
Spectral Preconditioning for Gradient Methods on Graded Non-convex Functions

Nikita Doikov¹ Sebastian U. Stich² Martin Jaggi¹

Abstract

The performance of optimization methods is often tied to the spectrum of the objective Hessian. Yet, conventional assumptions, such as smoothness, do often not enable us to make finely-grained convergence statements—particularly not for non-convex problems. Striving for a more intricate characterization of complexity, we introduce a unique concept termed *graded non-convexity*. This allows to partition the class of non-convex problems into a nested chain of subclasses. Interestingly, many traditional non-convex objectives, including partially convex problems, matrix factorizations, and neural networks, fall within these subclasses. As a second contribution, we propose gradient methods with spectral preconditioning, which employ inexact top eigenvectors of the Hessian to address the ill-conditioning of the problem, contingent on the grade. Our analysis reveals that these new methods provide provably superior convergence rates compared to basic gradient descent on applicable problem classes, particularly when large gaps exist between the top eigenvalues of the Hessian. Our theory is validated by numerical experiments executed on multiple practical machine learning problems.

1. Introduction

Motivation. The gradient method is an important and attractive tool for solving large-scale optimization problems. It has very cheap cost of every iteration and well established convergence guarantees, that hold starting from an arbitrary initial point and for a wide family of problem classes, including convex and non-convex problems. However, the major

¹Machine Learning and Optimization Laboratory (MLO), EPFL, Lausanne, Switzerland ²CISPA Helmholtz Center for Information Security, Saarbrücken, Germany. Correspondence to: Nikita Doikov <nikita.doikov@epfl.ch>, Sebastian U. Stich <stich@cispa.de>, Martin Jaggi <martin.jaggi@epfl.ch>.

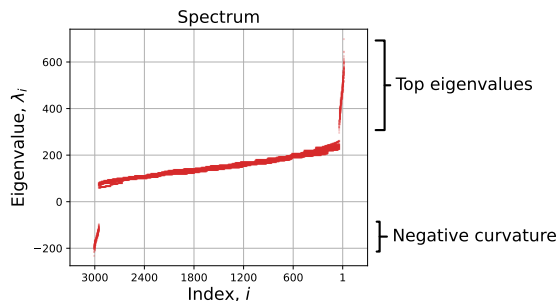


Figure 1. Spectrum of the Hessians for the matrix factorization problem (Example 2.10) of the optimization dimension $n = 3000$, for 10 random objectives.

drawback of the gradient method remains to be its slow rate of convergence: for solving modern optimization problems up to a reasonable accuracy level, it is often required to do a lot of gradient steps, due to *ill-conditioning* of the problem.

In order to improve the gradient direction, we can multiply the gradient by a specifically crafted matrix called *preconditioner*, which should adjust the method to the right geometry of the problem. However, finding a good preconditioner with strong theoretical guarantees is not easy, especially for non-convex problems. In this work, we propose a new family of *spectral preconditioners* that rely on an additional refined information about the function class. As a by-product, we establish convergence rates, that are provably better than those of the gradient methods.

Optimization theory suggests that the main complexity parameters that affect the rate of convergence for gradient methods are the spectrum of the Hessian $\nabla^2 f$ and its extremal characteristics, such as the bound for the *maximal eigenvalue* (the Lipschitz constant of the gradient) or the *condition number* (the ratio of the largest and smallest eigenvalues) (Nemirovski & Yudin, 1983; Nesterov, 2018). Moreover, some of the fundamental properties of the objective function, such as *convexity*, *weak convexity*, or *strong convexity*, that distinguish between all optimization problems and globally solvable ones, can be defined in terms of lower bounds on the spectrum. Thus (for a twice differentiable function),

$$f \text{ is convex} \Leftrightarrow \nabla^2 f \succeq \mathbf{0}. \quad (1)$$

Algorithm	Preconditioning	Non-convex Complexity	Strongly Convex	Arithmetic Cost
The Gradient Method	$\mathbf{H} := \mathbf{I}, \tau = 0$	$\mathcal{O}\left(\frac{\lambda_1}{\varepsilon^2}\right) \cdot F_0$	$\tilde{\mathcal{O}}\left(\frac{\lambda_1}{\lambda_n}\right)$	$\mathcal{O}(n)$
Spectral Preconditioning (ours)	$\mathbf{H} \approx \nabla_\tau^2 f, 1 \leq \tau < n$	$\mathcal{O}\left(\frac{\lambda_{\tau+1}}{\varepsilon^2} + \frac{L^{1/2}}{\varepsilon^{3/2}}\right) \cdot F_0$	$\tilde{\mathcal{O}}\left(\frac{\lambda_{\tau+1}}{\lambda_n} + MD\right)$	$\mathcal{O}(\tau^2 n + \tau^3)$
Newton's Method	$\mathbf{H} := \nabla^2 f, \tau = n$	$\mathcal{O}\left(\frac{L^{1/2}}{\varepsilon^{3/2}}\right) \cdot F_0$	$\tilde{\mathcal{O}}(MD)$	$\mathcal{O}(n^3)$

Table 1. Global complexity bounds for our method as compared to the classic gradient method and the regularized Newton methods. We denote by λ_i a uniform bound for the i th eigenvalue of the Hessian, sorted in a non-ascending order: $\lambda_1 \geq \dots \geq \lambda_n$. We denote by L the Lipschitz constant of the Hessian and, for convex objectives, M is the constant of quasi-self-concordance (see Section 6). $F_0 := f(\mathbf{x}_0) - f^*$, and D is the diameter of the initial sublevel set. We see that using the spectral preconditioner of order τ , we cut the top τ eigenvalues of the spectrum, which is the most significant complexity factor. We present the state-of-the-art global complexities from (Nesterov & Polyak, 2006; Doikov, 2023) for Newton's Method with cubic and gradient regularizations.

At the same time, for problems with specific structure, the worst-case guarantees for convergence of gradient methods obtained considering only the the largest and smallest eigenvalue of the Hessian can be too pessimistic. Indeed, we see that in practice, the distribution of eigenvalues can be quite specific, with a relatively small amount of top eigenvalues that are much larger than the others (Fig. 1). Consequently, any a priori information on the structure of the Hessian can be significant from the optimization perspective. Ultimately, we want to have algorithms that are able to benefit from this knowledge, achieving faster rates when the distribution of eigenvalues is far from uniform.

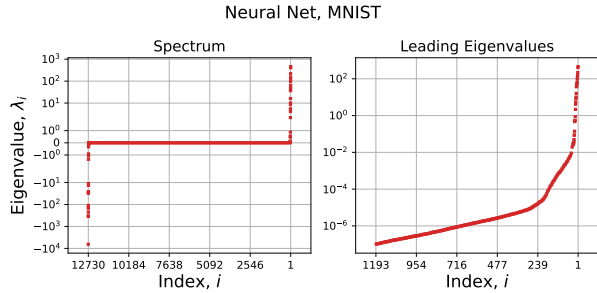


Figure 2. Left: spectrum of the Hessians for a two-layer fully connected neural network trained on MNIST dataset. Right: zoomed top eigenvalues. The total number of parameters is $n = 12730$ (the first layer: 12560, the second layer: 170). We see that the dimension τ of the subspace with positive eigenvalues is much bigger than the dimension of the last layer. However, there are only a few eigenvalues that are significantly larger than the others.

Contributions. In this work, we develop spectral preconditioning for the gradient methods that is able to tackle non-convex problems with *highly non-uniform* and clustered spectrum of the Hessian, that we often observe in practice. For that, we propose to make a step back from the common dichotomy between convex and non-convex problems. We introduce a new notion of *graded non-convex* functions, which granulate the class of all non-convex problems into nested family of subclasses. We say, for some

integer $\tau \geq 1$

$$f \text{ is non-convex of grade } \tau \Leftrightarrow \nabla_\tau^2 f \succeq \mathbf{0}, \quad (2)$$

where $\nabla_\tau^2 f$ is composed by the top τ eigenvectors of the full Hessian (see (6) for the formal definition). For $\tau = n$, where n is dimension of the problem, we obtain the entire Hessian $\nabla^2 f$ and inequality (2) means the standard convexity (1). It appears that for many practical non-convex problems, this condition is satisfied at least for some small τ . For example, for deep neural networks with convex loss it is satisfied at least for $\tau \geq d$, where d is the dimension of the last layer (however, in practice we observe that the actual values of τ can be much bigger, see Fig. 2).

Inequality (2) provides us with a certain *convex subspace* at each point (such that our function is locally convex in this subspace, see Fig. 3), which looks very attractive from the optimization standpoint. Furthermore, it contains the leading eigenvalues, which constitute the primary computational burden for first-order methods.

Based on our problem class, we propose to use a positive definite matrix $\mathbf{H} \approx \nabla_\tau^2 f(\mathbf{x})$ as a natural preconditioner for our method, that we call *spectral preconditioning*:

$$\mathbf{x}^+ = \mathbf{x} - (\mathbf{H} + \alpha \mathbf{I})^{-1} \nabla f(\mathbf{x}), \quad (3)$$

where and $\alpha \geq 0$ is a regularization parameter. For these iterations, we establish strong convergence guarantees, that improve with increasing the parameter τ (see Table 1).

Therefore, the method with any $\tau \geq 1$ works *provably better* than the basic gradient method ($\tau = 0$) in terms of the rate of convergence. For $\tau = n$ (the full Hessian), our algorithm becomes the regularized Newton method with the best global complexity bounds known in the literature. However, the most efficient version of our method corresponds to the case of small $\tau \approx 10^0 - 10^2$, when estimating the matrix \mathbf{H} can be done efficiently with the hot-start *power method*.

Related Work. Preconditioning is an important tool in numerical analysis and optimization (Nocedal & Wright,

2006). The basic example is preconditioning of the conjugate gradient method (Hestenes & Stiefel, 1952) for solving a system of linear equations. The choice of the right preconditioner is a difficult task and it often depends on an a-priori knowledge on the problem structure, as, for example, the Laplacian preconditioning for the graph-induced problems (Spielman & Teng, 2004; Vaidya, 1991) or for the systems involving partial differential equations (Mardal & Winther, 2011).

In optimization, a powerful approach for preconditioning that works for general problems is called the Newton method (Polyak, 2007). It uses the Hessian matrix $\nabla^2 f$ to alleviate the impact of ill-conditioned characteristics of the Hessian spectrum. The modern versions of this method provide us with the global rates of convergence that are significantly better than those of the first-order algorithms (Nesterov & Polyak, 2006; Cartis et al., 2011a;b; Grapiglia & Nesterov, 2017; Karimireddy et al., 2018; Doikov et al., 2022). However, computing and inverting the Hessian matrix is prohibitively expensive in terms of the arithmetic operations and memory usage to store big matrices.

The common approximate technique for improving the arithmetic cost of the methods with the full Hessian is called the *quasi-Newton methods*, such as SR1, DFP, BFGS, L-BFGS and others (Nocedal & Wright, 2006). Despite these methods show an outstanding performance on practical problems of a moderate dimension, it is a significant challenge to establish a rigorous theory of convergence for *quasi-second-order methods*, which would provably benefit globally from inexact Hessian information, while the local non-asymptotic theory of convergence for the classical quasi-Newton methods has emerged only very recently (Rodomanov & Nesterov, 2021a;b; Jin & Mokhtari, 2023). Note that employing a naive and straightforward line-search or damping approach in the basic Newton method can have as slow convergence as the plain gradient descent (Cartis et al., 2013). See also (Hanzely et al., 2022; Kamzolov et al., 2023; Jiang & Mokhtari, 2023) for the recent advancement in the global analysis of the damped and quasi-Newton methods for convex objectives. Some of the modern scalable techniques for the Newton method are *block-coordinate updates*, *stochastic subspaces* and *sketching* (Hanzely et al., 2020; Nutini et al., 2022; Hanzely, 2023), *distributed* and *lazy* computations (Qian et al., 2021; Safaryan et al., 2021; Islamov et al., 2021; 2022; Doikov et al., 2023; Doikov & Grapiglia, 2023) and *stochastic preconditioning* (Pasechnyuk et al., 2022; 2023), for convex and non-convex problems.

Our analysis is based on the new concept of *graded non-convexity*, which is related to studies of other generalized notions of convexity (Vial, 1983; Hörmander, 2007). At the same time, more refined specifications of the problem class that investigate the distribution of the eigenvalues were

considered recently in (Kovalev et al., 2018; Scieur & Pedregosa, 2020; Cunha et al., 2022; Goujaud et al., 2022). The motivation of our work to cut large gaps between the leading eigenvalues is closely related to recently proposed coordinate methods with *volume sampling* (Rodomanov & Kropotov, 2020) and *polynomial preconditioning* technique (Doikov & Rodomanov, 2023), analysing the convex problems with a specific structure. Another similar and relevant work is (Kunstner et al., 2024), which proposes a backtracking procedure for constructing a provable preconditioning for convex problems.

Additionally, several recent works study efficient preconditioning of gradient methods for specific classes of *non-convex* problems using a similar idea of employing the spectrum of the objective Hessian. These include overparameterized non-convex *matrix factorization* (Zhang et al., 2021; 2023; Ma et al., 2023) and *kernel ridge regression* problems (Ma & Belkin, 2017). Compared to these works, our preconditioning is designed to solve general smooth and possibly non-convex problems without assumptions on the intrinsic structure of the objective, but rather on the distribution of the Hessian spectrum.

Contents. The rest of the paper is organized as follows. In Section 2 we introduce the notion of *graded non-convexity*. We study its main properties and provide several examples. Section 3 contains our main algorithm (14) the Gradient Method with Spectral Preconditioning. We prove fast convergence rates for this method, showing an improvement when increasing the preconditioning order (Theorem 4.2). In Section 5 we show a simple modification of our method that allows to remove the dependency on the negative part of the spectrum (Theorem 5.1). In Section 6 we show improved rates of convergence for the special case of convex functions. In Section 7 we discuss the efficient implementation of our algorithms. Section 8 presents numerical experiments. Missing proofs are provided in the appendix.

Notation. We are interested in solving the following unconstrained optimization problem¹

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad (4)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a several times differentiable target function, which can be *non-convex*. We denote its global lower bound by $f^* := \inf_{\mathbf{x}} f(\mathbf{x})$, which we assume to be finite. In non-convex optimization, our goal is to find an approximate *stationary point* $\bar{\mathbf{x}}$ to (4), ensuring $\|\nabla f(\bar{\mathbf{x}})\| \leq \varepsilon$, for some $\varepsilon > 0$. We use the standard Euclidean norm for vectors: $\|\mathbf{x}\| := \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$, $\mathbf{x} \in \mathbb{R}^n$, and the corresponding operator norm for matrices: $\|\mathbf{A}\| := \max_{\mathbf{x}: \|\mathbf{x}\| \leq 1} \|\mathbf{A}\mathbf{x}\|$.

¹Our results can be also generalized to the *composite* formulation of optimization problems, see Section C in the appendix.

We denote the gradient of f at point $\mathbf{x} \in \mathbb{R}^n$ by $\nabla f(\mathbf{x}) \in \mathbb{R}^n$, and the Hessian matrix by $\nabla^2 f(\mathbf{x}) \in \mathbb{R}^{n \times n}$. Note that the Hessian is a symmetric matrix. Hence, for any point $\mathbf{x} \in \mathbb{R}^n$ we can introduce the following spectral decomposition:

$$\nabla^2 f(\mathbf{x}) = \sum_{i=1}^n \lambda_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x})^\top, \quad (5)$$

where the eigenvalues are sorted in a *non-ascending* order: $\lambda_1(\mathbf{x}) \geq \lambda_2(\mathbf{x}) \geq \dots \geq \lambda_n(\mathbf{x})$, and $\mathbf{u}_1(\mathbf{x}), \dots, \mathbf{u}_n(\mathbf{x}) \in \mathbb{R}^n$ are the corresponding orthonormal eigenvectors. There are always several possible choices for the eigenvector basis, hence, we assume that a specific selection has been fixed. Our results remain independent of the particular selection. For a fixed spectral decomposition (5), we denote by $\nabla_\tau^2 f(\mathbf{x})$, $1 \leq \tau \leq n$, the *Hessian of spectral order* τ :

$$\nabla_\tau^2 f(\mathbf{x}) := \sum_{i=1}^{\tau} \lambda_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x}) \mathbf{u}_i(\mathbf{x})^\top \in \mathbb{R}^{n \times n}. \quad (6)$$

For convenience, we set $\nabla_0^2 f(\mathbf{x}) \equiv \mathbf{0}$. Thus, $\text{rank}(\nabla_\tau^2 f(\mathbf{x})) \leq \tau$, and for $\tau = n$ we obtain the full Hessian. Of course, the decomposition of the form (6) is not unique, especially if $\lambda_\tau(\mathbf{x}) = \lambda_{\tau+1}(\mathbf{x})$ for a certain \mathbf{x} . However, we always assume that a unique choice has been fixed for ease of notation. We denote by $\nabla^3 f(\mathbf{x})$ the third derivative of f at point \mathbf{x} , which is a trilinear symmetric form. The action of this form onto some fixed directions $\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3 \in \mathbb{R}^n$ is denoted by $\nabla^3 f(\mathbf{x})[\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3] \in \mathbb{R}$.

2. Problem Classes

A standard assumption in non-convex optimization is to assume that the objective function is smooth, i.e. that its gradients are Lipschitz continuous. For twice-differentiable functions, this is equivalent to assuming that the norm of the Hessian, $\|\nabla^2 f(\mathbf{x})\|$, is uniformly bounded. In this section, we introduce a new concept that allows us to capture the distribution of the eigenvalues in more detail.

2.1. Grade of Non-convexity

We start with a formal definition of our problem class.

Definition 2.1. For a twice continuously differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and convex set $Q \subseteq \mathbb{R}^n$, we say that f is *non-convex of grade* τ , if

$$\nabla_\tau^2 f(\mathbf{x}) \succeq \mathbf{0}, \quad \forall \mathbf{x} \in Q. \quad (7)$$

In other words, (7) means that the top τ eigenvalues of the Hessian are non-negative everywhere on Q :

$$\lambda_\tau(\mathbf{x}) \geq 0, \quad \forall \mathbf{x} \in Q. \quad (8)$$

In our analysis, we will mostly use for simplicity $Q \equiv \mathbb{R}^n$ (the whole space). However, it can only refine our problem

class if some localization of a solution $\mathbf{x}^* \in Q$ is available (see also our discussion in Section 9). Definition 2.1 implies a certain restriction on a surface structure of the objective function. In differential geometry, condition (8) on the curvatures leads to the notion of *τ -convex surface* (Gromov, 1991).

For $0 \leq \tau \leq n$, we denote by \mathcal{F}_τ the set of functions $f \in \mathcal{F}_\tau$ that satisfies (7). By our definition, $\mathcal{F}_0 = C^2(\mathbb{R}^n)$ is the set of *all* twice continuously differentiable functions² and \mathcal{F}_n consists only of convex objectives. These classes are closed under multiplications by a non-negative scalar: $f \in \mathcal{F}_\tau \Rightarrow \alpha f \in \mathcal{F}_\tau$ for any $\alpha \geq 0$. Hence, we obtain the *nested family of functional cones*:

$$\boxed{\mathcal{F}_0 \supset \mathcal{F}_1 \supset \dots \supset \mathcal{F}_{n-1} \supset \mathcal{F}_n}$$

all functions convex functions

Intuitively, functions with larger grades should be easier to minimize. At the same time, a method that works for a certain \mathcal{F}_τ can also tackle all problems from \mathcal{F}_i for $i \geq \tau$.

2.2. Main Properties

Let us study some of the most basic properties that follow from our definition. First, we have the following important *grading rule*, that equips our sets with additional structure.

Proposition 2.2. Let $f \in \mathcal{F}_i$ and $g \in \mathcal{F}_j$, for some $0 \leq i, j \leq n$ such that $i + j \geq n$. Then, it holds:

$$\begin{aligned} f + g &\in \mathcal{F}_{i+j-n}, \\ \text{smax}(f, g) &\in \mathcal{F}_{i+j-n}, \end{aligned}$$

where $\text{smax}(f, g)(\mathbf{x}) \stackrel{\text{def}}{=} \ln(e^{f(\mathbf{x})} + e^{g(\mathbf{x})})$ is the *soft maximum of two functions*.

In particular, the summation of a function $f \in \mathcal{F}_\tau$ with *any convex* function $g \in \mathcal{F}_n$ cannot decrease its grade τ . The same holds for (soft) maximum operations³. Next, we observe that the grade is preserved under affine substitutions.

Proposition 2.3. Let $f \in \mathcal{F}_\tau(\mathbb{R}^n)$ be non-convex of grade τ , and let $\mathbf{A} \in \mathbb{R}^{n \times m}$, $\mathbf{b} \in \mathbb{R}^n$, with $m + \tau \geq n$. Denote $g(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b})$. Then, $g \in \mathcal{F}_{m-n+\tau}(\mathbb{R}^m)$ is non-convex of grade $m - n + \tau$.

For $\tau \geq 1$, our Definition 2.1 implies that the Hessian is not negative definite at any point: $\nabla^2 f(\mathbf{x}) \not\preceq \mathbf{0}$. This condition means that the function *cannot have strict local maxima*. This fact can be formalized as follows.

²Note that definition of our classes \mathcal{F}_τ , $0 \leq \tau \leq n$ depends on set Q and on the problem dimension n . Thus, it would be more formal to use notation $\mathcal{F}_\tau^n(Q)$. However, we omit extra indices since they should always be clear from the context.

³It holds $\max(f, g) = \lim_{\mu \rightarrow 0} \mu \cdot \text{smax}(f/\mu, g/\mu)$. We prefer to work with soft max operation to keep all functions in the smooth class.

Proposition 2.4. Let $f \in \mathcal{F}_\tau(Q)$ for $\tau \geq 1$. Then, the weak maximum principle holds: for any compact $K \subset Q$, the maximum is always achieved at the boundary,

$$\max_{\mathbf{x} \in K} f(\mathbf{x}) = \max_{\mathbf{x} \in \partial K} f(\mathbf{x}).$$

Finally, let us mention an intuitive geometric description of the surface of function f of a certain grade τ .

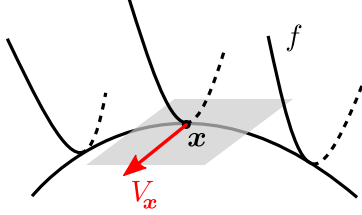


Figure 3. A surface of a non-convex function f . At point \mathbf{x} there is a subspace $V_{\mathbf{x}}$ where f is convex.

Proposition 2.5. Let for any \mathbf{x} there exists a vector subspace $V_{\mathbf{x}} \subseteq \mathbb{R}^n$ with $\dim(V_{\mathbf{x}}) \geq \tau$ such that

$$f(\mathbf{x} + \mathbf{h}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{h} \rangle, \quad \forall \mathbf{h} \in V_{\mathbf{x}}. \quad (9)$$

Then $f \in \mathcal{F}_\tau$ is non-convex of grade τ .

This inequality is stronger than our Definition 2.1: (9) is a sufficient condition for (7). Geometrically it means that for every point \mathbf{x} there exists a subspace $V_{\mathbf{x}}$ of the tangent space to the surface such that restriction of f onto this subspace is convex (see Fig. 3).

2.3. Examples

In this section, we provide several examples of non-convex objective functions of a non-trivial grade of non-convexity.

Example 2.6 (Quadratic Functions). Let $f(\mathbf{x}) = \frac{1}{2} \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle$ for some $\mathbf{A} = \mathbf{A}^\top \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$. Let the top τ eigenvalues of \mathbf{A} be positive: $\lambda_1(\mathbf{A}) \geq \dots \geq \lambda_\tau(\mathbf{A}) \geq 0$. Then $f \in \mathcal{F}_\tau$.

Consequently, adding a power of the Euclidean norm as a simple regularizer to a non-convex quadratic function as in Example 2.6, we obtain for $p > 2$ and some $\sigma > 0$:

$$f(\mathbf{x}) = \frac{1}{2} \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle + \frac{\sigma}{p} \|\mathbf{x}\|^p, \quad (10)$$

which describes a simple family of non-convex objectives that can realize all possible grades $f \in \mathcal{F}_\tau$, $0 \leq \tau \leq n$, while a global solution $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$ always exists, since f has bounded sublevel sets. The problems of the form (10) are important in applications to regularized second-order and high-order methods (Grapiglia & Nesterov, 2017; Nesterov, 2019; 2022).

Example 2.7 (Low-rank Vector Fields). Let $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ be a univariate (possibly non-convex) function, and $\mathbf{u} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a general differentiable mapping. Set

$$f(\mathbf{x}) = \varphi(\langle \mathbf{u}(\mathbf{x}), \mathbf{x} \rangle), \quad \mathbf{x} \in \mathbb{R}^n. \quad (11)$$

Assume that there exists $r \leq n - 1$ such that for any \mathbf{x} , it holds $r \geq \operatorname{rank}(\nabla \mathbf{u}(\mathbf{x}) + \nabla^2 \mathbf{u}(\mathbf{x})\mathbf{x})$. Then, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is non-convex of grade $n - r - 1$.

A particular case is when $\mathbf{u}(\mathbf{x}) \equiv \mathbf{u} \in \mathbb{R}^n$ is a constant vector field. Then, we obtain that the function (11) is non-convex of grade $n - 1$ (see Fig. 4).

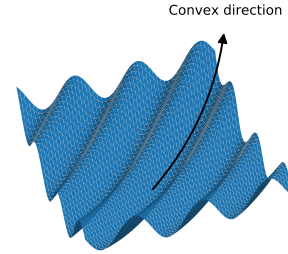


Figure 4. The graph of two-dimensional function $f(x, y) = \sin(x + y) + q(x, y)$, where q is a convex quadratic. The non-convex component has the structure of (11) with $\mathbf{u} \equiv (1, 1)^\top$.

Example 2.8 (Partial Convexity). Let $f(\mathbf{x}, \mathbf{y})$ depend on two groups of variables $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in \mathbb{R}^m$. Assume that for any fixed \mathbf{y} , the function of the first argument

$$f(\cdot, \mathbf{y}) : \mathbb{R}^n \rightarrow \mathbb{R}$$

is convex. Then $f : \mathbb{R}^{n+m} \rightarrow \mathbb{R}$ is non-convex of grade n .

We see that if the function is convex with respect to *some* of the variables, then, as a function of all variables, it is non-convex of the corresponding grade. Note, however, that the actual structure of the subspace when the Hessian is positive can be quite complicated. As a direct consequence, we obtain that matrix factorizations and deep neural network models with convex losses satisfy our assumption.

Example 2.9 (Diagonal Neural Networks). For a given $\mathbf{c} \in \mathbb{R}^n$, consider

$$f(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{x} \circ \mathbf{y} - \mathbf{c}\|^2, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad (12)$$

where \circ is the element-wise product. At every point $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{2n}$ the Hessian has the following set of $2n$ eigenvalues, a pair for each $1 \leq i \leq n$: $\lambda_i = \frac{1}{2} [x_i^2 + y_i^2 \pm \sqrt{(x_i^2 - y_i^2)^2 + 4(2x_i y_i - c_i)^2}]$, with at least one being always non-negative. Therefore, $f : \mathbb{R}^{2n} \rightarrow \mathbb{R}$ is non-convex of grade n .

The objective of the form (12) is a good model for studying the dynamics of gradient methods in deep learning (Woodworth et al., 2020; Pesme et al., 2021; Even et al., 2023; Pesme & Flammarion, 2023).

Example 2.10 (Matrix Factorizations). *For a given target matrix $\mathbf{C} \in \mathbb{R}^{n \times m}$, and some $r > 0$, consider*

$$f(\mathbf{X}, \mathbf{Y}) = \frac{1}{2} \|\mathbf{X}\mathbf{Y} - \mathbf{C}\|_F^2, \quad (13)$$

where $\mathbf{X} \in \mathbb{R}^{n \times r}$, $\mathbf{Y} \in \mathbb{R}^{r \times m}$. Then, $f : \mathbb{R}^{n \times r + r \times m} \rightarrow \mathbb{R}$ is non-convex of grade $\max\{m, n\} \times r$. More generally, the function $f(\mathbf{X}_1, \dots, \mathbf{X}_d) = \frac{1}{2} \|\mathbf{X}_1 \mathbf{X}_2 \cdots \mathbf{X}_d - \mathbf{C}\|_F^2$, $\mathbf{X}_i \in \mathbb{R}^{n_i \times m_i}$, is non-convex of grade $\max_{1 \leq i \leq d} [n_i \times m_i]$.

We observe that in many practical scenarios, it is very common to encounter objectives that have a subspace with positive eigenvalues of the Hessian. Indeed, the opposite seems rather quite rare and indicates that the target objective is *purely concave*. On the contrary, for deep neural network models with convex losses we always ensure the existence of a subspace with *positive curvature*, which serves as the main computational burden for a method to converge to a stationary point, when the problem is ill-conditioned.

In the next sections, we propose the spectral preconditioning technique for gradient methods in order to tackle ill-conditioning of the positive part of the spectrum. At the same time, as we will show, the impact of the negative curvature to optimization methods can easily be alleviated.

3. Spectral Preconditioning

In this section we present our proposed algorithm. The method aims to exploit the (possibly) convex-like structure of the function. When there is no such structure ($\tau = 0$) the method becomes equal to the standard gradient method.

At every iteration of our method, we choose a matrix $\mathbf{H} = \mathbf{H}^\top \succeq 0$ and perform the following preconditioned gradient step, for a given point $\mathbf{x} \in \mathbb{R}^n$ and gradient vector $\mathbf{g} \in \mathbb{R}^n$:

$$\text{GradStep}_{\mathbf{H}, \alpha}(\mathbf{x}, \mathbf{g}) := \mathbf{x} - (\mathbf{H} + \alpha \mathbf{I})^{-1} \mathbf{g},$$

where $\alpha \geq 0$ is some regularization parameter. Hence, the matrix \mathbf{H} plays the role of a *preconditioner*. We want to choose it as an approximation of the Hessian of a certain spectral order: $\mathbf{H} \approx \nabla_\tau^2 f(\mathbf{x})$. When $\tau = 0$, we have $\nabla_\tau^2 f(\mathbf{x}) \equiv \mathbf{0}$ and thus we do a step of plain gradient descent. In contrast, for $\tau = n$ we approximate the full Newton step. Let us present the method in algorithmic form.

Gradient Method with Spectral Preconditioning

Choose $\mathbf{x}_0 \in \mathbb{R}^n$ and $0 \leq \tau \leq n$.

For $k \geq 0$ iterate:

1. Estimate $\mathbf{H}_k \approx \nabla_\tau^2 f(\mathbf{x}_k) \in \mathbb{R}^{n \times n}$
2. Perform the gradient step, for some $\alpha_k \geq 0$:

$$\mathbf{x}_{k+1} = \text{GradStep}_{\mathbf{H}_k, \alpha_k}(\mathbf{x}_k, \nabla f(\mathbf{x}_k))$$

Up to now, we do not specify explicitly how we estimate $\nabla_\tau^2 f(\mathbf{x}_k)$. Our matrix \mathbf{H}_k should be easily computable and we aim to maintain a low rank representation,

$$\mathbf{H}_k = \sum_{i=1}^{\tau} a_{k,i} \mathbf{v}_{k,i} \mathbf{v}_{k,i}^\top, \quad (15)$$

for a set of positive numbers $(a_{k,i})_{i=1}^{\tau}$ and orthonormal vectors $(\mathbf{v}_{k,i})_{i=1}^{\tau}$, so that we can perform iterations of algorithm (14) cheaply. We discuss details on efficient implementation of every step in Section 7. To quantify the approximation errors, we denote

$$\delta_k := \|\mathbf{H}_k - \nabla_\tau^2 f(\mathbf{x}_k)\|, \quad \delta := \max_k \delta_k. \quad (16)$$

Thus, if $\delta = 0$, we use the exact Hessian of spectral order τ .

4. Global Convergence

Let us show the main convergence result for algorithm (14). We establish fast convergence rates to a stationary point of our objective (4), starting from an arbitrary initial point \mathbf{x}_0 . These rates become better when increasing the spectral order parameter τ used in our method.

We introduce the following characteristics of smoothness for our problem classes. We assume that the Hessian is Lipschitz continuous, with parameter $L \geq 0$, for all \mathbf{x}, \mathbf{y} :

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|. \quad (17)$$

To estimate the error of using the Hessian of a spectral order τ , we also introduce the following system of parameters:

$$\sigma_\tau := \sup_{\mathbf{x} \in \mathbb{R}^n} \|\nabla^2 f(\mathbf{x}) - \nabla_\tau^2 f(\mathbf{x})\|. \quad (18)$$

Thus, $\sigma_0 := \sup_{\mathbf{x}} \|\nabla^2 f(\mathbf{x})\|$ is simply the (best) Lipschitz constant of the gradient of f , and $\sigma_n = 0$. In general, the value of σ_τ is the uniform bound for both the $(\tau + 1)$ th eigenvalue of the Hessian, and its negative part, for $\tau < n$:

$$\sigma_\tau \geq \max\{\lambda_{\tau+1}(\mathbf{x}), -\lambda_n(\mathbf{x})\}, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (19)$$

Under these conditions, and using a sufficiently big regularization parameter α_k (a second-order ‘‘step size’’), we can prove the following progress for one iteration $\mathbf{x}_k \mapsto \mathbf{x}_{k+1}$ of our method.

Lemma 4.1. *Let $\alpha_k \geq \sqrt{\frac{L \|\nabla f(\mathbf{x}_k)\|}{2}} + \sigma_\tau + \delta$. Then*

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq \frac{1}{8\alpha_k} \|\nabla f(\mathbf{x}_{k+1})\|^2. \quad (20)$$

Note that we do not need to know the exact values of the parameters L , σ_τ , and δ_k . In practice, we can use an adaptive search to ensure sufficient progress (20). This lemma leads us to the basic global convergence result.

Theorem 4.2. *Let $f \in \mathcal{F}_\tau$ be non-convex of grade τ , where $0 \leq \tau \leq n$ is fixed. Let f have a Lipschitz Hessian with constant L and bounded parameter σ_τ . Consider iterations $\{\mathbf{x}_k\}_{k \geq 0}$ of algorithm (14) with*

$$\alpha_k = \sqrt{\frac{L \|\nabla f(\mathbf{x}_k)\|}{2}} + \sigma_\tau + \delta. \quad (21)$$

Then, for any $\varepsilon > 0$, it is enough to do $K =$

$$\left\lceil 8(f(\mathbf{x}_0) - f^*) \cdot \left(\sqrt{\frac{L}{2}} \frac{1}{\varepsilon^{3/2}} + \frac{\sigma_\tau + \delta}{\varepsilon^2} \right) + 2 \ln \frac{\|\nabla f(\mathbf{x}_0)\|}{\varepsilon} \right\rceil$$

steps to ensure $\min_{1 \leq i \leq K} \|\nabla f(\mathbf{x}_i)\| \leq \varepsilon$.

For $\tau = n$ (the full Newton), we have $\sigma_n = 0$ and the complexity becomes $\mathcal{O}(\frac{1}{\varepsilon^{3/2}})$ up to an additive logarithmic term. It can also be established for the cubic regularization of the Newton method (Nesterov & Polyak, 2006). For $\tau = 0$ we recover the rate of the gradient descent on non-convex problems (Nesterov, 2018). Note that the rule (21) for choosing α_k is very simple. It is inspired by the gradient regularization technique developed initially for convex optimization (Mishchenko, 2023; Doikov & Nesterov, 2023).

5. Cutting the Negative Spectrum

In our previous result, we saw that the complexity of the method depends on the parameter σ_τ , which becomes smaller when increasing the spectral order τ of the method. However, due to (19), it includes a bound for the absolute value of the negative curvature, which can be constantly big. In this section, we propose a simple modification of our step-size rule, which alleviates this issue.

Let us denote the *positive part* of the Hessian, $\nabla_+^2 f(\mathbf{x}) := \nabla_\xi^2 f(\mathbf{x})$, where $\xi = \operatorname{argmax}\{\tau : \nabla_\tau^2 f(\mathbf{x}) \succeq \mathbf{0}\}$. Correspondingly, the *negative part* is $\nabla_-^2 f(\mathbf{x}) := \nabla_+^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$. We introduce the following system of parameters (compare with (18)):

$$\sigma_\tau^+ := \sup_{\mathbf{x} \in \mathbb{R}^n} \|\nabla_+^2 f(\mathbf{x}) - \nabla_\tau^2 f(\mathbf{x})\|. \quad (22)$$

So, σ_τ^+ is the uniform bound for the positive tail of the spectrum, and it is no longer affected by the negative curvature:

$$\sigma_\tau^+ \geq \max\{\lambda_{\tau+1}(\mathbf{x}), 0\}, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (23)$$

Our new step-size rule is based on the cubic regularization technique. Namely, at every iteration $k \geq 0$, we compute the regularization parameter as the maximum of the following univariate concave function:

$$\alpha_k^* = \operatorname{argmax}_{\alpha > 0} \left[-\frac{1}{2} \langle (\mathbf{H}_k + (\alpha + \eta)\mathbf{I})^{-1} \mathbf{g}_k, \mathbf{g}_k \rangle - \frac{2\alpha^3}{3L^2} \right],$$

where $\mathbf{g}_k := \nabla f(\mathbf{x}_k)$ is the current gradient, $\eta \geq 0$ is a balancing term, and $L > 0$ is the Lipschitz constant of the

Hessian. This subproblem is well defined and has a unique global maximum, which can be found by using binary search or univariate Newton iterations (Conn et al., 2000; Nesterov & Polyak, 2006). Then, to eliminate the negative curvature, for the sequence generated by our method we use an extra sequence of test points $\{\mathbf{y}_k\}_{k \geq 1}$ defined by $\mathbf{y}_{k+1} := \mathbf{x}_k + \mathbf{P}_k(\mathbf{x}_{k+1} - \mathbf{x}_k)$, where \mathbf{P}_k is a projector which preserves the image of \mathbf{H}_k : $\mathbf{H}_k \mathbf{P}_k = \mathbf{H}_k$, but has a small intersection with the negative part of the Hessian,

$$\|\nabla_-^2 f(\mathbf{x}_k) \mathbf{P}_k\| = \|(\nabla^2 f(\mathbf{x}_k) - \nabla_+^2 f(\mathbf{x}_k)) \mathbf{P}_k\| \leq \delta_-,$$

for some $\delta_- > 0$. It can be built directly from the low-rank representation (15) of \mathbf{H}_k , as follows: $\mathbf{P}_k = \sum_{i=1}^\tau \mathbf{v}_{k,i} \mathbf{v}_{k,i}^\top$. Therefore, this matrix comes with no extra cost. We have $\delta_- = 0$ when $\mathbf{v}_{k,i}$ are orthogonal to all eigenvectors of $\nabla_-^2 f(\mathbf{x}_k)$. Note that the tolerance parameters δ and δ_- are induced by the approximation error $\mathbf{H}_k \approx \nabla_\tau^2 f(\mathbf{x}_k)$, and, contrary to σ_τ and σ_τ^+ , they do not describe our problem class. Using a high accuracy approximations for matrices \mathbf{H}_k (see Section 7), we can make both δ and δ_- arbitrarily close to zero. We are ready to state the better complexity result for algorithm (14) with a convergence rate that is independent of the negative curvature.

Theorem 5.1. *Let $f \in \mathcal{F}_\tau$ be non-convex of grade τ , where $0 \leq \tau \leq n$ is fixed. Let f have a Lipschitz Hessian with constant L and bounded parameter σ_τ^+ . Consider iterations $\{\mathbf{x}_k\}_{k \geq 0}$ of algorithm (14) with*

$$\alpha_k = \alpha_k^* + \eta, \quad (24)$$

where $\eta := \sigma_\tau^+ + \delta + \delta_-$ is the balancing term. Then, for any $\varepsilon > 0$, it is enough to do

$$K = \left\lceil 2(f(\mathbf{x}_0) - f^*) \left(\frac{3\sqrt{2L}}{\varepsilon^{3/2}} + \frac{16\eta}{\varepsilon^2} \right) \right\rceil$$

steps to ensure $\min_{1 \leq i \leq K} \|\nabla f(\mathbf{y}_i)\| \leq \varepsilon$.

This complexity bound is similar to that one of Theorem 4.2, where we substituted $\sigma_\tau \mapsto \sigma_\tau^+$. It can be much better in case of large negative eigenvalues of the Hessian, since then σ_τ^+ is much smaller than σ_τ . The cost of such an improvement is a more complex rule (24), involving α_k^* , and the guarantee is given for the auxiliary points $\{\mathbf{y}_k\}_{k \geq 1}$.

6. Convex Convergence Analysis

In this section, we provide the analysis of our method for a specific case of $f \in \mathcal{F}_n$ (convex objectives). In this case, we do not have the negative curvature part: $\nabla_-^2 f(\mathbf{x}) \equiv \mathbf{0}$, and $\sigma_\tau^+ \equiv \sigma_\tau$. Since $\mathcal{F}_n \subset \mathcal{F}_\tau$ for any τ , the results of the previous sections can be applied directly for the convex case. However, as we show, using a more refined technique, we can prove much better convergence rates for convex

and strongly convex problems. We denote by $\mu \geq 0$ the parameter of strong convexity, such that $\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}, \forall \mathbf{x}$.

When the Hessian is positive semidefinite, we can naturally use it to define the *local norm* at point \mathbf{x} by $\|\mathbf{u}\|_{\mathbf{x}} := \langle \nabla^2 f(\mathbf{x}) \mathbf{u}, \mathbf{u} \rangle^{1/2}$. The local norm becomes more appropriate for describing the right second-order geometry of the objective (Nesterov & Nemirovski, 1994). With its help, it is possible to characterize smoothness more accurately. Now, we assume that f is *quasi-self-concordant*⁴ with parameter $M \geq 0$, for all $\mathbf{x}, \mathbf{u}, \mathbf{v} \in \mathbb{R}^n$:

$$\nabla^3 f(\mathbf{x})[\mathbf{u}, \mathbf{u}, \mathbf{v}] \leq M \|\mathbf{u}\|_{\mathbf{x}}^2 \|\mathbf{v}\|. \quad (25)$$

This condition was considered in (Bach, 2010; Karimireddy et al., 2018; Sun & Tran-Dinh, 2019; Doikov, 2023). It appears that many practical problems satisfy this assumption, including *Logistic Regression*, *Soft Maximum*, *Matrix Balancing*, and *Matrix Scaling* problems. We adapt algorithm (14) to this problem class, establishing fast convergence rates that improve with increase of the spectral order τ . We denote by $D := \sup\{\|\mathbf{x} - \mathbf{x}^*\| : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ the diameter of the initial sublevel set which we assume to be bounded. We establish the following result.

Theorem 6.1. *Let $f \in \mathcal{F}_n$ be strongly convex with $\mu > 0$ and quasi-self-concordant with constant M . Consider iterations $\{\mathbf{x}_k\}_{k \geq 0}$ of algorithm (14) with*

$$\alpha_k = M \|\nabla f(\mathbf{x}_k)\| + \sigma_\tau + \delta. \quad (26)$$

Then, for any $\varepsilon > 0$, it is enough to do $K =$

$$4 \left\lceil \left(MD + \frac{\sigma_\tau + \delta}{2\mu} \right) \ln \frac{f(\mathbf{x}_0) - f^*}{\varepsilon} + \ln \frac{\|\nabla f(\mathbf{x}_0)\| D}{\varepsilon} \right\rceil$$

steps to ensure $f(\mathbf{x}_K) - f^ \leq \varepsilon$.*

Remark 6.2. Note that by adding the quadratic regularizer to our objective with small $\mu := \frac{\varepsilon}{D^2}$, we can turn any convex problem into strongly convex one. Hence, we obtain the following complexity for the general convex case:

$$\mathcal{O}\left(\frac{(\sigma_\tau + \delta)D^2}{\varepsilon} + MD \ln \frac{f(\mathbf{x}_0) - f^*}{\varepsilon} + \ln \frac{\|\nabla f(\mathbf{x}_0)\| D}{\varepsilon}\right).$$

Let us consider the simplest example of using only the one top eigenvector, $\tau = 1$. Then, employing the basic power method (see Section 7 and, e.g. (Golub & Van Loan, 2013)), the complexity of computing the preconditioner is $\mathcal{O}\left(\frac{\lambda_1}{\lambda_1 - \lambda_2} \ln \frac{1}{\delta}\right)$, where $\delta > 0$ is our target accuracy for computing the approximate eigenvector. Note that δ enters *under logarithm*. Hence, it seems reasonable to assume that we can make δ sufficiently small with only a small amount of extra effort. Combining the outer complexity and the

⁴Note that for the functions with Lipschitz continuous Hessian (17), we can bound the third derivative with the product of Euclidean norms: $\nabla^3 f(\mathbf{x})[\mathbf{u}, \mathbf{u}, \mathbf{v}] \leq L \|\mathbf{u}\|^2 \|\mathbf{v}\|$, which is in most cases less accurate than (25).

inner complexity of the power method, we observe that in some cases, we gain a provable benefit from using the preconditioner. Thus, in the convex case, the total complexity (taking into account the iterations of the power method) becomes $\tilde{\mathcal{O}}\left(\left[\frac{\lambda_2}{\lambda_n} + MD\right] \times \left[\frac{\lambda_1}{\lambda_1 - \lambda_2}\right]\right)$, where $\tilde{\mathcal{O}}(\cdot)$ hides logarithmic factors. It is much better than the complexity $\tilde{\mathcal{O}}\left(\frac{\lambda_1}{\lambda_n}\right)$ of the pure gradient descent when $\lambda_1 \gg \lambda_2$. These observations can be generalized to an arbitrary τ . However, even for $\tau = 1$, we can observe a significant improvement, which is also justified in our numerical examples (Fig. 11).

7. Efficient Implementation

The method (14) relies on an approximation $\mathbf{H}_k \approx \nabla^2 f(\mathbf{x}_k)$. We acknowledge that it might be costly to find such an approximation in general. Below, we outline a method that iteratively constructs low rank approximations by employing $T_k \geq 0$ steps of the power method. Here T_k is our parameter. In our experiments we use $T_k = 1$ (combined with hot-start from \mathbf{H}_{k-1}). However, other options are to choose $T_k = 0$ for some of the iterations (no update), or to spend more iteration in the initialization phase. We leave the explorations of such schedules to future work.

Power Method. For our experiments, we use a generalization of the classical power method for finding the top τ eigenvectors of a given matrix, which is also known as *orthogonal iterations*. We can write our matrix as $\mathbf{H}_k = \mathbf{V}_k \text{Diag}(\mathbf{a}_k) \mathbf{V}_k^\top$, where $\mathbf{a}_k \in \mathbb{R}_{\geq 0}^\tau$ and $\mathbf{V}_k \in \mathbb{R}^{n \times \tau}$ consists of orthonormal columns $\{\mathbf{v}_{k,i}\}_{i=1}^\tau$. We denote by $[\mathbf{A}]_Q := \mathbf{Q}$ the orthogonal matrix from the resulting QR-factorization of $\mathbf{A} = \mathbf{Q}\mathbf{R}$. It can be implemented as the standard Gram-Schmidt orthogonalization process with arithmetic complexity $\mathcal{O}(\tau^2 n)$ operations. Then, for each $k \geq 0$ we use the following procedure to update matrix \mathbf{V}_k :

Power Method
Set $\mathbf{V}_k := \mathbf{V}_{k-1}$ if $k \geq 1$ else random init.
For $t = 1 \dots T_k$ iterate:
Update $\mathbf{V}_k := [\nabla^2 f(\mathbf{x}_k) \mathbf{V}_k]_Q$
Set $\mathbf{a}_{k,i} := \langle \nabla^2 f(\mathbf{x}_k) \mathbf{v}_{k,i}, \mathbf{v}_{k,i} \rangle$.

(27)

This procedure converges with linear rate, with the main complexity factor proportional to the τ -th spectral gap (see, e.g., Theorem 8.2.2 from (Golub & Van Loan, 2013)). More advanced approaches include Oja's and Lanczos iterations (Kuczyński & Woźniakowski, 1992).

Low-Rank Updates. The Woodbury matrix identity provides us with the following formula:

$$(\mathbf{H}_k + \alpha_k \mathbf{I})^{-1} = \frac{1}{\alpha_k} [\mathbf{I} - \mathbf{V}_k (\mathbf{I} + \alpha_k \text{Diag}(\mathbf{a}_k))^{-1} \mathbf{V}_k^\top].$$

Therefore, we need only $O(\tau^2 n)$ arithmetical operations to perform the step, which is linear with respect to n and can be implemented very efficiently for small τ .

8. Experiments

We present illustrative numerical experiments on several machine learning problems. See Section A in the appendix for the details of our experiments and for extra plots.

Matrix Factorization. In Fig. 5, we show the convergence of the outer iterations of algorithm (14), solving an ill-conditioned instance of non-convex matrix factorization problems. We study two different strategies of choosing the regularization parameter: constant selection of parameters (left), and using the adaptive search from Section A.2 (right).

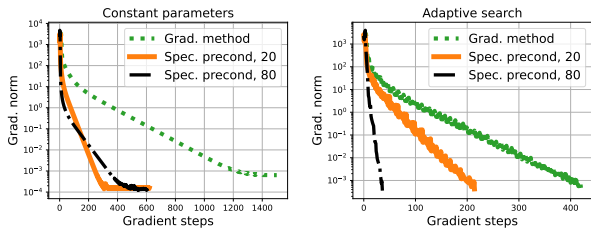


Figure 5. Solving an ill-conditioned matrix factorization. Left: best constant parameters, right: adaptive search. Spectral preconditioning helps the methods to converge faster in terms of the gradient steps.

We see that employing the spectral preconditioning helps the method converge faster, in terms of iterations (gradient steps). The use of the adaptive search accelerates the method even further, ensuring an automatic selection of the regularization parameter.

Logistic Regression. In the following experiments, we train a convex logistic regression model on several machine learning datasets, using the gradient method with spectral preconditioning. We also compare its performance with quasi-Newton methods: BFGS and the limited memory BFGS (L-BFGS) (Nocedal & Wright, 2006). The results are shown in Fig. 6. We see that employing spectral preconditioning even, for $\tau = 1$ and $\tau = 2$, can significantly accelerate the convergence of the gradient method in terms of both gradient steps and the total number of arithmetic operations. Its overall performance can be comparable to that of the quasi-Newton methods.

9. Conclusion

In this work, we propose using inexact top eigenvectors of the Hessian as a preconditioning for gradient methods.

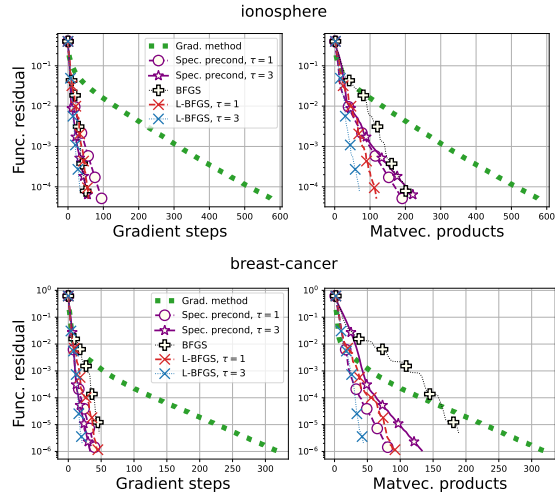


Figure 6. Training the logistic regression model. The method with spectral preconditioning works comparably to the BFGS algorithm.

We introduce the notion of *graded non-convexity*, which provides our problem classes with a uniform guarantee on the number of positive eigenvalues. We show that using a preconditioner of order $\tau \geq 1$ provably improves the rate of convergence, cutting the gap between the top τ eigenvalues.

There are several potential directions for future research. First, it could be interesting to discover whether it is possible to use a Hessian approximation H_k that has a higher rank than τ , where τ is the grade of non-convexity, and contains information on directions with negative curvature. In such a case, it seems necessary to use the cubic regularization of trust-region approach to properly bound the length of the method step. At the same time, it poses an interesting open question: can negative curvature be helpful for the method? It is also important to study a behavior of spectral preconditioning in a neighborhood of saddle points.

Another direction is the local analysis of the method. Indeed, besides using $Q = \mathbb{R}^n$ in (7), several other options exist. For example, the initial sublevel set $Q = \{x : f(x) \leq f(x_0)\}$, since, due to Lemma 4.1, all iterations belong to it. Alternatively, explicit regularization can be used with $Q = \{x : \|x\| \leq R\}$. Also, as in the classic analysis of Newton’s method, we can assume that the starting point is in a neighborhood Q of a local non-degenerate solution. In this case, the local perspective on graded non-convexity can enhance analysis of quasi-Newton methods.

Finally, it can be interesting to combine our results with recently proposed methods that leverage negative curvature (Carmon et al., 2018). In contrast to the pure first-order methods that typically depend on the leading eigenvalue λ_1 of the Hessian, in our work we allow for the elimination of this dependence.

Acknowledgements

This work was supported by the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number 22.00133.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Bach, F. Self-concordant analysis for logistic regression. 2010.
- Carmon, Y., Duchi, J. C., Hinder, O., and Sidford, A. Accelerated methods for nonconvex optimization. *SIAM Journal on Optimization*, 28(2):1751–1772, 2018.
- Cartis, C., Gould, N. I., and Toint, P. L. Adaptive cubic regularisation methods for unconstrained optimization. Part I: motivation, convergence and numerical results. *Mathematical Programming*, 127(2):245–295, 2011a.
- Cartis, C., Gould, N. I., and Toint, P. L. Adaptive cubic regularisation methods for unconstrained optimization. Part II: worst-case function-and derivative-evaluation complexity. *Mathematical programming*, 130(2):295–319, 2011b.
- Cartis, C., Gould, N. I., and Toint, P. L. An example of slow convergence for Newton’s method on a function with globally Lipschitz continuous Hessian. Technical report, Technical report, ERGO 13-008, School of Mathematics, Edinburgh University, 2013.
- Conn, A. R., Gould, N. I., and Toint, P. L. *Trust region methods*. SIAM, 2000.
- Cunha, L., Gidel, G., Pedregosa, F., Scieur, D., and Paquette, C. Only tails matter: Average-case universality and robustness in the convex regime. In *International Conference on Machine Learning*, pp. 4474–4491. PMLR, 2022.
- Doikov, N. Minimizing quasi-self-concordant functions by gradient regularization of Newton method. *arXiv preprint arXiv:2308.14742*, 2023.
- Doikov, N. and Grapiglia, G. N. First and zeroth-order implementations of the regularized Newton method with lazy approximated Hessians. *arXiv preprint arXiv:2309.02412*, 2023.
- Doikov, N. and Nesterov, Y. Gradient regularization of Newton method with Bregman distances. *Mathematical Programming*, pp. 1–25, 2023.
- Doikov, N. and Rodomanov, A. Polynomial preconditioning for gradient methods. In *40th International Conference on Machine Learning (ICML 2023)*, number 40, 2023.
- Doikov, N., Mishchenko, K., and Nesterov, Y. Super-universal regularized Newton method. *arXiv preprint arXiv:2208.05888*, 2022.
- Doikov, N., Chayti, E. M., and Jaggi, M. Second-order optimization with lazy Hessians. In *International Conference on Machine Learning*. PMLR, 2023.
- Even, M., Pesme, S., Gunasekar, S., and Flammarion, N. (s) gd over diagonal linear networks: Implicit regularisation, large stepsizes and edge of stability. *arXiv preprint arXiv:2302.08982*, 2023.
- Golub, G. H. and Van Loan, C. F. *Matrix computations*. JHU press, 2013.
- Goujaud, B., Scieur, D., Dieuleveut, A., Taylor, A. B., and Pedregosa, F. Super-acceleration with cyclical step-sizes. In *International Conference on Artificial Intelligence and Statistics*, pp. 3028–3065. PMLR, 2022.
- Grapiglia, G. N. and Nesterov, Y. Regularized Newton methods for minimizing functions with Hölder continuous Hessians. *SIAM Journal on Optimization*, 27(1): 478–506, 2017.
- Gromov, M. Sign and geometric meaning of curvature. *Rendiconti del Seminario Matematico e Fisico di Milano*, 61:9–123, 1991.
- Hanzely, F., Doikov, N., Nesterov, Y., and Richtarik, P. Stochastic subspace cubic Newton method. In *International Conference on Machine Learning*, pp. 4027–4038. PMLR, 2020.
- Hanzely, S. Sketch-and-project meets Newton method: Global $O(1/k^2)$ convergence with low-rank updates. *arXiv preprint arXiv:2305.13082*, 2023.
- Hanzely, S., Kamzolov, D., Pasechnyuk, D., Gasnikov, A., Richtárik, P., and Takác, M. A damped Newton method achieves global $O(\frac{1}{k^2})$ and local quadratic convergence rate. *Advances in Neural Information Processing Systems*, 35:25320–25334, 2022.
- Hestenes, M. R. and Stiefel, E. Methods of conjugate gradients for solving linear systems. *Journal of research of the National Bureau of Standards*, 49(6):409–436, 1952.
- Hörmander, L. *Notions of convexity*. Springer Science & Business Media, 2007.
- Islamov, R., Qian, X., and Richtárik, P. Distributed second order methods with fast rates and compressed communication. In *International conference on machine learning*, pp. 4617–4628. PMLR, 2021.

- Islamov, R., Qian, X., Hanzely, S., Safaryan, M., and Richtárik, P. Distributed Newton-type methods with communication compression and Bernoulli aggregation. *arXiv preprint arXiv:2206.03588*, 2022.
- Jiang, R. and Mokhtari, A. Accelerated quasi-newton proximal extragradient: Faster rate for smooth convex optimization. *arXiv preprint arXiv:2306.02212*, 2023.
- Jin, Q. and Mokhtari, A. Non-asymptotic superlinear convergence of standard quasi-Newton methods. *Mathematical Programming*, 200(1):425–473, 2023.
- Kamzolov, D., Ziu, K., Agafonov, A., and Takáč, M. Cubic regularized quasi-newton methods. *arXiv preprint arXiv:2302.04987*, 2023.
- Karimireddy, S. P., Stich, S. U., and Jaggi, M. Global linear convergence of Newton’s method without strong-convexity or Lipschitz gradients. *arXiv preprint arXiv:1806.00413*, 2018.
- Kovalev, D., Richtarik, P., Gorbunov, E., and Gasanov, E. Stochastic spectral and conjugate descent methods. *Advances in Neural Information Processing Systems*, 31, 2018.
- Kuczyński, J. and Woźniakowski, H. Estimating the largest eigenvalue by the power and Lanczos algorithms with a random start. *SIAM journal on matrix analysis and applications*, 13(4):1094–1122, 1992.
- Kunstner, F., Sanches Portella, V., Schmidt, M., and Harvey, N. Searching for optimal per-coordinate step-sizes with multidimensional backtracking. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ma, C., Xu, X., Tong, T., and Chi, Y. Provably accelerating ill-conditioned low-rank estimation via scaled gradient descent, even with overparameterization. *arXiv preprint arXiv:2310.06159*, 2023.
- Ma, S. and Belkin, M. Diving into the shallows: a computational perspective on large-scale shallow learning. *Advances in neural information processing systems*, 30, 2017.
- Mardal, K.-A. and Winther, R. Preconditioning discretizations of systems of partial differential equations. *Numerical Linear Algebra with Applications*, 18(1):1–40, 2011.
- Mishchenko, K. Regularized Newton method with global $\mathcal{O}(1/k^2)$ convergence. *SIAM Journal on Optimization*, 33(3):1440–1462, 2023.
- Nemirovski, A. and Yudin, D. Problem complexity and method efficiency in optimization. 1983.
- Nesterov, Y. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Nesterov, Y. Implementable tensor methods in unconstrained convex optimization. *Mathematical Programming*, pp. 1–27, 2019.
- Nesterov, Y. Quartic regularity. *arXiv preprint arXiv:2201.04852*, 2022.
- Nesterov, Y. and Nemirovski, A. *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.
- Nesterov, Y. and Polyak, B. T. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- Nocedal, J. and Wright, S. *Numerical optimization*. Springer Science & Business Media, 2006.
- Nutini, J., Laradji, I., and Schmidt, M. Let’s make block coordinate descent converge faster: faster greedy rules, message-passing, active-set complexity, and superlinear convergence. *Journal of Machine Learning Research*, 23(131):1–74, 2022.
- Pasechnyuk, D. A., Gasnikov, A., and Takáč, M. Effects of momentum scaling for sgd. *arXiv preprint arXiv:2210.11869*, 2022.
- Pasechnyuk, D. A., Gasnikov, A., and Takáč, M. Convergence analysis of stochastic gradient descent with adaptive preconditioning for non-convex and convex functions. *arXiv preprint arXiv:2308.14192*, 2023.
- Pesme, S. and Flammarion, N. Saddle-to-saddle dynamics in diagonal linear networks. *arXiv preprint arXiv:2304.00488*, 2023.
- Pesme, S., Pillaud-Vivien, L., and Flammarion, N. Implicit bias of sgd for diagonal linear networks: a provable benefit of stochasticity. *Advances in Neural Information Processing Systems*, 34:29218–29230, 2021.
- Polyak, B. T. Newton’s method and its use in optimization. *European Journal of Operational Research*, 181(3):1086–1096, 2007.
- Qian, X., Islamov, R., Safaryan, M., and Richtárik, P. Basis matters: better communication-efficient second order methods for federated learning. *arXiv preprint arXiv:2111.01847*, 2021.
- Rodomanov, A. and Kropotov, D. A randomized coordinate descent method with volume sampling. *SIAM Journal on Optimization*, 30(3):1878–1904, 2020.

- Rodomanov, A. and Nesterov, Y. New results on superlinear convergence of classical quasi-Newton methods. *Journal of optimization theory and applications*, 188:744–769, 2021a.
- Rodomanov, A. and Nesterov, Y. Rates of superlinear convergence for classical quasi-Newton methods. *Mathematical Programming*, pp. 1–32, 2021b.
- Safaryan, M., Islamov, R., Qian, X., and Richtárik, P. Fednl: Making Newton-type methods applicable to federated learning. *arXiv preprint arXiv:2106.02969*, 2021.
- Scieur, D. and Pedregosa, F. Universal average-case optimality of Polyak momentum. In *International conference on machine learning*, pp. 8565–8572. PMLR, 2020.
- Spielman, D. A. and Teng, S.-H. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pp. 81–90, 2004.
- Sun, T. and Tran-Dinh, Q. Generalized self-concordant functions: a recipe for Newton-type methods. *Mathematical Programming*, 178(1-2):145–213, 2019.
- Vaidya, P. M. Solving linear equations with symmetric diagonally dominant matrices by constructing good preconditioners. *A talk based on this manuscript*, 2(3.4):2–4, 1991.
- Vial, J.-P. Strong and weak convexity of sets and functions. *Mathematics of Operations Research*, 8(2):231–259, 1983.
- Woodworth, B., Gunasekar, S., Lee, J. D., Moroshko, E., Savarese, P., Golan, I., Soudry, D., and Srebro, N. Kernel and rich regimes in overparametrized models. In Abernethy, J. and Agarwal, S. (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 3635–3673. PMLR, 09–12 Jul 2020. URL <https://proceedings.mlr.press/v125/woodworth20a.html>.
- Zhang, G., Fattahi, S., and Zhang, R. Y. Preconditioned gradient descent for overparameterized nonconvex Burer–Monteiro factorization with global optimality certification. *Journal of Machine Learning Research*, 24(163): 1–55, 2023.
- Zhang, J., Fattahi, S., and Zhang, R. Y. Preconditioned gradient descent for over-parameterized nonconvex matrix factorization. *Advances in Neural Information Processing Systems*, 34:5985–5996, 2021.

Appendix

A. Extra Experiments and Details

In this section, let us provide the details of our experiments and additional numerical results, comparing the performance of the spectral preconditioning.

A.1. Matrix Factorization

We consider the following non-convex optimization problem, which has applications in dimensionality reduction and clustering:

$$\min_{\substack{\mathbf{X} \in \mathbb{R}^{n \times r} \\ \mathbf{Y} \in \mathbb{R}^{r \times m}} \left[f(\mathbf{X}, \mathbf{Y}) := \frac{1}{2} \|\mathbf{X}\mathbf{Y} - \mathbf{C}\|^2 \right], \quad (28)$$

where $\mathbf{C} \in \mathbb{R}^{n \times m}$ is a give target matrix and the norm is the standard Frobenius: $\|\mathbf{A}\| := \text{Trace}(\mathbf{A}^\top \mathbf{A})^{1/2}$ for $\mathbf{A} \in \mathbb{R}^{n, m}$.

In the following experiment, we fix $n = m = 40$ and generate matrix \mathbf{C} by the following procedure: we choose $r^* := 10$ (the target rank) and set

$$\mathbf{C} := \sum_{i=1}^{r^*} \mathbf{u}_i \mathbf{v}_i^\top,$$

where $\{\mathbf{u}_i\}_{i=1}^{r^*}$ and $\{\mathbf{v}_i\}_{i=1}^{r^*}$ are generated randomly and uniformly from the unit sphere. Therefore, matrix \mathbf{C} has r^* singular values that are close to 1, and all others are close to 0. In all our methods we use the initial point $(\mathbf{X}_0, \mathbf{Y}_0)$ generated randomly with entries from the standard normal distribution. A typical distribution of the spectrum of the Hessian of f at the initial point is shown in Fig. 7. We observe that there exist a group of eigenvalues which are significantly larger than the others. Therefore, in such cases, we expect to see an accelerated of the methods by using our spectral preconditioning technique.

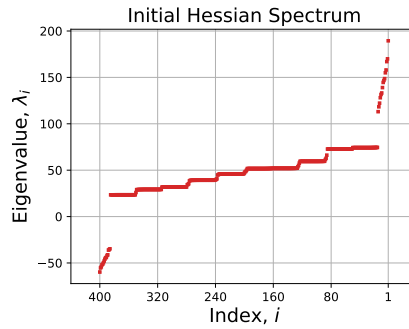


Figure 7. Spectrum of the Hessian at the initial point, $r = 5$ and $r^* = 10$.

We compare the classic gradient descent with the spectral preconditioning. In the gradient descent, we use the standard Armijo-type line search at each iteration to choose the step size. For the spectral preconditioning, we fix the regularization parameter, according to our theory, at iteration $k \geq 0$:

$$\alpha_k := \sqrt{L \|\nabla f(\mathbf{X}_k, \mathbf{Y}_k)\|} + \beta_k,$$

where we fix $L := 1$ and β_k is fitted using a simple adaptive search, that ensures sufficient progress of every step (see Section A.2). Namely, we start with an initial value of $\beta_0 := 0.05$. In each iteration, we increase the previous estimate by two until the sufficient condition is satisfied. After performing the method step, we decrease the found value by two. Thus, we allow parameters β_k to both increase and decrease with each iteration.

We use one iteration of the hot start power method to update our estimate of $\nabla_{\tau}^2 f$ (see Section 7).

- $r^* = 1$. In this experiment, we compare the performance of spectral preconditioning recovering the matrix \mathbf{C} of rank 1, using $r = 5$ components. The results are shown in Fig. 8.

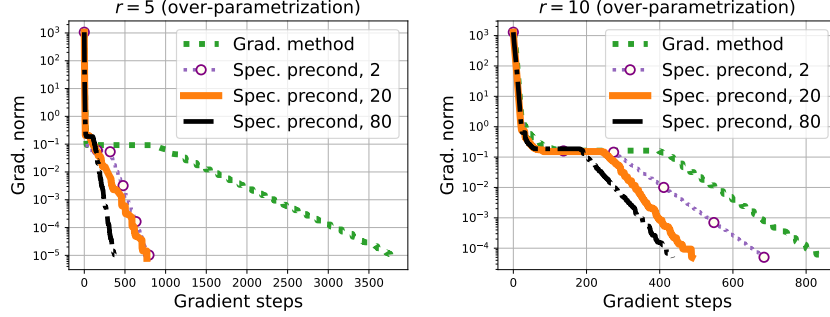


Figure 8. Convergence of the methods on the matrix factorization problem, recovering the matrix of rank $r^* = 1$.

- $r^* = 10$. In the following experiment, we test three different settings for the rank parameter: $r < r^*$ (under-parametrization), $r = r^*$ (true value), and $r > r^*$ (over-parametrization). The results are shown in Fig. 9.

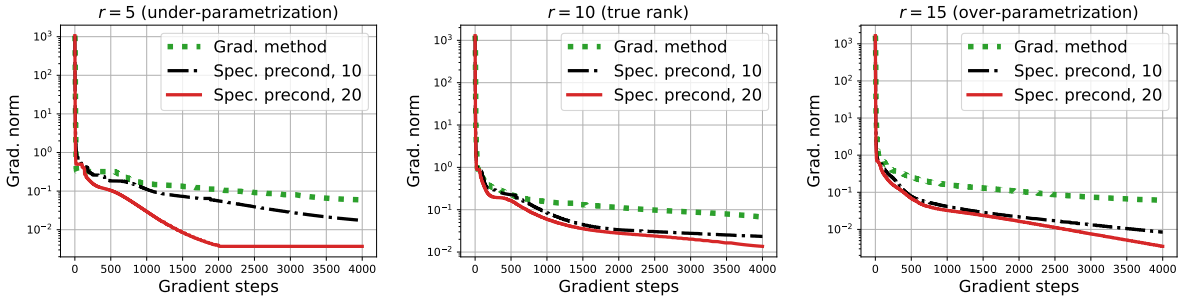


Figure 9. Convergence of the methods on the matrix factorization problem, for different values of r (rank estimate), recovering $r^* = 10$.

We see that it becomes easier to solve the problem, when parameter r is increasing. Note that the dimension of problem (28) also growth with r . In all cases, using the spectral preconditioning with a small value of τ accelerates convergence of the method, as predicted by our theory.

A.2. Selection of the Regularization Parameter

In this section, let us compare two different strategies for choosing the regularization parameter in the spectral preconditioning: a constant choice and adaptive search.

For solving optimization problem $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$, we consider the following rule (as in Theorem 4.2), for each iteration $k \geq 0$:

$$\alpha_k = \sqrt{L \|\nabla f(\mathbf{x}_k)\|} + \beta_k.$$

Then, one step of the method is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\mathbf{H}_k + \alpha_k \mathbf{I})^{-1} \nabla f(\mathbf{x}_k), \quad \text{where} \quad \mathbf{H}_k \approx \nabla_{\tau}^2 f(\mathbf{x}_k).$$

Our theory shows that choosing L and β_k sufficiently big, we guarantee a significant progress for each iteration (Lemma 4.1), that leads to good global convergence rates. The parameter L is the Lipschitz constant of the Hessian, while β_k depends on the τ th spectral gap and the accuracy of approximating the spectral Hessian. Hence, estimating L seems to be easier

Algorithm 1 Adaptive Gradient Method with Spectral Preconditioning

```

1: Input:  $\mathbf{x}_0 \in \mathbb{R}^n$ . Fix some  $L \geq 0$  and  $\beta_0 > 0$ .
2: for  $k = 0, 1, \dots$  do
3:   Maintain a low-rank approximation  $\mathbf{H}_k \approx \nabla_\tau^2 f(\mathbf{x}_k)$ 
4:   repeat
5:     Update  $\beta_k = 2 \cdot \beta_k$ 
6:     Set  $\alpha_k = \sqrt{L \|\nabla f(\mathbf{x}_k)\|} + \beta_k$ 
7:     Perform the gradient step:  $\mathbf{x}_{k+1} = \mathbf{x}_k - (\mathbf{H}_k + \alpha_k \mathbf{I})^{-1} \nabla f(\mathbf{x}_k)$ 
8:     until  $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geq \frac{1}{8\alpha_k} \|\nabla f(\mathbf{x}_{k+1})\|^2$ .
9:     Set  $\beta_{k+1} = \frac{1}{2} \cdot \beta_k$ 
10:  end for

```

than β_k (indeed, for some problems, there exist available bounds for the Lipschitz constants, as for example in the logistic regression or in the soft maximum problem). At the same time, estimating parameter β_k can be more difficult. Therefore, we use the following adaptive procedure, which fixes the value of L and chooses β_k adaptively in each iteration.

This adaptive search resembles a standard Armijo-type procedure, which is popular for the first-order and second-order methods (Nocedal & Wright, 2006; Nesterov, 2018). Importantly, it does not provide the method with any significant overhead, but makes us free from choosing the true value of parameter β_k .

Note that according to Lemma 4.1, the stopping condition in the inner loop will be satisfied at least when $\beta_k \geq \sigma_\tau + \delta$. Hence, the method is well-defined.

Let us denote by β_k^{init} the first value of parameter β_k after its initialization (in line 1 for β_0 and line 9 for β_{k+1}). Thus, the number of times the inner loop in lines 4-8 is performed at iteration $k \geq 0$ is $1 + \log_2(\beta_{k+1}^{\text{init}}/\beta_k^{\text{init}})$ and the total number T_k of the gradient steps during the first k iterations of the algorithm is estimated as

$$\begin{aligned}
 T_k &= \sum_{i=0}^{k-1} (1 + \log_2 \frac{\beta_{i+1}^{\text{init}}}{\beta_i^{\text{init}}}) = k + \log_2 \prod_{i=0}^{k-1} \frac{\beta_{i+1}^{\text{init}}}{\beta_i^{\text{init}}} = k + \log_2 \beta_k^{\text{init}} - \log \beta_0^{\text{init}} \\
 &\leq k + 1 + \log_2 \frac{\sigma_\tau + \delta}{\beta_0^{\text{init}}},
 \end{aligned}$$

where we used in the last bound that $\beta_k \leq \max\{2(\sigma_\tau + \delta), \beta_0^{\text{init}}\}$, $\forall k \geq 0$, due to Lemma 4.1. Therefore, up to an additional logarithmic term, the adaptive search method performs the same number of the gradient steps as the method with a fixed (known) value of the regularization constant. In average, it requires only one extra gradient computation (line 8) per iteration.

In the following graphs, let us compare two strategies for choosing the regularization parameter:

- Constant choice (fitted separately for each method and each problem);
- Adaptive search, provided by Algorithm 1.

We test the matrix factorization problem as in the previous experiments:

$$\min_{\substack{\mathbf{X} \in \mathbb{R}^{n \times r} \\ \mathbf{Y} \in \mathbb{R}^{r \times m}} \left[f(\mathbf{X}, \mathbf{Y}) := \frac{1}{2} \|\mathbf{X}\mathbf{Y} - \mathbf{C}\|^2 \right], \tag{29}$$

where $\mathbf{C} \in \mathbb{R}^{n \times m}$ is a give target matrix, generated randomly, of a fixed rank r^* . We set $n = m = 40$ and use the initial points $(\mathbf{X}_0, \mathbf{Y}_0)$ generated randomly with entries from the standard normal distribution.

Ill-conditioned problem. We consider an ill-conditioned problem of rank $r^* = 1$. The corresponding singular value of matrix \mathbf{C} is 400.

In the following graphs, we see convergence of the methods for the constant selection of parameters (left), and using the adaptive search (right). We observe that not all values of the gradient norm are monotonically decreasing. This does

not contradict our theory, since the guarantees given by Theorem 4.2 and Theorem 5.1 are given in terms of the *minimal gradient norm*. In Figure 5, we show the convergence in terms of the *current gradient norm* $\|\nabla f(\mathbf{x}_k)\|$, and in Figure 10, we demonstrate the convergence of the same methods for the *minimal gradient norm* $\min_{0 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|$, for each iteration $k \geq 0$.

We see that using the adaptive search might increase oscillations of the gradient norm. However, it also helps to reach the desired value of the minimal gradient norm faster. We find that using adaptive search is always helpful, as the method becomes independent of the parameters σ_τ and δ .

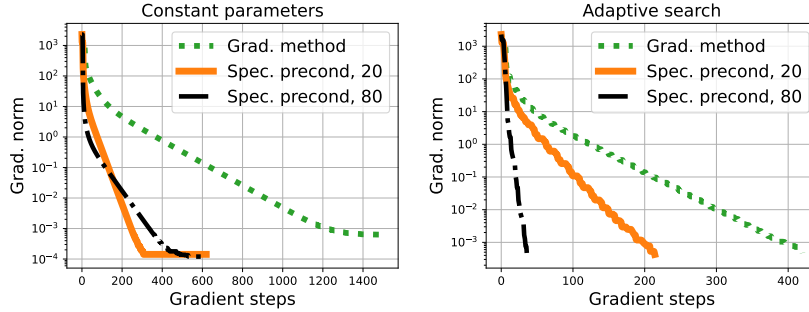


Figure 10. The graphs from Figure 5 for the minimal gradient norm $\min_{0 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|$.

A.3. Convex Models

In the following numerical experiments we train a logistic regression model on several machine learning datasets⁵. Thus, our optimization problem has the following form:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \left[f(\mathbf{x}) := \frac{1}{m} \sum_{i=1}^m \ln(1 + e^{\langle \mathbf{a}_i, \mathbf{x} \rangle}) \right], \tag{30}$$

where vectors $\{\mathbf{a}_i\}_{i=1}^m$ are determined by the dataset.

The results are shown in Fig. 11. We see that the spectral preconditioning works significantly better than the basic gradient descent and comparable with the powerful BFGS and the limited memory BFGS (L-BFGS) methods (Nocedal & Wright, 2006). The spectral preconditioning usually outperforms BFGS algorithm, for a low accuracy level, and works slightly worse than L-BFGS, using the same number of τ directions (the memory size).

The classical quasi-Newton methods could achieve better practical performance in some cases because of their superlinear local convergence (Rodomanov & Nesterov, 2021a). In contrast, the spectral preconditioning is designed to provide strong global convergence guarantees, including those for non-convex problems. Therefore, developing a low-rank spectral preconditioning technique equipped with local superlinear rates could be an interesting direction for future research.

⁵<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

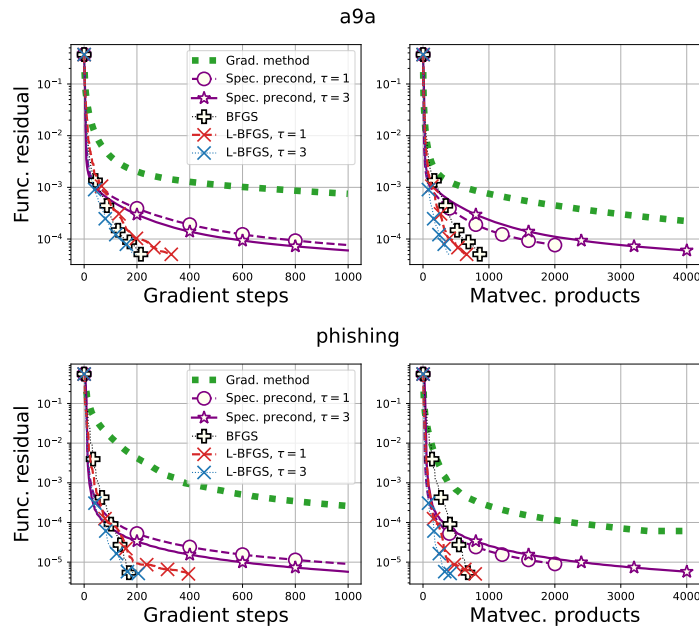


Figure 11. Training logistic regression using the spectral preconditioning with $\tau = 1$ and $\tau = 3$. We see that the spectral preconditioning works significantly better than the basic gradient descent and comparable with the powerful BFGS algorithm.

B. Proofs for Section 2

B.1. Proof of Proposition 2.2

Proof. For any two symmetric matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, and $1 \leq i, j \leq n$ such that $i + j \geq n + 1$, we have the following Weyl's inequality:

$$\lambda_{i+j-n}(\mathbf{A} + \mathbf{B}) \geq \lambda_i(\mathbf{A}) + \lambda_j(\mathbf{B}), \quad (31)$$

which immediately proves the statement about the sum. To prove the statement for soft maximum, let us denote

$$f(\mathbf{x}) := \text{smax}(f_1(\mathbf{x}), f_2(\mathbf{x})) = \ln(e^{f_1(\mathbf{x})} + e^{f_2(\mathbf{x})}),$$

and compute the derivatives. We obtain, for any $\mathbf{h} \in \mathbb{R}^n$:

$$\langle \nabla f(\mathbf{x}), \mathbf{h} \rangle = \sigma(f_1(\mathbf{x}) - f_2(\mathbf{x})) \cdot \langle \nabla f_1(\mathbf{x}), \mathbf{h} \rangle + \sigma(f_2(\mathbf{x}) - f_1(\mathbf{x})) \cdot \langle \nabla f_2(\mathbf{x}), \mathbf{h} \rangle,$$

where $\sigma(t) \stackrel{\text{def}}{=} \frac{1}{1+e^{-t}}$. Note that $\sigma'(t) = \sigma(t) \cdot \sigma(-t)$. Denoting $\sigma_1 \stackrel{\text{def}}{=} \sigma(f_1(\mathbf{x}) - f_2(\mathbf{x}))$, $\sigma_2 \stackrel{\text{def}}{=} \sigma(f_2(\mathbf{x}) - f_1(\mathbf{x}))$, we get

$$\begin{aligned} \langle \nabla^2 f(\mathbf{x}) \mathbf{h}, \mathbf{h} \rangle &= \sigma_1 \sigma_2 \left[\langle \nabla f_1(\mathbf{x}), \mathbf{h} \rangle - \langle \nabla^2 f(\mathbf{x}), \mathbf{h} \rangle \right]^2 + \sigma_1 \langle \nabla^2 f_1(\mathbf{x}) \mathbf{h}, \mathbf{h} \rangle + \sigma_2 \langle \nabla^2 f_2(\mathbf{x}) \mathbf{h}, \mathbf{h} \rangle \\ &\geq \sigma_1 \langle \nabla^2 f_1(\mathbf{x}) \mathbf{h}, \mathbf{h} \rangle + \sigma_2 \langle \nabla^2 f_2(\mathbf{x}) \mathbf{h}, \mathbf{h} \rangle. \end{aligned}$$

Using (31) completes the proof. \square

B.2. Proof of Proposition 2.3

Proof. The Hessian of g is given by

$$\nabla^2 g(\mathbf{x}) = \mathbf{A}^\top \nabla^2 f(L(\mathbf{x})) \mathbf{A} \in \mathbb{R}^{m \times m}.$$

The linear map $\mathbf{A} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ induces the isomorphism: $\mathbb{R}^m / \text{Ker}(\mathbf{A}) \xrightarrow{\varphi} \text{Im}(\mathbf{A})$ and thus $m = \dim \text{Ker}(\mathbf{A}) + \dim \text{Im}(\mathbf{A})$ (the classical Rank-nullity theorem).

By our assumption, $\nabla_\tau^2 f(L(\mathbf{x})) \succeq \mathbf{0}$. Let us denote by V_+ the subspace spanned by top τ eigenvectors of $\nabla^2 f(L(\mathbf{x}))$, and V_- is spanned by the rest. Hence, $\mathbb{R}^n = V_+ \oplus V_-$. Then

$$\begin{aligned} m &= \dim \text{Ker}(\mathbf{A}) + \dim \text{Im}(\mathbf{A}) \cap V_+ + \dim \text{Im}(\mathbf{A}) \cap V_- \\ &\leq \dim \text{Ker}(\mathbf{A}) + \dim \text{Im}(\mathbf{A}) \cap V_+ + n - \tau. \end{aligned}$$

Rearranging the terms, we get

$$\dim \text{Ker}(\mathbf{A}) + \dim \text{Im}(\mathbf{A}) \cap V_+ \geq m - n + \tau.$$

Therefore, we conclude that the linear subspace $U = \text{Ker}(\mathbf{A}) \oplus \varphi^{-1}(\text{Im}(\mathbf{A}) \cap V_+)$ has dimension $\dim U \geq m - n + \tau$, and for any $\mathbf{h} \in U$:

$$\langle \nabla^2 g(\mathbf{x}) \mathbf{h}, \mathbf{h} \rangle \geq 0,$$

which proves the required bound. \square

B.3. Proof of Proposition 2.4

Proof. For a small $\varepsilon > 0$, we consider the regularized objective

$$g(\mathbf{x}) := f(\mathbf{x}) + \frac{\varepsilon}{2} \|\mathbf{x}\|^2.$$

Assume that g achieves its maximum $\mathbf{x}_g^* = \arg \max_{\mathbf{x} \in K} g(\mathbf{x})$ in the interior, $\mathbf{x}_g^* \in \text{int } K$. Then, the second-order stationary condition implies that

$$\nabla^2 g(\mathbf{x}_g^*) \preceq \mathbf{0},$$

which is impossible due to $\lambda_1(\nabla^2 g(\mathbf{x}_g^*)) = \lambda_1(\nabla^2 f(\mathbf{x}_g^*)) + \varepsilon > 0$. Hence, $\mathbf{x}_g^* \in \partial K$ and we have that

$$\max_{\mathbf{x} \in K} f(\mathbf{x}) \leq \max_{\mathbf{x} \in K} g(\mathbf{x}) = \max_{\mathbf{x} \in \partial K} g(\mathbf{x}) \leq \max_{\mathbf{x} \in \partial K} f(\mathbf{x}) + \frac{\varepsilon}{2} R^2,$$

where R is the radius of a ball containing K . Tending ε to 0 completes the proof. \square

B.4. Proof of Proposition 2.5

Proof. Indeed, by the Taylor theorem, we have that for any $\mathbf{h} \in V_{\mathbf{x}}$ and $\lambda > 0$ there exists $\xi \in [0, 1]$ such that

$$\begin{aligned} 0 &\leq f(\mathbf{x} + \lambda\mathbf{h}) - f(\mathbf{x}) - \lambda\langle \nabla f(\mathbf{x}), \mathbf{h} \rangle \\ &= \frac{\lambda^2}{2} \langle \nabla^2 f(\mathbf{x} + \lambda\xi\mathbf{h})\mathbf{h}, \mathbf{h} \rangle. \end{aligned}$$

Dividing both sides by λ^2 and taking the limit $\lambda \rightarrow 0$, we obtain that

$$\langle \nabla^2 f(\mathbf{x})\mathbf{h}, \mathbf{h} \rangle \geq 0, \quad \forall \mathbf{h} \in V_{\mathbf{x}}.$$

Therefore, since $\dim(V_{\mathbf{x}}) \geq \tau$, we have that $\lambda_{\tau}(\mathbf{x}) \geq 0$. □

C. Composite Optimization and Proofs for Section 3

C.1. Composite Formulation

The main results of our work can be generalized to a more broad family of Composite Optimization Problems of the following form:

$$\min_{\mathbf{x} \in Q} [F(\mathbf{x}) := f(\mathbf{x}) + \psi(\mathbf{x})], \quad (32)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is our smooth objective (which can be non-convex), and $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a *simple* closed convex function (e.g., indicator of a given closed convex set, or a regularizer). We set $Q := \text{dom } \psi \subseteq \mathbb{R}^n$. Thus, the main properties of f (the grade of non-convexity and the level of smoothness can be defined with respect to this convex set Q).

In case of presence of the composite part $\psi(\cdot)$ in our problem, iterations of our method should be modified. For a given point $\mathbf{x} \in Q$, gradient vector $\mathbf{g} \in \mathbb{R}^n$ and matrix $\mathbf{H} = \mathbf{H}^T \succeq 0$, we define the composite gradient step, as follows:

$$\text{CompositeStep}_{\mathbf{H},\alpha}(\mathbf{x}, \mathbf{g}) := \underset{\mathbf{y} \in Q}{\text{argmin}} \left\{ \langle \mathbf{g}, \mathbf{y} \rangle + \psi(\mathbf{y}) + \frac{1}{2} \langle (\mathbf{H} + \alpha\mathbf{I})(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \right\}, \quad (33)$$

where $\alpha \geq 0$ is the regularization parameter. In general, due to the presence of regularization term in (33), for $\alpha > 0$, the composite subproblem is *strongly convex*, and we can employ the fast linear convergence of gradient methods as applied to (33), for computing this step inexactly.

With these modifications, we are ready to present a composite version of our algorithm for solving (32):

Composite Gradient Method with Spectral Preconditioning
<p>Choose $\mathbf{x}_0 \in Q$ and $0 \leq \tau \leq n$.</p> <p>For $k \geq 0$ iterate:</p> <ol style="list-style-type: none"> 1. Estimate $\mathbf{H}_k \approx \nabla_{\tau}^2 f(\mathbf{x}_k) \in \mathbb{R}^{n \times n}$ 2. Perform the composite gradient step, for some $\alpha_k \geq 0$: $\mathbf{x}_{k+1} = \text{CompositeStep}_{\mathbf{H}_k, \alpha_k}(\mathbf{x}_k, \nabla f(\mathbf{x}_k))$

(34)

Note that when $\psi \equiv 0$, composite step (33) coincides with our basic preconditioned gradient step defined earlier:

$$\text{CompositeStep}_{\mathbf{H},\alpha}(\mathbf{x}, \mathbf{g}) = \text{GradStep}_{\mathbf{H},\alpha}(\mathbf{x}, \mathbf{g}) = \mathbf{x} - (\mathbf{H} + \alpha\mathbf{I})^{-1} \mathbf{g},$$

and method (34) is identical to (14).

C.2. Convergence Analysis

In this section, we provide the proofs of our main convergence results from Section 3. We study the more general composite formulation (34) of our method, which covers the basic case when $\psi \equiv 0$.

Let us consider one step of our method: $\mathbf{x}^+ = \text{CompositeStep}_{\mathbf{H}, \alpha}(\mathbf{x}, \nabla f(\mathbf{x}))$, for some $\alpha > 0$, and establish its key properties. The new point \mathbf{x}^+ satisfies the following optimality condition (see, e.g., Theorem 3.1.23 in (Nesterov, 2018)):

$$\langle \nabla f(\mathbf{x}) + (\mathbf{H} + \alpha \mathbf{I})(\mathbf{x}^+ - \mathbf{x}), \mathbf{y} - \mathbf{x}^+ \rangle + \psi(\mathbf{y}) \geq \psi(\mathbf{x}^+), \quad \forall \mathbf{y} \in Q. \quad (35)$$

In other words, the vector $\psi'(\mathbf{x}^+) := -\nabla f(\mathbf{x}) - (\mathbf{H} + \alpha \mathbf{I})(\mathbf{x}^+ - \mathbf{x})$ belongs to the subdifferential of ψ at new point:

$$\psi'(\mathbf{x}^+) \in \partial\psi(\mathbf{x}^+).$$

We denote

$$F'(\mathbf{x}^+) := \nabla f(\mathbf{x}^+) + \psi'(\mathbf{x}^+), \quad (36)$$

that is the main object for which we prove the convergence of our method. Utilizing positive semi-definiteness of $\mathbf{H} = \mathbf{H}^\top \succeq 0$, we can bound the length of our displacement $r \stackrel{\text{def}}{=} \|\mathbf{x}^+ - \mathbf{x}\|$.

Lemma C.1. *For any $\mathbf{s} \in \partial\psi(\mathbf{x})$, we have*

$$r \leq \frac{\|\nabla f(\mathbf{x}) + \mathbf{s}\|}{\alpha}, \quad (37)$$

and

$$\langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \leq r \|\nabla f(\mathbf{x}) + \mathbf{s}\|. \quad (38)$$

Proof. Indeed, using convexity of ψ , we obtain, for any $\mathbf{s} \in \partial\psi(\mathbf{x})$ and our specific $\psi'(\mathbf{x}^+) \in \partial\psi(\mathbf{x}^+)$:

$$\begin{aligned} 0 &\leq \langle \psi'(\mathbf{x}^+) - \mathbf{s}, \mathbf{x}^+ - \mathbf{x} \rangle \\ &= \langle -\nabla f(\mathbf{x}) - (\mathbf{H} + \alpha \mathbf{I})(\mathbf{x}^+ - \mathbf{x}) - \mathbf{s}, \mathbf{x}^+ - \mathbf{x} \rangle \\ &= \langle \nabla f(\mathbf{x}) + \mathbf{s}, \mathbf{x} - \mathbf{x}^+ \rangle - \langle (\mathbf{H} + \alpha \mathbf{I})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle. \end{aligned}$$

Rearranging the terms and using Cauchy-Schwartz inequality, we get

$$r \|\nabla f(\mathbf{x}) + \mathbf{s}\| \geq \langle \nabla f(\mathbf{x}) + \mathbf{s}, \mathbf{x} - \mathbf{x}^+ \rangle \geq \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \alpha r^2.$$

Taking into account that $\mathbf{H} \succeq \mathbf{0}$ completes the proof. \square

Therefore, by choosing regularization parameter α appropriately, we can control the length of steps for our algorithm. When combined with smoothness properties (17), (18) of the objective, we can establish the global progress in terms of the objective function value. We denote

$$\delta := \|\mathbf{H} - \nabla_\tau^2 f(\mathbf{x})\|.$$

Lemma C.2. *Let $\mathbf{s} \in \partial\psi(\mathbf{x})$ be some subgradient of ψ at current point \mathbf{x} and let $\alpha \geq \sqrt{\frac{L\|\nabla f(\mathbf{x}) + \mathbf{s}\|}{3}} + \sigma_\tau + \delta$. Then*

$$F(\mathbf{x}) - F(\mathbf{x}^+) \geq \frac{\alpha}{2} r^2. \quad (39)$$

Proof. By Lipschitz continuity of the Hessian, we have

$$\begin{aligned} F(\mathbf{x}^+) &= f(\mathbf{x}^+) + \psi(\mathbf{x}^+) \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{6} r^3 + \psi(\mathbf{x}^+) \\ &\stackrel{(37)}{\leq} f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L\|\nabla f(\mathbf{x}) + \mathbf{s}\|}{6\alpha} r^2 + \psi(\mathbf{x}^+) \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} (\sigma_\tau + \delta + \frac{L\|\nabla f(\mathbf{x}) + \mathbf{s}\|}{3\alpha}) r^2 + \psi(\mathbf{x}^+), \end{aligned}$$

where $\sigma_\tau := \sup_{\mathbf{x}} \|\nabla^2 f(\mathbf{x}) - \nabla_\tau^2 f(\mathbf{x})\|$ as defined in (18). According to (35), it holds

$$\psi(\mathbf{x}^+) \leq \psi(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle - \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle - \alpha r^2.$$

Hence, substituting this inequality into the previous one, we can continue as follows:

$$\begin{aligned} F(\mathbf{x}^+) &\leq F(\mathbf{x}) - \frac{\alpha}{2} r^2 - \frac{1}{2} \left(\alpha - \sigma_\tau - \delta - \frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{3\alpha} \right) r^2 - \frac{1}{2} \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \\ &\leq F(\mathbf{x}) - \frac{\alpha}{2} r^2 - \frac{1}{2} \left(\alpha - \sigma_\tau - \delta - \frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{3\alpha} \right) r^2. \end{aligned}$$

To prove the result, it suffices to check $\alpha \geq \sigma_\tau + \delta + \frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{3\alpha}$, which is ensured by our choice. Indeed, denoting $B := \sigma_\tau + \delta$ and $C := \frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{2}$, the inequality which we need to ensure for $\alpha > 0$ is

$$\alpha^2 \geq B + \frac{C}{\alpha} \Leftrightarrow \alpha^2 - \alpha B - C \geq 0 \Leftrightarrow \alpha \geq \frac{B + \sqrt{B^2 + 4C}}{2}. \quad (40)$$

It is immediate to check that $\alpha := B + \sqrt{C}$ satisfies (40), hence it holds for any $\alpha \geq B + \sqrt{C}$, which is our choice in the condition of the lemma. \square

Now, let us relate the length r of the step with the norm of vector $F'(\mathbf{x}^+) := \nabla f(\mathbf{x}^+) + \psi'(\mathbf{x}^+)$.

Lemma C.3. *Let $\mathbf{s} \in \partial\psi(\mathbf{x})$ be some subgradient of ψ at current point \mathbf{x} and let $\alpha \geq \sqrt{\frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{2}} + \sigma_\tau + \delta$. Then*

$$r \geq \frac{1}{2\alpha} \|F'(\mathbf{x}^+)\|. \quad (41)$$

Proof. Using the definition of $\psi'(\mathbf{x}^+) := -\nabla f(\mathbf{x}) - (\mathbf{H} + \alpha\mathbf{I})(\mathbf{x}^+ - \mathbf{x})$ and Lipschitz continuity of the Hessian, we have

$$\begin{aligned} \|F'(\mathbf{x}^+)\| &= \|\nabla f(\mathbf{x}^+) + \psi'(\mathbf{x}^+)\| \\ &= \|\nabla f(\mathbf{x}^+) - \nabla f(\mathbf{x}) - (\mathbf{H} + \alpha\mathbf{I})(\mathbf{x}^+ - \mathbf{x})\| \\ &\leq \|\nabla f(\mathbf{x}^+) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x})\| + (\alpha + \sigma_\tau + \delta)r \\ &\stackrel{(17)}{\leq} \frac{L}{2} r^2 + (\alpha + \sigma_\tau + \delta)r \\ &\stackrel{(37)}{\leq} \left(\alpha + \sigma_\tau + \delta + \frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{2\alpha} \right) r \leq 2\alpha r, \end{aligned}$$

where the last inequality holds due to $\alpha \geq \sigma_\tau + \delta + \frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{2\alpha}$, which is ensured by our choice (see the end of the proof of Lemma C.2). \square

Therefore, combining these two Lemmas together, we obtain the following bound for one step of the method.

Corollary C.4. *Let $\alpha \geq \sqrt{\frac{L \|\nabla f(\mathbf{x}) + \mathbf{s}\|}{2}} + \sigma_\tau + \delta$ for some $\mathbf{s} \in \partial\psi(\mathbf{x})$. Then*

$$F(\mathbf{x}) - F(\mathbf{x}^+) \stackrel{(39),(41)}{\geq} \frac{1}{8\alpha} \|F'(\mathbf{x}^+)\|^2. \quad (42)$$

We are ready to prove the global complexity bound for convergence of our method.

Theorem C.5. *Let $f \in \mathcal{F}_\tau$ be non-convex of grade τ , where $0 \leq \tau \leq n$ is fixed. Let f have a Lipschitz Hessian with constant L and bounded parameter σ_τ . Consider iterations $\{\mathbf{x}_k\}_{k \geq 0}$ of algorithm (34) with*

$$\alpha_k = \sqrt{\frac{L \|F'(\mathbf{x}_k)\|}{2}} + \sigma_\tau + \delta,$$

where $F'(\mathbf{x}_k)$ is defined by (36) for $k \geq 1$, and $F'(\mathbf{x}_0) := \nabla f(\mathbf{x}_0) + \mathbf{s}$, for an arbitrary initial subgradient $\mathbf{s} \in \partial\psi(\mathbf{x}_0)$. Then, for any $\varepsilon > 0$, it is enough to do

$$K = \left\lceil 8(F(\mathbf{x}_0) - F^*) \cdot \left(\sqrt{\frac{L}{2}} \frac{1}{\varepsilon^{3/2}} + \frac{\sigma_\tau + \delta}{\varepsilon^2} \right) + 2 \ln \frac{\|F'(\mathbf{x}_0)\|}{\varepsilon} \right\rceil. \quad (43)$$

iterations to have $\min_{1 \leq i \leq K} \|F'(\mathbf{x}_i)\| \leq \varepsilon$.

Proof. Let us fix some $k \geq 0$ and assume that for any $0 \leq i \leq k$ we have $g_i \stackrel{\text{def}}{=} \|F'(\mathbf{x}_i)\| \geq \varepsilon$.

According to (42), we obtain, for $0 \leq i \leq k-1$:

$$\begin{aligned} F(\mathbf{x}_i) - F(\mathbf{x}_{i+1}) &\geq \frac{1}{8\alpha_k} g_{i+1}^2 = \frac{1}{8} \left[\frac{g_{i+1}}{g_i} \right]^2 \cdot \left(\sqrt{\frac{L}{2}} \frac{1}{g_{i+1}^{3/2}} + \frac{\sigma_\tau + \delta}{g_{i+1}^2} \right)^{-1} \\ &\geq \frac{1}{8} \left[\frac{g_{i+1}}{g_i} \right]^2 \cdot \left(\sqrt{\frac{L}{2}} \frac{1}{\varepsilon^{3/2}} + \frac{\sigma_\tau + \delta}{\varepsilon^2} \right)^{-1}. \end{aligned} \quad (44)$$

Denote

$$c := \frac{1}{8} \left(\sqrt{\frac{L}{2}} \frac{1}{\varepsilon^{3/2}} + \frac{\sigma_\tau + \delta}{\varepsilon^2} \right)^{-1}.$$

Then, telescoping bound (44) and using the inequality between arithmetic and geometric means, we get

$$\begin{aligned} F(\mathbf{x}_0) - F^* &\geq F(\mathbf{x}_0) - F(\mathbf{x}_k) \stackrel{(44)}{\geq} c \sum_{i=0}^{k-1} \left[\frac{g_{i+1}}{g_i} \right]^2 \geq ck \left[\prod_{i=0}^{k-1} \frac{g_{i+1}}{g_i} \right]^{2/k} \\ &= ck \left[\frac{g_k}{g_0} \right]^{2/k} \geq ck \left[\frac{\varepsilon}{g_0} \right]^{2/k} = ck \exp \left[-\frac{2}{k} \ln \frac{g_0}{\varepsilon} \right] \\ &\geq ck \left[1 - \frac{2}{k} \ln \frac{g_0}{\varepsilon} \right] = ck - 2c \ln \frac{g_0}{\varepsilon}. \end{aligned}$$

Hence, we obtain

$$k \leq \frac{F(\mathbf{x}_0) - F^*}{c} + 2 \ln \frac{g_0}{\varepsilon}.$$

Substituting the value of c completes the proof. \square

C.3. Proof of Theorem 4.2

Proof. It follows immediately from Theorem C.5 by substituting the non-composite case $\psi \equiv 0$. \square

D. Proofs for Section 5

The results of this section are applied to a basic unconstrained minimization problem (4). Let us consider one iteration $\mathbf{x} \mapsto \mathbf{x}^+$ of our method, which satisfies the following stationary condition

$$\nabla f(\mathbf{x}) + \mathbf{H}(\mathbf{x}^+ - \mathbf{x}) + \alpha \mathbf{I} = \mathbf{0}, \quad (45)$$

where regularization parameter is chosen as

$$\alpha := \alpha^* + \eta, \quad (46)$$

with

$$\alpha^* := \operatorname{argmax}_{\alpha > 0} \left[-\frac{1}{2} \langle (\mathbf{H} + (\alpha + \eta)\mathbf{I})^{-1} \nabla f(\mathbf{x}), \nabla f(\mathbf{x}) \rangle - \frac{2\alpha^3}{3L} \right], \quad (47)$$

and $\eta := \sigma_\tau^+ + \delta + \delta_- \geq 0$ is the balancing term to control the errors. Considering the decomposition of the Hessian onto positive and negative components:

$$\nabla^2 f(\mathbf{x}) \equiv \nabla_+^2 f(\mathbf{x}) - \nabla_-^2 f(\mathbf{x}), \quad \nabla_+^2 f(\mathbf{x}), \nabla_-^2 f(\mathbf{x}) \succeq \mathbf{0}, \quad (48)$$

we have that

$$\begin{aligned} \|\nabla_+^2 f(\mathbf{x}) - \nabla_\tau^2 f(\mathbf{x})\| &\leq \sigma_\tau^+ \\ \|\mathbf{H} - \nabla_\tau^2 f(\mathbf{x})\| &\leq \delta, \\ \|\nabla_-^2 f(\mathbf{x})\mathbf{P}\| &\leq \delta_-, \end{aligned} \quad (49)$$

where \mathbf{P} is a given projector onto image of \mathbf{H} , that satisfies

$$\mathbf{HP} = \mathbf{H}. \quad (50)$$

First, let us provide another description of regularization parameter α^* , that is more suitable for our analysis. First-order stationary condition for (47) gives

$$\frac{1}{2} \langle (\mathbf{H} + \alpha \mathbf{I})^{-1} \nabla f(\mathbf{x}), (\mathbf{H} + \alpha \mathbf{I})^{-1} \nabla f(\mathbf{x}) \rangle - \frac{2(\alpha^*)^2}{L} = 0,$$

which is equivalent to

$$\alpha^* = \frac{L}{2} \|(\mathbf{H} + \alpha \mathbf{I})^{-1} \nabla f(\mathbf{x})\| \stackrel{(45)}{=} \frac{L}{2} \|\mathbf{x}^+ - \mathbf{x}\| \equiv \frac{Lr}{2}, \quad (51)$$

where we use $r := \|\mathbf{x}^+ - \mathbf{x}\|$ as previously. Therefore, from (51) we see that α^* plays the role of the Cubic Regularization (Nesterov & Polyak, 2006) of our model.

Let us establish the main inequalities on the progress of each step.

Lemma D.1. *For the functional residual, it holds*

$$f(\mathbf{x}) - f(\mathbf{x}^+) \geq \frac{L}{3} r^3 + \frac{\eta}{2} r^2 \quad (52)$$

Proof. Indeed, using Lipschitz continuity of the Hessian and definition of our step, we obtain

$$\begin{aligned} f(\mathbf{x}^+) &\stackrel{(17),(48)}{\leq} f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle \nabla_+^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle - \frac{1}{2} \langle \nabla_-^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{6} r^3 \\ &\leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{1}{2} \langle \nabla_+^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{6} r^3 \\ &\stackrel{(35)}{=} f(\mathbf{x}) - \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle - \alpha r^2 + \frac{1}{2} \langle \nabla_+^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle + \frac{L}{6} r^3 \\ &\leq -\alpha r^2 + \frac{\sigma_r^+ + \delta}{2} r^2 + \frac{L}{6} r^3. \end{aligned}$$

Hence, rearranging the terms and using the definition of α , we obtain

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}^+) &\geq \alpha r^2 - \frac{\sigma_r^+ + \delta}{2} r^2 - \frac{L}{6} r^3 \\ &\stackrel{(46),(51)}{=} \frac{L}{2} r^3 + \eta r^2 - \frac{\sigma_r^+ + \delta}{2} r^2 - \frac{L}{6} r^3 \\ &\geq \frac{L}{3} r^3 + \frac{\eta}{2} r^2, \end{aligned}$$

which is the required bound. \square

Lemma D.2. *Let $\mathbf{y}^+ := \mathbf{x} + \mathbf{P}(\mathbf{x}^+ - \mathbf{x})$. Then, we can relate the gradient at \mathbf{y}^+ and the length r of the step, as follows*

$$\|\nabla f(\mathbf{y}^+)\| \leq 2\eta r + Lr^2. \quad (53)$$

Proof. By Lipschitz continuity of the Hessian, we have

$$\begin{aligned} \frac{L}{2} r^2 &\geq \frac{L}{2} \|\mathbf{y}^+ - \mathbf{x}\|^2 \\ &\stackrel{(17)}{\geq} \|\nabla f(\mathbf{y}^+) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{y}^+ - \mathbf{x})\| \\ &= \|\nabla f(\mathbf{y}^+) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})\mathbf{P}(\mathbf{x}^+ - \mathbf{x})\| \\ &\stackrel{(35)}{=} \|\nabla f(\mathbf{y}^+) + (\mathbf{H} - \nabla^2 f(\mathbf{x}))\mathbf{P}(\mathbf{x}^+ - \mathbf{x}) + \alpha(\mathbf{x}^+ - \mathbf{x})\| \\ &\geq \|\nabla f(\mathbf{y}^+)\| - \|(\mathbf{H} - \nabla^2 f(\mathbf{x}))\mathbf{P}(\mathbf{x}^+ - \mathbf{x})\| - \alpha r. \end{aligned}$$

Hence, rearranging the terms, we obtain

$$\begin{aligned} \|\nabla f(\mathbf{y}^+)\| &\leq \|(\mathbf{H} - \nabla^2 f(\mathbf{x}))\mathbf{P}(\mathbf{x}^+ - \mathbf{x})\| + \alpha r + \frac{L}{2}r^2 \\ &\stackrel{(48)}{=} \|(\mathbf{H} - \nabla_+^2 f(\mathbf{x}) + \nabla_-^2 f(\mathbf{x}))\mathbf{P}(\mathbf{x}^+ - \mathbf{x})\| + \alpha r + \frac{L}{2}r^2 \\ &\leq 2\eta r + Lr^2, \end{aligned}$$

which is the required bound. \square

Combining inequalities (52) and (53) together, we obtain the following bound on the progress.

Corollary D.3. *For one step of the method, it holds:*

$$f(\mathbf{x}) - f(\mathbf{x}^+) \geq \min\left\{\frac{1}{6\sqrt{2L}}\|\nabla f(\mathbf{y}^+)\|^3, \frac{1}{32\eta}\|\nabla f(\mathbf{y}^+)\|^2\right\}. \quad (54)$$

D.1. Proof of Theorem 5.1

Proof. Let us fix some $k \geq 1$ and assume that for any $1 \leq i \leq k$ we have $\|\nabla f(\mathbf{y}_i)\| \geq \varepsilon$.

According to (54), we obtain, for $0 \leq i \leq k-1$:

$$\begin{aligned} f(\mathbf{x}_i) - f(\mathbf{x}_{i+1}) &\geq \min\left\{\frac{1}{6\sqrt{2L}}\|\nabla f(\mathbf{y}_{i+1})\|^3, \frac{1}{32\eta}\|\nabla f(\mathbf{y}_{i+1})\|^2\right\} \\ &\geq \min\left\{\frac{1}{6\sqrt{2L}}\varepsilon^3, \frac{1}{32\eta}\varepsilon^2\right\}. \end{aligned}$$

Telescoping this bound for the first k iterations, we get

$$f(\mathbf{x}_0) - f^* \geq f(\mathbf{x}_0) - f(\mathbf{x}_k) \geq k \min\left\{\frac{1}{6\sqrt{2L}}\varepsilon^3, \frac{1}{32\eta}\varepsilon^2\right\},$$

which leads to the required complexity. \square

E. Proofs for Section 6

In this section, we provide a general analysis for the *composite* version of our method (34), when the target objective is convex: $f \in \mathcal{F}_n$ and quasi-Self-Concordant (25) with parameter $M > 0$.

Let us establish the progress for one step of our algorithm $\mathbf{x}^+ = \text{CompositeStep}_{\mathbf{H}, \alpha}(\mathbf{x}, \nabla f(\mathbf{x}))$ under this refined smoothness condition.

Lemma E.1. *Let $\alpha \geq M\|F'(\mathbf{x})\| + \sigma_\tau + \delta$. Then*

$$\langle F'(\mathbf{x}^+), \mathbf{x} - \mathbf{x}^+ \rangle \geq \frac{1}{2\alpha}\|F'(\mathbf{x}^+)\|^2. \quad (55)$$

Proof. By using basic properties of quasi-Self-Concordant functions (see Lemma 2.7 in (Doikov, 2023)), we have, for any two points $\mathbf{x}, \mathbf{x}^+ \in \mathbb{R}^n$:

$$\|\nabla f(\mathbf{x}^+) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x})\| \leq \frac{M}{2}\langle \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \cdot \varphi(Mr),$$

where $\varphi(t) \stackrel{\text{def}}{=} \frac{e^t - t - 1}{t^2}$ is a convex monotone univariate function. Let us show how we can control the right hand side. From (37), we have

$$\varphi(Mr) \leq \varphi\left(\frac{M\|F'(\mathbf{x})\|}{\alpha}\right) \stackrel{(*)}{\leq} 1,$$

where in $(*)$ we use our choice of the regularization coefficient: $\alpha \geq M\|F'(\mathbf{x})\|$.

For the Hessian, we can use

$$\begin{aligned} \langle \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle &\stackrel{(18)}{\leq} \sigma_\tau r^2 + \langle \mathbf{H}(\mathbf{x}^+ - \mathbf{x}), \mathbf{x}^+ - \mathbf{x} \rangle \stackrel{(38)}{\leq} \sigma_\tau r^2 + r \|F'(\mathbf{x})\| \\ &\stackrel{(37)}{\leq} \left(\frac{\sigma_\tau}{\alpha} + 1\right) \cdot r \|F'(\mathbf{x})\| \stackrel{(**)}{\leq} 2r \|F'(\mathbf{x})\|, \end{aligned}$$

where we used in (**) that $\alpha \geq \sigma_\tau$.

Therefore, we obtain

$$\|\nabla f(\mathbf{x}^+) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x})\| \leq Mr \|F'(\mathbf{x})\|. \quad (56)$$

Thus, using the definition of new subgradient $F'(\mathbf{x}^+)$ from the stationary condition (35), we have

$$\begin{aligned} \|F'(\mathbf{x}^+) + \alpha(\mathbf{x}^+ - \mathbf{x})\| &= \|\nabla f(\mathbf{x}^+) - \nabla f(\mathbf{x}) - \mathbf{H}(\mathbf{x}^+ - \mathbf{x})\| \\ &\leq \|\nabla f(\mathbf{x}^+) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{x}^+ - \mathbf{x})\| \\ &\quad + r \|\mathbf{H} - \nabla^2 f(\mathbf{x})\| \\ &\stackrel{(56)}{\leq} r(\sigma_\tau + \delta + M \|F'(\mathbf{x})\|). \end{aligned}$$

Taking square of both sides and rearranging the terms, we obtain

$$\langle F'(\mathbf{x}^+), \mathbf{x} - \mathbf{x}^+ \rangle \geq \frac{1}{2\alpha} \|F'(\mathbf{x}^+)\|^2 + \frac{r^2}{2\alpha} (\alpha^2 - (\sigma_\tau + \delta + M \|\nabla f(\mathbf{x})\|)^2).$$

Taking into account our choice of α completes the proof. \square

We denote by D the diameter of the initial sublevel set which we assume to be bounded:

$$D := \sup_{\mathbf{x}} \left\{ \|\mathbf{x} - \mathbf{x}^*\| : F(\mathbf{x}) \leq F(\mathbf{x}_0) \right\} < +\infty.$$

We prove the following result.

Theorem E.2. *Let $f \in \mathcal{F}_n$ be convex and quasi-Self-Concordant with constant M . Consider iterations $\{\mathbf{x}_k\}_{k \geq 0}$ of algorithm (34) with*

$$\alpha_k = M \|F'(\mathbf{x}_k)\| + \sigma_\tau + \delta,$$

where $F'(\mathbf{x}_k)$ is defined by (36) for $k \geq 1$, and $F'(\mathbf{x}_0) := \nabla f(\mathbf{x}_0) + \mathbf{s}$, for an arbitrary initial subgradient $\mathbf{s} \in \partial\psi(\mathbf{x}_0)$. Then, for any $\varepsilon > 0$, it is enough to do

$$K = 4 \left\lceil \left(MD + \frac{\sigma_\tau + \delta}{2\mu} \right) \ln \frac{F(\mathbf{x}_0) - F^*}{\varepsilon} + \ln \frac{\|F'(\mathbf{x}_0)\| D}{\varepsilon} \right\rceil \quad (57)$$

steps to ensure $F(\mathbf{x}_k) - F^* \leq \varepsilon$.

Proof. Let us prove the following rate of convergence, for the iterations of our method,

$$F(\mathbf{x}_k) - F^* \leq \exp\left(-\frac{k}{4} \left[\frac{\sigma_\tau + \delta}{2\mu} + MD \right]^{-1}\right) (F(\mathbf{x}_0) - F^*) + \exp\left(-\frac{k}{4}\right) \cdot \|F'(\mathbf{x}_0)\| D, \quad (58)$$

which immediately leads to the complexity bound (57).

We denote $g_k = \|F'(\mathbf{x}_k)\|$. Then, by convexity, we have for every iteration $k \geq 0$:

$$\begin{aligned} F(\mathbf{x}_k) - F(\mathbf{x}_{k+1}) &\geq \langle F'(\mathbf{x}_{k+1}), \mathbf{x}_k - \mathbf{x}_{k+1} \rangle \\ &\stackrel{(55)}{\geq} \frac{1}{2(\sigma_\tau + \delta + M g_k)} \left(\frac{g_{k+1}}{g_k}\right)^2 g_k^2 \\ &= \frac{1}{2} \left(\frac{g_{k+1}}{g_k}\right)^2 \cdot \left(\frac{\sigma_\tau + \delta}{g_k^2} + \frac{M}{g_k}\right)^{-1}. \end{aligned} \quad (59)$$

Denoting the functional residual by $F_k \stackrel{\text{def}}{=} F(\mathbf{x}_k) - F^*$, we have by convexity and strong convexity:

$$g_k \geq \frac{F_k}{D}, \quad g_k^2 \geq 2\mu F_k.$$

Substituting these bounds into (59), we obtain

$$F_k - F_{k+1} \geq \frac{1}{2} \left(\frac{g_{k+1}}{g_k} \right)^2 \cdot \left(\frac{\sigma_\tau + \delta}{2\mu F_k} + \frac{MD}{F_k} \right)^{-1} = \frac{q}{2} \left(\frac{g_{k+1}}{g_k} \right)^2 F_k, \quad (60)$$

where condition number q is defined by

$$q := \left(\frac{\sigma_\tau + \delta}{2\mu} + MD \right)^{-1}.$$

To show the desired rate, we use concavity of logarithm,

$$\ln \frac{a}{b} = \ln a - \ln b \geq \frac{1}{a}(a - b), \quad \forall a, b > 0,$$

and conclude that

$$\ln \frac{F_k}{F_{k+1}} \geq \frac{1}{F_k}(F_k - F_{k+1}) \stackrel{(60)}{\geq} \frac{q}{2} \left(\frac{g_{k+1}}{g_k} \right)^2.$$

Telescoping this bound and using inequality between arithmetic and geometric means, we get

$$\begin{aligned} \ln \frac{F_0}{F_k} &\geq \frac{q}{2} \sum_{i=0}^{k-1} \left[\frac{g_{i+1}}{g_i} \right]^2 \geq \frac{kq}{2} \left[\prod_{i=0}^{k-1} \frac{g_{i+1}}{g_i} \right]^{2/k} = \frac{kq}{2} \left[\frac{g_k}{g_0} \right]^{2/k} \\ &= \frac{kq}{2} \exp \left[\frac{2}{k} \ln \frac{g_k}{g_0} \right] \geq \frac{kq}{2} \left(1 + \frac{2}{k} \ln \frac{g_k}{g_0} \right) \geq \frac{kq}{2} \left(1 + \frac{2}{k} \ln \frac{F_k}{g_0 D} \right). \end{aligned}$$

The last inequality leads to the required bound (58). □

E.1. Proof of Theorem 6.1

Proof. It follows immediately from Theorem E.2 by substituting the non-composite case $\psi \equiv 0$. □