

Design Criteria for SGD Preconditioners: Local Conditioning, Noise Floors, and Basin Stability

Mitchell Scott

*Department of Mathematics
Emory University*

mitchell.scott@emory.edu

Tianshi Xu

*Department of Mathematics
Emory University*

tianshi.xu@emory.edu

Ziyuan Tang

*Department of Computer Science
University of Minnesota Twin Cities*

tang0389@umn.edu

Alexandra Pichette-Emmons

*Department of Mathematics
University of Kentucky*

hugolarochelle@google.com

Qiang Ye

*Department of Mathematics
University of Kentucky*

qye3@uky.edu

Yousef Saad

*Department of Computer Science
University of Minnesota Twin Cities*

saad@umn.edu

Yuanzhe Xi

*Department of Mathematics
Emory University*

yuanzhe.xi@emory.edu

Reviewed on OpenReview: <https://openreview.net/forum?id=vo8F0Bt6f6>

Abstract

1 Stochastic Gradient Descent (SGD) often slows in the late stage of training due to anisotropic
2 curvature and gradient noise. We analyze preconditioned SGD in the geometry induced
3 by a symmetric positive definite matrix \mathbf{M} . Our bounds make explicit how both the
4 convergence rate and the stochastic noise floor depend on \mathbf{M} . For nonconvex objectives, we
5 establish a basin-stability guarantee in a local \mathbf{M} -metric neighborhood around a minimizer
6 set: under local smoothness and a local PL condition, we give an explicit lower bound on the
7 probability that the iterates remain in the basin up to a time horizon. This perspective is
8 particularly relevant in Scientific Machine Learning (SciML), where reaching small training
9 losses under stochastic updates is closely tied to physical fidelity, numerical stability, and
10 constraint satisfaction. Our framework covers both diagonal/adaptive and curvature-aware
11 preconditioners and yields a practical criterion: choose \mathbf{M} to improve local conditioning
12 while attenuating noise in the \mathbf{M}^{-1} -norm. Experiments on a quadratic diagnostic and three
13 SciML benchmarks support the predicted rate-floor behavior.

1 Introduction

Stochastic Gradient Descent (SGD) has long been the workhorse of large-scale machine learning. Since its early application to multilayer perceptrons in the 1960s (Amari, 1967), its simplicity, scalability, and low per-iteration cost have made it a popular optimizer for deep learning models (Bottou et al., 2018). Classical convergence theory for SGD under noisy gradients typically guarantees a sublinear rate of $\mathcal{O}(1/k)$ under convexity and smoothness assumptions (Robbins & Monro, 1951; Blum, 1954). The theory for SGD convergence under various combinations of conditions is well studied and documented in Garrigos & Gower (2024); Khaled & Richtárik (2023), and Bach (2024).

Recent theoretical developments have established *linear convergence* for SGD under stronger conditions, such as strong convexity, smoothness, and bounded noise (Bottou et al., 2018). When the loss F is c -strongly convex, has L -Lipschitz gradients, and the learning rate α satisfies $\alpha \leq \mu/(LK_G)$, the iterates \mathbf{w}_k satisfy

$$\mathbb{E}[F(\mathbf{w}_k) - F_*] \leq (1 - \alpha c \mu)^{k-1} \left(F(\mathbf{w}_1) - F_* - \frac{\alpha LK}{2c\mu} \right) + \frac{\alpha LK}{2c\mu}, \quad (1)$$

where μ , K , and K_G are constants associated with the stochastic gradients (defined in Assumptions 9–11), and let \mathbf{w}^* denote the unique minimizer and $F_* := F(\mathbf{w}^*)$ the optimal value. Eq. (1) highlights two late-stage drivers: a linear contraction factor $1 - \alpha c \mu$ and a stochastic error floor

$$\frac{\alpha LK}{2c\mu} = \frac{\alpha}{2\mu} \kappa K,$$

where $\kappa := \frac{L}{c}$ is the (Euclidean) condition number associated with curvature. For any admissible α , the floor scales with κ and K , while the contraction depends on the product $\alpha c \mu$.

Many successful optimizers can be viewed as *preconditioned variants of SGD*. Adaptive methods such as Adagrad (Duchi et al., 2011), Adam (Kingma & Ba, 2017), and RMSProp (Hinton, 2014), structured second-order approaches including Shampoo (Gupta et al., 2018), natural gradient descent (Amari, 1998), K-FAC (Martens & Grosse, 2015; Ishikawa & Karakida, 2024), and Sophia (Liu et al., 2024), as well as quasi-Newton methods like L-BFGS (Liu & Nocedal, 1989; Chen et al., 2014), all apply a linear transformation to the gradient that reshapes both curvature and gradient noise. From this perspective, their empirical effectiveness indicates that late-stage optimization is influenced not only by the choice of learning rates, but also by how the preconditioning alters local conditioning and the geometry of stochastic noise. Despite their widespread use, however, there is still no unified theoretical framework that identifies which properties of a preconditioner determine the late-stage convergence rate and the attainable noise floor.

Motivated by this perspective, we study the preconditioned SGD update in the following form

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k), \quad (2)$$

where $\mathbf{M} \succ 0$ is a symmetric positive definite (SPD) matrix that defines the geometry in which both curvature and noise are measured, $g(\mathbf{w}_k, \boldsymbol{\xi}_k) = \nabla_w F_k(\mathbf{w})$ is the stochastic gradient, α_k is the learning rate, $\boldsymbol{\xi}_k$ is an i.i.d. sample drawn at iteration k . The standard (vanilla) SGD update is recovered when $\mathbf{M} = \mathbf{I}$. Our goal is not to propose a new optimizer, but to provide a principled framework to analyze and compare preconditioners in the late stage of training.

Main contributions We investigate how preconditioning influences the late-stage behavior of SGD within a well-behaved basin of the loss surface. By analyzing preconditioned SGD in the \mathbf{M} -induced geometry, we show how rescaling the gradient affects both the convergence rate and the attainable noise floor, and we derive criteria that clarify which properties of a preconditioner matter in the late stage of training.

- 1. Preconditioned SGD in the strongly convex baseline.** We extend the classical “linear rate + noise floor” theory for SGD to updates preconditioned by a fixed SPD matrix \mathbf{M} . The resulting bounds show that late-stage behavior is controlled by (i) an effective conditioning in the \mathbf{M} -geometry and (ii) the preconditioned gradient-noise level; the attainable error floor scales with their *product*. Since admissible constant stepsizes are limited by \mathbf{M} -smoothness, improved conditioning allows larger stepsizes and hence faster contraction. With diminishing stepsizes, we obtain an $\mathcal{O}(1/k)$ rate.

55 **2. Local nonconvex regime with basin stability.** Under a local \mathbf{M} -PL condition and local
 56 smoothness, we establish late-stage convergence guarantees inside a well-behaved basin around a
 57 minimizer set, again with an explicit rate–floor structure. In addition, we provide a basin-stability
 58 bound that lower-bounds the probability of remaining in the basin up to a horizon.

59 **3. Design criteria and empirical evidence.** Our theory yields a simple design principle: choose \mathbf{M}
 60 to improve local conditioning while attenuating noise in the \mathbf{M}^{-1} -norm; the attainable late-stage
 61 floor tracks their product. We validate this mechanism on (i) a quadratic diagnostic where the
 62 relevant constants can be computed in closed form, and (ii) three SciML benchmarks where late-stage
 63 behavior is strongly tied to final accuracy.

64 While late-stage convergence is broadly relevant, it is especially important in SciML. Here, training losses
 65 encode physically meaningful quantities (e.g., PDE residuals, boundary conditions, stability). Unlike standard
 66 ML tasks where moderate error may still be acceptable, small reductions in the final loss can determine whether
 67 solutions conserve invariants, remain stable over long horizons, or meet scientific accuracy requirements. In
 68 this setting, the optimizer’s asymptotic behavior—and particularly the final noise floor—directly governs
 69 physical fidelity (Rathore et al., 2024).

70 2 Related work

71 Recent work has advanced the theoretical understanding of preconditioned and adaptive variants of SGD
 72 under various structural and noise assumptions. Koren et al. (2022) showed that preconditioned SGD
 73 achieves a rate of $\mathcal{O}(1/\sqrt{k})$ for general stochastic convex optimization, though convergence can stagnate in
 74 the presence of persistent gradient noise. Faw et al. (2022) further established that adaptive SGD attains
 75 an order-optimal $\tilde{\mathcal{O}}(1/\sqrt{k})$ rate for nonconvex smooth objectives under affine variance conditions, without
 76 requiring bounded gradients or finely tuned learning rates. More recently, Attia & Koren (2023) derived
 77 high-probability guarantees of $\tilde{\mathcal{O}}(1/k + \sigma_0/\sqrt{k})$ for adaptive methods in both convex and nonconvex settings,
 78 relaxing the need for strong smoothness or prior parameter knowledge.

79 These results primarily address *global* convergence behavior across general problem classes. In contrast, our
 80 analysis focuses on the *asymptotic regime*—the late stage of training where iterates lie within a well-behaved
 81 basin around a local minimizer and optimization progress is limited by curvature anisotropy and gradient
 82 noise. In this regime, we show that both the convergence rate and the noise floor of the preconditioned SGD
 83 are determined by curvature and variance quantities measured in the preconditioned geometry. This local,
 84 geometry-aware viewpoint clarifies why curvature-informed preconditioners and adaptive algorithms yield
 85 faster and more stable late-stage convergence.

86 Other techniques such as batch normalization (Lange et al., 2022) and weight decay (Loshchilov & Hutter,
 87 2017; Barrett & Dherin, 2020) can also be interpreted as implicit forms of preconditioning, though they operate
 88 through different regularization mechanisms. For comprehensive surveys of explicit preconditioned SGD
 89 and related adaptive methods, we refer the reader to Ye (2024). Beyond convergence rates, preconditioning
 90 has also been studied as an implicit regularization that may affect generalization (Amari et al., 2021). Our
 91 paper, however, focuses on optimization and training loss rather than test error or generalization. This
 92 emphasis is deliberate in many SciML problems, where the training objective often directly measures physically
 93 meaningful quantities such as PDE residuals and boundary-condition violations, so driving the training loss
 94 low is itself important. At the same time, preconditioning changes the optimization trajectory and therefore
 95 the algorithm’s implicit bias, so it may also affect generalization. Understanding how the metric-dependent
 96 quantities in our analysis interact with out-of-sample accuracy is an important direction for future work.

97 3 Preconditioned SGD convergence analysis

98 We first analyze the globally strongly convex case as a *baseline* to make the role of the preconditioned geometry
 99 explicit. Although this setting is rarely realized in deep learning, it reveals the essential mechanism through
 100 which preconditioning affects convergence. The analysis shows how curvature and noise floor transform

under a change of metric, providing a principled way to compare different choices of \mathbf{M} . This also lays the groundwork for the local nonconvex analysis in Section 3.2, where \mathbf{M} influences both basin size and stability.

3.1 Convergence in the globally strongly convex setting

We establish convergence guarantees for preconditioned SGD when the objective is globally strongly convex. This simplified setting allows for a transparent analysis of how a preconditioner reshapes both the effective curvature and the gradient noise. While the derivations parallel the Euclidean case, expressing them in the \mathbf{M} -induced geometry makes the dependence on the preconditioner explicit and lays the groundwork for the more general nonconvex results to follow.

Curvature assumptions. Preconditioning redefines smoothness and strong convexity through effective constants (\hat{L}, \hat{c}) measured in the \mathbf{M} -induced norm.

Assumption 1 (\mathbf{M} -strong convexity). $F: \mathbb{R}^d \rightarrow \mathbb{R}$ is \mathbf{M} -strongly convex: there exists $\hat{c} > 0$ such that

$$F(\bar{\mathbf{w}}) \geq F(\mathbf{w}) + \nabla F(\mathbf{w})^\top (\bar{\mathbf{w}} - \mathbf{w}) + \frac{1}{2} \hat{c} \|\bar{\mathbf{w}} - \mathbf{w}\|_{\mathbf{M}}^2, \quad \forall \bar{\mathbf{w}}, \mathbf{w} \in \mathbb{R}^d.$$

Assumption 2 (\mathbf{M} -Lipschitz gradient). ∇F is \mathbf{M} -Lipschitz with constant $\hat{L} > 0$:

$$\|\nabla F(\bar{\mathbf{w}}) - \nabla F(\mathbf{w})\|_{\mathbf{M}^{-1}} \leq \hat{L} \|\bar{\mathbf{w}} - \mathbf{w}\|_{\mathbf{M}}, \quad \forall \bar{\mathbf{w}}, \mathbf{w} \in \mathbb{R}^d.$$

These conditions are direct analogues of the Euclidean definitions. Writing $\mathbf{M}^{-1} = \mathbf{P}\mathbf{P}^\top$ gives the spectral characterization:

Lemma 3.1. Let F be twice differentiable and $\mathbf{M}^{-1} = \mathbf{P}\mathbf{P}^\top$. Then: (i) ∇F is \mathbf{M} -Lipschitz with constant $\hat{L} \iff$ all eigenvalues of $\mathbf{P}^\top \nabla^2 F(\mathbf{w}) \mathbf{P}$ are $\leq \hat{L}$; (ii) F is \mathbf{M} -strongly convex with constant $\hat{c} \iff$ all eigenvalues of $\mathbf{P}^\top \nabla^2 F(\mathbf{w}) \mathbf{P}$ are $\geq \hat{c}$.

Hence, preconditioning improves the effective condition number whenever $\hat{L}/\hat{c} < L/c$.

Noise assumptions. We measure the first and second moments of the stochastic gradient in the \mathbf{M}^{-1} -norm. Specifically, holding \mathbf{w}_k fixed, we define the variance with respect to the sampling of $\boldsymbol{\xi}_k$ by

$$\mathbb{V}_{\boldsymbol{\xi}_k}[g(\mathbf{w}_k, \boldsymbol{\xi}_k), \|\cdot\|_{\mathbf{M}^{-1}}] := \mathbb{E}_{\boldsymbol{\xi}_k}[\|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2] - \|\mathbb{E}_{\boldsymbol{\xi}_k}[g(\mathbf{w}_k, \boldsymbol{\xi}_k)]\|_{\mathbf{M}^{-1}}^2. \quad (3)$$

Assumption 3 (Moment bounds in \mathbf{M}^{-1}). For the iterates of (2), there exist constants $\mu_G \geq \mu > 0$, $K \geq 0$, and $K_V \geq 0$ such that, for all k ,

$$\langle \nabla F(\mathbf{w}_k), \mathbb{E}_{\boldsymbol{\xi}_k}[g(\mathbf{w}_k, \boldsymbol{\xi}_k)] \rangle_{\mathbf{M}^{-1}} \geq \mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2, \quad (4)$$

$$\|\mathbb{E}_{\boldsymbol{\xi}_k}[g(\mathbf{w}_k, \boldsymbol{\xi}_k)]\|_{\mathbf{M}^{-1}} \leq \mu_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}, \quad (5)$$

$$\mathbb{V}_{\boldsymbol{\xi}_k}[g(\mathbf{w}_k, \boldsymbol{\xi}_k), \|\cdot\|_{\mathbf{M}^{-1}}] \leq K + K_V \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2. \quad (6)$$

We call K the *preconditioned noise level* because the variance in the \mathbf{M}^{-1} -norm satisfies

$$\mathbb{V}_{\boldsymbol{\xi}}[g(\mathbf{w}, \boldsymbol{\xi}), \|\cdot\|_{\mathbf{M}^{-1}}] = \text{tr}(\mathbf{M}^{-1} \boldsymbol{\Sigma}(\mathbf{w})),$$

where $\boldsymbol{\Sigma}(\mathbf{w}) := \text{Cov}(g(\mathbf{w}, \boldsymbol{\xi}) \mid \mathbf{w})$. In the stationary case $\boldsymbol{\Sigma}(\mathbf{w}) \equiv \boldsymbol{\Sigma}$, we have the fixed $\text{tr}(\mathbf{M}^{-1} \boldsymbol{\Sigma})$. More generally, on a region containing the iterates it is natural to choose $K \geq \sup_{\mathbf{w}} \text{tr}(\mathbf{M}^{-1} \boldsymbol{\Sigma}(\mathbf{w}))$, so K is a uniform baseline for the preconditioned noise.

Under these assumptions we obtain the usual linear and sublinear rates, but with constants that depend explicitly on the preconditioned geometry.

Theorem 3.2. Under Assumptions 1–3 (with $F_{\min} = F_*$), suppose (2) uses a fixed learning rate $\alpha_k = \bar{\alpha}$ with

$$0 < \bar{\alpha} \leq \frac{\mu}{\hat{L} K_G} \quad \text{where } K_G = K_V + \mu_G^2 \geq \mu^2 > 0.$$

Then, for all $k \in \mathbb{N}$,

$$\mathbb{E}[F(\mathbf{w}_k) - F_*] \leq \frac{\bar{\alpha} \hat{L} K}{2 \hat{c} \mu} + (1 - \bar{\alpha} \hat{c} \mu)^{k-1} \left(F(\mathbf{w}_1) - F_* - \frac{\bar{\alpha} \hat{L} K}{2 \hat{c} \mu} \right) \xrightarrow{k \rightarrow \infty} \frac{\bar{\alpha} \hat{L} K}{2 \hat{c} \mu}. \quad (7)$$

Theorem 3.2 shows that, with a fixed learning rate $\bar{\alpha}$, preconditioned SGD contracts linearly with factor $1 - \bar{\alpha} \hat{c} \mu$ and converges to an asymptotic floor

$$\frac{\bar{\alpha} \hat{L} K}{2 \hat{c} \mu} = \frac{\bar{\alpha}}{2 \mu} \left(\frac{\hat{L}}{\hat{c}} \right) K.$$

Thus, the floor factorizes into an *effective condition number* \hat{L}/\hat{c} and a *preconditioned noise level* K . In the late stage of training, we have $F(\mathbf{w}_k) - F_* = \mathcal{O}(\bar{\alpha} K)$ and $\|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 = \mathcal{O}(\bar{\alpha} K)$. Substituting into the variance bound (6) gives

$$\mathbb{V}_{\xi_k}[g(\mathbf{w}_k, \xi_k), \|\cdot\|_{\mathbf{M}^{-1}}] \leq K + \mathcal{O}(\bar{\alpha} K),$$

so for small $\bar{\alpha}$ the variance is dominated by the baseline K term.

Moreover, since $\mathbb{V}_{\xi}[g(\mathbf{w}, \xi), \|\cdot\|_{\mathbf{M}^{-1}}] = \text{tr}(\mathbf{M}^{-1} \Sigma(\mathbf{w}))$, we may view K as an upper baseline for the preconditioned noise $\text{tr}(\mathbf{M}^{-1} \Sigma(\mathbf{w}))$ along the late-stage trajectory. Preconditioning reduces this baseline through its effect on $\text{tr}(\mathbf{M}^{-1} \Sigma(\mathbf{w}))$; choosing \mathbf{M} to attenuate high-variance directions lowers this trace and thus lowers the effective noise floor.

Theorem 3.3. *Under Assumptions 1–3 (with $F_{\min} = F_*$), suppose (2) uses $\alpha_k = \beta/(\gamma + k)$ with $\beta > \frac{1}{\hat{c} \mu}$ and $\gamma > 0$ chosen so that $\alpha_1 \leq \mu/(\hat{L} K_G)$. Then, for all $k \in \mathbb{N}$,*

$$\mathbb{E}[F(\mathbf{w}_k) - F_*] \leq \frac{\nu}{\gamma + k}, \quad \nu := \max \left\{ \frac{\beta^2 \hat{L} K}{2(\beta \hat{c} \mu - 1)}, (\gamma + 1)(F(\mathbf{w}_1) - F_*) \right\}. \quad (8)$$

With diminishing learning rates, the noise floor vanishes and Theorem 3.3 shows that preconditioned SGD attains the optimal $\mathcal{O}(1/k)$ rate. Preconditioning no longer changes the rate itself—it always decays like $1/k$ —but it directly influences the leading constant ν which has the same structure as the fixed-learning-rate floor: an effective condition number \hat{L}/\hat{c} multiplied by the preconditioned noise level K . Thus even when the noise floor disappears, late-stage performance is still governed by the same metric-dependent quantities (\hat{L}, \hat{c}, K) . Consequently, effective preconditioners must again balance curvature alignment (to reduce \hat{L}/\hat{c}) with noise attenuation (to reduce K), improving both the asymptotic constants in the $\mathcal{O}(1/k)$ regime.

3.2 Local convergence in the nonconvex setting

The empirical loss $F(\mathbf{w})$ over network parameters is typically *nonconvex*, and its local geometry near minimizers is rarely strictly convex. Empirical studies show that trained models often converge to regions that are flat in many directions and exhibit highly degenerate curvature—manifested as a cluster of very small or near-zero eigenvalues in the Hessian—arising from overparameterization, symmetries, and parameter non-identifiability (Sagun et al., 2018; Ghorbani et al., 2019). Despite this degeneracy, the optimization dynamics remain structured: iterates contract along directions with significant curvature while the loss changes little along flat directions. To describe this late-stage regime without assuming strong convexity, we impose a *local Polyak–Lojasiewicz (PL)* condition (Chan, 1979; Karimi et al., 2016) in the \mathbf{M} -geometry, which enforces gradient domination only in informative directions and tolerates flat or weakly curved subspaces. This flat-tolerant formulation provides a natural framework to study how preconditioning reshapes local curvature and noise, governing contraction rates, asymptotic error floors, and stability during the final phase of optimization.

Additional local assumptions. Fix an SPD matrix \mathbf{M} and an open neighborhood $\mathcal{U} \subset \mathbb{R}^d$. Assume the local minimizer set

$$\mathcal{S} := \arg \min_{\mathbf{w} \in \mathcal{U}} F(\mathbf{w}) \neq \emptyset, \quad F_* := \min_{\mathbf{w} \in \mathcal{U}} F(\mathbf{w}) = F(\mathbf{s}) \text{ for any } \mathbf{s} \in \mathcal{S}.$$

162 Write $\|x\|_{\mathbf{M}} := (x^\top \mathbf{M}x)^{1/2}$ and $\text{dist}_{\mathbf{M}}(\mathbf{w}, \mathcal{S}) := \inf_{\mathbf{s} \in \mathcal{S}} \|\mathbf{w} - \mathbf{s}\|_{\mathbf{M}}$. For radii $0 < r < r_+$, define the \mathbf{M} -metric
163 neighborhoods

$$\mathcal{N}_r := \{\mathbf{w} : \text{dist}_{\mathbf{M}}(\mathbf{w}, \mathcal{S}) \leq r\}, \quad \mathcal{N}_{r_+} := \{\mathbf{w} : \text{dist}_{\mathbf{M}}(\mathbf{w}, \mathcal{S}) \leq r_+\} \subseteq \mathcal{U}.$$

164 We assume the following conditions hold on \mathcal{N}_r (for the iterates) and on \mathcal{N}_{r_+} (for the exit bound).

165 **Assumption 4** (Local \mathbf{M} -PL on \mathcal{N}_r). *There exists $\hat{\mu}_{\text{PL}} > 0$ such that, for all $\mathbf{w} \in \mathcal{N}_r$,*

$$2\hat{\mu}_{\text{PL}}(F(\mathbf{w}) - F_*) \leq \|\nabla F(\mathbf{w})\|_{\mathbf{M}^{-1}}^2.$$

166 **Assumption 5** (Local \mathbf{M} -Lipschitz gradient on a convex neighborhood of \mathcal{N}_{r_+}). *There exists an open convex
167 set \mathcal{V} with $\mathcal{N}_{r_+} \subset \mathcal{V} \subseteq \mathcal{U}$ and a constant $\hat{L} > 0$ such that, for all $\bar{\mathbf{w}}, \mathbf{w} \in \mathcal{V}$,*

$$\|\nabla F(\bar{\mathbf{w}}) - \nabla F(\mathbf{w})\|_{\mathbf{M}^{-1}} \leq \hat{L} \|\bar{\mathbf{w}} - \mathbf{w}\|_{\mathbf{M}}.$$

168 **Assumption 6** (Local stochastic gradient conditions on \mathcal{N}_r). *Let (\mathcal{F}_k) denote the natural filtration and set
169 $g_k := g(\mathbf{w}_k, \boldsymbol{\xi}_k)$. There exist constants $\mu \in (0, 1]$, $K_G \geq 0$, and $K \geq 0$ such that, for every k with $\mathbf{w}_k \in \mathcal{N}_r$,*

$$\langle \nabla F(\mathbf{w}_k), \mathbb{E}[g_k | \mathcal{F}_k] \rangle_{\mathbf{M}^{-1}} \geq \mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2, \quad \mathbb{E}[\|g_k\|_{\mathbf{M}^{-1}}^2 | \mathcal{F}_k] \leq K_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + K.$$

170 **Assumption 7** (Local quadratic growth (QG) on \mathcal{N}_{r_+}). *There exists $\alpha_{\text{QG}} > 0$ such that, for all $\mathbf{w} \in \mathcal{N}_{r_+}$,*

$$F(\mathbf{w}) - F_* \geq \frac{\alpha_{\text{QG}}}{2} \text{dist}_{\mathbf{M}}(\mathbf{w}, \mathcal{S})^2.$$

171 **Assumption 8** (Controlled one-step overshoot on \mathcal{N}_r). *Fix radii $0 < r < r_+$ and a horizon $T \geq 1$, and set
172 $\Delta := r_+ - r$. There exist deterministic numbers $(\delta_k)_{k=1}^{T-1}$ with $\delta_k \in [0, 1)$ such that for every $k \leq T - 1$,*

$$\mathbf{1}_{\{\mathbf{w}_k \in \mathcal{N}_r\}} \alpha_k^2 \mathbb{E}[\|g_k\|_{\mathbf{M}^{-1}}^2 | \mathcal{F}_k] \leq \delta_k \Delta^2 \quad \text{a.s.}$$

173 Lemma 3.4 gives the one-step containment probability implied by Assumption 8.

174 **Lemma 3.4** (Containment probability implied by Assumption 8). *Under Assumption 8, for every $k \leq T - 1$,*

$$\mathbf{w}_k \in \mathcal{N}_r \implies \mathbb{P}(\mathbf{w}_{k+1} \in \mathcal{N}_{r_+} | \mathcal{F}_k) \geq 1 - \delta_k.$$

175 These local assumptions are the basin-restricted analogue of the global conditions in Section 3.1. The local
176 \mathbf{M} -PL condition replaces global strong convexity by a *gradient-domination* inequality in the \mathbf{M} -metric: it
177 enforces curvature only in directions that drive descent while permitting flat or weakly curved directions.
178 The local \mathbf{M} -Lipschitz gradient assumption on a convex neighborhood $\mathcal{V} \supset \mathcal{N}_{r_+}$ provides a quadratic upper
179 model along any update segment that stays in \mathcal{V} :

$$F(\bar{\mathbf{w}}) \leq F(\mathbf{w}) + \nabla F(\mathbf{w})^\top (\bar{\mathbf{w}} - \mathbf{w}) + \frac{\hat{L}}{2} \|\bar{\mathbf{w}} - \mathbf{w}\|_{\mathbf{M}}^2.$$

180 In our finite-horizon analysis, this condition is invoked only on trajectories for which the iterates (and hence
181 the corresponding update segments, by convexity) remain inside \mathcal{V} up to time T .

182 The local stochastic gradient condition (Assumption 6) mirrors the global moment bounds in Assumption 3,
183 but is only required to hold when $\mathbf{w}_k \in \mathcal{N}_r$. It imposes a first-moment alignment condition and a *second-*
184 *moment* bound in the \mathbf{M}^{-1} -norm, which is the natural scale for preconditioned updates. The local QG
185 condition ensures that the objective grows at least quadratically with $\text{dist}_{\mathbf{M}}(\mathbf{w}, \mathcal{S})$ near the basin boundary—a
186 property that holds, for example, when curvature is positive in normal directions—and it supplies the barrier
187 needed in the optional-stopping/exit-time argument.

188 Assumption 8 controls rare one-step overshoots from the inner basin \mathcal{N}_r to outside the enlarged neighborhood
189 \mathcal{N}_{r_+} . When $\mathbf{w}_k \in \mathcal{N}_r$, the preconditioned update moves a distance

$$\|\mathbf{w}_{k+1} - \mathbf{w}_k\|_{\mathbf{M}} = \alpha_k \|g_k\|_{\mathbf{M}^{-1}}.$$

190 Since $\text{dist}_{\mathbf{M}}(\mathbf{w}_k, \mathcal{S}) \leq r$ on \mathcal{N}_r , the triangle inequality implies that $\mathbf{w}_{k+1} \notin \mathcal{N}_{r_+}$ can occur only if $\alpha_k \|g_k\|_{\mathbf{M}^{-1}} >$
 191 Δ with $\Delta := r_+ - r$. Assumption 8 bounds the conditional second moment of $\|g_k\|_{\mathbf{M}^{-1}}$ relative to Δ ; therefore,
 192 by Markov's inequality,

$$\mathbb{P}(\mathbf{w}_{k+1} \notin \mathcal{N}_{r_+} \mid \mathcal{F}_k) \leq \delta_k \quad \text{whenever } \mathbf{w}_k \in \mathcal{N}_r.$$

193 Together, these assumptions describe a local regime that accommodates moderate nonconvexity and flatness
 194 while still providing sufficient structure for quantitative finite-horizon convergence and stability guarantees
 195 under stochastic gradients.

196 **Theorem 3.5** (Convergence in a local basin up to a finite horizon). *Fix radii $0 < r < r_+$ and a horizon*
 197 *$T \geq 1$, and let*

$$\tau := \inf\{k \geq 1 : \mathbf{w}_k \notin \mathcal{N}_r\}, \quad \Omega_T := \{\tau > T\}.$$

198 *Assume: (i) Assumptions 4 and 6 hold on \mathcal{N}_r ; (ii) Assumption 5 holds on a convex set \mathcal{V} with $\mathcal{N}_{r_+} \subset \mathcal{V} \subseteq \mathcal{U}$;*
 199 *(iii) Assumption 7 holds on \mathcal{N}_{r_+} ; (iv) Assumption 8 holds with horizon T and failure probabilities $(\delta_k)_{k=1}^{T-1}$;*
 200 *and (v) the conditional-moment version of Assumption 6 holds on Ω_T (i.e., the first/second-moment bounds*
 201 *are valid when conditioning on $(\mathcal{F}_k, \Omega_T)$ for $k \leq T - 1$).*

202 *Suppose $\mathbf{w}_1 \in \mathcal{N}_r$ and use a constant stepsize $\alpha_k = \bar{\alpha}$ such that*

$$0 < \bar{\alpha} \leq \frac{\mu}{\hat{L}K_G} \quad (\text{if } K_G > 0), \quad \text{and} \quad 0 < \bar{\alpha} < \frac{1}{\mu\hat{\mu}_{\text{PL}}}.$$

203 *Define*

$$\rho := \bar{\alpha} \hat{\mu}_{\text{PL}} \mu \in (0, 1), \quad C := \frac{\bar{\alpha} \hat{L} K}{2 \hat{\mu}_{\text{PL}} \mu}, \quad B := \frac{\alpha_{\text{QG}}}{2} r^2.$$

204 *For all $1 \leq k \leq T$,*

$$\mathbb{E}[F(\mathbf{w}_k) - F_* \mid \tau > T] \leq C + (1 - \rho)^{k-1} (F(\mathbf{w}_1) - F_* - C).$$

205 *The probability of remaining in \mathcal{N}_r up to time T satisfies*

$$\mathbb{P}(\tau > T) \geq \left[1 - \frac{F(\mathbf{w}_1) - F_* + \frac{\hat{L}}{2} \bar{\alpha}^2 K (T - 1)}{B} - \sum_{k=1}^{T-1} \delta_k \right]_+,$$

206 *where $[x]_+ := \max\{0, x\}$.*

207 **Theorem 3.6** (Diminishing learning rate in a local basin up to a finite horizon). *Fix radii $0 < r < r_+$ and a*
 208 *horizon $T \geq 1$, and let*

$$\tau := \inf\{k \geq 1 : \mathbf{w}_k \notin \mathcal{N}_r\}, \quad \Omega_T := \{\tau > T\}.$$

209 *Assume: (i) Assumptions 4 and 6 hold on \mathcal{N}_r ; (ii) Assumption 5 holds on a convex set \mathcal{V} with $\mathcal{N}_{r_+} \subset \mathcal{V} \subseteq \mathcal{U}$;*
 210 *(iii) Assumption 7 holds on \mathcal{N}_{r_+} ; (iv) Assumption 8 holds with horizon T and failure probabilities $(\delta_k)_{k=1}^{T-1}$;*
 211 *and (v) the conditional-moment version of Assumption 6 holds on Ω_T .*

212 *Suppose $\mathbf{w}_1 \in \mathcal{N}_r$ and use harmonic stepsizes*

$$\alpha_k = \frac{\beta}{\gamma + k}, \quad \gamma > 0,$$

213 *with*

$$0 < \alpha_1 = \frac{\beta}{\gamma + 1} \leq \frac{\mu}{\hat{L}K_G} \quad (\text{if } K_G > 0), \quad \text{and} \quad \beta > \frac{1}{\mu\hat{\mu}_{\text{PL}}} \quad (\text{equivalently } a := \beta\mu\hat{\mu}_{\text{PL}} > 1).$$

214 *Define*

$$m := \mu\hat{\mu}_{\text{PL}}, \quad c := \frac{\hat{L}K}{2}, \quad B := \frac{\alpha_{\text{QG}}}{2} r^2, \quad \nu := \max \left\{ \frac{c\beta^2}{\beta m - 1}, (\gamma + 1)[F(\mathbf{w}_1) - F_*] \right\}.$$

215 For all $1 \leq k \leq T$,

$$\mathbb{E}[F(\mathbf{w}_k) - F_* | \tau > T] \leq \frac{\nu}{\gamma + k}.$$

216 The probability of remaining in \mathcal{N}_r up to time T satisfies

$$\mathbb{P}(\tau > T) \geq \left[1 - \frac{F(\mathbf{w}_1) - F_* + c \sum_{k=1}^{T-1} \alpha_k^2}{B} - \sum_{k=1}^{T-1} \delta_k \right]_+,$$

217 where $[x]_+ := \max\{0, x\}$.

218 Theorem 3.5 (fixed stepsize) and Theorem 3.6 (harmonic stepsizes) characterize late-stage optimization *after*
 219 the iterates have entered a well-behaved local basin \mathcal{N}_r . Both results are stated on the finite-horizon survival
 220 event

$$\Omega_T := \{\tau > T\}, \quad \tau := \inf\{k \geq 1 : \mathbf{w}_k \notin \mathcal{N}_r\},$$

221 so that along Ω_T the local \mathbf{M} -smoothness and local \mathbf{M} -PL inequalities apply to the entire trajectory up
 222 to time T and yield explicit descent recursions. With a constant stepsize $\bar{\alpha}$, Theorem 3.5 gives conditional
 223 geometric contraction to the noise floor $C = \frac{\bar{\alpha} \hat{L} K}{2 \hat{\mu}_{\text{PL}} \mu}$, whereas with harmonic stepsizes $\alpha_k = \beta/(\gamma + k)$,
 224 Theorem 3.6 yields the conditional $\mathcal{O}(1/k)$ rate. In both cases, the constants are *local* and expressed in
 225 the \mathbf{M} -geometry. Unlike global strongly convex analyses, no global curvature or global variance control is
 226 required; the bounds depend only on the basin actually explored by the iterates.

227 The basin-stability guarantees are also local, and they make two distinct failure mechanisms explicit. The
 228 first is an objective barrier controlled by the local QG constant α_{QG} and the basin radius r through

$$B := \frac{\alpha_{\text{QG}}}{2} r^2,$$

229 which quantifies the minimum objective increase needed to reach the boundary $\mathcal{N}_{r_+} \setminus \mathcal{N}_r$. The second is one-
 230 step overshoot: Assumption 8 allows rare updates that jump from \mathcal{N}_r to outside the enlarged neighborhood
 231 \mathcal{N}_{r_+} , with conditional failure probabilities δ_k . Here, $\sum_{k=1}^{T-1} \delta_k$ quantifies the accumulated overshoot risk: if
 232 the tails/second moments are large, or if the basin margin $\Delta = r_+ - r$ is small, then δ_k may be large, and
 233 the stability bound becomes conservative.

234 Because all constants in the local bounds are \mathbf{M} -dependent, a well-chosen preconditioner \mathbf{M} can improve late-
 235 stage behavior by: (i) enhancing local conditioning (increasing $\hat{\mu}_{\text{PL}}$ and/or decreasing \hat{L} , thereby strengthening
 236 contraction); (ii) reducing the preconditioned noise level K ; and (iii) improving stability by reducing the
 237 overshoot probabilities δ_k (e.g., via smaller $\mathbb{E}[\|g_k\|_{\mathbf{M}^{-1}}^2]$ and/or a larger margin $\Delta = r_+ - r$) and, when
 238 aligned with normal-space curvature, by increasing the barrier parameter $B = \frac{\alpha_{\text{QG}}}{2} r^2$.

239 3.3 Practical preconditioners for SGD

240 A wide range of preconditioning strategies are used in modern machine learning. On the first-order side,
 241 adaptive methods such as Adam (Kingma & Ba, 2017), AMSGrad (Reddi et al., 2018), PAdam (Chen
 242 et al., 2020), and Yogi (Zaheer et al., 2018) implicitly apply *diagonal* preconditioners by rescaling gradients
 243 with running estimates of coordinatewise second moments. On the second-order side, *curvature-aware*
 244 preconditioners exploit Hessian or Fisher Information Matrix (FIM) structure, including the empirical
 245 FIM (Schraudolph, 2002), full or mini-batch Hessians (Fletcher, 2013; Garg et al., 2024), mini-batch
 246 quasi-Newton updates (Griffin et al., 2022), and Kronecker-factored FIM (K-FAC) (Martens & Grosse,
 247 2015). Classical schemes such as L-BFGS (Liu & Nocedal, 1989; Chen et al., 2014) can also be viewed as
 248 low-rank, history-based preconditioners. A particularly important special case is *natural gradient descent*,
 249 obtained by choosing \mathbf{M} as the Fisher information matrix (Amari, 1998). In that case, the update follows
 250 the local information geometry of the model rather than the Euclidean geometry of parameter space. For
 251 exponential-family and least-squares settings, the Fisher matrix is closely related to, and in many cases
 252 coincides with, the generalized Gauss-Newton matrix (Schraudolph, 2002; Martens, 2020). Our present theory
 253 treats a fixed SPD metric \mathbf{M} , so it does not directly analyze the fully time-varying choice $\mathbf{M}_k = \mathbf{F}(\mathbf{w}_k)$.

254 Nevertheless, it provides a local lens on natural-gradient-type methods. Appendix B.2 summarizes these
 255 approaches and their computational trade-offs.

256 The convergence analysis in Sections 3.1–3.2 suggests two practical mechanisms through which preconditioners
 257 shape late-stage behavior:

- 258 • *Local conditioning.* Curvature-aware preconditioners (e.g., Fisher, Gauss–Newton, Hessian, K-FAC)
 259 tend to reduce the metric–smoothness constant \hat{L} and can increase the local PL constant $\hat{\mu}_{\text{PL}}$. In our
 260 bounds, this improves the effective local condition number $\hat{L}/\hat{\mu}_{\text{PL}}$, permits larger admissible fixed
 261 learning rates $\alpha \leq \mu/(\hat{L}K_G)$, and reduces the leading constant under diminishing learning rates.
- 262 • *Noise attenuation.* Preconditioners aligned with the gradient-noise structure reduce the preconditioned
 263 noise level K in the late-stage regime. Together with improved conditioning (smaller \hat{L}/\hat{c} or $\hat{L}/\hat{\mu}_{\text{PL}}$),
 264 this lowers the noise floor, which scales with their *product*. Fisher-based and related methods are
 265 especially effective because they explicitly incorporate gradient statistics.

266 These two mechanisms—improved conditioning and reduced preconditioned noise—match the behavior
 267 observed in Section 4. Curvature-matched preconditioners (Fisher, Gauss–Newton, K-FAC, Hessian) typically
 268 yield faster late-stage contraction by reducing \hat{L} and, in some cases, increasing $\hat{\mu}_{\text{PL}}$, while their use of gradient
 269 second-moment information tends to reduce K . Adaptive/diagonal methods likewise lower K by damping
 270 high-variance coordinates, though their alignment with curvature is typically weaker. Recent theory further
 271 suggests that in anisotropic settings, Kronecker-structured preconditioning can be statistically necessary for
 272 efficient feature learning, whereas entry-wise/diagonal scaling offers only partial improvements (Zhang et al.,
 273 2025).

274 4 Numerical results

275 Many machine-learning benchmarks illustrate the benefits of preconditioned SGD (e.g., Schmidt et al. (2021);
 276 Schneider et al. (2019)), but our emphasis is on SciML, where driving the loss to very small values is
 277 tightly linked to physical fidelity, numerical stability, and constraint satisfaction. We therefore structure the
 278 experiments in two parts.

279 First, we analyze a *diagnostic quadratic model* in which all the quantities in our theory— \hat{L} , $\hat{\mu}_{\text{PL}}$, and the
 280 preconditioned noise level K —admit closed forms. This allows us to directly compute the geometry– and
 281 noise–dependent metrics from Sections 3.1–3.2 and verify their influence on rate and floor.

282 Second, we examine three representative SciML problems: noisy Franke surface regression (Franke, 1979), a
 283 Poisson–type PINN, and Green’s–function learning for diffusion and convection–diffusion (Zhang et al., 2024;
 284 Rathore et al., 2024; Hao et al., 2024; Xu et al., 2025), to see how the theoretical mechanisms are reflected in
 285 practical settings.

286 4.1 Diagnostic quadratic model

287 To isolate the effects predicted by the theory, we consider the quadratic objective

$$F(\mathbf{w}) = \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*) + F_*, \quad \mathbf{H} \succeq 0,$$

288 here \mathbf{H} specifies curvature. We test two simple, analytically tractable choices: Euclidean SGD ($\mathbf{M} = \mathbf{I}$)
 289 and a low-rank curvature-aware preconditioner $\mathbf{M} = \mathbf{I} + \mathbf{U}_s(\tilde{\mathbf{\Lambda}}_s - \mathbf{I})\mathbf{U}_s^\top$, where \mathbf{U}_s contains the top (or
 290 bottom) s eigenvectors of \mathbf{H} and $\tilde{\mathbf{\Lambda}}_s$ is a diagonal matrix. This model captures the essential effect of curvature
 291 information. We used a fixed learning rate.

292 Instead of forming a dataset, we synthesize unbiased mini–batch gradients

$$g_k = \nabla F(\mathbf{w}_k) + \zeta_k, \quad \mathbb{E}[\zeta_k] = 0, \quad \text{Cov}(\zeta_k) = \frac{1}{B}\mathbf{\Sigma}.$$

293 To match the second-order statistics of least-squares problems near \mathbf{w}^* , we set $\Sigma = \sigma^2 \mathbf{H}$, giving $K =$
 294 $\frac{\sigma^2}{B} \text{tr}(\mathbf{M}^{-1} \mathbf{H})$. We choose $d = 100$ and construct $\mathbf{H} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ with $\mathbf{\Lambda}$ log-uniform grid on $[10^{-2}, 10^2]$ and
 295 \mathbf{U} Haar-distributed. We set $\mathbf{w}^* = 0$, $F_* = 0$, and initialize $\mathbf{w}_1 \sim \mathcal{N}(0, 10^{-4} \mathbf{I})$, and report averages over 30
 296 runs. To illustrate how individual eigenvalues affect constants $(\hat{L}, \hat{\mu}_{\text{PL}}, K)$, we design three groups of tests
 297 targeting different part of the spectrum of \mathbf{H} .

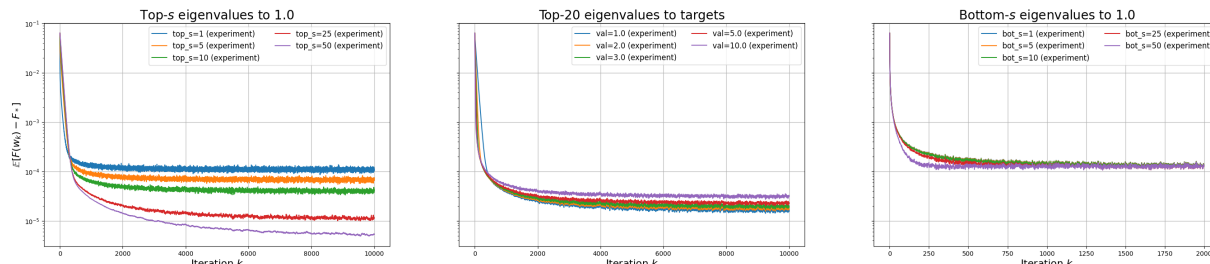


Figure 1: Convergence behavior under different deflation-based preconditioners. Left: deflating the largest s eigenvalues ($s \in \{1, 5, 10, 25, 50\}$). Middle: deflating the top 20 eigenvalues to target values $1.0, 2.0, 3.0, 5.0, 10.0$. Right: deflating the smallest s eigenvalues ($s \in \{1, 5, 10, 25, 50\}$).

298 Figure 1 shows how deflating different parts of the spectrum of \mathbf{H} affects the key theoretical constants.
 299 Denote the eigenpairs of \mathbf{H} as $(\lambda_i, \mathbf{u}_i)$, and let \mathbf{U}_s contain the selected eigenvectors. We construct a spectral
 300 preconditioner of the form $\mathbf{M} = \mathbf{I} + \mathbf{U}_s (\tilde{\mathbf{\Lambda}}_s - \mathbf{I}) \mathbf{U}_s^\top$, where $\tilde{\mathbf{\Lambda}}_s = \text{diag}(\tau_1, \dots, \tau_s)$ assigns a target value τ_i
 301 to the i -th chosen eigendirection. Deflating the largest s eigenvalues (left panel)—i.e., setting $\tau_i = \lambda_i$ so
 302 that these preconditioned eigenvalues become 1—reduces the smoothness constant \hat{L} and the noise level
 303 $K = \frac{\sigma^2}{B} \text{tr}(\mathbf{M}^{-1} \mathbf{H})$ while leaving $\hat{\mu}_{\text{PL}}$ unchanged, yielding a lower noise floor.

304 To isolate the effect of the noise term, the middle panel fixes $\hat{\mu}_{\text{PL}}$. It deflates the top 20 eigenvalues into a
 305 common value v lying between λ_{21} and λ_d by setting $\tau_i = \lambda_i/v$, so that \hat{L} and $\hat{\mu}_{\text{PL}}$ remain unchanged while
 306 K varies. The resulting steady-state losses track this change in K , in line with the predicted noise-floor
 307 scaling. Deflating the smallest s eigenvalues (right panel)—that is, selecting the bottom eigenvectors and
 308 assigning target values τ_i equal to these smallest eigenvalues so that the preconditioned eigenvalues λ_i/τ_i
 309 move to 1—does increase $\hat{\mu}_{\text{PL}}$, but it simultaneously enlarges K . The two effects counterbalance each other,
 310 yielding only modest overall gains, consistent with the predicted noise-floor behavior.

311 4.2 SciML problems

312 We then briefly summarize the three SciML tasks used to evaluate late-stage optimization behavior under
 313 different preconditioners.

314 **Noisy Franke surface regression.** The Franke function is a classical multiscale benchmark consisting
 315 of several Gaussian peaks with heterogeneous length scales. We sample 256 points uniformly in $[0, 1]^2$
 316 and corrupt the values with Gaussian noise $\varepsilon \sim \mathcal{N}(0, 10^{-4})$. The combination of multiscale structure and
 317 observational noise yields a loss landscape with varying curvature, making it well suited for evaluating how
 318 preconditioning affects convergence in practice. The target surface is

$$f(x, y) = 0.75e^{-\frac{(9x-2)^2 + (9y-2)^2}{4}} + 0.75e^{-\frac{(9x+1)^2}{49} - \frac{9y+1}{10}} + 0.5e^{-\frac{(9x-7)^2 + (9y-3)^2}{4}} - 0.2e^{-(9x-4)^2 - (9y-7)^2}.$$

319 **Physics-informed neural networks (PINNs).** We train a PINN to solve the 2D Poisson problem

$$-\Delta u = f(x, y) = 8\pi^2 \sin(2\pi x) \sin(2\pi y) \quad \text{in } (0, 1)^2, \quad u = 0 \text{ on } \partial[0, 1]^2,$$

320 whose exact solution is $u(x, y) = \sin(2\pi x) \sin(2\pi y)$. The training set includes 1,000 interior residual points
 321 and 200 boundary points. The weighted loss (PDE residual weight 1.0, boundary weight 100.0) produces a
 322 challenging composite landscape known to stress first-order methods (Krishnapriyan et al., 2021). The right
 323 panel of Fig. 3 visualizes the source term $f(x, y)$.

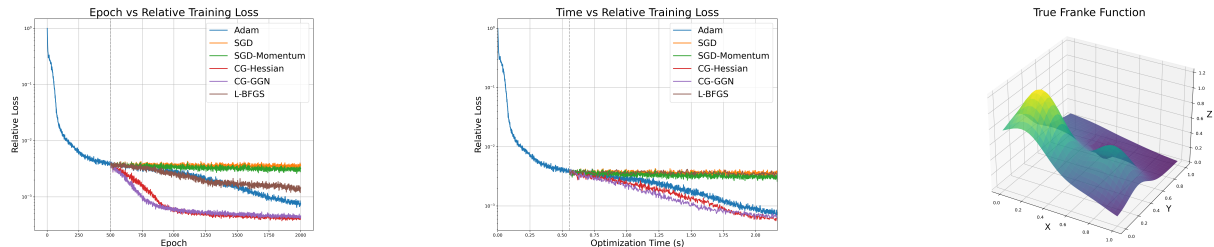


Figure 2: Franke-function regression (mean over 5 runs). Left: relative training loss vs. epochs with the switch to Phase II at epoch 500. Center: relative training loss vs. wall-clock time. Here, relative training loss denotes the training objective normalized by its value at the first plotted epoch, so the initial plotted value is 1.0.

Right: Franke surface.

324 **Green’s-function learning.** We learn Green’s functions for the 1D convection–diffusion operator

$$\mathcal{L}u := -\nu u'' + \beta u', \quad u(0) = u(1) = 0,$$

325 under two regimes: (i) diffusion-dominated ($\nu = 1.0, \beta = 0$) and (ii) convection-dominated ($\nu = 0.1, \beta = 1.0$).
 326 The Green’s function satisfies $\mathcal{L}G(x, y) = \delta(x - y)$, where we approximate the delta distribution by a narrow
 327 Gaussian with width $\sigma = 0.01$. Training uses: (a) 1,000 uniformly sampled (x, y) pairs for PDE residuals, (b)
 328 500 near-diagonal samples ($|x - y|$ small) to capture the near-singularity, and (c) 200 boundary samples. This
 329 produces a highly multiscale and stiffness–dominated optimization problem, ideal for testing curvature-aware
 330 preconditioners.

331 **Baselines and protocol.** Across all SciML tasks, we compare vanilla SGD, momentum, Adam, L–BFGS,
 332 and curvature-aware preconditioners (CG–Hessian and CG–GGN/Fisher). Matrix–free CG with a fixed
 333 iteration budget is used to apply Hessian or Gauss–Newton/Fisher updates. Following standard SciML
 334 practice, we adopt a two-phase schedule: Phase I uses Adam to reach a comparable local basin; Phase II
 335 switches to the target optimizer to isolate late-stage behavior. Because our nonconvex theory is local, the
 336 basin reached at the end of Phase I can influence the local constants $(\hat{L}, \hat{\mu}_{PL}, K)$ encountered in Phase II
 337 and hence may affect which optimizer performs best after the switch. We therefore use the same Adam
 338 warm start, switch point, architecture, and seed protocol across all methods to control for basin selection
 339 and interpret the Phase II results as comparisons conditional on entering a comparable basin rather than
 340 fully basin-agnostic rankings. For all loss-versus-epoch and loss-versus-time plots, we report the *relative*
 341 *training loss*, i.e., the task-specific training objective normalized so that the first plotted epoch has value 1.0.
 342 Thus values below 1 indicate reduction relative to the initial training loss. We report loss vs. epochs and
 343 wall–clock time, with all architectural and implementation details in Appendix D. All implementations use
 344 JAX (Bradbury et al., 2018); code and data are available in the supplemental material.

345 4.3 Noisy data regression

346 After the Adam warm start (Phase I), Phase II separates the methods (Fig. 2): *Adam*, *L–BFGS*, *CG–*
 347 *GGN*, and *CG–Hessian* descend faster than *SGD* and *SGD+Momentum*. The two curvature-aware variants,
 348 *CG–Hessian* and *CG–GGN*, track one another closely–showing similar contraction and reaching essentially
 349 the same loss floor. The similar performance of *CG–Hessian* and *CG–GGN* suggests that both methods
 350 provide comparable normal-space curvature and covariance matrix structure approximation. Adam’s diagonal
 351 rescaling and *L–BFGS*’s low-rank curvature information also mitigate anisotropy and stabilize noisy directions,
 352 which explains their advantage over *SGD*. In wall–clock time, the faster descent of curvature-aware methods
 353 compensates for their higher per-step cost.

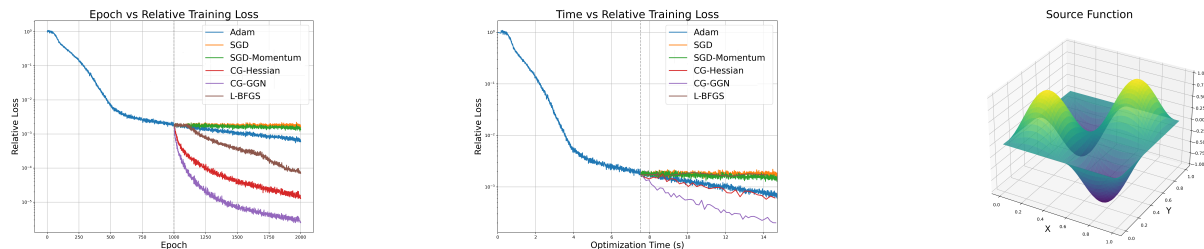


Figure 3: PINN for a Poisson-type PDE (mean over 5 runs). Left: relative training loss vs. epochs with Phase I \rightarrow Phase II at epoch 1,000. Center: relative training loss vs. wall-clock time. Right: source term.

354 4.4 Physics-informed neural networks (PINNs)

355 With the same two-phase protocol, Phase II shows a consistent ranking (Fig. 3). At the bottom, *Adam*
 356 and *SGD/SGD+Momentum* lack explicit curvature information and progress slowly. *L-BFGS* achieves
 357 intermediate performance: it captures limited curvature through its low-rank approximation and line search,
 358 but the memory constraint prevents it from matching the full curvature captured by the two *CG* methods.
 359 At the top tier, *CG-GGN* and *CG-Hessian* both achieve better performance as curvature-aware methods,
 360 with *CG-GGN* showing a slight advantage.

361 For PINNs, which minimize weighted least-squares residuals, the Gauss-Newton approximation $\mathbf{J}^\top \mathbf{J}$ is
 362 naturally aligned with the gradient covariance structure and thus provides more effective noise attenua-
 363 tion—consistent with our theory, where the preconditioned noise level is governed by $\text{tr}(\mathbf{M}^{-1} \Sigma(\mathbf{w}))$ in the late
 364 stage. The Hessian approximation, by contrast, can introduce negative curvature and additional anisotropy.
 365 In wall-clock time, *CG-GGN* achieves the best accuracy within a comparable time budget, despite its higher
 366 per-step cost.

367 4.5 Green’s function learning

368 After Phase I, Phase II again shows a clear separation of methods (Figs. 4–5). In both the diffusion- and
 369 convection-dominated cases, *CG-GGN* continues to drive the loss down, whereas *CG-Hessian*, *L-BFGS*,
 370 *Adam*, *SGD*, and *SGD+Momentum* quickly form a tight cluster and improve only marginally. Compared with
 371 the earlier PINNs experiment, the Green’s-function tasks are more near-singular due to the smoothed-delta
 372 forcing, leading to a more challenging, highly anisotropic optimization problem.

373 Although we did not directly measure the local constants $(\hat{L}, \hat{\mu}_{\text{PL}}, K)$ on this run, the observed advantage
 374 of *CG-GGN* is consistent with the structure of PINN objectives. First, for squared-residual losses, the
 375 Gauss-Newton/Fisher matrix is positive semidefinite, avoiding the negative-curvature directions introduced
 376 by second-derivative terms in the exact Hessian. This makes the preconditioner more stable and better
 377 suited to CG. Second, Fisher-type preconditioners are built from gradient second moments and therefore tend
 378 to *whiten* gradient noise, reducing the preconditioned noise level K . In contrast, a Hessian preconditioner
 379 includes second-order terms that are often misaligned with the gradient-noise covariance, and the damping
 380 needed to handle indefiniteness diminishes curvature gains while weakening noise attenuation.

381 These two effects—better alignment with useful curvature and more effective noise whitening—explain why
 382 *CG-GGN* reaches lower losses within comparable wall-clock time, despite its higher per-step cost.

383 The right panels of Figs. 4 and 5 display the learned Green’s functions $G(x, y)$ at three representative source
 384 locations y together with simple operator and boundary checks for *CG-GGN*. The kernels are localized
 385 around the source locations and decay toward the Dirichlet boundaries, and the corresponding operator
 386 evaluations produce narrow spikes at $x = y$, in line with the smoothed-delta forcing used in the training
 387 loss. This suggests that the lower training losses achieved by *CG-GGN* reflect a reasonable Green’s-function
 388 approximation rather than a purely numerical artifact.

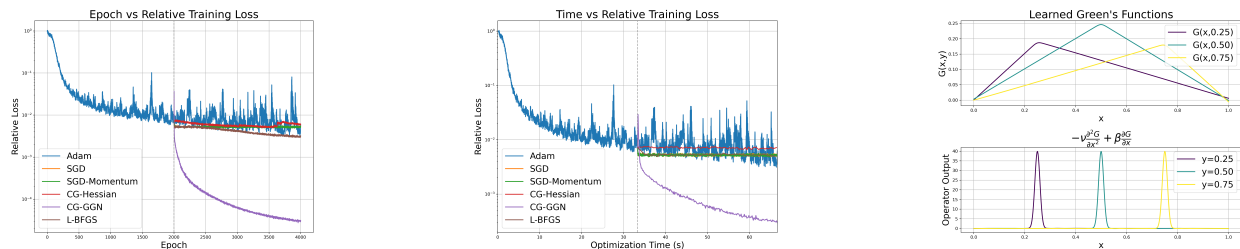


Figure 4: Laplacian Green’s function learning (mean over 5 runs). Left: relative training loss vs. epochs with Phase I \rightarrow Phase II at epoch 2,000. Center: relative training loss vs. wall-clock time. Right: learned $G(x, y)$ for three source locations and operator checks.

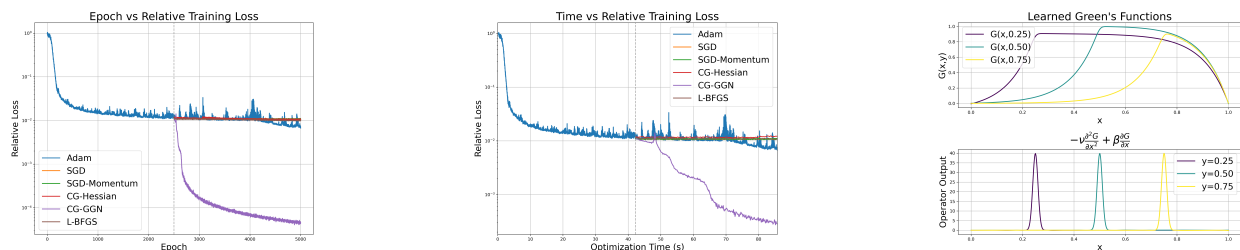


Figure 5: Convection–diffusion Green’s function learning (mean over 5 runs). Left: relative training loss vs. epochs with Phase I \rightarrow Phase II at epoch 2,500. Center: relative training loss vs. wall-clock time. Right: learned $G(x, y)$ for three source locations and operator checks.

389 We conclude the numerical experiments by connecting the CG–GGN preconditioner to the theoretical
 390 convergence framework developed in this paper. We empirically examine the quantities L and K that govern
 391 the convergence of preconditioned SGD for the PINNs problem and two Green’s function learning problems.
 392 Because a CG-based preconditioner with only a few iterations typically does not significantly alter the
 393 cluster of near-zero eigenvalues, we treat the \mathbf{M} –PL constant as unchanged and attribute the quality of the
 394 preconditioner primarily to its effect on L and K . For these three problems, we fix the random seed to 42
 395 and analyze the network parameters at epoch 250 in Phase II. After preconditioning, the L value reduced by
 396 factors of 78x, 3710x, and 1923x, respectively. We additionally quantify the impact of preconditioning on the
 397 noise level K . Using the same network parameters \mathbf{w} , we sample 100 independent mini-batches, construct the
 398 preconditioner \mathbf{M}^{-1} from the first batch, and observe that after preconditioning the estimated trace of the
 399 gradient-noise covariance matrix is reduced by factors of 12x, 1505x, and 203x, respectively. This substantial
 400 reduction demonstrates that the CG–GGN preconditioner effectively attenuates gradient noise. Consistent
 401 with our theory, the combined improvements in conditioning and noise reduction yield both faster linear
 402 convergence and a significantly lower asymptotic noise floor.

403 5 Conclusion

404 We developed a local, geometry-aware theory for preconditioned SGD that makes two effects explicit: (1) the
 405 rate inside a basin is controlled by a preconditioner-dependent condition number in the \mathbf{M} –metric, and (2)
 406 the noise floor is governed by the preconditioned noise. We additionally obtained a basin-stability guarantee,
 407 giving an explicit probability that iterates remain in a region where these local properties hold. Together,
 408 the results motivate a simple rule: choose \mathbf{M} to improve local conditioning while suppressing noise in the
 409 \mathbf{M}^{-1} –norm.

410 A key next direction is *covariance-aware* preconditioning. Our bounds suggest that effective design should
 411 jointly target conditioning and noise attenuation, motivating structured covariance models and adaptive
 412 schemes that update curvature and noise statistics simultaneously. Extending basin-stability guarantees

⁴¹³ to nonstationary noise and developing online diagnostics for the local constants would move toward fully
⁴¹⁴ adaptive, geometry- and noise-aware SGD.

References

- 415
- 416 Shunichi Amari. A theory of adaptive pattern classifiers. *IEEE Transactions on Electronic Computers*, EC-16
417 (3):299–307, 1967. ISSN 0367-7508. doi: 10.1109/PGEC.1967.264666. URL <https://ieeexplore.ieee.org/document/4039068>. Conference Name: IEEE Transactions on Electronic Computers.
- 418
- 419 Shunichi Amari. Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):251–276,
420 February 1998. ISSN 0899-7667. doi: 10.1162/089976698300017746. URL <https://doi.org/10.1162/089976698300017746>.
- 421
- 422 Shunichi Amari, Jimmy Ba, Roger Baker Grosse, Xuechen Li, Atsushi Nitanda, Taiji Suzuki, Denny Wu, and
423 Ji Xu. When does preconditioning help or hurt generalization? In *International Conference on Learning*
424 *Representations*, 2021. URL https://openreview.net/forum?id=S724o4_WB3.
- 425 Amit Attia and Tomer Koren. Sgd with adagrad stepsizes: Full adaptivity with high probability to unknown
426 parameters, unbounded gradients and affine variance. In *International Conference on Machine Learning*,
427 pp. 1147–1171. PMLR, 2023.
- 428 Francis Bach. *Learning Theory from First Principles*. The MIT Press, 2024. ISBN 9780262049443.
- 429 David GT Barrett and Benoit Dherin. Implicit gradient regularization. *arXiv preprint arXiv:2009.11162*,
430 2020.
- 431 Albert S. Berahas, Jorge Nocedal, and Martin Takáč. A multi-batch l-bfgs method for machine learning,
432 2016. URL <https://arxiv.org/abs/1605.06049>.
- 433 Julius R. Blum. Approximation methods which converge with probability one. *The Annals of Mathematical*
434 *Statistics*, 25(2):382–386, 1954. ISSN 0003-4851, 2168-8990. doi: 10.1214/aoms/1177728794. URL
435 [https://projecteuclid.org/journals/annals-of-mathematical-statistics/volume-25/issue-2/](https://projecteuclid.org/journals/annals-of-mathematical-statistics/volume-25/issue-2/Approximation-Methods-which-Converge-with-Probability-one/10.1214/aoms/1177728794.full)
436 [Approximation-Methods-which-Converge-with-Probability-one/10.1214/aoms/1177728794.full](https://projecteuclid.org/journals/annals-of-mathematical-statistics/volume-25/issue-2/Approximation-Methods-which-Converge-with-Probability-one/10.1214/aoms/1177728794.full).
437 Publisher: Institute of Mathematical Statistics.
- 438 Raghu Bollapragada, Dheevatsa Mudigere, Jorge Nocedal, Hao-Jun Michael Shi, and Ping Tak Peter Tang. A
439 progressive batching l-bfgs method for machine learning, 2018. URL <https://arxiv.org/abs/1802.05374>.
- 440 Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learn-
441 ing. *SIAM Review*, 60(2):223–311, 2018. doi: 10.1137/16M1080173. URL [https://doi.org/10.1137/](https://doi.org/10.1137/16M1080173)
442 [16M1080173](https://doi.org/10.1137/16M1080173).
- 443 James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin,
444 George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable
445 transformations of Python+NumPy programs, 2018. URL <http://github.com/jax-ml/jax>.
- 446 Han Chan. Some methods of unconstrained minimization. *USSR Computational Mathematics and Mathe-*
447 *matical Physics*, 19(2):31–44, 1979. ISSN 0041-5553. doi: [https://doi.org/10.1016/0041-5553\(79\)90004-1](https://doi.org/10.1016/0041-5553(79)90004-1).
448 URL <https://www.sciencedirect.com/science/article/pii/0041555379900041>.
- 449 Jinghui Chen, Dongruo Zhou, Yiqi Tang, Ziyang Yang, Yuan Cao, and Quanquan Gu. Closing the generalization
450 gap of adaptive gradient methods in training deep neural networks. In Christian Bessiere (ed.), *Proceedings*
451 *of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pp. 3267–3275.
452 International Joint Conferences on Artificial Intelligence Organization, 7 2020. doi: 10.24963/ijcai.2020/452.
453 URL <https://doi.org/10.24963/ijcai.2020/452>. Main track.
- 454 Weizhu Chen, Zhenghao Wang, and Jingren Zhou. Large-scale l-bfgs using mapreduce. In *Proceedings of*
455 *the 28th International Conference on Neural Information Processing Systems - Volume 1, NIPS’14*, pp.
456 1332–1340, Cambridge, MA, USA, 2014. MIT Press.
- 457 John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic
458 optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.

- 459 Matthew Faw, Isidoros Tziotis, Constantine Caramanis, Aryan Mokhtari, Sanjay Shakkottai, and Rachel
460 Ward. The power of adaptivity in sgd: Self-tuning step sizes with unbounded gradients and affine variance.
461 In *Conference on Learning Theory*, pp. 313–355. PMLR, 2022.
- 462 R. Fletcher. *Practical Methods of Optimization*. Wiley, 2013. ISBN 978-1-118-72318-0. URL https://books.google.com/books?id=_WuAvIx0EE4C.
- 464 Richard Franke. A Critical Comparison of Some Methods for Interpolation of Scattered Data. Technical
465 report, Graduate School of Operational and Information Sciences (GSOIS), 1979. URL <https://apps.dtic.mil/sti/citations/ADA081688>.
- 467 Sachin Garg, Albert S. Berahas, and Michał Dereziński. Second-order information promotes mini-batch
468 robustness in variance-reduced gradients, 2024. URL <https://arxiv.org/abs/2404.14758>.
- 469 Guillaume Garrigos and Robert M. Gower. Handbook of convergence theorems for (stochastic) gradient
470 methods, 2024. URL <http://arxiv.org/abs/2301.11235>.
- 471 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization via
472 hessian eigenvalue density. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*
473 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp.
474 2232–2241. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/ghorbani19b.html>.
- 475 Joshua D Griffin, Majid Jahani, Martin Takáč, Seyedalireza Yektamaram, and Wenwen Zhou. A minibatch
476 stochastic Quasi-Newton method adapted for nonconvex deep learning problems. *Optimization Online*,
477 2022.
- 478 Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization.
479 In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine*
480 *Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1842–1850. PMLR, 10–15 Jul 2018.
481 URL <https://proceedings.mlr.press/v80/gupta18a.html>.
- 482 Wenrui Hao, Rui Peng Li, Yuanzhe Xi, Tianshi Xu, and Yahong Yang. Multiscale neural networks for
483 approximating green’s functions, 2024. URL <https://arxiv.org/abs/2410.18439>.
- 484 Geoffrey Hinton. Neural networks for machine learning lecture 6e: rmsprop: Divide the gradient by a running
485 average of its recent magnitude, 2014. URL https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
- 487 Satoki Ishikawa and Ryo Karakida. On the parameterization of second-order optimization effective towards
488 the infinite width. In *The Twelfth International Conference on Learning Representations*, 2024. URL
489 <https://openreview.net/forum?id=g8sGBSjYk>.
- 490 Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient
491 methods under the polyak-łojasiewicz condition. In Paolo Frasconi, Niels Landwehr, Giuseppe Manco, and
492 Jilles Vreeken (eds.), *Machine Learning and Knowledge Discovery in Databases*, pp. 795–811, Cham, 2016.
493 Springer International Publishing. ISBN 978-3-319-46128-1.
- 494 Ahmed Khaled and Peter Richtárik. Better theory for SGD in the nonconvex world. *Transactions on Machine*
495 *Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=AU4qHN2Vks>. Survey
496 Certification.
- 497 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <http://arxiv.org/abs/1412.6980>.
- 499 Tomer Koren, Roi Livni, Yishay Mansour, and Uri Sherman. Benign underfitting of stochastic gradient
500 descent. *Advances in Neural Information Processing Systems*, 35:19605–19617, 2022.

- 501 Aditi Krishnapriyan, Amir Gholami, Shandian Zhe, Robert Kirby, and Michael W Mahoney. Characterizing
502 possible failure modes in physics-informed neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin,
503 P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34,
504 pp. 26548–26560. Curran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/
505 paper/2021/file/df438e5206f31600e6ae4af72f2725f1-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/df438e5206f31600e6ae4af72f2725f1-Paper.pdf).
- 506 Frederik Kunstner, Lukas Balles, and Philipp Hennig. Limitations of the empirical fisher approximation
507 for natural gradient descent. In *Proceedings of the 33rd International Conference on Neural Information
508 Processing Systems*, Red Hook, NY, USA, 2019. Curran Associates Inc.
- 509 Susanna Lange, Kyle Helfrich, and Qiang Ye. Batch normalization preconditioning for neural network training.
510 *Journal of Machine Learning Research*, 23(72):1–41, 2022.
- 511 Xi-Lin Li. Preconditioned stochastic gradient descent. *IEEE Transactions on Neural Networks and Learning
512 Systems*, 29(5):1454–1466, 2018. doi: 10.1109/TNNLS.2017.2672978.
- 513 Dong C. Liu and Jorge Nocedal. On the limited memory BFGS method for large scale optimization.
514 *Mathematical Programming*, 45(1):503–528, August 1989. ISSN 1436-4646. doi: 10.1007/BF01589116. URL
515 <https://doi.org/10.1007/BF01589116>.
- 516 Hong Liu, Zhiyuan Li, David Leo Wright Hall, Percy Liang, and Tengyu Ma. Sophia: A scalable stochastic
517 second-order optimizer for language model pre-training. In *The Twelfth International Conference on
518 Learning Representations*, 2024. URL <https://openreview.net/forum?id=3xHDeA8Noi>.
- 519 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*,
520 2017.
- 521 James Martens. New insights and perspectives on the natural gradient method. *Journal of Machine Learning
522 Research*, 21(146):1–76, 2020. URL <http://jmlr.org/papers/v21/17-678.html>.
- 523 James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature.
524 In *Proceedings of the 32nd International Conference on International Conference on Machine Learning -
525 Volume 37*, ICML’15, pp. 2408–2417. JMLR.org, 2015.
- 526 Pratik Rathore, Weimu Lei, Zachary Frangella, Lu Lu, and Madeleine Udell. Challenges in training pinns:
527 a loss landscape perspective. In *Proceedings of the 41st International Conference on Machine Learning*,
528 ICML’24. JMLR.org, 2024.
- 529 Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International
530 Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=ryQu7f-RZ>.
- 531 Herbert Robbins and Sutton Monro. A stochastic approximation method. *The Annals of Math-
532 ematical Statistics*, 22(3):400–407, 1951. ISSN 0003-4851, 2168-8990. doi: 10.1214/aoms/
533 1177729586. URL [https://projecteuclid.org/journals/annals-of-mathematical-statistics/
534 volume-22/issue-3/A-Stochastic-Approximation-Method/10.1214/aoms/1177729586.full](https://projecteuclid.org/journals/annals-of-mathematical-statistics/volume-22/issue-3/A-Stochastic-Approximation-Method/10.1214/aoms/1177729586.full). Pub-
535 lisher: Institute of Mathematical Statistics.
- 536 Levent Sagun, Utku Evci, V. Ugur Guney, Yann Dauphin, and Leon Bottou. Empirical analysis of the hessian
537 of over-parametrized neural networks, 2018. URL <https://openreview.net/forum?id=rJrTwxCb>.
- 538 Robin M Schmidt, Frank Schneider, and Philipp Hennig. Descending through a crowded valley - benchmarking
539 deep learning optimizers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International
540 Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9367–9376.
541 PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/schmidt21a.html>.
- 542 Frank Schneider, Lukas Balles, and Philipp Hennig. DeepOBS: A deep learning optimizer benchmark suite.
543 In *International Conference on Learning Representations*, 2019. URL [https://openreview.net/forum?
544 id=rJg6ssC5Y7](https://openreview.net/forum?id=rJg6ssC5Y7).

- 545 Nicol N. Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. *Neural*
546 *Computation*, 14(7):1723–1738, 2002. ISSN 0899-7667, 1530-888X. doi: 10.1162/08997660260028683. URL
547 <https://direct.mit.edu/neco/article/14/7/1723-1738/6626>.
- 548 Tianshi Xu, Rui Peng Li, and Yuanzhe Xi. Neural approximate inverse preconditioners, 2025. URL
549 <https://arxiv.org/abs/2510.13034>.
- 550 Qiang Ye. Preconditioning for accelerated gradient descent optimization and regularization, 2024. URL
551 <https://arxiv.org/abs/2410.00232>.
- 552 Manzil Zaheer, Sashank Reddi, Devendra Sachan, Satyen Kale, and Sanjiv Kumar. Adaptive meth-
553 ods for nonconvex optimization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-
554 Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Cur-
555 ran Associates, Inc., 2018. URL [https://proceedings.neurips.cc/paper_files/paper/2018/file/](https://proceedings.neurips.cc/paper_files/paper/2018/file/90365351ccc7437a1309dc64e4db32a3-Paper.pdf)
556 [90365351ccc7437a1309dc64e4db32a3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/90365351ccc7437a1309dc64e4db32a3-Paper.pdf).
- 557 Handi Zhang, Langchen Liu, and Lu Lu. Federated scientific machine learning for approximating functions and
558 solving differential equations with data heterogeneity, 2024. URL <https://arxiv.org/abs/2410.13141>.
- 559 Thomas TCK Zhang, Behrad Moniri, Ansh Nagwekar, Faraz Rahman, Anton Xue, Hamed Hassani, and
560 Nikolai Matni. On the concurrence of layer-wise preconditioning methods and provable feature learning.
561 In *Forty-second International Conference on Machine Learning*, 2025. URL [https://openreview.net/](https://openreview.net/forum?id=aRUUFFycNh)
562 [forum?id=aRUUFFycNh](https://openreview.net/forum?id=aRUUFFycNh).

Appendix

A Notation used in paper

In general, capital bold letters are matrices (\mathbf{A}), lower case bold letters are vectors (\mathbf{v}), and lower case Greek or Latin letters are constants (ν, c). Moreover, there are some notation that is used consistently throughout the paper. A reference table for these symbols is given in Table 1.

Table 1: Reference for recurring notation in the paper.

Symbol	Definition
k	iteration counter
\mathbf{w}	Model parameters
$F(\mathbf{w})$	Objective function at point w
$F_* := F(\mathbf{w}^*)$	minimum function value at minimizer
α	learning rate
α_k	learning rate scheduler/ learning rate at epoch k
$\bar{\alpha}$	fixed learning rate
c, L	strong convexity, Lipschitz constant for $\ \cdot\ _2 = \ \cdot\ _{\mathbf{I}}$
\hat{c}, \hat{L}	strong convexity, Lipschitz constant for preconditioned case: $\ \cdot\ _{\mathbf{M}}$
$\hat{\mu}_{PL}$	PL constant for preconditioned case: $\ \cdot\ _{\mathbf{M}}$
\mathcal{B}	mini-batch of the dataset
\mathbf{M}	generic preconditioner where \mathbf{M}^{-1} is applied to a vector
$g(\cdot, \cdot)$	gradient vector
$\kappa(\mathbf{M})$	Condition number of \mathbf{M} (always based on $\ \cdot\ _2$)
μ, μ_G	lower and upper bound constants on the first moment of the gradient
K, K_V	constant and scaling values of the affine bound on the gradient's variance
K_G	Constant needed for learning rate upper bound, dependent on $K_V + \mu_G^2 > 0$.
$\mathbb{E}_{\xi}, \mathbb{V}_{\xi}$	Expectation and Variance of gradient with random realization ξ
β, γ	constants affecting the lower and upper bound on α_k for diminishing learning rate proofs
ν	convergence constant in $\mathcal{O}((\gamma + k)^{-1})$
r	radius of convex basin around local minimum
$\mathcal{N}_r, \mathcal{N}_{r+}$	local neighborhood around minimizer, slightly larger local neighborhood for containment
τ	smallest iteration number where $\mathbf{w}_k \notin \mathcal{N}_r$.
C	The stochastic noise floor defined $\bar{\alpha}\hat{L}K/(2\hat{c}\mu)$
$\mathcal{N}_{\mathbf{M}}(\mathbf{w})$	instantaneous preconditioned noise $\text{tr}(\mathbf{M}^{-1}\Sigma(\mathbf{w}))$
K	uniform baseline for $\mathcal{N}_{\mathbf{M}}(\mathbf{w})$ on the analysis region (noise floor constant)
α_{QG}	quadratic growth constant of locally convex basin a distance from the minimizer

B Mathematical preliminaries

B.1 Preconditioning

The condition number from a linear equation $\mathbf{Ax} = \mathbf{b}$ bounds the accuracy of the solution \mathbf{x} , and is defined as

$$\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|,$$

where if not stated $\|\cdot\| = \|\cdot\|_2$. If \mathbf{A} is ill-conditioned, i.e. has a large condition number, then a small perturbation in \mathbf{b} can result in a large perturbation of the solution \mathbf{x} . In addition to the accuracy of the solution, the convergence rate of iterative methods, such as conjugate gradient, depends on $r = \frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}$.

It is easy to see that $r < 1$, but if $\kappa \gg 1$, then convergence will be extremely slow as $r \rightarrow 1$. This motivates the need for ways to reduce the condition number, through a technique called *preconditioning*. Throughout

575 this paper, we assume that \mathbf{M} is the preconditioner, and we only have access to the action of \mathbf{M}^{-1} onto a
 576 vector. More technically, we say \mathbf{M} is an efficient preconditioner to the matrix \mathbf{A} such that

$$\kappa(\mathbf{M}^{-1}\mathbf{A}) < \kappa(\mathbf{A}).$$

577 For clarity, even though we call \mathbf{M} the preconditioner, we don't explicitly form it. Additionally, we don't
 578 form \mathbf{M}^{-1} either but just observe the action of the preconditioner on a vector, $\mathbf{M}^{-1}\mathbf{v}$.

There are different ways we can utilize the preconditioner \mathbf{M} . First, assume \mathbf{M}^{-1} exists, then the *left* preconditioned system is

$$\mathbf{M}^{-1}(\mathbf{A}\mathbf{x} - \mathbf{b}) = 0.$$

Both the original linear system and the left-preconditioned system give the same solution. Additionally, we could solve the right preconditioned system

$$\mathbf{A}\mathbf{M}^{-1}(\mathbf{M}\mathbf{x}) = \mathbf{b}.$$

579 This requires us to solve $\mathbf{A}\mathbf{M}^{-1}\mathbf{y} = \mathbf{b}$ for \mathbf{y} , and then to recover the original solution, we would need to do
 580 another linear system solve $\mathbf{M}\mathbf{x} = \mathbf{y}$ for \mathbf{x} .

581 These two techniques can be combined to perform *split* preconditioning. If we employ \mathbf{M} as the right
 582 preconditioner, and \mathbf{N} as the left preconditioner, we compute

$$\mathbf{N}\mathbf{A}\mathbf{M}^{-1}(\mathbf{M}\mathbf{x}) = \mathbf{N}\mathbf{b}.$$

583 This is beneficial if one would like to scale the rows and columns of \mathbf{A} differently. Additionally, observe that
 584 if \mathbf{A} is symmetric and $\mathbf{N}^\top = \mathbf{M}^{-1}$, then $\mathbf{N}\mathbf{A}\mathbf{M}^{-1}$ is also symmetric.

585 In the preconditioned version of CG (PCG), one solves the equivalent system $\mathbf{M}^{-1}\mathbf{A}\mathbf{x} = \mathbf{M}^{-1}\mathbf{b}$ using a similar
 586 three-term recurrence, but applied to the transformed system. The key requirement is that the preconditioner
 587 \mathbf{M} be symmetric positive definite and chosen so that $\mathbf{M}^{-1}\mathbf{A}$ has a significantly smaller condition number
 588 than \mathbf{A} itself. For practical purposes, PCG is used in matrix-free settings where only the action $\mathbf{M}^{-1}\mathbf{v}$ is
 589 required, not the explicit matrix \mathbf{M}^{-1} .

590 B.2 Preconditioners for SGD

591 In this section, we briefly review several preconditioners commonly used in the ML literature. First, if we
 592 define \mathbf{g}_k to be the sum of the squared gradients up until iteration k , we arrive at AdaGrad (Duchi et al.,
 593 2011)

$$\mathbf{M}_{\text{AdaGrad}} = \text{diag}(\sqrt{\mathbf{g}_k} + \varepsilon).$$

594 The issues with this is the gradient squared will only increase, leading to premature stopping. To coun-
 595 teract that, exponentially moving weighted averages are widely used in diagonal preconditioners such as
 596 Adam (Kingma & Ba, 2017) and its momentum-less counterpart RMSProp (Hinton, 2014):

$$\mathbf{M}_{\text{Adam}} = \text{diag}(\sqrt{\mathbf{s}_k} + \varepsilon),$$

597 where here \mathbf{s}_k is an exponential moving average of squared gradients, and $\varepsilon > 0$ is a small constant added for
 598 numerical stability. While computationally efficient and robust to scaling, such diagonal preconditioners fail
 599 to capture cross-parameter curvature, which may lead to suboptimal convergence in ill-conditioned problems.

600 The Hessian matrix of the loss function,

$$\mathbf{H}(\mathbf{w}) = \nabla^2 \mathcal{L}(\mathbf{w}),$$

601 captures the exact second-order structure of the problem and provides the most complete curvature information.
 602 However, computing or storing the full Hessian is typically infeasible in high-dimensional neural network
 603 (NN) models. Moreover, it is not guaranteed to be positive definite in nonconvex settings, which complicates
 604 its direct use as a preconditioner.

605 To reduce computational cost, one can approximate the Hessian using a single mini-batch, \mathcal{B} :

$$\mathbf{H}_{\mathcal{B}}(\mathbf{w}) = \nabla^2 \mathcal{L}_{\mathcal{B}}(\mathbf{w}).$$

606 This matrix is cheaper to compute and can be updated online, but suffers from high variance and may not
 607 preserve important curvature directions observed over the full dataset. While the Newton and quasi-Newton
 608 methods work well for deterministic optimization, many have provided a distinction between these and other
 609 methods for designing preconditioners in the stochastic setting (Li, 2018; Bottou et al., 2018).

610 As opposed to constructing the Hessian, an alternative is the Gauss-Newton Hessian approximation, which
 611 assumes the difference between the model and label is small in a least-squares norm. This idea was further
 612 generalized to loss functions of the form $\ell(\theta) = \sum_n a_n(b_n(\theta))$ in Schraudolph (2002). This generalized
 613 Gauss-Newton matrix (GGN), which ignores second order information of b_n , is SPD when a_n is convex even
 614 when the true Hessian is indefinite.

615 Another alternate method is the FIM defined as

$$\mathbf{F}(\mathbf{w}) = \mathbb{E}_{x,y} [\nabla_{\mathbf{w}} \log p_{\mathbf{w}}(y | x) \nabla_{\mathbf{w}} \log p_{\mathbf{w}}(y | x)^{\top}],$$

616 which is guaranteed to be SPD under mild regularity conditions. For models trained with exponential-family
 617 losses, the FIM coincides with the GGN (Martens, 2020; Schraudolph, 2002). Its structure allows for stable
 618 and curvature-aware preconditioning. Using the FIM as the preconditioner in the stochastic gradient descent
 619 algorithm yields Natural Gradient Descent from online learning Amari (1998).

620 The empirical FIM estimates the expectation in the FIM using a finite mini-batch:

$$\mathbf{F}_{\text{emp}}(\mathbf{w}) = \frac{1}{|\mathcal{B}|} \sum_{(x,y) \in \mathcal{B}} \nabla_{\mathbf{w}} \log p_{\mathbf{w}}(y | x) \nabla_{\mathbf{w}} \log p_{\mathbf{w}}(y | x)^{\top}.$$

621 It is symmetric and positive semidefinite, and is often used in practice due to its lower computational overhead
 622 compared to the full FIM. However, it may introduce bias depending on the mini-batch size and model
 623 quality (Kunstner et al., 2019).

624 Finally, the L-BFGS algorithm is a popular quasi-Newton method that builds a low-rank approximation
 625 to the inverse Hessian using a history of gradients and iterates. It is well-suited to medium-scale problems
 626 and has seen empirical success in ML (Bottou et al., 2018). Additional variants of L-BFGS have also been
 627 proposed (Berahas et al., 2016; Bollapragada et al., 2018). While not traditionally framed as a preconditioner,
 628 L-BFGS can be interpreted as implicitly applying a data-driven curvature approximation.

629 C Assumptions and proofs of theorems

630 C.1 Assumptions

631 **Assumption 9** (Strong Convexity). *The objective function $F: \mathbb{R}^d \rightarrow \mathbb{R}$ is strongly convex in that there exists*
 632 *a constant $c > 0$ such that*

$$F(\bar{\mathbf{w}}) \geq F(\mathbf{w}) + \nabla F(\mathbf{w})^{\top} (\bar{\mathbf{w}} - \mathbf{w}) + \frac{1}{2} c \|\bar{\mathbf{w}} - \mathbf{w}\|_2^2, \quad \forall (\bar{\mathbf{w}}, \mathbf{w}) \in \mathbb{R}^d \times \mathbb{R}^d$$

633 From elementary optimization, this assumption is equivalent to F having a unique minimizer $\mathbf{w}^* \in \mathbb{R}^d$. We
 634 define $F_* := F(\mathbf{w}^*)$.

635 **Assumption 10** (Lipschitz continuity of gradient). *The objective function $F: \mathbb{R}^d \rightarrow \mathbb{R}$ is continuously*
 636 *differentiable and the gradient function of F , $\nabla F: \mathbb{R}^d \rightarrow \mathbb{R}^d$, is Lipschitz continuous with Lipschitz constant*
 637 *$L > 0$, i.e.*

$$\|\nabla F(\mathbf{w}) - \nabla F(\bar{\mathbf{w}})\|_2 \leq L \|\mathbf{w} - \bar{\mathbf{w}}\|_2$$

638 for all $\{\mathbf{w}, \bar{\mathbf{w}}\} \subset \mathbb{R}^d$.

639 **Remark 1.** If F is continuously twice differentiable, then ∇F is Lipschitz continuous with Lipschitz constant
640 L if and only if the eigenvalues of the matrix $\nabla^2 F(\mathbf{w})$ are bounded above by L for all w . F is strongly convex
641 with constant c if and only if the eigenvalues of the matrix $\nabla^2 F(\mathbf{w})$ is bounded below by c for all w . Therefore,
642 L/c is an upper bound of the condition number of $\nabla^2 F(\mathbf{w})$.

643 Lipschitz continuity of gradient is an assumption made in nearly all convergence analyses of gradient-based
644 methods (Khaled & Richtárik, 2023).

645 **Assumption 11** (Bounds on First and Second Moments of Gradient). Assume

646 1. There exist scalars $\mu_G \geq \mu > 0$ such that, for all $k \in \mathbb{N}$,

$$\nabla F(\mathbf{w}_k)^\top \mathbb{E}_{\xi_k} [g(\mathbf{w}_k, \xi_k)] \geq \mu \|\nabla F(\mathbf{w}_k)\|_2^2 \quad (9)$$

$$\|\mathbb{E}_{\xi_k} [g(\mathbf{w}_k, \xi_k)]\|_2 \leq \mu_G \|\nabla F(\mathbf{w}_k)\|_2 \quad (10)$$

648 2. There exist scalars $K \geq 0$ and $K_V \geq 0$ such that, for all $k \in \mathbb{N}$,

$$\mathbb{V}_{\xi_k} [g(\mathbf{w}_k, \xi_k)] \leq K + K_V \|\nabla F(\mathbf{w}_k)\|_2^2 \quad (11)$$

649 where $\mathbb{V}_{\xi_k} [g(\mathbf{w}_k, \xi_k)] := \mathbb{E}_{\xi_k} [\|g(\mathbf{w}_k, \xi_k)\|_2^2] - \|\mathbb{E}_{\xi_k} [g(\mathbf{w}_k, \xi_k)]\|_2^2$.

650 **Theorem C.1** (Strongly convex objective function, fixed learning rate (Bottou et al., 2018)). Under
651 Assumptions 9,10, 11, suppose that the SGD algorithm is run with fixed learning rates, $\alpha_k = \bar{\alpha}$ for all $k \in \mathbb{N}$
652 where

$$0 < \bar{\alpha} \leq \frac{\mu}{LK_G} \quad \text{and} \quad K_G := K_V + \mu_G^2 \geq \mu^2 > 0.$$

653 Then, the expected optimality gap satisfies the following for all $k \in \mathbb{N}$:

$$\mathbb{E}[F(\mathbf{w}_k) - F_*] \leq \frac{\bar{\alpha}LK}{2c\mu} + (1 - \bar{\alpha}c\mu)^{k-1} \left(F(\mathbf{w}_1) - F_* - \frac{\bar{\alpha}LK}{2c\mu} \right) \xrightarrow{k \rightarrow \infty} \frac{\bar{\alpha}LK}{2c\mu} \quad (12)$$

654 Note that it follows from (10) and (11) that $\mathbb{E}_{\xi_k} [\|g(\mathbf{w}_k, \xi_k)\|_2^2] \leq K + K_G \|\nabla F(\mathbf{w}_k)\|_2^2$ with $K_G := K_V + \mu_G^2 \geq$
655 $\mu^2 > 0$.

656 **Theorem C.2** (Strongly convex objective function, diminishing learning rates (Bottou et al., 2018)). Under
657 the same assumptions as Theorem C.1, suppose that the SGD algorithm is run with a learning rate sequence
658 such that, for all $k \in \mathbb{N}$,

$$\alpha_k = \frac{\beta}{\gamma + k} \text{ for some } \beta > \frac{1}{c\mu} \text{ and } \gamma > 0 \text{ such that } \alpha_1 \leq \frac{\mu}{LK_G}$$

659 Then, the expected optimality gap satisfies the following for all $k \in \mathbb{N}$:

$$\mathbb{E}[F(\mathbf{w}_k) - F_*] \leq \frac{\nu}{\gamma + k} \quad (13)$$

660 where

$$\nu := \max \left\{ \frac{\beta^2 LK}{2(\beta c\mu - 1)}, (\gamma + 1)(F(\mathbf{w}_1) - F_*) \right\} \quad (14)$$

661 Under the assumption of strong convexity, the optimality gap can be bounded at any point by the 2-norm
662 squared of the gradient of the objective function at that particular point. That is,

$$2c(F(\mathbf{w}) - F_*) \leq \|\nabla F(\mathbf{w})\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d$$

663 As before, F has a unique minimizer, denoted as $\mathbf{w}^* \in \mathbb{R}^d$ with $F_* := F(\mathbf{w}^*)$.

664 Previously, the optimality gap was bounded by the 2-norm of the gradient of the objective function squared.
 665 Here, however, the optimality gap is bounded by the \mathbf{M} -norm of the gradient of the objective function squared.
 666 That is,

$$2\hat{c}(F(\mathbf{w}) - F(\mathbf{w}_*)) \leq \|\nabla F(\mathbf{w})\|_{\mathbf{M}^{-1}}^2$$

667 This result is used several times in the upcoming proofs. We repeat Lemma 3.1 here for convenience below:

668 **Lemma C.3.** *Let F be twice differentiable and $\mathbf{M}^{-1} = \mathbf{P}\mathbf{P}^\top$. Then: (i) ∇F is \mathbf{M} -Lipschitz with constant $\hat{L} \iff$ all eigenvalues of $\mathbf{P}^\top \nabla^2 F(\mathbf{w})\mathbf{P}$ are $\leq \hat{L}$; (ii) F is \mathbf{M} -strongly convex with constant $\hat{c} \iff$ all eigenvalues of $\mathbf{P}^\top \nabla^2 F(\mathbf{w})\mathbf{P}$ are $\geq \hat{c}$.*

671 *Proof.* We consider a change of parameter as used in preconditioning. Let $\mathbf{w} = \mathbf{P}\mathbf{z}$ and $\bar{\mathbf{w}} = \mathbf{P}\bar{\mathbf{z}}$. Then
 672 $\mathbf{w} - \bar{\mathbf{w}} = \mathbf{P}(\mathbf{z} - \bar{\mathbf{z}})$ which gives $\mathbf{P}^{-1}(\mathbf{w} - \bar{\mathbf{w}}) = \mathbf{z} - \bar{\mathbf{z}}$. Define $f(\mathbf{z}) = F(\mathbf{P}\mathbf{z})$. Then $\nabla_{\mathbf{z}} f(\mathbf{z}) = \mathbf{P}^\top \nabla_{\mathbf{w}} F(\mathbf{w})$ and
 673 $\nabla_{\mathbf{z}}^2 f(\mathbf{z}) = \mathbf{P}^\top \nabla_{\mathbf{w}}^2 F(\mathbf{w})\mathbf{P}$. Hence

$$\|\nabla f(\mathbf{z}) - \nabla f(\bar{\mathbf{z}})\|_2 = \|\mathbf{P}^\top \nabla_{\mathbf{w}} F(\mathbf{w}) - \mathbf{P}^\top \nabla_{\mathbf{w}} F(\bar{\mathbf{w}})\|_2 = \|\nabla_{\mathbf{w}} F(\mathbf{w}) - \nabla_{\mathbf{w}} F(\bar{\mathbf{w}})\|_{\mathbf{M}^{-1}}.$$

674 Therefore, the \mathbf{M} -Lipschitz continuity of the gradient for F is equivalent to the Lipschitz continuity of the
 675 gradient for f , which is equivalent to that $\nabla_{\mathbf{z}}^2 f(\mathbf{z})$, i.e. $\mathbf{P}^\top \nabla_{\mathbf{w}}^2 F(\mathbf{w})\mathbf{P}$, has eigenvalues bounded above by \hat{L} .
 676 Similarly, the statement on \mathbf{M} -strong convexity follows from

$$F(\mathbf{w}) + \nabla F(\mathbf{w})^\top (\bar{\mathbf{w}} - \mathbf{w}) + \frac{1}{2} \hat{c} \|\bar{\mathbf{w}} - \mathbf{w}\|_{\mathbf{M}}^2 = f(\mathbf{z}) + \nabla_{\mathbf{z}} f(\mathbf{z})^\top (\bar{\mathbf{z}} - \mathbf{z}) + \frac{1}{2} \hat{c} \|\bar{\mathbf{z}} - \mathbf{z}\|_2^2.$$

677

□

678 We may assume \hat{L} and \hat{c} are respectively the maximum and the minimum of the eigenvalues of $\mathbf{P}^\top \nabla^2 F(\mathbf{w})\mathbf{P}$
 679 for all \mathbf{w} . So $\frac{\hat{L}}{\hat{c}}$ plays the role of the condition number of the preconditioned matrix $\mathbf{P}^\top \nabla^2 F(\mathbf{w})\mathbf{P}$. If we
 680 assume $\mathbf{M}^{-1} = \mathbf{P}\mathbf{P}^\top$ is such that $\frac{\hat{L}}{\hat{c}}$ is smaller than $\frac{L}{c}$, it basically reduces the condition number. We will
 681 demonstrate that this accelerates the speed of convergence.

682 An important lemma comes directly from this assumption.

683 **Lemma C.4.** *Under the assumption of \mathbf{M} -Lipschitz continuity of gradient,*

$$F(\mathbf{w}) \leq F(\bar{\mathbf{w}}) + \nabla F(\bar{\mathbf{w}})^\top (\mathbf{w} - \bar{\mathbf{w}}) + \frac{1}{2} \hat{L} \|\mathbf{w} - \bar{\mathbf{w}}\|_{\mathbf{M}}^2 \quad (15)$$

Proof. Consider the following,

$$\begin{aligned} F(\mathbf{w}) &= F(\bar{\mathbf{w}}) + \int_0^1 (\nabla F(\bar{\mathbf{w}} + t(\mathbf{w} - \bar{\mathbf{w}})))^\top \mathbf{P}\mathbf{P}^{-1}(\mathbf{w} - \bar{\mathbf{w}}) dt \\ &= F(\bar{\mathbf{w}}) + \nabla F(\bar{\mathbf{w}})^\top (\mathbf{w} - \bar{\mathbf{w}}) + \int_0^1 (\nabla F(\bar{\mathbf{w}} + t(\mathbf{w} - \bar{\mathbf{w}})) - \nabla F(\bar{\mathbf{w}}))^\top \mathbf{P}\mathbf{P}^{-1}(\mathbf{w} - \bar{\mathbf{w}}) dt \\ &\leq F(\bar{\mathbf{w}}) + \nabla F(\bar{\mathbf{w}})^\top (\mathbf{w} - \bar{\mathbf{w}}) + \int_0^1 \hat{L} \|t(\mathbf{w} - \bar{\mathbf{w}})\|_{\mathbf{M}} \|\mathbf{w} - \bar{\mathbf{w}}\|_{\mathbf{M}} dt \end{aligned}$$

684 which gives us our consequence that was to be shown. □

685 Notice that combining the variance definition (Eq. 3) with Assumption 3, we have the following

$$\mathbb{E}_{\boldsymbol{\xi}_k} [\|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2] \leq K_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + K \text{ with } K_G := K_V + \mu_G^2 \geq \mu^2 > 0 \quad (16)$$

686 The proof for the two theorems relies on the following lemmas.

687 **Lemma C.5.** *Under Assumption 1, the iterates of Eq. 2 satisfy the following inequality for all $k \in \mathbb{N}$:*

$$\mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\alpha_k \nabla F(\mathbf{w}_k)^\top \mathbb{E}_{\boldsymbol{\xi}_k} [g(\mathbf{w}_k, \boldsymbol{\xi}_k)] + \frac{1}{2} \alpha_k^2 \hat{L} \mathbb{E}_{\boldsymbol{\xi}_k} [\|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2] \quad (17)$$

688 *Proof.* Let $\mathbf{w} = \mathbf{w}_{k+1}$ and $\bar{\mathbf{w}} = \mathbf{w}_k$. Then, by Assumption 1,

$$F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k) \leq \nabla F(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{1}{2} \hat{L} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_{\mathbf{M}}^2$$

Recalling that Eq. 2 gives $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k)$, we then have,

$$\begin{aligned} F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k) &\leq \nabla F(\mathbf{w}_k)^\top (-\alpha_k \mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k)) + \frac{1}{2} \hat{L} \|-\alpha_k \mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}}^2 \\ &\leq -\alpha_k \nabla F(\mathbf{w}_k)^\top \mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k) + \frac{1}{2} \alpha_k^2 \hat{L} \|\mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}}^2 \\ &\leq -\alpha_k \nabla F(\mathbf{w}_k)^\top \mathbf{M}^{-1} g(\mathbf{w}_k, \boldsymbol{\xi}_k) + \frac{1}{2} \alpha_k^2 \hat{L} \|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2 \end{aligned}$$

689 Take the expectation of both sides

$$\mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\alpha_k \nabla F(\mathbf{w}_k)^\top \mathbf{M}^{-1} \mathbb{E}_{\boldsymbol{\xi}_k} [g(\mathbf{w}_k, \boldsymbol{\xi}_k)] + \frac{1}{2} \alpha_k^2 \hat{L} \mathbb{E}_{\boldsymbol{\xi}_k} [\|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2]$$

690 Thus, the desired result is achieved. \square

Lemma C.6. *Under Assumptions 1 and 2, the iterates of Eq. 2 satisfy the following inequalities for all $k \in \mathbb{N}$:*

$$\mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\mu \alpha_k \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{1}{2} \alpha_k^2 \hat{L} \mathbb{E}_{\boldsymbol{\xi}_k} [\|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2] \quad (18)$$

$$\leq -\left(\mu - \frac{1}{2} \alpha_k \hat{L} K_G\right) \alpha_k \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{1}{2} \alpha_k^2 \hat{L} K \quad (19)$$

Proof. By Lemma C.5 and Assumption 2, it follows that

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) &\leq -\alpha_k \mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{1}{2} \alpha_k^2 \hat{L} \mathbb{E}_{\boldsymbol{\xi}_k} [\|g(\mathbf{w}_k, \boldsymbol{\xi}_k)\|_{\mathbf{M}^{-1}}^2] \\ &\leq -\alpha_k \mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{1}{2} \alpha_k^2 \hat{L} (K_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + K) \\ &\leq -\left(\mu - \frac{1}{2} \alpha_k \hat{L} K_G\right) \alpha_k \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{1}{2} \alpha_k^2 \hat{L} K \end{aligned}$$

691 Hence, we have the desired inequalities. \square

692 The final lemma necessary is as follows.

693 **Lemma C.7.** *Under assumptions 1, 2, and 3 (with F_* being the minimum of F), suppose Eq. 2 is run with*
694 *a learning rate sequence such that for all $k \in \mathbb{N}$, assume $\alpha_k \leq \frac{\mu}{\hat{L} K_G}$. (Note that α_k could be constant for all*
695 *$k \in \mathbb{N}$). Then the following inequality holds*

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F_*] \leq (1 - \alpha_k \hat{c} \mu) \mathbb{E}[F(\mathbf{w}_k) - F_*] + \frac{1}{2} \alpha_k^2 \hat{L} K \quad (20)$$

Proof. Given the assumptions and using Lemma C.6, we have $\mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) \leq -\hat{c} \alpha_k \mu (F(\mathbf{w}_k) - F_*) + \frac{1}{2} \alpha_k^2 \hat{L} K$. Subtract F_* from both sides and take the total expectation. We denote this total expectation as $\mathbb{E}[\cdot]$, which represents the expected value taken with respect to all random variables. That is, $\mathbb{E}[F(\mathbf{w}_k)] = \mathbb{E}_{\boldsymbol{\xi}_1} \mathbb{E}_{\boldsymbol{\xi}_2} \dots \mathbb{E}_{\boldsymbol{\xi}_{k-1}} [F(\mathbf{w}_k)]$.

$$\begin{aligned} \mathbb{E}[\mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) - F_*] &\leq \mathbb{E} \left[-\hat{c} \alpha_k \mu (F(\mathbf{w}_k) - F_*) + \frac{1}{2} \alpha_k^2 \hat{L} K - F_* \right] \\ \mathbb{E}[\mathbb{E}_{\boldsymbol{\xi}_k} [F(\mathbf{w}_{k+1})] - F_*] &\leq \mathbb{E} [-\hat{c} \alpha_k \mu (F(\mathbf{w}_k) - F_*) - F(\mathbf{w}_k) - F_*] + \frac{1}{2} \alpha_k^2 \hat{L} K \\ &\leq \mathbb{E} [-\hat{c} \alpha_k \mu F(\mathbf{w}_k) + \hat{c} \alpha_k \mu F_* + F(\mathbf{w}_k) - F_*] + \frac{1}{2} \alpha_k^2 \hat{L} K \\ &\leq (1 - \hat{c} \alpha_k \mu) \mathbb{E}[F(\mathbf{w}_k) - F_*] + \frac{1}{2} \alpha_k^2 \hat{L} K \end{aligned}$$

696 which is our desired inequality (20). □

697 C.2 Proofs of main theorems

698 C.2.1 Proof of Theorem 3.2

Proof. Using Lemma C.6, we have for all $k \in \mathbb{N}$:

$$\begin{aligned}
\mathbb{E}_{\xi_k}[F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) &\leq -\left(\mu - \frac{1}{2}\bar{\alpha}\hat{L}K_G\right)\bar{\alpha}\|\nabla F(\mathbf{w}_k)\|_{\mathbb{M}^{-1}}^2 + \frac{1}{2}\bar{\alpha}^2\hat{L}K \\
&\leq -\left(\mu - \frac{1}{2}\left(\frac{\mu}{\hat{L}K_G}\right)\hat{L}K_G\right)\bar{\alpha}\|\nabla F(\mathbf{w}_k)\|_{\mathbb{M}^{-1}}^2 + \frac{1}{2}\bar{\alpha}^2\hat{L}K \\
&= -\frac{1}{2}\bar{\alpha}\mu\|\nabla F(\mathbf{w}_k)\|_{\mathbb{M}^{-1}}^2 + \frac{1}{2}\bar{\alpha}^2\hat{L}K \\
&\leq -\frac{1}{2}\bar{\alpha}\mu[2\hat{c}(F(\mathbf{w}_k) - F(\mathbf{w}_*))] + \frac{1}{2}\bar{\alpha}^2\hat{L}K \\
&\leq -\bar{\alpha}\hat{c}\mu(F(\mathbf{w}_k) - F_*) + \frac{1}{2}\bar{\alpha}^2\hat{L}K
\end{aligned}$$

Now, subtract the constant $\frac{\bar{\alpha}\hat{L}K}{2\hat{c}\mu}$ from both sides of inequality (Eq. 20)

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F_*] - \frac{\bar{\alpha}\hat{L}K}{2\hat{c}\mu} \leq (1 - \bar{\alpha}\hat{c}\mu)\mathbb{E}[F(\mathbf{w}_k) - F_*] + \frac{1}{2}\bar{\alpha}\hat{L}K - \frac{\bar{\alpha}\hat{L}K}{2\hat{c}\mu} \quad (21)$$

$$= (1 - \bar{\alpha}\hat{c}\mu) \left(\mathbb{E}[F(\mathbf{w}_k) - F_*] - \frac{\bar{\alpha}\hat{L}K}{2\hat{c}\mu} \right) \quad (22)$$

699 We must now notice the following chain of inequalities.

$$0 < \bar{\alpha}\hat{c}\mu \leq \frac{\hat{c}\mu^2}{\hat{L}K_G}$$

700 This inequality holds by the theorem assumption that $0 < \bar{\alpha} \leq \frac{\mu}{\hat{L}K_G}$.

$$\frac{\hat{c}\mu^2}{\hat{L}K_G} \leq \frac{\hat{c}\mu^2}{\hat{L}\mu^2} = \frac{\hat{c}}{\hat{L}}$$

701 This inequality holds by (16) from Assumption 3.

702 Now, note that since $\hat{c} \leq \hat{L}$, it follows that $\frac{\hat{c}}{\hat{L}} \leq 1$. The result thus follows by applying C.6 repeatedly through
703 iteration $k \in \mathbb{N}$. □

704 **Corollary C.7.1.** *If $g(\mathbf{w}_k, \xi_k)$ is an unbiased estimate of $\nabla F(w_k)$, and the variance of $g(\mathbf{w}_k, \xi_k)$ is bounded
705 by a constant K independent of $\nabla F(\mathbf{w}_k)$, Then for a fixed learning rate bounded by $\frac{K_G}{\hat{L}K_G}$, $\mathbb{E}[F(\mathbf{w}_k) - F_*]$
706 decreases to below $\frac{\bar{\alpha}\hat{L}K}{2\hat{c}\mu}$ at the rate of $\frac{\hat{c}}{\hat{L}}$.*

707 C.2.2 Proof of Theorem 3.3

Proof. Since the learning rates are diminishing and by the theorem statement, we have $\alpha_k\hat{L}K_G \leq \alpha_1\hat{L}K_G \leq \mu$
for all $k \in \mathbb{N}$. By Lemma C.6 and Assumption 3,

$$\begin{aligned}
\mathbb{E}_{\xi_k}[F(\mathbf{w}_{k+1})] - F(\mathbf{w}_k) &\leq -\left(\mu - \frac{1}{2}\alpha_k\hat{L}K_G\right)\alpha_k\|\nabla F(\mathbf{w}_k)\|_{\mathbb{M}^{-1}}^2 + \frac{1}{2}\alpha_k^2\hat{L}K \\
&\leq -\left(\mu - \frac{1}{2}\mu\right)\alpha_k\|\nabla F(\mathbf{w}_k)\|_{\mathbb{M}^{-1}}^2 + \frac{1}{2}\alpha_k^2\hat{L}K \\
&\leq -\alpha_k\mu\hat{c}(F(\mathbf{w}_k) - F_*) + \frac{1}{2}\alpha_k^2\hat{L}K
\end{aligned}$$

By Lemma C.7, using (20), we have

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F_*] \leq (1 - \alpha_k \hat{c}\mu) \mathbb{E}[F(\mathbf{w}_k) - F_*] + \frac{1}{2} \alpha_k^2 \hat{L}K$$

708 Now, we prove the convergence result via induction. Consider the base case, $k = 1$.

709 Since $\nu \geq (\gamma + 1)(F(\mathbf{w}_1) - F_*)$ and $\nu \geq \frac{\beta^2 \hat{L}K}{2(\beta \hat{c}\mu - 1)}$, it follows that $\mathbb{E}[F(\mathbf{w}_1) - F_*] \leq \frac{\nu}{\gamma + 1}$.

Now, we assume that (8) holds for some $k \geq 1$. Thus

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}_{k+1}) - F_*] &\leq (1 - \alpha_k \hat{c}\mu) \mathbb{E}[F(\mathbf{w}_k) - F_*] + \frac{1}{2} \alpha_k^2 \hat{L}K \\ &\leq (1 - \alpha_k \hat{c}\mu) \frac{\nu}{\gamma + k} + \frac{1}{2} \alpha_k^2 \hat{L}K \\ &= \left(1 - \frac{\beta}{\gamma + k} \hat{c}\mu\right) \frac{\nu}{\gamma + k} + \frac{1}{2} \left(\frac{\beta}{\gamma + k}\right)^2 \hat{L}K \\ &= \left(1 - \frac{\beta \hat{c}\mu}{\tilde{k}}\right) \frac{\nu}{\tilde{k}} + \frac{\beta^2 \hat{L}K}{2\tilde{k}^2} \\ &= \left(\frac{\tilde{k} - 1}{\tilde{k}^2}\right) \nu - \left(\frac{\beta \hat{c}\mu - 1}{\tilde{k}^2}\right) \nu + \frac{\beta^2 \hat{L}K}{2\tilde{k}^2} \end{aligned}$$

710 where $\tilde{k} := \gamma + k$. Note that $\left(\frac{\beta \hat{c}\mu - 1}{\tilde{k}^2}\right) \nu - \frac{\beta^2 \hat{L}K}{2\tilde{k}^2} \geq 0$ since $\nu \geq \frac{\beta^2 \hat{L}K}{2(\beta \hat{c}\mu - 1)}$.

711 Thus,

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F_*] \leq \left(\frac{\tilde{k} - 1}{\tilde{k}^2}\right) \nu - \left(\frac{\beta \hat{c}\mu - 1}{\tilde{k}^2}\right) \nu + \frac{\beta \hat{L}K}{2\tilde{k}^2} \stackrel{\dagger}{\leq} \frac{\nu}{\tilde{k} + 1}$$

712 where (\dagger) follows since $\tilde{k}^2 \geq (\tilde{k} + 1)(\tilde{k} - 1)$. □

713 C.2.3 Proof of Lemma 3.4

714 *Proof.* Fix $k \leq T - 1$ and assume $\mathbf{w}_k \in \mathcal{N}_r$, i.e. $\text{dist}_{\mathbf{M}}(\mathbf{w}_k, \mathcal{S}) \leq r$. If $\mathbf{w}_{k+1} \notin \mathcal{N}_{r_+}$ then $\text{dist}_{\mathbf{M}}(\mathbf{w}_{k+1}, \mathcal{S}) >$
715 $r_+ = r + \Delta$. By the triangle inequality,

$$\text{dist}_{\mathbf{M}}(\mathbf{w}_{k+1}, \mathcal{S}) \leq \text{dist}_{\mathbf{M}}(\mathbf{w}_k, \mathcal{S}) + \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_{\mathbf{M}} \leq r + \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_{\mathbf{M}},$$

716 hence $\|\mathbf{w}_{k+1} - \mathbf{w}_k\|_{\mathbf{M}} > \Delta$. Using $\mathbf{w}_{k+1} - \mathbf{w}_k = -\alpha_k \mathbf{M}^{-1} g_k$ we have $\|\mathbf{w}_{k+1} - \mathbf{w}_k\|_{\mathbf{M}} = \alpha_k \|g_k\|_{\mathbf{M}^{-1}}$, so

$$\mathbb{P}(\mathbf{w}_{k+1} \notin \mathcal{N}_{r_+} \mid \mathcal{F}_k) \leq \mathbb{P}(\alpha_k \|g_k\|_{\mathbf{M}^{-1}} > \Delta \mid \mathcal{F}_k).$$

717 Markov's inequality and Assumption 8 yield

$$\mathbb{P}(\alpha_k \|g_k\|_{\mathbf{M}^{-1}} > \Delta \mid \mathcal{F}_k) \leq \frac{\alpha_k^2 \mathbb{E}[\|g_k\|_{\mathbf{M}^{-1}}^2 \mid \mathcal{F}_k]}{\Delta^2} \leq \delta_k.$$

718 □

719 C.2.4 Proof of Theorem 3.5

720 *Proof.* Fix $\alpha_k = \bar{\alpha}$ and let $\mathcal{F}_k := \sigma(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{k-1})$. Write $g_k := g(\mathbf{w}_k, \boldsymbol{\xi}_k)$ and define

$$\tau := \inf\{k \geq 1 : \mathbf{w}_k \notin \mathcal{N}_r\}, \quad \Omega_T := \{\tau > T\}.$$

721 Fix $k \leq T - 1$ and work on Ω_T . Then $\mathbf{w}_k, \mathbf{w}_{k+1} \in \mathcal{N}_r \subset \mathcal{N}_{r_+} \subset \mathcal{V}$. By convexity of \mathcal{V} , the segment
722 $[\mathbf{w}_k, \mathbf{w}_{k+1}] \subset \mathcal{V}$, and by Assumption 5 (local \mathbf{M} -smoothness),

$$F(\mathbf{w}_{k+1}) \leq F(\mathbf{w}_k) - \bar{\alpha} \nabla F(\mathbf{w}_k)^\top \mathbf{M}^{-1} g_k + \frac{\hat{L}}{2} \bar{\alpha}^2 \|g_k\|_{\mathbf{M}^{-1}}^2 \quad \text{on } \Omega_T. \quad (23)$$

Taking conditional expectation given $(\mathcal{F}_k, \Omega_T)$ and using the conditional-moment version of Assumption 6 on Ω_T yields

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}_{k+1}) - F_* | \mathcal{F}_k, \Omega_T] &\leq (F(\mathbf{w}_k) - F_*) - \bar{\alpha} \mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 \\ &\quad + \frac{\hat{L}}{2} \bar{\alpha}^2 \left(K_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + K \right). \end{aligned}$$

723 Using $\bar{\alpha} \leq \mu / (\hat{L} K_G)$ gives $\bar{\alpha} \mu - \frac{\hat{L}}{2} \bar{\alpha}^2 K_G \geq \frac{\mu}{2} \bar{\alpha}$, hence

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F_* | \mathcal{F}_k, \Omega_T] \leq (F(\mathbf{w}_k) - F_*) - \frac{\mu}{2} \bar{\alpha} \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{\hat{L}}{2} \bar{\alpha}^2 K. \quad (24)$$

724 On Ω_T we have $\mathbf{w}_k \in \mathcal{N}_r$, so Assumption 4 implies $\|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 \geq 2\hat{\mu}_{\text{PL}}(F(\mathbf{w}_k) - F_*)$. Substituting into
725 (24) gives

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F_* | \mathcal{F}_k, \Omega_T] \leq (1 - \rho)(F(\mathbf{w}_k) - F_*) + \rho C,$$

726 with $\rho := \bar{\alpha} \hat{\mu}_{\text{PL}} \mu \in (0, 1)$ and $C := \frac{\bar{\alpha} \hat{L} K}{2\hat{\mu}_{\text{PL}} \mu}$. Taking expectations under $\mathbb{P}(\cdot | \Omega_T)$ and defining $x_k :=$
727 $\mathbb{E}[F(\mathbf{w}_k) - F_* | \Omega_T]$ yields for $k \leq T - 1$,

$$x_{k+1} \leq (1 - \rho)x_k + \rho C.$$

728 Iterating gives, for all $1 \leq k \leq T$,

$$x_k \leq C + (1 - \rho)^{k-1}(F(\mathbf{w}_1) - F_* - C),$$

729 which is the desired conditional geometric bound.

730 Define overshoot events

$$A_k := \{\mathbf{w}_k \in \mathcal{N}_r, \mathbf{w}_{k+1} \notin \mathcal{N}_{r+}\}, \quad k = 1, \dots, T - 1,$$

731 and the no-overshoot event $\mathcal{E}_T := \bigcap_{k=1}^{T-1} A_k^c$. By Lemma 3.4, $\mathbb{P}(A_k) \leq \delta_k$, hence by the union bound

$$\mathbb{P}(\mathcal{E}_T^c) \leq \sum_{k=1}^{T-1} \delta_k. \quad (25)$$

732 Let $\sigma := \tau \wedge T$. On $\mathcal{E}_T \cap \{\tau \leq T\}$ we have $\mathbf{w}_\tau \in \mathcal{N}_{r+} \setminus \mathcal{N}_r$, hence by Assumption 7,

$$F(\mathbf{w}_\tau) - F_* \geq B := \frac{\alpha_{\text{QG}}}{2} r^2.$$

733 Since $\mathbf{w}_\sigma = \mathbf{w}_\tau$ on $\{\tau \leq T\}$,

$$B \mathbf{1}_{\{\tau \leq T\}} \mathbf{1}_{\mathcal{E}_T} \leq (F(\mathbf{w}_\sigma) - F_*) \mathbf{1}_{\mathcal{E}_T}.$$

734 Taking expectations gives

$$B \mathbb{P}(\tau \leq T, \mathcal{E}_T) \leq \mathbb{E}[(F(\mathbf{w}_\sigma) - F_*) \mathbf{1}_{\mathcal{E}_T}]. \quad (26)$$

735 We upper bound the RHS of (26). For each $k = 1, \dots, T - 1$, define the prefix no-overshoot event

$$\mathcal{E}_{k+1} := \bigcap_{j=1}^k A_j^c,$$

736 so that $\mathcal{E}_{k+1} \in \mathcal{F}_{k+1}$ and $\mathcal{E}_T \subseteq \mathcal{E}_{k+1}$. On $\mathcal{E}_{k+1} \cap \{k < \tau\}$ we have $\mathbf{w}_k \in \mathcal{N}_r$ and $\mathbf{w}_{k+1} \in \mathcal{N}_{r+} \subset \mathcal{V}$, so by
737 smoothness,

$$F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k) \leq -\bar{\alpha} \nabla F(\mathbf{w}_k)^\top \mathbf{M}^{-1} g_k + \frac{\hat{L}}{2} \bar{\alpha}^2 \|g_k\|_{\mathbf{M}^{-1}}^2 \quad \text{on } \mathcal{E}_{k+1} \cap \{k < \tau\}.$$

738 Taking conditional expectation given \mathcal{F}_k and using Assumption 6 (valid on $\{k < \tau\}$ since then $\mathbf{w}_k \in \mathcal{N}_r$)
739 yields

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k) | \mathcal{F}_k] \leq -\bar{\alpha}\mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{\hat{L}}{2}\bar{\alpha}^2(K_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + K) \leq \frac{\hat{L}}{2}\bar{\alpha}^2 K,$$

740 where the last inequality uses that the first term is nonpositive and we drop it.

Now note that $F(\mathbf{w}_\sigma) - F(\mathbf{w}_1) = \sum_{k=1}^{T-1} (F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k)) \mathbf{1}_{\{k < \tau\}}$ and that on \mathcal{E}_T we have $\mathcal{E}_T \subseteq \mathcal{E}_{k+1}$, hence the above bound applies on $\mathcal{E}_T \cap \{k < \tau\}$ for every $k \leq T-1$. Therefore,

$$\begin{aligned} \mathbb{E}[(F(\mathbf{w}_\sigma) - F(\mathbf{w}_1)) \mathbf{1}_{\mathcal{E}_T}] &= \sum_{k=1}^{T-1} \mathbb{E}[(F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k)) \mathbf{1}_{\mathcal{E}_T} \mathbf{1}_{\{k < \tau\}}] \\ &= \sum_{k=1}^{T-1} \mathbb{E}[\mathbf{1}_{\mathcal{E}_T} \mathbf{1}_{\{k < \tau\}} \mathbb{E}[F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k) | \mathcal{F}_k]] \\ &\leq \sum_{k=1}^{T-1} \frac{\hat{L}}{2} \bar{\alpha}^2 K = \frac{\hat{L}}{2} \bar{\alpha}^2 K (T-1), \end{aligned}$$

741 which implies

$$\mathbb{E}[(F(\mathbf{w}_\sigma) - F_*) \mathbf{1}_{\mathcal{E}_T}] \leq (F(\mathbf{w}_1) - F_*) + \frac{\hat{L}}{2} \bar{\alpha}^2 K (T-1). \quad (27)$$

742 Combining (26) and (27) yields

$$\mathbb{P}(\tau \leq T, \mathcal{E}_T) \leq \frac{F(\mathbf{w}_1) - F_* + \frac{\hat{L}}{2} \bar{\alpha}^2 K (T-1)}{B}.$$

743 Finally, using (25),

$$\mathbb{P}(\tau \leq T) \leq \mathbb{P}(\tau \leq T, \mathcal{E}_T) + \mathbb{P}(\mathcal{E}_T^c) \leq \frac{F(\mathbf{w}_1) - F_* + \frac{\hat{L}}{2} \bar{\alpha}^2 K (T-1)}{B} + \sum_{k=1}^{T-1} \delta_k,$$

744 and rearranging gives the stated lower bound on $\mathbb{P}(\tau > T)$ (with truncation at 0). \square

745 C.2.5 Proof of Theorem 3.6

746 *Proof.* Let $\mathcal{F}_k := \sigma(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{k-1})$, set $\alpha_k = \beta/(\gamma + k)$, and write $g_k := g(\mathbf{w}_k, \boldsymbol{\xi}_k)$. Define $\tau := \inf\{k \geq 1 : \mathbf{w}_k \notin \mathcal{N}_r\}$, $\Omega_T := \{\tau > T\}$, and $S_k := F(\mathbf{w}_k) - F_*$.

748 Fix $k \leq T-1$ and work on Ω_T . Then $\mathbf{w}_k, \mathbf{w}_{k+1} \in \mathcal{N}_r \subset \mathcal{N}_{r_+} \subset \mathcal{V}$. Since \mathcal{V} is convex, $[\mathbf{w}_k, \mathbf{w}_{k+1}] \subset \mathcal{V}$ and
749 Assumption 5 implies the \mathbf{M} -smoothness inequality:

$$F(\mathbf{w}_{k+1}) \leq F(\mathbf{w}_k) - \alpha_k \nabla F(\mathbf{w}_k)^\top \mathbf{M}^{-1} g_k + \frac{\hat{L}}{2} \alpha_k^2 \|g_k\|_{\mathbf{M}^{-1}}^2 \quad \text{on } \Omega_T.$$

Take conditional expectation given \mathcal{F}_k and using Assumption 6 (valid on $\{k < \tau\}$ since then $\mathbf{w}_k \in \mathcal{N}_r$) yields:

$$\mathbb{E}[S_{k+1} | \mathcal{F}_k, \Omega_T] \leq S_k - \alpha_k \mu \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{\hat{L}}{2} \alpha_k^2 (K_G \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + K).$$

750 Because $\alpha_k \leq \alpha_1 = \beta/(\gamma + 1) \leq \mu/(\hat{L}K_G)$, we have $\mu\alpha_k - \frac{\hat{L}}{2}\alpha_k^2 K_G \geq \frac{\mu}{2}\alpha_k$, hence

$$\mathbb{E}[S_{k+1} | \mathcal{F}_k, \Omega_T] \leq S_k - \frac{\mu}{2} \alpha_k \|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 + \frac{\hat{L}}{2} \alpha_k^2 K.$$

751 On Ω_T we have $\mathbf{w}_k \in \mathcal{N}_r$, so Assumption 4 yields $\|\nabla F(\mathbf{w}_k)\|_{\mathbf{M}^{-1}}^2 \geq 2\hat{\mu}_{\text{PL}} S_k$. Therefore, with $m := \mu\hat{\mu}_{\text{PL}}$ and
752 $c := \hat{L}K/2$,

$$\mathbb{E}[S_{k+1} | \mathcal{F}_k, \Omega_T] \leq (1 - m\alpha_k) S_k + c\alpha_k^2.$$

753 Now take expectation under $\mathbb{P}(\cdot \mid \Omega_T)$ and define $x_k := \mathbb{E}[S_k \mid \Omega_T]$. Then for all $k \leq T-1$,

$$x_{k+1} \leq (1 - m\alpha_k)x_k + c\alpha_k^2.$$

754 Substituting $\alpha_k = \beta/(\gamma + k)$ gives

$$x_{k+1} \leq \left(1 - \frac{a}{\gamma + k}\right)x_k + \frac{b}{(\gamma + k)^2}, \quad a := \beta m, \quad b := c\beta^2.$$

755 Since $\beta > 2/(\hat{\mu}_{\text{PL}}\mu)$, we have $a > 1$. Let

$$\nu := \max\left\{\frac{b}{a-1}, (\gamma+1)x_1\right\}, \quad x_1 = F(\mathbf{w}_1) - F_*.$$

756 We prove by induction that $x_k \leq \nu/(\gamma + k)$ for $1 \leq k \leq T$. The base case holds because $x_1 \leq \nu/(\gamma + 1)$ by
757 definition of ν . Assuming $x_k \leq \nu/(\gamma + k)$, we obtain

$$x_{k+1} \leq \left(1 - \frac{a}{\gamma + k}\right)\frac{\nu}{\gamma + k} + \frac{b}{(\gamma + k)^2} = \frac{\nu}{\gamma + k} + \frac{b - a\nu}{(\gamma + k)^2}.$$

758 Using $\nu \geq b/(a-1)$ implies $b - a\nu \leq -\nu$, hence

$$x_{k+1} \leq \frac{\nu}{\gamma + k} - \frac{\nu}{(\gamma + k)^2} \leq \frac{\nu}{\gamma + k} - \frac{\nu}{(\gamma + k)(\gamma + k + 1)} = \frac{\nu}{\gamma + k + 1}.$$

759 Thus $x_k \leq \nu/(\gamma + k)$ for all $1 \leq k \leq T$, i.e.

$$\mathbb{E}[F(\mathbf{w}_k) - F_* \mid \Omega_T] \leq \frac{\nu}{\gamma + k}, \quad 1 \leq k \leq T.$$

760 Define overshoot events $A_k := \{\mathbf{w}_k \in \mathcal{N}_r, \mathbf{w}_{k+1} \notin \mathcal{N}_{r+}\}$ for $k = 1, \dots, T-1$ and $\mathcal{E}_T := \bigcap_{k=1}^{T-1} A_k^c$. By
761 Lemma 3.4, $\mathbb{P}(A_k) \leq \delta_k$, hence

$$\mathbb{P}(\mathcal{E}_T^c) \leq \sum_{k=1}^{T-1} \delta_k.$$

762 Let $\sigma := \tau \wedge T$. On $\mathcal{E}_T \cap \{\tau \leq T\}$ we have $\mathbf{w}_\tau \in \mathcal{N}_{r+} \setminus \mathcal{N}_r$, so Assumption 7 yields

$$F(\mathbf{w}_\tau) - F_* \geq B := \frac{\alpha_{\text{QG}}}{2} r^2.$$

763 Since $\mathbf{w}_\sigma = \mathbf{w}_\tau$ on $\{\tau \leq T\}$, it follows that

$$B \mathbf{1}_{\{\tau \leq T\}} \mathbf{1}_{\mathcal{E}_T} \leq (F(\mathbf{w}_\sigma) - F_*) \mathbf{1}_{\mathcal{E}_T}.$$

764 Taking expectations gives

$$B \mathbb{P}(\tau \leq T, \mathcal{E}_T) \leq \mathbb{E}[(F(\mathbf{w}_\sigma) - F_*) \mathbf{1}_{\mathcal{E}_T}].$$

765 We upper bound the right-hand side by telescoping. For $k = 1, \dots, T-1$, define the prefix event $\mathcal{E}_k := \bigcap_{j=1}^{k-1} A_j^c$
766 (so $\mathcal{E}_k \in \mathcal{F}_k$ and $\mathcal{E}_T \subseteq \mathcal{E}_k$). On $\mathcal{E}_k \cap \{k < \tau\}$ we have $\mathbf{w}_k \in \mathcal{N}_r$ and $\mathbf{w}_{k+1} \in \mathcal{N}_{r+} \subset \mathcal{V}$, so the smoothness
767 inequality and Assumption 6 imply

$$\mathbb{E}[F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k) \mid \mathcal{F}_k] \leq c\alpha_k^2 \quad \text{on } \mathcal{E}_k \cap \{k < \tau\},$$

768 using again $\alpha_k \leq \mu/(\hat{L}K_G)$ to drop the (nonpositive) gradient-dependent part. Multiplying by $\mathbf{1}_{\mathcal{E}_T} \mathbf{1}_{\{k < \tau\}}$
769 and taking expectations yields

$$\mathbb{E}[(F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k)) \mathbf{1}_{\mathcal{E}_T} \mathbf{1}_{\{k < \tau\}}] \leq c\alpha_k^2.$$

770 Summing over $k = 1, \dots, T - 1$ and using $F(\mathbf{w}_\sigma) - F(\mathbf{w}_1) = \sum_{k=1}^{T-1} (F(\mathbf{w}_{k+1}) - F(\mathbf{w}_k)) \mathbf{1}_{\{k < \tau\}}$ gives

$$\mathbb{E}[(F(\mathbf{w}_\sigma) - F(\mathbf{w}_1)) \mathbf{1}_{\mathcal{E}_T}] \leq c \sum_{k=1}^{T-1} \alpha_k^2,$$

771 hence

$$\mathbb{E}[(F(\mathbf{w}_\sigma) - F_*) \mathbf{1}_{\mathcal{E}_T}] \leq (F(\mathbf{w}_1) - F_*) + c \sum_{k=1}^{T-1} \alpha_k^2.$$

772 Therefore,

$$\mathbb{P}(\tau \leq T, \mathcal{E}_T) \leq \frac{F(\mathbf{w}_1) - F_* + c \sum_{k=1}^{T-1} \alpha_k^2}{B}.$$

773 Finally,

$$\mathbb{P}(\tau > T) \geq 1 - \mathbb{P}(\tau \leq T, \mathcal{E}_T) - \mathbb{P}(\mathcal{E}_T^c) \geq 1 - \frac{F(\mathbf{w}_1) - F_* + c \sum_{k=1}^{T-1} \alpha_k^2}{B} - \sum_{k=1}^{T-1} \delta_k,$$

774 and truncation gives the $\max\{0, \cdot\}$ form. □

775 D Numerical experiments

776 D.1 Implementation details

777 The algorithms in this paper were implemented in Python using `jax` (version 0.5.0), `flax` (version 0.10.0),
 778 and `optax` (version 0.2.4). All timing results reported in Section 4 were measured on a consistent hardware
 779 platform running Ubuntu 24.04.2 LTS, equipped with an Intel(R) Core(TM) i7-12700K CPU (8 Performance-
 780 cores @ 3.60 GHz and 4 Efficient-cores @ 2.70 GHz), and 64 GB of system memory. All experiments were
 781 executed in double precision arithmetic to ensure numerical stability for the challenging SciML problems.

782 D.2 Baseline methods and experimental setting

783 Our experiments evaluated several optimization algorithms to validate our theoretical analysis of precon-
 784 ditioning effects. We implemented vanilla SGD, SGD with momentum ($\beta = 0.9$), and the preconditioned
 785 methods using GGN and Hessian approximations. The goal was to minimize the mean square error loss.
 786 For plotting purposes, we normalized the training losses so that the first loss was recorded as 1.0. The
 787 preconditioned methods employ conjugate gradient to efficiently approximate matrix-vector products with the
 788 inverse preconditioner, avoiding the prohibitive cost of explicitly forming and inverting the full matrices. This
 789 approach provides a computationally tractable way to incorporate curvature information into the optimization
 790 process. For Adam (with $\beta_1 = 0.9$, $\beta_2 = 0.999$) and L-BFGS (with memory size 100 and maximum line
 791 search of 100 steps), we utilized the implementations available in the `optax` library.

792 Our experimental protocol employed a structured two-phase optimization strategy. Phase I uses Adam to
 793 reach a comparable local basin; Phase II switches to the target optimizer to isolate late-stage behavior.
 794 Because our nonconvex theory is local, the basin reached at the end of Phase I can influence the local
 795 constants (\hat{L} , $\hat{\mu}_{\text{PL}}$, K) encountered in Phase II and hence may affect which optimizer performs best after the
 796 switch. We therefore use the same Adam warm start, switch point, architecture, and seed protocol across
 797 all methods to control for basin selection and interpret the Phase II results as comparisons conditional on
 798 entering a comparable basin rather than fully basin-agnostic rankings. This established a common starting
 799 point in the optimization landscape and helped navigate past initial high-gradient regions. In Phase II, we
 800 transitioned to the respective optimization methods for direct performance comparison. The specific duration
 801 of each phase varied by task complexity and is detailed in the respective experimental sections.

802 We individually optimized learning rates for each method-task combination through grid search, deliberately
 803 omitting learning rate schedulers to isolate the inherent convergence properties of each optimizer. For
 804 Adam, we searched within the range $\{0.001, 0.0005, 0.0002, 0.0001, \dots, 0.00001\}$. The preconditioned methods
 805 required different learning rate ranges due to their curvature properties: CG-Hessian and CG-GGN used

806 $\{1.0, 0.5, \dots, 0.001\}$. This difference reflects our theoretical analysis that effective preconditioning can support
 807 larger learning rates when operating near local minima. For vanilla SGD and momentum SGD, we initially
 808 explored the same ranges as Adam and expanded to wider intervals when necessary to ensure optimal
 809 performance. This methodology ensured a fair comparison by allowing each optimizer to operate at its most
 810 effective learning rate for each specific task.

811 To ensure robust experimental results, we conducted each experiment five times using different random
 812 seeds (42 to 46 for Phase I and 43 to 47 for Phase II). This approach accounts for the inherent stochasticity
 813 in neural network training processes and allows us to report mean performance metrics. For our timing
 814 analysis, we implemented a precise measurement protocol that isolates the computational efficiency of the
 815 optimization methods themselves. Specifically, we excluded all data generation and preprocessing overhead,
 816 capturing only the cumulative duration of the actual training iterations on identical hardware configurations.
 817 This methodology provides an equitable assessment of computational efficiency, particularly important
 818 when comparing methods with substantially different per-iteration costs, such as first-order methods versus
 819 preconditioned approaches that require conjugate gradient iterations.

820 D.3 Noisy data regression

821 For the Franke function regression experiment, we used a neural network with two hidden layers of 50 neurons
 822 each and ReLU activation functions. We resampled the dataset every epoch, generating 256 points with
 823 additive Gaussian noise as described in Section 4.2 and illustrated in the left panel of Figure 6. For the
 824 preconditioned methods, we employed 5 conjugate gradient iterations. The right panel of Figure 6 extends
 825 our main results by displaying not only the mean performance across 5 independent runs but also the variance
 826 bands for each optimization method.

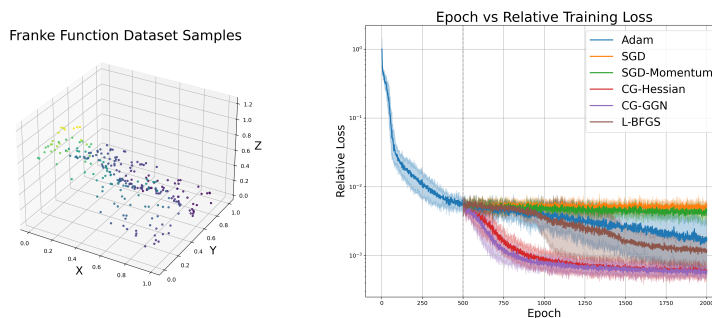


Figure 6: Left: Visualization of the Franke function dataset sampling. Right: Franke function regression performance averaged over 5 independent runs. Left: Training loss versus epochs with Phase I transitioning to Phase II at epoch 500 with variance.

827 D.4 Physics-informed neural networks

828 For solving the Poisson equation with PINNs, we used a neural network with two hidden layers of 50 neurons
 829 each and tanh activation functions. We resampled the dataset every epoch, generating 1,000 points within
 830 the domain and 200 points on the boundary, as described in Appendix 4.2 and illustrated in the left panel of
 831 Figure 7. For the preconditioned methods, we employed 20 conjugate gradient iterations. The right panel of
 832 Figure 7 shows that the mean loss trajectory is accompanied by a tight variance envelope across 5 independent
 833 runs.

834 D.5 Green's function learning

835 For both cases in the Green's function experiments, we used a neural network with five hidden layers of
 836 20 neurons each and tanh activation functions. We resampled the dataset every epoch, generating 1,000
 837 points within the domain, 500 points such that x is close to y , and 200 points on the boundary. For the

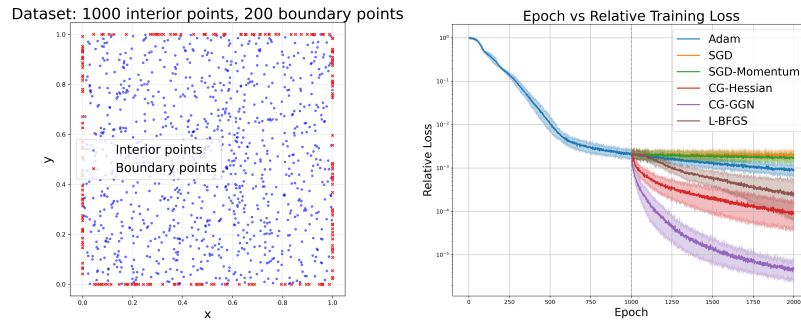


Figure 7: Left: Visualization of the sampling strategy for the 2D Poisson equation PINNs. The plot shows the distribution of 1,000 collocation points within the domain (blue) and 200 points along the boundary (red) used for enforcing the PDE and boundary conditions respectively. Right: Poisson equation PINNs performance averaged over 5 independent runs. Training loss versus epochs with Phase I transitioning to Phase II at epoch 1,000 with variance.

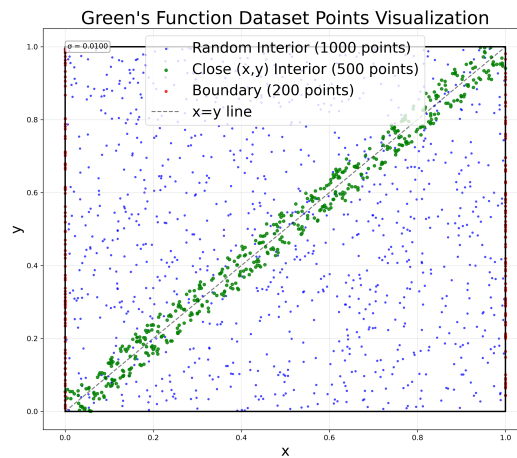


Figure 8: Visualization of the sampling strategy for Green's function learning. The plot shows three categories of training points: randomly distributed interior points (blue, 1,000 points), points concentrated near the diagonal where x is close to y (green, 500 points) to capture the near-singularity behavior characteristic of Green's functions, and boundary points (red, 200 points) used to enforce homogeneous Dirichlet boundary conditions.

838 preconditioned methods, we employed 20 conjugate gradient iterations. Figure 9 extends our main results by
839 displaying not only the mean performance across 5 independent runs but also the variance bands for each
840 optimization method.

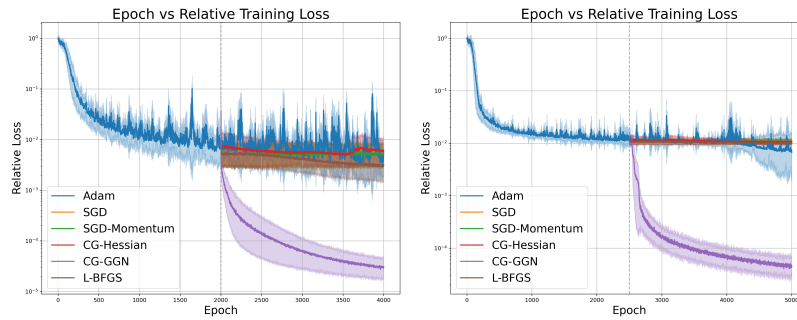


Figure 9: Green’s function learning performance averaged over 5 independent runs. Left: Training loss versus epochs with Phase I transitioning to Phase II at epoch 2,000 with variance for Laplacian. Right: Training loss versus epochs with Phase I transitioning to Phase II at epoch 2,500 with variance for convection-diffusion.