwencke & Furtlehner, 2025, Definition 9: generalized Green's function). Let  $\mathcal{H}$  be an Hilbert space,  $D : \mathcal{H} \rightarrow L^2(\Omega, \mu)$  be a linear differential operator,  $\mathcal{H}_0 \subset \mathcal{H}$  a subspace isometrically embedded in  $\mathcal{H}$  and  $f \in L^2(\Omega, \mu)$ . A generalized Green's function of  $D$  on  $\mathcal{H}_0$  is then any kernel function  $g : \Omega \times \Omega \rightarrow \mathbb{R}$  such that the operator:

$$R_{\mathcal{H}_0} : \begin{cases} L^2(\Omega \rightarrow \mathbb{R}, \mu) & \rightarrow \mathcal{H} \\ f & \mapsto \left( x \in \Omega \mapsto \int_{\Omega} g(x, s) f(s) \mu(ds) \right) \end{cases},$$

verifies the equation:

$$D \circ R_{\mathcal{H}_0} = \Pi_{D[\mathcal{H}_0]}^\perp \tag{61}$$

**Proposition 2** (Schwencke & Furtlehner, 2025, Proposition 3). *Let  $D : \mathcal{H} \rightarrow L^2(\Omega, \mu)$  be a linear differential operator, and  $\mathcal{H}_0 := \text{Span}(u_p : 1 \leq p \leq P) \subset \mathcal{H}$  a subspace isometrically embedded in  $\mathcal{H}$ . Then the generalized Green's function of  $D$  on  $\mathcal{H}_0$  is given by: for all  $x, y \in \Omega$*

$$g_{\mathcal{H}_0}(x, y) := \sum_{1 \leq p, q \leq P} u_p(x) G_{p,q}^\dagger D[u_q](y), \quad (62)$$

with: for all  $1 \leq p, q \leq P$ ,

$$G_{p,q} := \langle D[u_p], D[u_q] \rangle_{L^2(\Omega \rightarrow \mathbb{R}, \mu)}. \quad (63)$$

**Our goal.** We aim to

- (i) generalize Schwencke & Furtlehner (2025, Proposition 3) to arbitrary Reproducing Kernel Hilbert Spaces;
- (ii) establish a direct connection to the regularization framework introduced earlier. This will provide a novel reinterpretation of the Green's function in the regularized operator setting.

**Operator framework.** Consider the operator  $D : \mathcal{H}_D \rightarrow L^2(\Omega, \mu)$  from Section G.1, and assume that there exists an RKHS  $\mathcal{H}_0$  isometrically embedded in  $\mathcal{H}_D$  (for instance, any finite-dimensional RKHS, see Schwencke & Furtlehner, 2025, Corollary 1). For Schwencke & Furtlehner (2025, Definition 9) to be well-posed, the range  $D[\mathcal{H}_0]$  must be a closed subspace of  $L^2(\Omega, \mu)$ . As argued earlier, this is guaranteed if  $D$  is continuously invertible: indeed, in this case

$$D[\mathcal{H}_0] = (D^{-1})^{-1}[\mathcal{H}_0], \quad (64)$$

and the inverse image of a closed subspace under a continuous operator is closed.

**Key observation.** Thus, to generalize Schwencke & Furtlehner (2025, Proposition 3), we require  $D$  to be continuously invertible. Conveniently, this is precisely the property enforced by the regularization schemes we introduced earlier.

In what follows, we first focus on the cutoff regularization, which offers the clearest interpretation in terms of Green's functions. We then briefly revisit the case of Ridge regression. Before delving further into our main goal, let us first establish two general facts.

**Lemma 1.** *Let  $(\mathcal{H}_0, \|\cdot\|_{\mathcal{H}_0})$  be an RKHS on a set  $X$  with reproducing kernel  $k$ . Suppose that  $\|\cdot\|_{bis}$  is a norm equivalent to  $\|\cdot\|_{\mathcal{H}_0}$ . Then  $(\mathcal{H}_0, \|\cdot\|_{bis})$  is also an RKHS.*

*Proof.* The key point is to show that there exists a reproducing kernel for the inner product  $\langle \cdot, \cdot \rangle_{bis}$  associated with  $\|\cdot\|_{bis}$ . Our argument follows the simple reasoning in Paulsen & Raghupathi (2016, Definitions 1–2).

Since, for every  $x \in X$ , the point evaluation functional

$$\delta_x : u \in \mathcal{H}_0 \mapsto u(x) \quad (65)$$

is continuous with respect to  $\|\cdot\|_{\mathcal{H}_0}$  by the definition of an RKHS, it is also continuous with respect to the equivalent norm  $\|\cdot\|_{bis}$ . Therefore, by the Riesz representation theorem, for each  $x \in X$ , there exists a unique element  $k_x^{bis} \in \mathcal{H}_0$  such that for all  $u \in \mathcal{H}_0$

$$\langle k_x^{bis}, u \rangle_{bis} = u(x). \quad (66)$$

In particular, this defines a reproducing kernel for the norm  $\|\cdot\|_{bis}$ , given by

$$k_{bis}(x, y) = \langle k_x^{bis}, k_y^{bis} \rangle_{bis} = k_x^{bis}(y), \quad \forall x, y \in X. \quad (67)$$

Hence  $(\mathcal{H}_0, \|\cdot\|_{bis})$  is indeed an RKHS.  $\square$

**Lemma 2.** *Let  $\mathcal{H}_A, \mathcal{H}_B$  be two Hilbert spaces. If  $U : \mathcal{H}_A \rightarrow \mathcal{H}_B$  is an isometry, then*

$$U^*U = I_{\mathcal{H}_A}, \quad UU^* = \Pi_{\text{Im } U}. \quad (68)$$

*In particular  $\text{Im } U$  is closed in  $\mathcal{H}_B$ .*

1674 *Proof.* The first identity follows from the fact that for all  $x, y \in \mathcal{H}_A$ ,  
 1675 
$$\langle x, U^*U[y] \rangle_{\mathcal{H}_A} = \langle U[x], U[y] \rangle_{\mathcal{H}_B} = \langle x, y \rangle_{\mathcal{H}_A}. \quad (69)$$

1676 Thus  $(U^*U(y) - y) \in \mathcal{H}_A^\perp$ , i.e.  $U^*U = I_{\mathcal{H}_A}$ . For the second, the key point is to show that  $\text{Im } U$  is  
 1677 closed, i.e.  $\text{Im } U = \overline{\text{Im } U}$ .  
 1678

1679 Let  $y \in \overline{\text{Im } U}$ , and  $(y_n) \in \text{Im } U^{\mathbb{N}}$  with  $y_n \rightarrow y$ . Since  $(y_n)$  is Cauchy, and  $y_n = U(x_n)$  for some  
 1680  $(x_n) \in \mathcal{H}_A^{\mathbb{N}}$ , we have

$$1681 \quad \|U(x_n) - U(x_m)\|_{\mathcal{H}_B} = \|x_n - x_m\|_{\mathcal{H}_A}, \quad (70)$$

1682 so  $(x_n)$  is also Cauchy and converges to  $x \in \mathcal{H}_A$ , since  $\mathcal{H}_A$  is complete. Since  $U$  is an isometry, we  
 1683 have for all  $x \in \mathcal{H}_A$

$$1684 \quad \|U(x)\|_{\mathcal{H}_B} = \|x\|_{\mathcal{H}_A}. \quad (71)$$

1685 In particular,  $U$  is bounded with operator norm  $\|U\| = 1$ , and hence continuous. Thus  $U(x) = y$ ,  
 1686 hence  $y \in \text{Im } U$ . We conclude that  $\text{Im } U$  is closed in  $\mathcal{H}_B$ . Finally:

1687  
 1688 • For  $y \in \text{Im } U$ , say  $y = U(x)$ , we have  
 1689 
$$UU^*(y) = U(U^*U)(x) = U(x) = y. \quad (72)$$

1690  
 1691 • For  $y \in (\text{Im } U)^\perp$ , we check that  $UU^*(y) = 0$ . Indeed, for any  $z \in \mathcal{H}_B$ ,  
 1692 
$$\langle z, UU^*(y) \rangle_{\mathcal{H}_B} = \langle UU^*(z), y \rangle_{\mathcal{H}_B} = 0, \quad (73)$$
  
 1693 since  $UU^*(z) \in \text{Im } U$ . Thus  $UU^*(y) \in \mathcal{H}_B^\perp$ , i.e.  $UU^*(y) = 0$ .  $\square$   
 1694

1695 We are interested in the restriction of  $D$  to the domain  $\mathcal{H}_0$ . Since the restriction  $D^*D : \mathcal{H}_D \rightarrow \mathcal{H}_D$   
 1696 does not, *a priori*, map  $\mathcal{H}_0$  into itself, we first need to adapt the setting in order to apply the spectral  
 1697 theorem of Theorem 2.

1698 Because  $\mathcal{H}_0 \subset \mathcal{H}_D$  isometrically, we have for all  $u, v \in \mathcal{H}_0$ :

$$1699 \quad \begin{aligned} \langle D[u], D[v] \rangle_{L^2(\Omega, \mu)} &= \langle D[\Pi_{\mathcal{H}_0} u], D[\Pi_{\mathcal{H}_0} v] \rangle_{L^2(\Omega, \mu)} \\ 1700 &= \langle \Pi_{\mathcal{H}_0} u, D^*D[\Pi_{\mathcal{H}_0} v] \rangle_{\mathcal{H}_D} \\ 1701 &= \langle u, (\Pi_{\mathcal{H}_0} D^*D \Pi_{\mathcal{H}_0})[v] \rangle_{\mathcal{H}_D}, \end{aligned} \quad (74)$$

1702 where we used in the last step that  $\Pi_{\mathcal{H}_0}$  is self-adjoint.  
 1703

1704 We can therefore apply the spectral theorem Theorem 2 to the bounded self-adjoint operator  
 1705  $\Pi_{\mathcal{H}_0} D^*D \Pi_{\mathcal{H}_0} : \mathcal{H}_0 \rightarrow \mathcal{H}_0$ , obtaining the analogue of the decomposition in Equation (54):

$$1706 \quad \Pi_{\mathcal{H}_0} D^*D \Pi_{\mathcal{H}_0} = \int_{\mathbb{R}_+} \lambda \pi_D^{\mathcal{H}_0}(d\lambda), \quad I_{\mathcal{H}_0} = \int_{\mathbb{R}_+} \pi_D^{\mathcal{H}_0}(d\lambda). \quad (75)$$

1707  
 1708 **Regularized spaces.** Fixing  $\alpha > 0$ , and analogously to Equations (55) and (56), we define the  
 1709 regularized projection and subspace:

$$1710 \quad \Pi_{D, \mathcal{H}_0}^\alpha := \int_{\alpha^2}^{+\infty} \pi_D^{\mathcal{H}_0}(d\lambda), \quad \mathcal{H}_{D, \mathcal{H}_0}^\alpha := \text{Im } \Pi_{D, \mathcal{H}_0}^\alpha \subset \mathcal{H}_0 \subset \mathcal{H}_D. \quad (76)$$

1711 Let  $k : \Omega \times \Omega \rightarrow \mathbb{R}$  be the reproducing kernel of  $\mathcal{H}_0$ . Then, by Paulsen & Raghupathi (2016,  
 1712 Theorem 2.5),  $\mathcal{H}_{D, \mathcal{H}_0}^\alpha$  is an RKHS with reproducing kernel

$$1713 \quad k_\alpha(x, y) := \Pi_{D, \mathcal{H}_0}^\alpha[k(x, \cdot)](y), \quad \forall x, y \in \Omega. \quad (77)$$

1714  
 1715 **Norm equivalence.** Since  $\mathcal{H}_{D, \mathcal{H}_0}^\alpha \subset \mathcal{H}_0 \subset \mathcal{H}_D$ , inequality in Equation (26) remains valid, i.e. for  
 1716 all  $u \in \mathcal{H}_{D, \mathcal{H}_0}^\alpha$ :

$$1717 \quad \|D[u]\|_{L^2(\Omega, \mu)} \leq \|u\|_{\mathcal{H}_D}. \quad (78)$$

1718 Furthermore, by an argument entirely analogous to Equation (57), we also have

$$1719 \quad \|D[u]\|_{L^2(\Omega, \mu)} \geq \alpha \|u\|_{\mathcal{H}_D}, \quad \forall u \in \mathcal{H}_{D, \mathcal{H}_0}^\alpha. \quad (79)$$

1720 In particular, the functional

$$1721 \quad \|\cdot\|_D : \begin{cases} \mathcal{H}_{D, \mathcal{H}_0}^\alpha & \rightarrow \mathbb{R} \\ u & \mapsto \|D[u]\|_{L^2(\Omega, \mu)} \end{cases} \quad (80)$$

1722 defines a norm equivalent to  $\|\cdot\|_{\mathcal{H}_D}$  on  $\mathcal{H}_{D, \mathcal{H}_0}^\alpha$ . By Lemma 1, the pair  $(\mathcal{H}_{D, \mathcal{H}_0}^\alpha, \|\cdot\|_D)$  is itself an  
 1723 RKHS with a reproducing kernel  $k_D$ .  
 1724

1728 **Isometry property.** The crucial observation is that  $D$  is an isometry with respect to this norm.  
 1729 Indeed, for all  $u, v \in \mathcal{H}_{D, \mathcal{H}_0}^\alpha$ ,

$$1730 \langle u, v \rangle_D = \langle D[u], D[v] \rangle_{L^2(\Omega, \mu)}. \quad (81)$$

1732 This allows us to characterize the associated Green’s function.

1733 **Theorem 1.** *The generalized Green’s function of the operator  $D$  in the regularized space  $\mathcal{H}_{D, \mathcal{H}_0}^\alpha$  is given, for all  $x, y \in \Omega$ , by*

$$1734 g_D(x, y) := D[k_D(x, \cdot)](y), \quad (12)$$

1737 *Proof.* For all  $f \in L^2(\Omega, \mu)$  and  $x \in \Omega$ ,

$$1738 \int_{\Omega} g_D(x, s) f(s) \mu(ds) = \langle g_D(x, \cdot), f \rangle_{L^2(\Omega, \mu)}$$

$$1739 = \langle D[k_D(x, \cdot)], f \rangle_{L^2(\Omega, \mu)} \quad (82)$$

$$1740 = \langle k_D(x, \cdot), D^* f \rangle_D$$

$$1741 = (D^* f)(x).$$

1742 Since  $D$  is an isometry, Lemma 2 gives  $DD^* = \Pi_{D[\mathcal{H}_{D, \mathcal{H}_0}^\alpha]}$ . Therefore,

$$1743 D \left[ x \mapsto \int_{\Omega} g_D(x, s) f(s) \mu(ds) \right] = D[D^* f] = \Pi_{D[\mathcal{H}_{D, \mathcal{H}_0}^\alpha]} f, \quad (83)$$

1744 which precisely shows that  $g_D$  is a generalized Green’s function. □

1752 The key insight of Theorem 1 is that, in the PINNs setting—and most notably in our algorithm—we implicitly construct the reproducing kernel  $k_D$  associated with the norm  $\|\cdot\|_D$  on the regularized tangent space  $T_\theta^\alpha \mathcal{M}$  of the neural network manifold  $\mathcal{M}$ , at cutoff level  $\alpha$ . This kernel is precisely the PINNs NNTK introduced by Schwencke & Furtlehner (2025).

1757 A crucial consequence is that the regularization of the Gram matrix is not merely a “numerical trick” to guarantee stability: it is the very mechanism that ensures the Green’s function is well defined.

1760 **Conceptual interpretation.** This perspective also offers a profound interpretation of the procedure: rather than attempting to invert the operator  $D$  directly, we build a kernel  $k_D$  whose associated metric makes  $D$  an isometry, and thus ensures that  $D^*$  acts as the generalized left-inverse of  $D$ . The magic of the kernel lies in the following facts:

1764 (i) We never need to compute  $D^*$  explicitly, since it is implicitly encoded in the relation

$$1765 \langle D[k_D(x, \cdot)], f \rangle_{L^2(\Omega, \mu)} = \langle k_D(x, \cdot), D^* f \rangle_D. \quad (84)$$

1768 (ii) The same formula allows us to directly evaluate the generalized solution  $D^* f$ : indeed, for all  $x \in \Omega$ , the reproducing property gives

$$1769 D^* f(x) = \langle k_D(x, \cdot), D^* f \rangle_D. \quad (85)$$

1772 **Comparison with Ridge regression.** An analogous analysis holds for Ridge regression. However, instead of inverting  $D$  “via isometry,” we invert the augmented operator  $(D, \alpha I_{\mathcal{H}_D})$ .

## 1775 G.5 CONVERGENCE OF AMSTRAMGRAM IN THE NTK-REGIME

1777 Let us consider we are in the so-called NTK-regime (Jacot et al., 2018) under infinite-width assumption. In this case the NTK converges to a fixed kernel  $k$ . Under mild assumptions(?), we may assume that  $k$  is  $2m$ -times differentiable. Let us denote  $\mathcal{H}(k)$  the Reproducing Kernel Hilbert Space associated to  $k$  with its associated norm  $\|\cdot\|_k$ .

1781 We have the following result.

1782 **Proposition 3.** Let  $(X, \|\cdot\|_X)$  be a finite-dimensional normed space and let  $k : X \times X \rightarrow \mathbb{R}$  be a  
 1783 positive semidefinite kernel of class  $\mathcal{C}^{2m}$  with respect to the product topology on  $X \times X$ .

1784  
 1785 Then every  $f \in \mathcal{H}(k)$  is of class  $\mathcal{C}^m$  on  $X$ . Furthermore, for every multi-index  $|\alpha| \leq m$  and  $x \in X$ ,  
 1786 there exist  $k_x^\alpha \in \mathcal{H}(k)$  such that for all  $f \in \mathcal{H}(k)$

$$1787 \partial^\alpha f(x) = \langle k_x^\alpha, f \rangle_{\mathcal{H}(k)}. \quad (86)$$

1788  
 1789 *Proof.* We begin by fixing some notation. For multi-indices  $\alpha, \beta$ , we denote by  $\partial^{(\alpha, \beta)} k$  the mixed  
 1790 partial derivatives of  $k$  with respect to its two kernel arguments. Concretely, for a simple index  $i$ , we  
 1791 set

$$1792 \partial^{(i, 0)} k(x, y) := \lim_{\varepsilon \rightarrow 0} \frac{k(x + \varepsilon e_i, y) - k(x, y)}{\varepsilon}, \quad (87)$$

1793 and for a simple index  $j$ ,

$$1794 \partial^{(i, j)} k(x, y) := \lim_{\varepsilon \rightarrow 0} \frac{\partial^{(i, 0)} k(x, y + \varepsilon e_j) - \partial^{(i, 0)} k(x, y)}{\varepsilon}. \quad (88)$$

1795 Note that  $\partial^{(\alpha, \beta)} k$  is a mixed partial derivative of total order  $|\alpha| + |\beta|$ . This hints at the requirement  
 1800 that  $k$  be of class  $\mathcal{C}^{2m}$  in order to control derivatives of functions in  $\mathcal{H}(k)$  up to order  $m$ , as we shall  
 1801 see.

1802 By the Schwarz–Clairaut–Young theorem on the symmetry of mixed derivatives, these partial  
 1803 derivatives may be computed in any order, even when alternating between the two kernel variables.

1804 Suppose now that such an element  $k_x^\alpha$ , as described in Equation (86), exists for  $x \in X$  and a  
 1805 multi-index  $\alpha$ , and moreover lies in  $\mathcal{H}(k)$ . Then for each  $y \in X$ , the reproducing property gives

$$1806 \langle k_x^\alpha, k_y \rangle = k_x^\alpha(y) \stackrel{\text{equation 86}}{=} \partial^\alpha k_y(x) = \partial^{(\alpha, 0)} k(x, y). \quad (89)$$

1807 This entirely characterizes  $k_x^\alpha$  as a function in  $\mathcal{F}(X \rightarrow \mathbb{R})$ . Hence, three things remain to be checked:

- 1808 (i) that for all  $x \in X$ ,  $\partial^{(\alpha, 0)} k(x, \cdot) \in \mathcal{H}(k)$ ,
- 1809 (ii) that Equation (86) indeed holds for all  $f \in \mathcal{H}(k)$ ,
- 1810 (iii) that these partial derivatives are continuous.

1811 We will establish these properties by induction on the order  $m = |\alpha|$  of the partial derivative.

1812  **$m = 1$ ,  $k$  of class  $\mathcal{C}^2$ .** Fix  $1 \leq i \leq \dim(X)$  and  $y \in X$ .

1813 (i)  $\partial^{(i, 0)} k(x, \cdot) \in \mathcal{H}(k)$ . Since  $X \times X$  is equipped with the product topology, for any ball  $B$   
 1814 centered at  $y$  and contained in  $X$ , the set  $B \times X$  is a neighborhood of  $(y, x)$  for every  $x \in X$ . Then  
 1815 for any  $n > N$ , with  $N$  large enough so that  $y + 2^{-n} e_i \in B$ , we have  $k_{y+2^{-n} e_i} \in \mathcal{H}$ , and we may  
 1816 define

$$1817 c_n^i := \frac{k_{y+2^{-n} e_i} - k_y}{2^{-n}} \in \mathcal{H}. \quad (90)$$

1818 Furthermore, for any  $x \in X$ ,

$$1819 \frac{k_{y+2^{-n} e_i}(x) - k_y(x)}{2^{-n}} = \frac{k((y, x) + (2^{-n} e_i, 0)) - k((y, x))}{2^{-n}}, \quad (91)$$

1820 and this converges to  $\partial^{(i, 0)} k(x, y)$  as  $n \rightarrow \infty$ . In other words, the sequence  $(c_n^i)_{n \geq N}$  converges  
 1821 pointwise to the desired function.

1822 It remains to show that the convergence also holds in  $\|\cdot\|_{\mathcal{H}(k)}$ . Since convergence in  $\|\cdot\|_{\mathcal{H}(k)}$  implies  
 1823 pointwise convergence, uniqueness of the limit will then ensure that  $\partial^{(i, 0)} k(x, \cdot) \in \mathcal{H}(k)$ .

1836 Since  $\mathcal{H}(k)$  is complete, it is enough to show that  $(c_n^i)$  is Cauchy. To lighten the notation, let us write  
 1837  $h_n := 2^{-n}$  and, with a slight abuse, use  $k(x + h_n, \cdot)$  instead of  $k(x + h_n e_i, \cdot)$ . For  $p, q \geq N$ , we  
 1838 compute:

$$\begin{aligned}
 1839 \langle c_p^i, c_q^i \rangle_{\mathcal{H}(k)} &= \left\langle \frac{k_{x+h_p} - k_x}{h_p}, \frac{k_{x+h_q} - k_x}{h_q} \right\rangle_{\mathcal{H}(k)} & (92) \\
 1840 &= \frac{1}{h_p h_q} \left[ k(x + h_p, x + h_q) + k(x, x) - k(x + h_p, x) - k(x + h_q, x) \right] \\
 1841 &= \frac{1}{h_p} \left[ \frac{k(x + h_p, x + h_q) - k(x + h_p, x)}{h_q} - \frac{k(x + h_q, x) - k(x, x)}{h_q} \right].
 \end{aligned}$$

1842 Taking the limit as  $q \rightarrow \infty$ , this gives

$$1843 \lim_{q \rightarrow \infty} \langle c_p^i, c_q^i \rangle_{\mathcal{H}(k)} = \frac{1}{h_p} \left[ \partial^{(i,0)} k(x + h_p, x) - \partial^{(i,0)} k(x, x) \right], \quad (93)$$

1844 and therefore

$$1845 \lim_{p \rightarrow \infty} \lim_{q \rightarrow \infty} \langle c_p^i, c_q^i \rangle_{\mathcal{H}(k)} = \partial^{((i,i),0)} k(x, x). \quad (94)$$

1846 Note that if  $p = q$ , then since  $k$  is of class  $\mathcal{C}^2$  on  $X \times X$  (and in particular at  $(x, x)$ ), we also have

$$1847 \lim_{p \rightarrow \infty} \langle c_p^i, c_p^i \rangle_{\mathcal{H}(k)} = \partial^{((i,i),0)} k(x, x). \quad (95)$$

1848 Finally, observe that

$$1849 \langle c_n^i - c_m^i, c_n^i - c_m^i \rangle_{\mathcal{H}(k)} = \langle c_n^i, c_n^i \rangle_{\mathcal{H}(k)} + \langle c_m^i, c_m^i \rangle_{\mathcal{H}(k)} - 2 \langle c_n^i, c_m^i \rangle_{\mathcal{H}(k)}. \quad (96)$$

1850 Hence,

$$1851 \lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} \langle c_n^i - c_m^i, c_n^i - c_m^i \rangle_{\mathcal{H}(k)} = 0, \quad (97)$$

1852 which shows that  $(c_n^i)$  is indeed Cauchy in  $\mathcal{H}(k)$ .

1853 *Remark 5.* Using the symmetry of  $k$  in Equation (93), we also obtain

$$1854 \lim_{p \rightarrow \infty} \lim_{q \rightarrow \infty} \langle c_p^i, c_q^i \rangle_{\mathcal{H}(k)} = \partial^{(i,i)} k(x, x), \quad (98)$$

1855 which shows that for all  $x \in X$ ,

$$1856 \partial^{(i,i)} k(x, x) = \partial^{((i,i),0)} k(x, x). \quad (99)$$

1857 In particular, this identity requires only that  $k$  be of class  $\mathcal{C}^1$  in each of its variables separately.  
 1858 However, this is not sufficient to ensure convergence of the diagonal terms

$$1859 \lim_{p \rightarrow \infty} \langle c_p^i, c_p^i \rangle_{\mathcal{H}(k)} = \partial^{((i,i),0)} k(x, x). \quad (100)$$

1860 This clarifies why the stronger assumption  $k \in \mathcal{C}^{2m}$  is imposed in order to recover  $\mathcal{C}^m$  regularity for  
 1861 functions in  $\mathcal{H}(k)$ .

1862 **(ii) Equation (86) holds.** Let  $f \in \mathcal{H}(k)$  and  $x \in X$ . Then for any  $n \geq N$ ,

$$1863 \frac{f(x) - f(x + 2^{-n} e_i)}{2^{-n}} = \langle c_n^i, f \rangle_{\mathcal{H}(k)}. \quad (101)$$

1864 Since  $c_n^i$  converges to  $\partial^{(i,0)} k_x$  in  $\mathcal{H}(k)$ , the right-hand side converges to  $\langle \partial^{(i,0)} k_x, f \rangle_{\mathcal{H}(k)}$  in  $\mathbb{R}$ .  
 1865 This proves both that  $\partial^i f(x)$  exists and that it is reproduced by  $\partial^{(i,0)} k_x$ .

1866 **(iii)  $\partial^i f(\cdot)$  is continuous.** This is essentially an adaptation of the proof of the ??, now applied to  
 1867  $\partial^{(i,0)} k_{(\cdot)}$ . Namely, given  $x \in X$  and a sequence  $(x_n) \in X^{\mathbb{N}}$  converging to  $x$ ,

$$1868 |\partial^i f(x_n) - \partial^i f(x)| \leq \left\| \partial^{(i,0)} k_{x_n} - \partial^{(i,0)} k_x \right\|_{\mathcal{H}(k)} \|f\|_{\mathcal{H}(k)}. \quad (102)$$

Moreover, applying the reproducing property of Equation (86) to  $\partial^{(\alpha,0)}k_x$  itself,

$$\left\| \partial^{(i,0)}k_{x_n} - \partial^{(i,0)}k_x \right\|_{\mathcal{H}(k)}^2 = \partial^{(i,i)}k(x, x) + \partial^{(i,i)}k(x_n, x_n) - 2\partial^{(i,i)}k(x, x_n), \quad (103)$$

which converges to 0 as  $n \rightarrow \infty$ , since  $k$  is of class  $\mathcal{C}^2$ .

**Induction step:  $k$  of class  $\mathcal{C}^{2(m+1)}$ .** Suppose that for every multi-index  $\alpha$  with  $|\alpha| \leq m$ , we have  $\partial^{(\alpha,0)}k_x \in \mathcal{H}(k)$  and that it satisfies Equation (86), thereby yielding continuous partial derivatives.

Fix now a multi-index  $\alpha$  such that  $|\alpha| = m$ , let  $1 \leq i \leq \dim(X)$ , and  $y \in X$ . We shall denote by  $[\alpha, i]$  the concatenation of the multi-index  $\alpha$  and the index  $i$ .

As in the case  $m = 1$ , introduce the sequence

$$c_n^{[\alpha, i]} := \frac{\partial^{(\alpha,0)}k_{y+2^{-n}e_i} - \partial^{(\alpha,0)}k_y}{2^{-n}}. \quad (104)$$

Every term of this sequence lies in  $\mathcal{H}(k)$ , since by the induction hypothesis both  $\partial^{(\alpha,0)}k_{y+2^{-n}e_i}$  and  $\partial^{(\alpha,0)}k_y$  belong to  $\mathcal{H}(k)$ .

By the same arguments as in the base case, we see that this sequence converges pointwise to  $\partial^{([\alpha, i], 0)}k_y$  and is Cauchy in  $\mathcal{H}(k)$ . Hence it converges in  $\|\cdot\|_{\mathcal{H}(k)}$  to  $\partial^{([\alpha, i], 0)}k_y$ , showing that

$$\partial^{([\alpha, i], 0)}k_y \in \mathcal{H}(k). \quad (105)$$

This proves point (i).

Points (ii) and (iii) follow by entirely analogous arguments, and we shall not repeat the details here.  $\square$

**Corollary 1.** *Assume  $\Omega$  is of finite measure for  $\mu$ . Let  $D$  be an operator of order at most  $m$ . Then  $D : \mathcal{H}(k) \rightarrow L^2(\Omega)$  is continuous.*

*Proof.* This is an immediate consequence of previous proposition.  $\square$

Note in particular, that there exist thus  $C_D > 0$  such that for all  $u \in \mathcal{H}(k)$

$$\|D[u]\|_{L^2(\Omega)} \leq C_D \|u\|_k \quad (106)$$

We can thus define the adjoint of  $D : \mathcal{H}(k) \rightarrow L^2(\Omega)$  and thus the gram operator  $D^*D$ . Equivalently this is given the Gram kernel:

$$\mathbb{G}(x, y) := \langle D[k(x, \cdot)], D[k(y, \cdot)] \rangle_{L^2(\Omega)}, \quad (107)$$

By the same mechanism with spectral theorem depicted in Section G.3, we may regularize this operator by cutoff at level  $\alpha$ , yielding an approximation space  $\mathcal{H}(k, \alpha) \subset \mathcal{H}(k)$  such that for all  $u \in \mathcal{H}(k, \alpha)$

$$\|D[u]\|_{L^2(\Omega)}^2 \geq \alpha^2 \|u\|_k^2. \quad (108)$$

Overall this shows that  $u \mapsto \|D[u]\|_{L^2(\Omega)}$  is a norm equivalent to  $\|\cdot\|_k$  in  $\mathcal{H}(k, \alpha)$ , we can then apply then results of Section G.4 to show that  $D$  has a Green's function on  $\mathcal{H}(k, \alpha)$ .

In particular, we can then apply standard kernel regression results with Green's functions (?) to show convergence of AMStramGRAM in the NTK-regime at this level of regularization  $\alpha$ .

The only missing step is then to show that AMStramGRAM yields a stable regularization level. Let  $(r_{\max, t})_{t \in \mathbb{N}}$  be the the maximum cutoff rank sequence and  $\epsilon$  the precision level. We the have the following lemma that concludes the convergence of AMStramGRAM.

**Lemma 3.** *The sequence  $\left( \max_{t \in \mathbb{N}} (\Delta_{r_{\max, t}}, \epsilon) \right)_{t \in \mathbb{N}}$  is convergent.*

*Proof.* By Line 11 in Algorithm 2, we know that  $(\Delta_{r_{\max, t}})$  is non-increasing. Furthermore we obviously have for all  $t \in \mathbb{N}$ ,  $\max_{t \in \mathbb{N}} (\Delta_{r_{\max, t}}, \epsilon) \geq \epsilon$ . Therefore there exist  $\alpha_\infty > \epsilon$  such that  $\lim_{t \rightarrow \infty} \Delta_{r_{\max, t}} = \alpha_\infty$ .  $\square$

---

## 1944 H PROOFS

### 1945 H.1 STATEMENT AND PROOF OF PROPOSITION 4

1946 We start by recalling the following statements from Schwencke & Furtlehner (2025).

1947 **Definition** (Schwencke & Furtlehner, 2025, Definition 4). A linear operator  $A : \mathcal{H} \rightarrow \mathcal{H}$  is an  
1948 integral operator given that there is  $k : \Omega \times \Omega \rightarrow \mathbb{K}, \mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ , such that: for all  $f \in \mathcal{H}$ , for all  $x \in \Omega$

$$1949 A(f)(x) = \langle k(x, \cdot), f \rangle_{\mathcal{H}}. \quad (109)$$

1950 **Lemma** (Schwencke & Furtlehner, 2025, Lemma 1). Let us be  $\mathcal{H}_0 := \text{Span}(u_p : 1 \leq p \leq P) \subset \mathcal{H}$   
1951 and consider the Gram matrix  $G_{pq} := \langle u_p, u_q \rangle_{\mathcal{H}}$  of  $(u_p)$  and its eigen-decomposition  $G = U \Delta^2 U^t$ .  
1952 Then:

$$1953 L_p := \sum_{1 \leq q \leq P} u_q U_{q,p} \Delta_p^\dagger, \quad (110)$$

1954 is an orthonormal basis of  $\mathcal{H}_0$ . In particular,  $\Pi_{\mathcal{H}_0}$  is an integral operator whose kernel is:

$$1955 k(x, y) = \sum_{1 \leq p, q \leq P} u_p(x) G_{p,q}^\dagger u_q(y). \quad (111)$$

1956 Furthermore  $L_p$  are the left-singular vector of the so-called **synthesis operator**<sup>3</sup>:

$$1957 \mathcal{T} : \begin{cases} \mathbb{R}^P & \rightarrow \mathcal{H}_0 \\ \alpha & \mapsto \sum_{1 \leq p \leq P} \alpha_p u_p \end{cases}. \quad (112)$$

1958 **Proposition 4.** Given the scalar loss

$$1959 \ell(\theta) := \mathcal{L}(u_\theta) \stackrel{(6)}{=} \frac{1}{2} \|u_\theta - f\|_{L^2(\Omega, \mu)}^2, \quad (113)$$

1960 the Natural Gradient update of Equation (8)

$$1961 u_{\theta_{t+1}} \leftarrow u_{\theta_t} - \eta \Pi_{T_{\theta_t} \mathcal{M}}(\nabla \mathcal{L}_{u_{\theta_t}}); \quad \theta_{t+1} \leftarrow \theta_t - \eta \text{d}u_{\theta_t}^\dagger(\Pi_{T_{\theta_t} \mathcal{M}}(\nabla \mathcal{L}_{u_{\theta_t}})) \quad (8)$$

1962 can be equivalently written as

$$1963 \theta_{t+1} \leftarrow \theta_t - \eta G_{\theta_t}^\dagger \nabla \ell(\theta_t); \quad G_{\theta_{t,p,q}} := \langle \partial_p u_{\theta_t}, \partial_q u_{\theta_t} \rangle_{L^2(\Omega, \mu)}. \quad (9)$$

1964 *Proof.* Since the tangent space  $T_\theta \mathcal{M}$  of Equation (7):

$$1965 T_\theta \mathcal{M} := \text{Im}(\text{d}u_\theta) = \text{Span}(\partial_p u_\theta : 1 \leq p \leq P) \subset \mathcal{H}, \quad (7)$$

1966 is finite-dimensional, we may invoke Schwencke & Furtlehner (2025, Lemma 1). This result shows  
1967 that the **Natural Neural Tangent Kernel (NNTK)**, given by

$$1968 NNTK_\theta(x, y) := \sum_{1 \leq p, q \leq P} (\partial_p u_\theta(x)) G_{\theta_{pq}}^\dagger (\partial_q u_\theta(y))^t, \quad G_{\theta_{p,q}} := \langle \partial_p u_\theta, \partial_q u_\theta \rangle_{\mathcal{H}}, \quad (114)$$

1969 is the kernel of the orthogonal projection  $\Pi_{T_\theta \mathcal{M}}^\perp$  onto  $T_\theta \mathcal{M}$ . Therefore, by Equation (109), for all  
1970  $x \in \Omega$ ,

$$1971 \Pi_{T_\theta \mathcal{M}}^\perp(\nabla \mathcal{L}_{|u_\theta})(x) = \langle NNTK_\theta(x, \cdot), \nabla \mathcal{L}_{|u_\theta} \rangle_{\mathcal{H}} \\ 1972 \stackrel{(114)}{=} \sum_{1 \leq p, q \leq P} \partial_p u_\theta(x) G_{\theta_{pq}}^\dagger \langle \partial_q u_\theta, \nabla \mathcal{L}_{|u_\theta} \rangle_{\mathcal{H}}. \quad (115)$$

1973 Next, note that

$$1974 \langle \partial_q u_\theta, \nabla \mathcal{L}_{|u_\theta} \rangle_{\mathcal{H}} = \text{d}\mathcal{L}_{|u_\theta}(\partial_q u_\theta) \stackrel{\text{chain rule}}{=} \partial_q \mathcal{L}(u_\theta) \stackrel{(113)}{=} \partial_q \ell(\theta). \quad (116)$$

1975 <sup>3</sup>Name and notation are taken from Adcock & Huybrechs (2019).

Therefore, by linearity of  $du_{\theta}^{\dagger}$ ,

$$du_{\theta}^{\dagger}(\Pi_{\mathcal{T}_{\theta}, \mathcal{M}}^{\perp}(\nabla \mathcal{L}|_{u_{\theta}})) \stackrel{(115), (116)}{=} \sum_{1 \leq p, q \leq P} du_{\theta}^{\dagger}(\partial_p u_{\theta}) G_{\theta pq}^{\dagger} \partial_q \ell(\theta). \quad (117)$$

Finally, observe that  $\partial_p u_{\theta} = du_{\theta}(e^{(p)})$ , where  $e^{(p)}$  is the  $p$ -th canonical basis vector of  $\mathbb{R}^P$ . If  $du_{\theta}$  were invertible, we would directly obtain

$$du_{\theta}^{\dagger}(\partial_p u_{\theta}) = e^{(p)}, \quad (118)$$

which would complete the argument. However, this invertibility does not hold in general.

To address this, recall that  $du_{\theta}$  can be identified with the synthesis operator  $\mathcal{T}$  introduced in Equation (112) of Schwencke & Furtlehner (2025, Lemma 1). From the final part of that lemma, we know that  $\text{Im } du_{\theta}^{\dagger} = \text{Im } G_{\theta}^{\dagger}$ . Consequently,

$$G_{\theta}^{\dagger} e^{(p)} = G_{\theta}^{\dagger} du_{\theta}^{\dagger}(\partial_p u_{\theta}). \quad (119)$$

Putting all pieces together yields the desired update rule, thereby completing the proof.  $\square$

## H.2 RIDGE-REGRESSION IMPLEMENTATION ANAGRAM

In the following, we show that a Ridge-regression can be implemented in ANaGRAM's update rule given by Equation (10).

**Proposition 5.** *A Ridge-regression can be implemented in the SVD-based update Equation (10) by replacing the pseudo-inverse  $\widehat{\Delta}^{\dagger}$  with*

$$\left( \frac{\widehat{\Delta}_{t,i}}{\widehat{\Delta}_{t,i}^2 + S\alpha} \right)_{1 \leq i \leq r_{svd}}. \quad (120)$$

*Proof.* As shown in (Schwencke & Furtlehner, 2025, Section E), the ANaGRAM's update of Equation (10):

$$\theta_{t+1} \leftarrow \theta_t - \eta \widehat{\phi}^{\dagger} \widehat{\nabla} \mathcal{L}_{\theta_t}; \quad \widehat{\phi}_{i,p} := \partial_p u_{\theta}(x_i); \quad \left( \widehat{\nabla} \mathcal{L}_{\theta} \right)_i := u_{\theta}(x_i) - f(x_i), \quad (10)$$

is equivalent to the update with the empirical matrix  $\widehat{\mathcal{G}}_{\theta}$ :

$$\theta_{t+1} \leftarrow \theta_t - \eta \widehat{\mathcal{G}}_{\theta_t}^{\dagger} \nabla \ell(\theta_t); \quad \widehat{\mathcal{G}}_{\theta_t} := \frac{1}{S} \widehat{\phi}_{\theta_t}^t \widehat{\phi}_{\theta_t}, \quad (121)$$

where  $\ell$  is defined in Equation (5):

$$\ell(\theta) := \frac{1}{2S} \sum_{i=1}^S (u_{\theta}(x_i) - f(x_i))^2. \quad (5)$$

Thus, we get immediately

$$\nabla \ell(\theta_t) = \frac{1}{S} \widehat{\phi}^t \widehat{\nabla} \mathcal{L}_{\theta} = \frac{1}{S} \widehat{U} \widehat{\Delta} \widehat{V}_{\theta}^t \widehat{\nabla} \mathcal{L}_{\theta}, \quad (122)$$

where we used the SVD decomposition of  $\widehat{\phi}$ :

$$\widehat{\phi} = \widehat{U} \widehat{\Delta} \widehat{V}_{\theta}^t. \quad (123)$$

Using Equation (123) again, we have

$$\widehat{\mathcal{G}}_{\theta} = \frac{1}{S} \widehat{U} \widehat{\Delta}_{\theta}^2 \widehat{U}_{\theta}^t, \quad (124)$$

thus for all  $\alpha > 0$

$$\widehat{\mathcal{G}}_{\theta} + \alpha I_d = \frac{1}{S} \widehat{U} \widehat{\Delta}_{\theta}^2 \widehat{U}_{\theta}^t + \alpha \widehat{U} \widehat{U}^t = \widehat{U} \left( \text{diag} \left( \frac{\widehat{\Delta}_{\theta_i}^2}{S} + \alpha \right)_{1 \leq i \leq r_{svd}} \right) \widehat{U}_{\theta}^t, \quad (125)$$

which implies

$$\left(\widehat{\mathcal{G}}_{\theta} + \alpha I_d\right)^{-1} = \widehat{U} \left( \text{diag} \left( \frac{S}{\widehat{\Delta}_{\theta_i}^2 + S\alpha} \right)_{1 \leq i \leq r_{\text{svd}}} \right) \widehat{U}_{\theta}^t. \quad (126)$$

This finally yields

$$\left(\widehat{\mathcal{G}}_{\theta} + \alpha I_d\right)^{-1} \nabla \ell(\theta_t) \stackrel{(122)}{=} \widehat{U} \left( \text{diag} \left( \frac{S}{\widehat{\Delta}_{\theta_i}^2 + S\alpha} \right)_{1 \leq i \leq r_{\text{svd}}} \right) \widehat{U}_{\theta}^t \frac{1}{S} \widehat{U} \widehat{\Delta} \widehat{V}_{\theta}^t \widehat{\nabla} \mathcal{L}_{\theta} \quad (127)$$

$$= \widehat{U} \left( \text{diag} \left( \frac{\widehat{\Delta}_{t,i}}{\widehat{\Delta}_{\theta_i}^2 + S\alpha} \right)_{1 \leq i \leq r_{\text{svd}}} \right) \widehat{V}_{\theta}^t \widehat{\nabla} \mathcal{L}_{\theta}, \quad (128)$$

which concludes the proof.  $\square$

### H.3 PROOF OF PROPOSITION 1

To prove Proposition 1, we need the following lemma:

**Lemma 4.** For  $1 \leq M \leq N \leq r_{\text{svd}}$ :

$$\left(\text{RCE}_M^S\right)^2 - \left(\text{RCE}_N^S\right)^2 = \frac{1}{S} \left\| \Pi_N^M \widehat{V}^T \widehat{\nabla} \mathcal{L} \right\|_{\mathbb{R}^S}^2. \quad (129)$$

*Proof.* Let us first recall the definition of the  $\text{RCE}_N^S$  in Equation (13), namely

$$\text{RCE}_N^S := \frac{1}{\sqrt{S}} \left\| \widehat{V} \Pi_N^0 \widehat{V}^T \widehat{\nabla} \mathcal{L} - \widehat{\nabla} \mathcal{L} \right\|_{\mathbb{R}^S}. \quad (13)$$

Fixing  $1 \leq N \leq M \leq r_{\text{svd}}$  and applying the same reasoning as in Equation (139) to  $\text{RCE}_M^S$  and  $\text{RCE}_N^S$  (see the proof of Proposition 1 in Section H.3), we get

$$S \left(\text{RCE}_M^S\right)^2 = \widehat{\nabla} \mathcal{L}_{\theta}^t \widehat{\nabla} \mathcal{L}_{\theta} - \sum_{p=1}^M \left( \widehat{V}_{\theta_p}^t \widehat{\nabla} \mathcal{L}_{\theta} \right)^2; \quad S \left(\text{RCE}_N^S\right)^2 = \widehat{\nabla} \mathcal{L}_{\theta}^t \widehat{\nabla} \mathcal{L}_{\theta} - \sum_{p=1}^N \left( \widehat{V}_{\theta_p}^t \widehat{\nabla} \mathcal{L}_{\theta} \right)^2, \quad (130)$$

and therefore

$$\begin{aligned} S \left( \left(\text{RCE}_N^S\right)^2 - \left(\text{RCE}_M^S\right)^2 \right) &= \sum_{p=1}^N \left( \widehat{V}_{\theta_p}^t \widehat{\nabla} \mathcal{L}_{\theta} \right)^2 - \sum_{p=1}^M \left( \widehat{V}_{\theta_p}^t \widehat{\nabla} \mathcal{L}_{\theta} \right)^2 \\ &\stackrel{M \leq N}{=} \sum_{p=M+1}^N \left( \widehat{V}_{\theta_p}^t \widehat{\nabla} \mathcal{L}_{\theta} \right)^2 = \sum_{p=M+1}^N \left( e^{(p)t} \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta} \right)^2 \\ &= \sum_{p=M+1}^N \left( \widehat{\nabla} \mathcal{L}_{\theta}^t \widehat{V} e^{(p)} \right) \left( e^{(p)t} \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta} \right) \\ &= \widehat{\nabla} \mathcal{L}_{\theta}^t \widehat{V} \left( \underbrace{\sum_{p=M+1}^N e^{(p)} e^{(p)t}}_{=\Pi_N^M \text{ by Equation (14)}} \right) \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta} \\ &= \left\langle \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta}, \Pi_N^M \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta} \right\rangle_{\mathbb{R}^S} \\ &\stackrel{\Pi_N^M = \Pi_N^M}{=} \left\langle \Pi_N^M \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta}, \Pi_N^M \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta} \right\rangle_{\mathbb{R}^S} \\ &= \left\| \Pi_N^M \widehat{V}^t \widehat{\nabla} \mathcal{L}_{\theta} \right\|_{\mathbb{R}^S}^2, \end{aligned} \quad (131)$$

where we use in the penultimate equality, the fact that  $\Pi_N^M$  is an orthogonal projection.  $\square$

2106 *Remark 6.* The above lemma provides an interesting property that gives a further understanding of  
 2107 what is happening during the flattening, *i.e.*  $\text{RCE}_M^S - \text{RCE}_N^S \approx 0$ . In particular, as  $(\text{RCE}_M^S)^2 -$   
 2108  $(\text{RCE}_N^S)^2 = (\text{RCE}_M^S - \text{RCE}_N^S)(\text{RCE}_M^S + \text{RCE}_N^S)$ , therefore flattening for the components in the  
 2109 range  $[N_{\text{flat}}, r_{\text{cutoff}}]$  means that  $\frac{1}{S} \left\| \Pi_N^M \widehat{V}^T \widehat{\nabla} \mathcal{L} \right\|_{\mathbb{R}^S}^2 \approx 0$ . In other words, the problem is "learned" for  
 2110 those components, as the projection of the functional gradient (which is proportional to the error) on  
 2111 their corresponding span is null. The proof of this lemma is provided in Section H.3.  
 2112

2113 **Proposition 1.**  $\text{RCE}_N^S$  is a non-increasing function of  $N$ , *i.e.* for all  $1 \leq M, N \leq r_{\text{svd}}$ :

$$2114 \quad M \leq N \implies \text{RCE}_M^S \geq \text{RCE}_N^S. \quad (15)$$

2115 Furthermore, assuming that  $(x_i)_{i=1}^S$  are *i.i.d* sampled from  $\mu$ , we have  $\mu$ -almost surely

$$2116 \quad \lim_{S \rightarrow \infty} \text{RCE}_N^S = \left\| \nabla \mathcal{L}_{u_\theta} - \Pi_{T_N^0 \mathcal{M}}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega, \mu)} = \left\| \Pi_{[T_N^0 \mathcal{M}]^\perp}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega, \mu)}, \quad (16)$$

2117 where  $T_N^M \mathcal{M} = \text{Span}(V_{t,i} : M \leq i \leq N)$ , while  $(V_{t,i})_{1 \leq i \leq r_{\text{svd}}}$  are the right singular-vectors of the  
 2118 differential  $du_\theta$  ordered in a decreasing order according to their associated singular values.  
 2119

2120 *Proof.* The first statement is a direct consequence of Lemma 4 proven above.

2121 Let us now show that the second statement takes place. Since  $\nabla \mathcal{L}_{u_\theta} \in L^2(\Omega, \mu)$  and  $\text{Im } du_\theta \subset$   
 2122  $L^2(\Omega, \mu)$ , the law of large numbers yields that for all  $1 \leq p, q \leq P$

$$2123 \quad \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{i=1}^S [\nabla \mathcal{L}_{u_\theta}(x_i)]^2 = \lim_{S \rightarrow \infty} \frac{1}{S} \widehat{\nabla} \mathcal{L}_\theta^t \widehat{\nabla} \mathcal{L}_\theta = \int_{\Omega} [\nabla \mathcal{L}_{u_\theta}(x)]^2 \mu(\mathrm{d}x) \quad a.s., \quad (132)$$

$$2124 \quad \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{i=1}^S \partial_p u_\theta(x_i) \nabla \mathcal{L}_{u_\theta}(x_i) = \lim_{S \rightarrow \infty} \frac{1}{S} \widehat{\phi}_{\theta_p}^t \widehat{\nabla} \mathcal{L}_\theta = \int_{\Omega} \partial_p u_\theta(x) \nabla \mathcal{L}_{u_\theta}(x) \mu(\mathrm{d}x) \quad a.s., \quad (133)$$

$$2125 \quad \lim_{S \rightarrow \infty} \frac{1}{S} \sum_{i=1}^S \partial_p u_\theta(x_i) \partial_q u_\theta(x_i) = \lim_{S \rightarrow \infty} \frac{1}{S} \widehat{\phi}_{\theta_p}^t \widehat{\phi}_{\theta_q} = \int_{\Omega} \partial_p u_\theta(x) \partial_q u_\theta(x) \mu(\mathrm{d}x) \quad a.s. \quad (134)$$

2126 In particular, this implies

$$2127 \quad \lim_{S \rightarrow \infty} \frac{1}{S} \widehat{\phi}^t \widehat{\phi} = G_\theta = U_\theta \Delta_\theta^2 U_\theta^t \quad a.s. \quad (135)$$

2128 Since the eigenvectors (and eigenvalues) are continuous functions of the matrix coefficients (by  
 2129 polynomial dependence through the characteristic polynomial) and taking into account that  $\frac{1}{S} \widehat{\phi}^t \widehat{\phi} =$   
 2130  $\frac{1}{S} \widehat{U} \Delta_\theta^2 \widehat{U}^t$ , this yields

$$2131 \quad \lim_{S \rightarrow \infty} \widehat{U} = U_\theta \quad a.s.; \quad \lim_{S \rightarrow \infty} \frac{1}{S} \widehat{\Delta}^2 = \Delta_\theta^2 \quad a.s. \quad (136)$$

2132 By continuity of the square root and the inverse on  $\mathbb{R}_+^*$ , we get that for all  $1 \leq p \leq P$  such that  
 2133  $\Delta_{\theta_p} > 0$

$$2134 \quad \lim_{S \rightarrow \infty} \sqrt{S} \widehat{\Delta}_{\theta_p}^{-1} = \Delta_{\theta_p}^{-1} \quad a.s., \quad (137)$$

2135 and thus for all  $1 \leq p \leq P$  such that  $\Delta_{\theta_p} > 0$ , we have *almost surely*

2160

2161

2162

2163

2164

2165

2166

2167

2168

2169

2170

2171

2172

2173

2174

2175

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

2187

2188

2189

2190

2191

2192

2193

2194

2195

2196

2197

2198

2199

2200

2201

2202

2203

2204

2205

2206

2207

2208

2209

2210

2211

2212

2213

$$\begin{aligned}
\lim_{S \rightarrow \infty} \frac{1}{\sqrt{S}} \widehat{V}_{\theta_p}^T \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} &= \lim_{S \rightarrow \infty} \sqrt{S} \widehat{\Delta}_{\theta_p}^{-1} \widehat{U}_{\theta_p}^t \left( \sum_{q=1}^P e^{(q)} e^{(q)t} \right) \frac{1}{S} \widehat{\phi}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \\
&= \sum_{q=1}^P \left( \lim_{S \rightarrow \infty} \sqrt{S} \widehat{\Delta}_{\theta_p}^{-1} \right) \left( \lim_{S \rightarrow \infty} \widehat{U}_{\theta_p}^t e^{(q)} \right) \left( \lim_{S \rightarrow \infty} \frac{1}{S} \widehat{\phi}_{\theta_q}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \right) \\
&= \sum_{q=1}^P \Delta_{\theta_p}^{-1} U_{\theta_p}^t e^{(q)} \int_{\Omega} \partial_q u_{\theta}(x) \nabla \mathcal{L}_{u_{\theta}}(x) \mu(\mathrm{d}x) \\
&= \int_{\Omega} \mathrm{d}u_{\theta} \left( U_{\theta_p} \Delta_{\theta_p}^{-1} \right) (x) \nabla \mathcal{L}_{u_{\theta}}(x) \mu(\mathrm{d}x) \\
&= \int_{\Omega} V_{\theta_p}(x) \nabla \mathcal{L}_{u_{\theta}}(x) \mu(\mathrm{d}x),
\end{aligned} \tag{138}$$

where we used in the last equality, the identification of the singular vectors of  $\mathrm{d}u_{\theta}$  in (Schwencke & Furtlehner, 2025, Lemma 1 p. 24, section C.2). Returning to the definition of the RCE $_N^S$  in Equation (13), namely

$$\text{RCE}_N^S := \frac{1}{\sqrt{S}} \left\| \widehat{V} \Pi_N^0 \widehat{V}^T \widehat{\nabla} \widehat{\mathcal{L}} - \widehat{\nabla} \widehat{\mathcal{L}} \right\|_{\mathbb{R}^S}, \tag{13}$$

we get

$$\begin{aligned}
S (\text{RCE}_N^S)^2 &= \left\langle \widehat{V} \Pi_N^0 \widehat{V}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} - \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}, \widehat{V} \Pi_N^0 \widehat{V}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} - \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \right\rangle_{\mathbb{R}^S} \\
&= \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} + \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{V} \overbrace{\Pi_N^0 \widehat{V}^t \widehat{V} \Pi_N^0}^{=I_d} \widehat{V}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} - 2 \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{V} \Pi_N^0 \widehat{V}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \\
&= \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} - \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{V} \Pi_N^0 \widehat{V}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \\
&= \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} - \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{V} \left( \sum_{p=1}^N e^{(p)} e^{(p)t} \right) \widehat{V}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \\
&= \widehat{\nabla} \widehat{\mathcal{L}}_{\theta}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} - \sum_{p=1}^N \left( \widehat{V}_{\theta_p}^t \widehat{\nabla} \widehat{\mathcal{L}}_{\theta} \right)^2,
\end{aligned} \tag{139}$$

where in the second equality, we use the fact that  $\widehat{V}$  is orthogonal and  $\Pi_N^0$  is a projection. Combining Equations (132) and (138), this yields

$$\lim_{S \rightarrow \infty} (\text{RCE}_N^S)^2 = \int_{\Omega} \nabla \mathcal{L}_{u_{\theta}}(x)^2 \mu(\mathrm{d}x) - \sum_{p=1}^N \left( \int_{\Omega} V_{\theta_p}(x) \nabla \mathcal{L}_{u_{\theta}}(x) \mu(\mathrm{d}x) \right)^2 \quad a.s. \tag{140}$$

By Fubini's theorem, we have *almost surely*

$$\begin{aligned}
\sum_{p=1}^N \left( \int_{\Omega} V_{\theta_p}(x) \nabla \mathcal{L}_{u_{\theta}}(x) \mu(\mathrm{d}x) \right)^2 &= \int_{\Omega^2} \nabla \mathcal{L}_{u_{\theta}}(x) \left( \sum_{p=1}^N V_{\theta_p}(x) V_{\theta_p}(y) \right) \nabla \mathcal{L}_{u_{\theta}}(y) \mu(\mathrm{d}x) \mu(\mathrm{d}y) \\
&= \int_{\Omega} \nabla \mathcal{L}_{u_{\theta}}(x) \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^{\perp} \nabla \mathcal{L}_{u_{\theta}}(x) \mu(\mathrm{d}x) \\
&= \left\| \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^{\perp} \nabla \mathcal{L}_{u_{\theta}} \right\|_{L^2(\Omega, \mu)}^2,
\end{aligned} \tag{141}$$

where in the second equality, we used (Schwencke & Furtlehner, 2025, Theorem 4 p. 23, section C.2) and the fact that  $\left( \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^{\perp} \right)^2 = \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^{\perp}$  in the third. Therefore, from Equation (140) and Equation (141)

$$\lim_{S \rightarrow \infty} (\text{RCE}_N^S)^2 = \left\| \nabla \mathcal{L}_{u_{\theta}} \right\|_{L^2(\Omega, \mu)}^2 - \left\| \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^{\perp} \nabla \mathcal{L}_{u_{\theta}} \right\|_{L^2(\Omega, \mu)}^2 \quad a.s., \tag{142}$$

$$= \left\| \nabla \mathcal{L}_{u_{\theta}} - \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^{\perp} \nabla \mathcal{L}_{u_{\theta}} \right\|_{L^2(\Omega, \mu)}^2 \quad a.s., \tag{143}$$

where in the second equality, we use in the reverse order a reasoning similar to Equation (139). Finally, we obtain

$$\left\| \nabla \mathcal{L}_{u_\theta} - \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega, \mu)}^2 = \left\| \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega, \mu)}^2, \quad (144)$$

which comes from the canonical decomposition in Hilbert spaces, *i.e.* using that  $\text{Span}(V_{\theta_i} : 1 \leq i \leq N)$  is a closed subspace and

$$\nabla \mathcal{L}_{u_\theta} = \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)}^\perp \nabla \mathcal{L}_{u_\theta} + \Pi_{\text{Span}(V_{\theta_i} : 1 \leq i \leq N)} \nabla \mathcal{L}_{u_\theta}. \quad (145)$$

This completes the proof.  $\square$

**Corollary 2.** For  $1 \leq M \leq N \leq r_{\text{svd}}$ :

$$\lim_{S \rightarrow \infty} (RCE_M^S)^2 - (RCE_N^S)^2 = \left\| \Pi_{T_N^M \mathcal{M}}^\perp \nabla \mathcal{L}_{u_\theta} \right\|_{L^2(\Omega)}^2 \quad (146)$$

*Proof.* Apply Proposition 1 to Equation (129) of Lemma 4.  $\square$

## I APPENDIX: EMPIRICAL ANALYSIS OF BATCHWISE TRUNCATED NATURAL GRADIENT

### I.1 EMPIRICAL ANALYSIS OF KERNEL ALIGNMENT

To quantify the stability of the natural gradient update under resampling, we introduce the *inverse kernel alignment* metric. This metric measures the similarity between the tangent spaces induced by the same kernel on different batches of training data. Formally, let  $B_1, B_2$  be two independent training batches and  $B_{\text{test}}$  a test batch. For each batch  $k \in \{1, 2\}$ , we compute the SVD of the feature matrix  $\hat{\phi}_k = \hat{V}_k \hat{\Delta}_k \hat{U}_k^\top$ . The truncated inverse Gram matrix at rank  $r$  is given by  $G_k^{\dagger(r)} = \sum_{i=1}^r \sigma_{k,i}^{-2} \mathbf{u}_{k,i} \mathbf{u}_{k,i}^\top$ , where  $\mathbf{u}_{k,i}$  is the  $i$ -th column of  $\hat{U}_k$ . We define the alignment  $\mathcal{A}(r)$  as the cosine similarity between  $G_1^{\dagger(r)}$  and  $G_2^{\dagger(r)}$  relative to the geometry of the test set:

$$\mathcal{A}(r) := \frac{\langle G_1^{\dagger(r)}, G_2^{\dagger(r)} \rangle_{\Sigma_{\text{test}}}}{\|G_1^{\dagger(r)}\|_{\Sigma_{\text{test}}} \|G_2^{\dagger(r)}\|_{\Sigma_{\text{test}}}}, \quad (147)$$

where  $\Sigma_{\text{test}} = \hat{\phi}_{\text{test}}^\top \hat{\phi}_{\text{test}}$  and  $\langle A, B \rangle_{\Sigma_{\text{test}}} = \text{Tr}(A \Sigma_{\text{test}} B \Sigma_{\text{test}})$ . This metric evaluates whether the optimization directions prescribed by different batches are consistent when applied to the test distribution. Computationally, this reduces to comparing the projected test features  $W_k = \hat{\phi}_{\text{test}}^\top \hat{U}_k$ : the numerator becomes  $\sum_{i,j=1}^r \sigma_{1,i}^{-2} \sigma_{2,j}^{-2} (\mathbf{w}_{1,i}^\top \mathbf{w}_{2,j})^2$ , where  $\mathbf{w}_{k,i}$  are columns of  $W_k$ .

Figure 20 and Figure 21 show the alignment curves for the heat equation. We observe a characteristic behavior:

- At low quantities of retained components (small  $r$ ), the kernels are very similar. This corresponds to the largest eigenvalues, which capture the dominant modes of the loss landscape and are less sensitive to changes in the training samples.
- At large quantities of retained components (large  $r$ ), the kernels again show high alignment, as they both describe almost the same space (approaching the full empirical tangent space).
- In the intermediate regime, the alignment drops, indicating that the subspaces spanned by the intermediate eigenvectors are more sensitive to the specific batch of data points.

This analysis confirms that truncating the spectrum (as done in AMStraMGRAM) or using the full spectrum is relatively stable, whereas intermediate truncations might be more sensitive to sampling noise.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

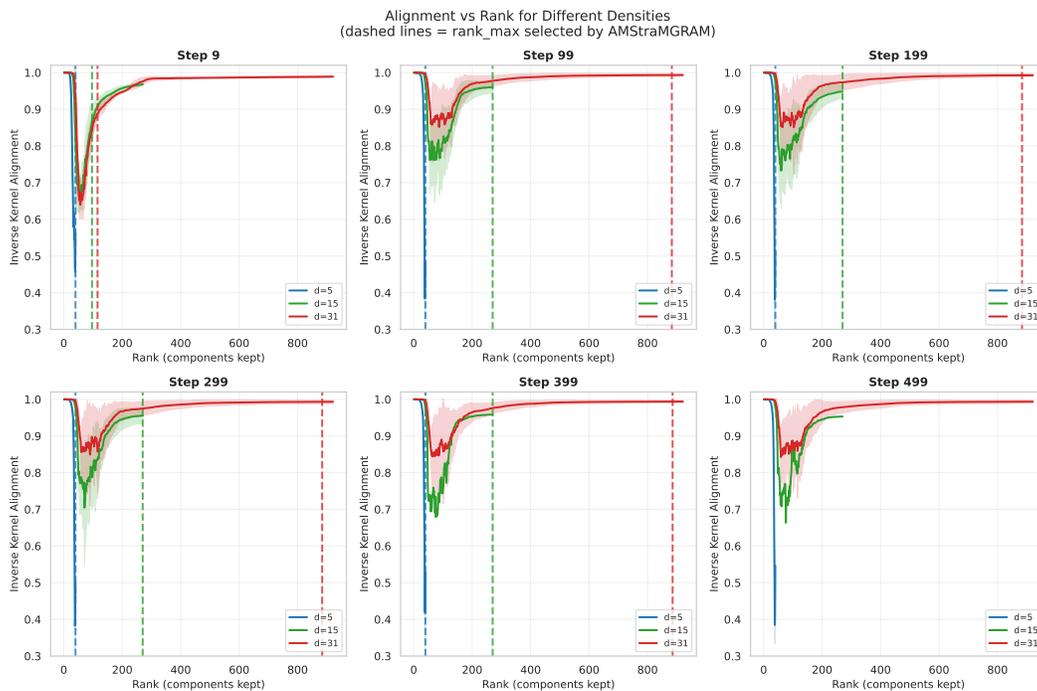


Figure 20: Inverse kernel alignment as a function of the rank  $r$  (number of retained components) for the heat equation: Impact of Sample Density. The dashed lines indicate the rank selected by AMStraMGRAM.

## J LAZY TRAINING FOR THE NATURAL GRADIENT (CLASSICAL VERSION)

We recall from (Chizat et al., 2019, Section 1.2) that lazy training is characterized by monitoring the following two quantities:

**Relative change of the loss** defined by the ratio

$$\Delta(\ell)(\theta, h) := \frac{|\ell(\theta + h) - \ell(\theta)|}{\ell(\theta)}. \quad (148)$$

**Relative change of the model differential** defined by the ratio

$$\Delta(du)(\theta, h) := \frac{\|du_{\theta+h} - du_{\theta}\|}{\|du_{\theta}\|}, \quad (149)$$

where  $\|\cdot\|$  denotes the operator norm, that is,

$$\|du_{\theta}\| = \sup_{\|\delta\|_{\mathbb{R}^P}=1} \|du_{\theta}(\delta)\|_{L^2(\Omega)}. \quad (150)$$

Lazy training occurs when  $\Delta(du)(\theta, h) \ll \Delta(\ell)(\theta, h)$ , meaning that the model differential remains essentially constant while the loss may vary significantly; *i.e.*, the training dynamics stay close to those of the linearized model.

Equivalently, this condition can be expressed by introducing the ratio

$$\kappa_{u_{\theta}}(\theta, h) := \frac{\Delta(du)(\theta, h)}{\Delta(\ell)(\theta, h)}, \quad (151)$$

and requiring that  $\kappa_{u_{\theta}}(\theta, h) \ll 1$ .

In Chizat et al. (2019), an explicit formula is derived in the setting of quadratic regression under standard gradient descent. However, the situation differs in our case because the optimization relies

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

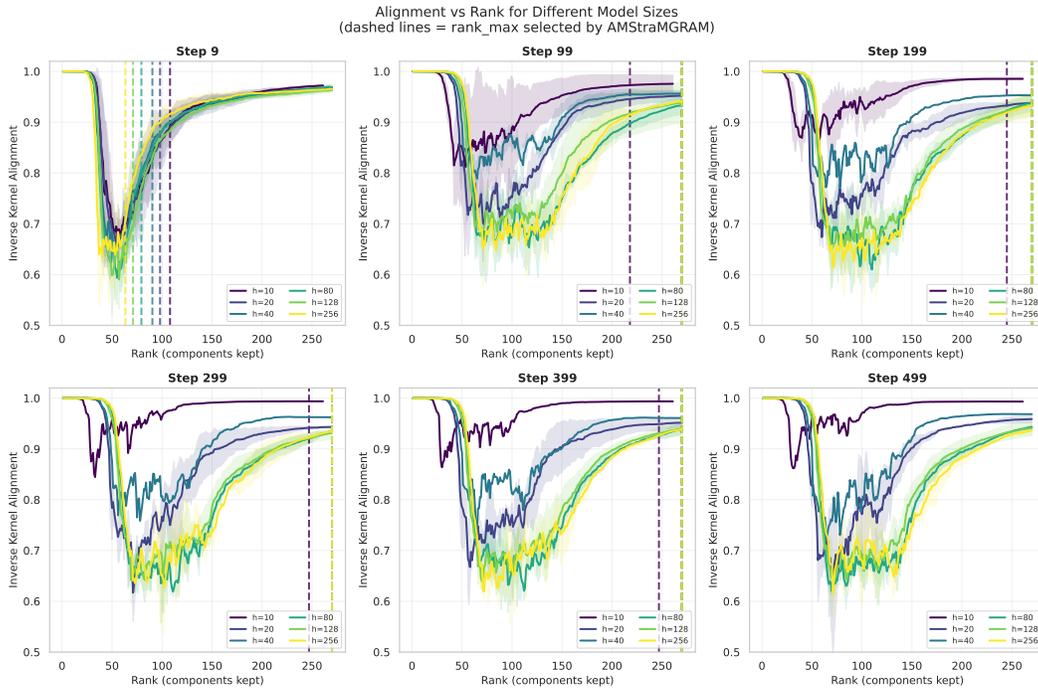


Figure 21: Inverse kernel alignment as a function of the rank  $r$  (number of retained components) for the heat equation: Impact of Model Size. The dashed lines indicate the rank selected by AMStramGRAM.

on the natural gradient. We therefore proceed step by step. As in Chizat et al. (2019), we start from the Taylor expansion

$$\begin{aligned}
\ell(\theta + h) &= \ell(\theta) + d\ell_\theta(h) + o(\|h\|_{\mathbb{R}^P}) \\
&= \mathcal{L}(u_\theta) + \langle du_\theta(h), \nabla \mathcal{L} \rangle_{L^2(\Omega)} + o(\|h\|_{\mathbb{R}^P}) \\
&= \frac{1}{2} \|u_\theta - f\|_{L^2(\Omega)}^2 + \langle du_\theta(h), u_\theta - f \rangle_{L^2(\Omega)} + o(\|h\|_{\mathbb{R}^P}),
\end{aligned} \tag{152}$$

from which we obtain, for  $\|h\|_{\mathbb{R}^P}$  sufficiently small,

$$|\ell(\theta + h) - \ell(\theta)| \simeq |d\ell_\theta(h)| = \left| \langle du_\theta(h), \nabla \mathcal{L} \rangle_{L^2(\Omega)} \right| = \left| \langle du_\theta(h), u_\theta - f \rangle_{L^2(\Omega)} \right|. \tag{153}$$

Since  $h$  follows the regularized natural gradient at level  $N$ , we have the identity

$$du_\theta(h) = \eta \Pi_N(\nabla \mathcal{L}) = \eta \Pi_N(f - u_\theta). \tag{154}$$

Substituting equation 154 into equation 153 gives

$$|\ell(\theta + h) - \ell(\theta)| \simeq \left| \langle \eta \Pi_N(f - u_\theta), u_\theta - f \rangle_{L^2(\Omega)} \right| = \eta \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}^2. \tag{155}$$

Thus, the quantity  $\Delta(\ell)$  in Equation (148) becomes, under the regularized natural gradient at level  $N$ , denoted  $\Delta(\ell)(\theta, N)$ ,

$$\Delta(\ell)(\theta, N) \simeq \frac{\eta \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}^2}{\frac{1}{2} \|f - u_\theta\|_{L^2(\Omega)}^2}. \tag{156}$$

Let us now focus on the relative change of the model differential  $\Delta(du)$ . Recall first that the SVD of  $du_\theta$  is given by

$$du_\theta = V_\theta \Delta_\theta U_\theta^\top. \tag{157}$$

2376 In particular,

$$2377 \quad \|du_\theta\| = \sup_{\|\delta\|_{\mathbb{R}^P}=1} \|du_\theta(\delta)\|_{L^2(\Omega)} = \Delta_{\theta_{\max}}. \quad (158)$$

2379 Next, for any  $\delta \in \mathbb{R}^P$ , another Taylor expansion yields

$$2381 \quad du_{\theta+h}(\delta) = du_\theta(\delta) + d^2u_\theta(\delta, h) + o(\|h\|_{\mathbb{R}^P}), \quad (159)$$

2382 from which we deduce, for  $\|h\|_{\mathbb{R}^P}$  sufficiently small,

$$2384 \quad \|du_{\theta+h} - du_\theta\| \simeq \|d^2u_\theta(\cdot, h)\| \leq C_2 \|h\|_{\mathbb{R}^P}, \quad (160)$$

2385 where

$$2387 \quad C_2 := \|d^2u_\theta\| = \sup_{\substack{\|\delta_1\|_{\mathbb{R}^P}=1 \\ \|\delta_2\|_{\mathbb{R}^P}=1}} \|d^2u_\theta(\delta_1, \delta_2)\|_{L^2(\Omega)}. \quad (161)$$

2390 Combining equation 158 and equation 160 yields

$$2391 \quad \Delta(du) \lesssim \frac{C_2 \|h\|_{\mathbb{R}^P}}{\Delta_{\theta_{\max}}}. \quad (162)$$

2394 Substituting this estimate into Equation (151) yields

$$2396 \quad \kappa_{u_\theta}(\theta, N) \lesssim \frac{C_2 \|h\|_{\mathbb{R}^P} \frac{1}{2} \|f - u_\theta\|_{L^2(\Omega)}^2}{\Delta_{\theta_{\max}} \eta \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}^2}. \quad (163)$$

2399 We now estimate  $\|h\|_{\mathbb{R}^P}$ . Since  $h$  follows the natural-gradient update of Equation (10),

$$2401 \quad \theta_{t+1} \leftarrow \theta_t - \eta \hat{\phi}^\dagger \widehat{\nabla \mathcal{L}}_{\theta_t}; \quad \hat{\phi}_{i,p} := \partial_p u_\theta(x_i); \quad \left(\widehat{\nabla \mathcal{L}}_{\theta_t}\right)_i := u_\theta(x_i) - f(x_i), \quad (10)$$

2403 we obtain, omitting the learning rate  $\eta$  for readability,

$$2404 \quad \frac{1}{\eta^2} \|h\|_{\mathbb{R}^P}^2 = \left\| \hat{\phi}^\dagger \widehat{\nabla \mathcal{L}}_{\theta_t} \right\|_{\mathbb{R}^P}^2. \quad (164)$$

2407 Using the SVD  $\hat{\phi}^\dagger = \widehat{U} \widehat{\Delta}^\dagger \widehat{V}^T$ , this becomes

$$2408 \quad \begin{aligned} 2409 \quad \frac{1}{\eta^2} \|h\|_{\mathbb{R}^P}^2 &= \left\| \widehat{U} \widehat{\Delta}^\dagger \widehat{V}^T \widehat{\nabla \mathcal{L}}_{\theta_t} \right\|_{\mathbb{R}^P}^2 \\ 2410 &= \widehat{\nabla \mathcal{L}}_{\theta_t}^T \widehat{V} \widehat{\Delta}^\dagger (\widehat{U}^T \widehat{U}) \widehat{\Delta}^\dagger \widehat{V}^T \widehat{\nabla \mathcal{L}}_{\theta_t} \\ 2411 &= \widehat{\nabla \mathcal{L}}_{\theta_t}^T \widehat{V} (\widehat{\Delta}^\dagger)^2 \widehat{V}^T \widehat{\nabla \mathcal{L}}_{\theta_t} \\ 2412 &= \sum_{i=1}^N \frac{\left\| \widehat{V}_i^T \widehat{\nabla \mathcal{L}}_{\theta_t} \right\|_{\mathbb{R}^P}^2}{\widehat{\Delta}_i^2} \leq \frac{1}{\widehat{\Delta}_N^2} \left\| \Pi_N \widehat{\nabla \mathcal{L}}_{\theta_t} \right\|_{\mathbb{R}^P}^2, \end{aligned} \quad (165)$$

2414 where  $N$  is the number of retained components.

2418 By Equation (138) (see the proof of Proposition 1), we have almost surely

$$2420 \quad \lim_{S \rightarrow \infty} \frac{1}{\widehat{\Delta}_N^2} \left\| \Pi_N \widehat{\nabla \mathcal{L}}_{\theta_t} \right\|_{\mathbb{R}^P}^2 = \frac{1}{\Delta_N} \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}^2. \quad (166)$$

2422 Substituting this into Equation (163), we obtain almost surely

$$2424 \quad \kappa_{u_\theta}(\theta, N) \lesssim \frac{C_2 \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)} \frac{1}{2} \|f - u_\theta\|_{L^2(\Omega)}^2}{\Delta_{\theta_{\max}} \Delta_N \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}^2}. \quad (167)$$

2428 We now use the decomposition

$$2429 \quad \|f - u_\theta\|_{L^2(\Omega)}^2 = \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}^2 + \|\Pi_N^\perp(f - u_\theta)\|_{L^2(\Omega)}^2. \quad (168)$$

---

2430 By construction of the flattening phase,  
2431

$$2432 \quad \|\Pi_N^\perp(f - u_\theta)\|_{L^2(\Omega)}^2 \leq \epsilon^2, \quad (169)$$

2433 where  $\epsilon$  is the target accuracy. Plugging into Equation (167) yields  
2434

$$2435 \quad \kappa_{u_\theta}(\theta, N) \lesssim \frac{C_2 \frac{1}{2} \left( \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)} + \frac{\epsilon^2}{\|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}} \right)}{\Delta_{\theta_{\max}} \Delta_N}. \quad (170)$$

2436 Since  $\|\Pi_N(f - u_\theta)\|_{L^2(\Omega)} \gtrsim \epsilon$ , the second term is dominated, and redefining  $C := \frac{(1+\epsilon)C_2}{2}$ , we  
2437 obtain  
2438

$$2439 \quad \kappa_{u_\theta}(\theta, N) \lesssim \frac{C \|\Pi_N(f - u_\theta)\|_{L^2(\Omega)}}{\Delta_{\theta_{\max}} \Delta_N} \leq \frac{C \|f - u_\theta\|_{L^2(\Omega)}}{\Delta_{\theta_{\max}} \Delta_N}. \quad (171)$$

2440 Finally, note that  $\Delta_{\theta_{\max}} = \Delta_0$  is typically several orders of magnitude larger than  $\|f - u_\theta\|_{L^2(\Omega)}$ , and  
2441 that  $\Delta_N$  decreases exponentially fast. Hence, as  $N \rightarrow 0$ , the ratio  $\kappa_{u_\theta}(\theta, N)$  decays exponentially to  
2442 zero, which is precisely the hallmark of the lazy-training regime.  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483