
Polynomial Preconditioning for Gradient Methods

Nikita Doikov¹ Anton Rodomanov²

Abstract

We study first-order methods with preconditioning for solving structured nonlinear convex optimization problems. We propose a new family of preconditioners generated by symmetric polynomials. They provide first-order optimization methods with a provable improvement of the condition number, cutting the gaps between highest eigenvalues, without explicit knowledge of the actual spectrum. We give a stochastic interpretation of this preconditioning in terms of coordinate volume sampling and compare it with other classical approaches, including the Chebyshev polynomials. We show how to incorporate a polynomial preconditioning into the Gradient and Fast Gradient Methods and establish the corresponding global complexity bounds. Finally, we propose a simple adaptive search procedure that automatically chooses the best possible polynomial preconditioning for the Gradient Method, minimizing the objective along a low-dimensional Krylov subspace. Numerical experiments confirm the efficiency of our preconditioning strategies for solving various machine learning problems.

1. Introduction

Motivation. Preconditioning is an important tool for improving the performance of numerical algorithms. The classical example is the preconditioned *Conjugate Gradient Method* (Hestenes & Stiefel, 1952) for solving a system of linear equations. It proposes to modify the initial system in a way to improve its eigenvalue distribution and thus to accelerate the convergence of the method. The question of choosing the right preconditioner heavily depends on the problem structure, and there exist many problem-specific recommendations which provide us with a good trade-off between computational cost and the spectrum properties of

¹EPFL, Switzerland ²UCLouvain, Belgium. Correspondence to: Nikita Doikov <nikita.doikov@epfl.ch>, Anton Rodomanov <anton.rodomanov@uclouvain.be>.

the new system. Some notable examples include *Jacobi* or the *diagonal* preconditioners, *symmetric successive over-relaxation*, the *incomplete Cholesky* factorization (Golub & Van Loan, 2013), *Laplacian* preconditioning for graph problems (Vaidya, 1991; Spielman & Teng, 2004), preconditioners for discretizations of system of *partial differential equations* (Mardal & Winther, 2011).

Another important class of numerical algorithms are the second-order methods or *Newton’s Method* (see, e.g. (Nesterov, 2018)), that aims to solve difficult ill-conditioned problems by using local curvature information (the Hessian matrix) as a preconditioner at every step. However, being a powerful algorithm, each iteration of Newton’s Method is very expensive. It requires to solve a system of linear equations with the Hessian matrix, and in case of quadratic objective it is equivalent to solving the original problem.

In this paper, our goal is to solve a general *nonlinear optimization problem* with a structured convex objective by the efficient first-order methods. Thus, in the case of unconstrained minimization of a smooth function: $\min_{\mathbf{x}} f(\mathbf{x})$, the simplest method that we study is as follows, for $k \geq 0$:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{P} \nabla f(\mathbf{x}_k), \quad (1)$$

where $\alpha_k > 0$ is a stepsize and \mathbf{P} is a fixed preconditioning matrix. $\mathbf{P} := \mathbf{I}$ corresponds to the classical gradient descent. Another natural choice is $\mathbf{P} := \mathbf{B}^{-1}$, where \mathbf{B} is a *curvature matrix* of our problem¹, which is directly available for the algorithm. That resembles the Newton-type direction, and the method with this preconditioner tends to converge much faster in practice (see Figure 1). However, computing \mathbf{B}^{-1} (or solving the corresponding linear system with \mathbf{B}) is a very expensive operation in the large scale setting.

Instead of using \mathbf{B}^{-1} , we propose a new family of *Symmetric Polynomial Preconditioners*, that provably improve the spectrum of the objective. The first member of our family is

$$\mathbf{P} := \text{tr}(\mathbf{B})\mathbf{I} - \mathbf{B}. \quad (2)$$

We prove that using preconditioner (2) within method (1), makes the condition number *insensitive to the gap between*

¹See the definition of \mathbf{B} in our Assumption 2.1 and the corresponding Examples 2.2, 2.3, 2.4 of different problems.

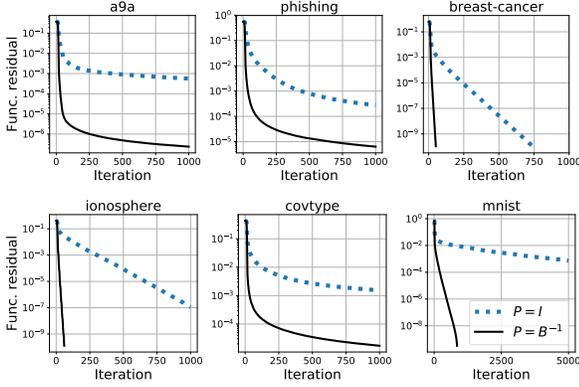


Figure 1: Training logistic regression with the standard gradient descent ($P = I$), and using the inverse of the curvature matrix ($P = B^{-1}$) as a preconditioner in (1). The latter method works much faster, while it can be very expensive to compute B^{-1} for large scale problems.

the top two eigenvalues of the curvature matrix. Since it is quite common for real data to have a highly nonuniform spectrum with several large gaps between the top eigenvalues (see Figure 2), our preconditioning can significantly improve the convergence of the first-order methods. At the same time, one step of the form (1),(2) is still cheap to compute. It involves just the standard matrix operations (trace and the matrix-vector product), without the need to solve linear systems with the curvature matrix as in Newton’s Method.

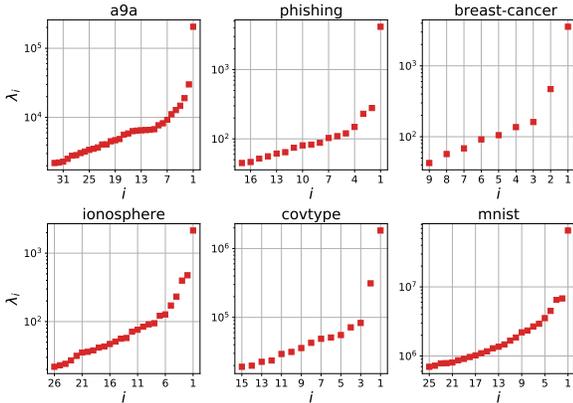


Figure 2: Leading eigenvalues (in the logarithmic scale) of the curvature matrix B , for several typical datasets². There are large gaps between the top eigenvalues.

This approach works for general structured nonlinear problems (not necessarily quadratics) and also for the problems with possible composite parts (e.g., constrained minimization or non-smooth regularization).

²<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

Our new family of Symmetric Polynomial Preconditioners gradually interpolate between the first preconditioner (2) and $P \propto B^{-1}$ as the other extreme case. We show that increasing the order of the preconditioner, we are able to cut off several top eigenvalues of the curvature matrix, without knowing the actual spectrum. We can incorporate these preconditioners both into the Gradient Method, as well as into the accelerated Fast Gradient Method (Nesterov, 1983), with a further provable improvement of the condition number.

Finally, we address the common question of choosing the best possible preconditioner. We propose a new adaptive strategy for the basic nonlinear Gradient Method based on the Krylov subspace minimization. In this approach, preconditioner P is defined as a general polynomial of the curvature matrix B of a fixed (small) degree τ :

$$P := a_0 I + a_1 B + \dots + a_\tau B^\tau,$$

where the vector of coefficients $\mathbf{a} \in \mathbb{R}^{\tau+1}$ is found by solving a certain linear system of size $\tau + 1$ in each iteration of the method. It has a plain interpretation of projecting the direction $B^{-1} \nabla f(\mathbf{x}_k)$ onto an affine set $\mathcal{K}_{\mathbf{x}_k}^\tau$, which is the Krylov subspace:

$$\mathcal{K}_{\mathbf{x}}^\tau \stackrel{\text{def}}{=} \text{span} \{ \nabla f(\mathbf{x}), B \nabla f(\mathbf{x}), \dots, B^\tau \nabla f(\mathbf{x}) \}. \quad (3)$$

In case of small τ , we can solve this linear system easily and obtain the best preconditioning guarantee for our method, which is adaptive for each iteration.

Related Work. It is widely known that the standard Conjugate Gradient Method is optimal in the class of the first-order algorithms for unconstrained minimization of convex quadratic functions (Nemirovski, 1995). The k th iteration of the Conjugate Gradients finds the full minimum of the objective over the k -dimensional Krylov subspace, and thus it provably solves the problem after $k = n$ iterations, where n is the dimension of the problem. Quadratic minimization is equivalent to solving a system of linear equations, therefore it is often referred as the linear case. Polynomial preconditioning for solving large linear systems has been extensively studied during the last several decades; see (Dubois et al., 1979; Johnson et al., 1983; Saad, 1985; Van Gijzen, 1995; Liu et al., 2015; Loe & Morgan, 2022) and references therein. See also Section 5.3 for the comparison of our preconditioning strategies with the linear Conjugate Gradient Method.

The situation with nonlinear problems is more difficult. Along with the basic Gradient Method, the classical approaches include the Nonlinear Conjugate Gradients and Quasi-Newton Methods (see, e.g. (Nocedal & Wright, 2006)), which typically demonstrate a decent practical performance, while replicating the standard Conjugate Gradients in the linear case. However, these methods lack of

Polynomial Preconditioning for Gradient Methods

	Preconditioner	Condition number, β/α	Methods	Cost
Classical Gradient Method	$\mathbf{P} = \mathbf{I}$	λ_1/λ_n	GM, FGM	cheap
"Full Preconditioning"	$\mathbf{P} = \mathbf{B}^{-1}$	1	GM, FGM	expensive
Symmetric Polynomial Preconditioning (ours)	$\mathbf{P} = \mathbf{P}_\tau$ (11)	$\lambda_1/\lambda_n \cdot \xi_\tau(\boldsymbol{\lambda})$	GM, FGM	cheap for small τ
Krylov Subspace Minimization (ours)	optimal polynomial.	$\lambda_{\tau+1}/\lambda_n$	GM	cheap for small τ

Table 1: The value β/α for different preconditioning strategies, $\boldsymbol{\lambda} = \boldsymbol{\lambda}(\mathbf{B})$. Note that $\xi_\tau(\boldsymbol{\lambda}) \leq 1$, and $\xi_\tau(\boldsymbol{\lambda}) \rightarrow 0$ in case of large spectral gaps, namely when $\lambda_1/\lambda_{\tau+1} \rightarrow \infty$ (see Section 4). For solving the problem with $\epsilon > 0$ accuracy, GM needs $k(\epsilon) = \mathcal{O}(\beta/\alpha \cdot 1/\epsilon)$ and $k(\epsilon) = \mathcal{O}(\beta/\alpha \cdot L/\mu \cdot \log 1/\epsilon)$ iterations for convex and strongly convex functions correspondingly. FGM needs only $\sqrt{k(\epsilon)}$ iterations (Theorems 3.1 and 3.2).

having any good *global complexity bounds*, and thus in the worst-case scenario they can actually perform even worse than the Gradient Method (Gupta et al., 2023). At the same time, the Fast Gradient Method developed by (Nesterov, 1983) is *optimal* for the class of nonlinear problems with a *uniformly bounded eigenvalues* of the Hessian (Nemirovski & Yudin, 1983). This assumption does not take into account the actual distribution of the spectrum. Hence, it can not distinguish the problems with large gaps between the top eigenvalues, as in Figure 2.

There have been several attempts to study more specific problem formulations, and so to gain a provable advantage for the optimization algorithms by leveraging the spectrum of the Hessian. Thus, the quadratic minimization problems were studied under the assumption of a particular *probability distribution* for the eigenvalues (Scieur & Pedregosa, 2020; Cunha et al., 2022), or assuming a certain fixed *spectral gap* (Goujaud et al., 2022), revealing the advantages of employing the Heavy-ball Method (Polyak, 1987) in these cases. Another example is the Stochastic Spectral Descent (Kovalev et al., 2018), which improves the condition number for quadratic problems if we know some of the eigenvectors.

In this work, we consider a refined smoothness characterization of the objective with the curvature matrix \mathbf{B} (Assumption 2.1). It is similar in spirit to that one used in Stochastic Dual Newton Ascent (Qu et al., 2016). An important particular instance of this class of algorithms is the Randomized Coordinate Descent with *Volume Sampling* (Rodomanov & Kropotov, 2020). In the latter method, it was proposed to select subsets of variables of certain size m proportionally to the determinants of principal submatrices of \mathbf{B} . While this approach was practically implementable only for the subsets of size $m = 1$ or 2 , it was shown that, in theory, the method is insensitive to the large spectral gap between the top $m - 1$ eigenvalues.

Surprisingly, our new family of the Symmetric Polynomial Preconditioners can be viewed as a *deterministic version* of the Volume Sampling technique (with $m = \tau + 1$ where τ is the degree of a preconditioning polynomial; preconditioner (2) corresponds to $\tau = 1$). Thus, we provide the Volume

Sampling with a novel deterministic interpretation, which also leads to new accelerated and composite optimization algorithms (see Section 4.3 for a detailed comparison).

Contributions. We propose several polynomial preconditioning strategies for first-order methods for solving a general composite convex optimization problem, and prove their better global complexity guarantees, specifically:

- We study the convergence of the basic Gradient Method (GM, Algorithm 1) and the accelerated Fast Gradient Method (FGM, Algorithm 2) with a general (arbitrarily fixed) preconditioning matrix. We introduce *two* condition numbers, that are designated to the different parts of the objective (L/μ for nonlinearity and β/α for the curvature matrix), and show that they serve as main complexity factors.
- We develop a new family of Symmetric Polynomial Preconditioners (Section 4). Combining them with the preconditioned Gradient Methods, we establish a significant improvement of the curvature condition number β/α in case of *large gaps* between the top eigenvalues of the matrix (see Table 1).
- Then, we propose a new adaptive procedure based on the Krylov subspace minimization (Algorithm 3) that achieves the *best polynomial* preconditioning. We present the guarantees we can get, including cutting off the top eigenvalues directly and by employing the Chebyshev polynomials, and compare this approach with the Symmetric Polynomial Preconditioning.
- Numerical experiments are provided.

2. Notation and Assumptions

We consider the following optimization problem given in the *composite* form:

$$F^* = \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ F(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{x}) + \psi(\mathbf{x}) \right\}, \quad (4)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable convex function which is the *main* part of the problem, and $\psi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper closed convex function that can be nondifferentiable but has a *simple* structure. For example, it can be an indicator of a convex set, or ℓ_1 -regularizer.

Additionally, we fix some symmetric positive-definite matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ (notation $\mathbf{B} = \mathbf{B}^\top \succ 0$). This matrix plays the key role in our characterization of the smoothness properties of f . Namely, we assume the following (considering for simplicity two-times differentiable functions):

Assumption 2.1. *The Hessian of f is uniformly bounded, for some constants $L \geq \mu \geq 0$:*

$$\mu \mathbf{B} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{B}, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (5)$$

Having fixed the matrix \mathbf{B} , we define the corresponding induced norm by $\|\mathbf{x}\|_{\mathbf{B}} \stackrel{\text{def}}{=} \langle \mathbf{B}\mathbf{x}, \mathbf{x} \rangle^{1/2}$, $\mathbf{x} \in \mathbb{R}^n$. Thus, matrix \mathbf{B} is responsible for fixing the coordinate system in the problem. Then, condition (5) can be rewritten in terms of the global lower and upper bound on the first-order approximation of f (Nesterov, 2018):

$$\begin{aligned} \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{B}}^2 &\leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \\ &\leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{B}}^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n. \end{aligned} \quad (6)$$

In what follows, we denote by $\boldsymbol{\lambda} = \boldsymbol{\lambda}(\mathbf{B}) \in \mathbb{R}^n$ the vector of eigenvalues for the matrix \mathbf{B} , sorted in a nonincreasing order: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$.

The classical example is $\mathbf{B} := \mathbf{I}$ (identity matrix). Then, condition (5) implies that the function f is (strongly) convex and has the Lipschitz continuous gradient. However, by choosing a specific \mathbf{B} , we tend to achieve a better granularity of the description of our problem class and thus to improve the convergence properties of the methods.

Example 2.2. *Let $\mathbf{a} \in \mathbb{R}^n$. Then, the quadratic function*

$$f(\mathbf{x}) = \frac{1}{2} \langle \mathbf{B}\mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{a}, \mathbf{x} \rangle,$$

satisfies condition (5) with $L = \mu = 1$.

We see that in this case, the so-called *condition number* L/μ is just 1, which means that preconditioning the Gradient Method (1) with the matrix $\mathbf{P} := \mathbf{B}^{-1}$ would give an immediate convergence to the solution. However, inverting the matrix is prohibitively expensive for large scale problems. Our aim is to find a suitable trade-off between improving the condition number and the arithmetic cost of algorithm steps. Let us consider the following important example which can be met in many practical applications.

Example 2.3. *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be a given data matrix, and $\mathbf{b} \in \mathbb{R}^m$ be a given vector. Denote,*

$$f(\mathbf{x}) = g(\mathbf{A}\mathbf{x} + \mathbf{b})$$

Then, the derivatives are as follows: $\nabla f(\mathbf{x}) = \mathbf{A}^\top \nabla g(\mathbf{A}\mathbf{x} + \mathbf{b})$ and $\nabla^2 f(\mathbf{x}) = \mathbf{A}^\top \nabla^2 g(\mathbf{A}\mathbf{x} + \mathbf{b}) \mathbf{A}$. Hence, assuming: $\mu \mathbf{I}_m \preceq \nabla^2 g(\mathbf{x}) \preceq L \mathbf{I}_m$, $\forall \mathbf{x}$, with some $L \geq \mu \geq 0$, condition (5) is satisfied³ with

$$\mathbf{B} := \mathbf{A}^\top \mathbf{A}.$$

At the same time, for $\mathbf{B} := \mathbf{I}_n$ (the standard Euclidean norm), the Lipschitz constant increases by the factor $\lambda_1(\mathbf{A}^\top \mathbf{A})$, which makes the problem extremely ill-conditioned.

A particular case of this example is *separable optimization*, or *generalized linear models* (Bishop, 2006), which covers the classical regression and classification models.

Example 2.4. *Let*

$$f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \phi(\langle \mathbf{a}_i, \mathbf{x} \rangle), \quad \mathbf{x} \in \mathbb{R}^n,$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a loss function satisfying: $\mu \leq \phi''(t) \leq L$, $\forall t \in \mathbb{R}$, with some $L \geq \mu \geq 0$. Then, forming the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ whose rows are $\mathbf{a}_1^\top, \dots, \mathbf{a}_m^\top$ and setting $\mathbf{B} := \mathbf{A}^\top \mathbf{A}$, condition (5) holds.

3. Preconditioned Gradient Methods

A natural intention would be to use the global upper bound (6) as a model for the smooth part of the objective. However, the direct minimization of such upper model requires to solve the linear system with the matrix \mathbf{B} , which can computationally unfeasible for large scale problems.

Instead, let us fix for our *preconditioner* some positive definite symmetric matrix $\mathbf{P} = \mathbf{P}^\top \succ 0$, which satisfies the following bound, for some $\alpha := \alpha(\mathbf{P})$ and $\beta := \beta(\mathbf{P}) \geq \alpha > 0$:

$$\alpha \mathbf{B}^{-1} \preceq \mathbf{P} \preceq \beta \mathbf{B}^{-1}. \quad (7)$$

We are going to use this matrix instead of \mathbf{B}^{-1} in our methods. For a fixed symmetric positive definite matrix \mathbf{P} and parameter $M > 0$, we denote the *gradient step* from a point $\mathbf{x} \in \text{dom } \psi$ along a *gradient direction* $\mathbf{g} \in \mathbb{R}^n$ by

$$\begin{aligned} \text{GradStep}_{M, \mathbf{P}}(\mathbf{x}, \mathbf{g}) \\ \stackrel{\text{def}}{=} \underset{\mathbf{y} \in \text{dom } \psi}{\text{argmin}} \left\{ \langle \mathbf{g}, \mathbf{y} \rangle + \psi(\mathbf{y}) + \frac{M}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathbf{P}^{-1}}^2 \right\}. \end{aligned}$$

This operation is well-defined since the objective function in the above minimization problem is strongly convex. We assume that both ψ and \mathbf{P} are reasonably simple so that the corresponding gradient step can be efficiently computed. An important case is $\psi = 0$ for which we have

³Here, we assume that $\mathbf{A}^\top \mathbf{A} \succ 0$ which is typically the case when $m \gg n$. Otherwise, we can reduce the dimensionality.

$\text{GradStep}_{M,\mathbf{P}}(\mathbf{x}, \mathbf{g}) = \mathbf{x} - \frac{1}{M}\mathbf{P}\mathbf{g}$. The latter expression can be efficiently computed whenever one can cheaply multiply the matrix \mathbf{P} by any vector.

3.1. Preconditioned Basic Gradient Method

First, we consider the basic first-order scheme shown in Algorithm 1 for solving the composite problem (4). For simplicity, in this section, we only present a version of this method with a fixed step size and assume that all necessary constants are known. An adaptive version of Algorithm 1 which does not have these limitations and is more efficient in practice can be found in Appendix C.

Algorithm 1 Preconditioned Basic Gradient Method

Input: $\mathbf{x}_0 \in \text{dom } \psi$, $\mathbf{P} = \mathbf{P}^\top \succ 0$, $M > 0$.
for $k = 0, 1, \dots$ **do**
 Compute $\mathbf{x}_{k+1} = \text{GradStep}_{M,\mathbf{P}}(\mathbf{x}_k, \nabla f(\mathbf{x}_k))$.
end for

For Algorithm 1, we can prove the following results.

Theorem 3.1. *Consider Algorithm 1 with $M = \beta L$. Then, at each iteration $k \geq 1$, we have*

$$F(\mathbf{x}_k) - F^* \leq \frac{\beta L \|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{B}}^2}{\alpha k}. \quad (8)$$

When $\mu > 0$, the convergence is linear: for all $k \geq 1$,

$$F(\mathbf{x}_k) - F^* \leq \left(1 - \frac{1}{4} \frac{\alpha \mu}{\beta L}\right)^k [F(\mathbf{x}_0) - F^*]. \quad (9)$$

We see that one of the principal complexity factors in the above estimates is the condition number β/α which depends on the choice of our preconditioner \mathbf{P} (see (7)). For the basic choice $\mathbf{P} = \mathbf{I}$, we have $\beta/\alpha = \lambda_1/\lambda_n$. However, as we show in the following sections, it is possible to use more efficient (and still quite cheap) preconditioners which improve this condition number.

3.2. Preconditioned Fast Gradient Method

Now let us consider an accelerated scheme shown in Algorithm 2. This algorithm is one of the standard variants of the Fast Gradient Method (FGM) known as the Method of Similar Triangles (see, e.g., Section 6.1.3 in (Nesterov, 2018)) but adapted to our assumptions (5) and (7).

As in other versions of FGM, to properly handle strongly convex problems, Algorithm 2 requires the knowledge of the strong convexity parameter α and μ (or, more precisely, their product $\rho = \alpha\mu$). For non-strongly convex problems, we can always choose $\alpha = \mu = 0$. See also Appendix C for a variant of Algorithm 2 which can automatically adjust the constant M in iterations.

Algorithm 2 Preconditioned Fast Gradient Method

Input: $\mathbf{x}_0 \in \text{dom } \psi$, $\mathbf{P} = \mathbf{P}^\top \succ 0$, $M > 0$, $\rho \geq 0$.
 Set $\mathbf{v}_0 = \mathbf{x}_0$, $A_0 = 0$.
for $k = 0, 1, \dots$ **do**
 Find a_{k+1} from eq. $\frac{M a_{k+1}^2}{A_k + a_{k+1}} = 1 + \rho(A_k + a_{k+1})$.
 Set $A_{k+1} = A_k + a_{k+1}$, $H_k = \frac{1 + \rho A_{k+1}}{a_{k+1}}$.
 Set $\theta_k = \frac{a_{k+1}}{A_{k+1}}$, $\omega_k = \frac{\rho}{H_k}$, $\gamma_k = \frac{\omega_k(1 - \theta_k)}{1 - \omega_k \theta_k}$.
 Set $\hat{\mathbf{v}}_k = (1 - \gamma_k)\mathbf{v}_k + \gamma_k \mathbf{x}_k$.
 Set $\mathbf{y}_k = (1 - \theta_k)\mathbf{x}_k + \theta_k \hat{\mathbf{v}}_k$.
 Compute $\mathbf{v}_{k+1} = \text{GradStep}_{H_k, \mathbf{P}}(\hat{\mathbf{v}}_k, \nabla f(\mathbf{y}_k))$.
 Set $\mathbf{x}_{k+1} = (1 - \theta_k)\mathbf{x}_k + \theta_k \mathbf{v}_{k+1}$.
end for

The convergence results for Algorithm 2 are as follows.

Theorem 3.2. *Consider Algorithm 2 with $M = \beta L$ and $\rho = \alpha\mu$. Then, at each iteration $k \geq 1$, we have*

$$F(\mathbf{x}_k) - F^* \leq 2 \frac{\beta L \|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{B}}^2}{\alpha k^2}. \quad (10)$$

When $\mu > 0$, the convergence is linear: for all $k \geq 1$,

$$F(\mathbf{x}_k) - F^* \leq \left(1 - \sqrt{\frac{\alpha \mu}{\beta L}}\right)^{k-1} \frac{\beta L}{\alpha} \|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{B}}^2.$$

Comparing these estimates with those from Theorem 3.1, we see that the accelerated scheme is much more efficient. For instance, to reach accuracy $\epsilon > 0$ in terms of the objective function in the non-strongly convex case, Algorithm 1 needs $k(\epsilon) = \frac{\beta L \|\mathbf{x}_0 - \mathbf{x}^*\|_{\mathbf{B}}^2}{\alpha \epsilon}$ iterations, while for Algorithm 2 this number is only $k_2(\epsilon) = \sqrt{2k(\epsilon)}$. Similar conclusions are valid in the strongly convex case.

Despite having much weaker dependency on the condition number β/α , Algorithm 2 is still quite sensitive to it. Thus, the proper choice of the preconditioner \mathbf{P} is important for both our methods.

4. Symmetric Polynomial Preconditioning

We would like to have a *family* of preconditioners \mathbf{P}_τ for our problem indexed by some parameter τ . Varying τ should provide us with a trade off between the spectral quality of approximation (7) of the inverse matrix and the arithmetical cost of computing the preconditioner.

Surprisingly, such a family of preconditioners can be built by using *symmetric polynomials*, the classical objects of Algebra. We prove that our preconditioning improves the condition number $\frac{\beta}{\alpha}$ of the problem, by automatically cutting off the large gaps between the top eigenvalues.

4.1. Definition and Basic Properties

We define the family of symmetric matrices $\{P_\tau\}_{0 \leq \tau \leq n-1}$ recursively. We start with identity matrix: $P_0 \stackrel{\text{def}}{=} I$. Then,

$$P_\tau \stackrel{\text{def}}{=} \frac{1}{\tau} \sum_{i=1}^{\tau} (-1)^{i-1} P_{\tau-i} U_i, \quad (11)$$

where $U_\tau \stackrel{\text{def}}{=} \text{tr}(B^\tau)I - B^\tau$ are the auxiliary matrices. It turns out that matrices (11) serve as a good approximation of the inverse matrix: $P_\tau \approx B^{-1}$, up to some multiplicative constant, and the quality of such approximation is gradually improving when increasing parameter τ . Let us look at several first members. Clearly,

$$P_1 = \text{tr}(B)I - B, \quad (12)$$

which is very easy to handle, by having computed the trace of the curvature matrix. Then, multiplying P_1 by any vector would require just one matrix-vector multiplication with our original B . Further,

$$\begin{aligned} P_2 &= \frac{1}{2} \text{tr}(P_1 B)I - P_1 B \\ &= \frac{1}{2} [\text{tr}(B)]^2 I - \text{tr}(B^2)I - \text{tr}(B)B + B^2, \end{aligned} \quad (13)$$

thus its use would cost just *two* matrix-vector products with B , having evaluated^{4,5} the numbers $\text{tr}(B)$ and $\text{tr}(B^2)$.

It is clear that in general $P_\tau = p_\tau(B)$, where p_τ is a polynomial of a fixed degree τ with coefficients that can be found recursively from (11). Let us give a useful interpretation for our family of preconditioners, that also explains their name. For $\mathbf{a} \in \mathbb{R}^{n-1}$, we denote by $\sigma_0(\mathbf{a}), \dots, \sigma_{n-1}(\mathbf{a})$ the *elementary symmetric polynomials* in $n-1$ variables⁶. It is known that every symmetric polynomial (that is invariant to any permutation of the variables) can be represented as a weighed sum of elementary symmetric polynomials (Dummit & Foote, 2004). We establish the following important characterization.

Lemma 4.1. *Let $B = Q \text{Diag}(\boldsymbol{\lambda}) Q^\top$ be the spectral decomposition. Then,*

$$P_\tau = Q \text{Diag}(\sigma_\tau(\boldsymbol{\lambda}_{-1}), \dots, \sigma_\tau(\boldsymbol{\lambda}_{-n})) Q^\top, \quad (14)$$

where $\boldsymbol{\lambda}_{-i} \in \mathbb{R}^{n-1}$ is the vector that contains all elements of $\boldsymbol{\lambda}$ except λ_i .

In particular, we justify $P_\tau \succ 0$. For $\tau = n-1$, we get

⁴Note that $\text{tr}(B^2) = \sum_{i=1}^n \|B[:, i]\|_2^2$, where $B[:, i] \in \mathbb{R}^n$ is the i th column of B .

⁵For general τ , we can also use a stochastic estimate of the trace: $\xi_\tau \stackrel{\text{def}}{=} n \langle B^\tau \mathbf{u}, \mathbf{u} \rangle$, where $\mathbf{u} \in \mathbb{R}^n$ is uniformly distributed on the unit sphere. It would give an unbiased estimate: $\mathbb{E}[\xi_\tau] = n \mathbb{E}[\text{tr}(\mathbf{u}^\top B^\tau \mathbf{u})] = n \text{tr}(\mathbb{E}[\mathbf{u} \mathbf{u}^\top] B^\tau) = \text{tr}(B^\tau)$.

⁶That is $\sigma_\tau(\mathbf{a}) \stackrel{\text{def}}{=} \sum_{1 \leq i_1 < \dots < i_\tau \leq n-1} a_{i_1} \dots a_{i_\tau}$.

$$P_{n-1} \stackrel{(4.1)}{=} \det(B) B^{-1} \stackrel{\text{def}}{=} \text{Adj}(B), \quad (15)$$

which gives us the true inverse matrix B^{-1} up to the constant factor $\det(B)$.

4.2. Approximation Quality

Now, let us show that the quality of approximation $P_\tau \approx B^{-1}$ and that the corresponding condition number $\frac{\beta}{\alpha}$ is improving when τ is increasing.

Theorem 4.2. *For any τ , we have*

$$\lambda_n \sigma_\tau(\boldsymbol{\lambda}_{-n}) B^{-1} \preceq P_\tau \preceq \lambda_1 \sigma_\tau(\boldsymbol{\lambda}_{-1}) B^{-1}. \quad (16)$$

Therefore, the condition number is bounded as

$$\frac{\beta}{\alpha} \stackrel{(16)}{=} \frac{\lambda_1}{\lambda_n} \cdot \xi_\tau(\boldsymbol{\lambda}), \quad \text{where} \quad \xi_\tau(\boldsymbol{\lambda}) \stackrel{\text{def}}{=} \frac{\sigma_\tau(\boldsymbol{\lambda}_{-1})}{\sigma_\tau(\boldsymbol{\lambda}_{-n})}.$$

Note that $\xi_\tau(\boldsymbol{\lambda}) \leq 1$. It is equal to 1 for $\tau = 0$ and can be much smaller for bigger values of τ . For example, for $\tau = 1$ (thus using preconditioner P_1 given by (12)), we get

$$\xi_1(\boldsymbol{\lambda}) = \frac{\sum_{i=2}^n \lambda_i}{\sum_{i=1}^{n-1} \lambda_i} \leq 1,$$

and it is much smaller than 1 when $\lambda_1 \gg \lambda_2$, which corresponds to the case when the highest eigenvalue is well separated from the others. Therefore, the methods with preconditioner P_1 achieve a *provable acceleration* in the case of large gap between λ_1 and λ_2 , *without explicit knowledge* of the spectrum of B . The price of using P_1 instead of P_0 is just one extra matrix-vector product per iteration. Let us summarize the main properties of $\xi_\tau(\boldsymbol{\lambda})$ (see also Figure 3).

Lemma 4.3. *It holds that $\xi_0(\boldsymbol{\lambda}) = 1$, $\xi_{n-1}(\boldsymbol{\lambda}) = \frac{\lambda_n}{\lambda_1}$, and $\xi_\tau(\boldsymbol{\lambda})$ monotonically decreases with τ . Moreover,*

$$\xi_\tau(\boldsymbol{\lambda}) \rightarrow 0 \quad \text{when} \quad \frac{\lambda_1}{\lambda_{\tau+1}} \rightarrow \infty.$$

A more explicit upper bound which implies the above limit is as follows.

Lemma 4.4. *For any $0 \leq \tau \leq n-1$, we have*

$$\xi_\tau(\boldsymbol{\lambda}) \leq \frac{\sum_{i=\tau+1}^n \lambda_i}{\lambda_1 + \sum_{i=\tau+1}^{n-1} \lambda_i}.$$

We see that τ interpolates P_τ between I and $\text{Adj}(B)$, while the condition number $\frac{\beta}{\alpha}$ changes from $\frac{\lambda_1}{\lambda_n}$ to 1. Therefore, we obtain an extra degree of freedom in our methods for choosing an appropriate small value of τ , that improves the spectrum of the problem by cutting off large gaps between λ_1 and $\lambda_{\tau+1}$.

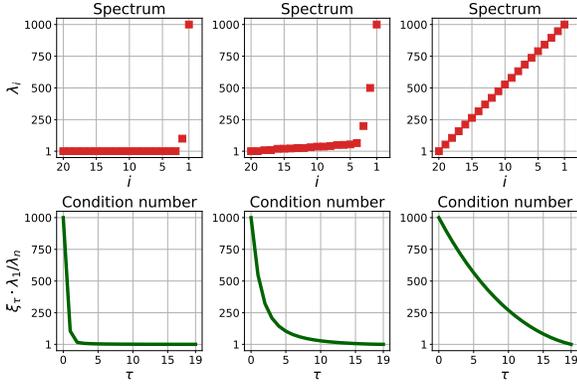


Figure 3: *Above*: different distributions of eigenvalues $\lambda(\mathbf{B})$. *Below*: the corresponding improvement of the condition number $\frac{\beta}{\alpha} = \xi_\tau(\lambda) \cdot \lambda_1/\lambda_n$ by using the preconditioner \mathbf{P}_τ of order τ .

4.3. Stochastic Representation

Let us provide another interesting interpretation for our family of preconditioners. One way of approximating the inverse matrix \mathbf{B}^{-1} could be to extract from \mathbf{B} a randomly selected principal submatrix of size $\tau + 1$, compute its inverse, put it back into the original “big matrix” and zero out all other elements outside the submatrix. It turns out that, in expectation, the result of this operation is exactly proportional to our preconditioner \mathbf{P}_τ if we pick the submatrix from a special *volume sampling* distribution (Deshpande et al., 2006).

Theorem 4.5. *For any $0 \leq \tau \leq n - 1$, it holds that*

$$\mathbf{P}_\tau \propto \mathbb{E}_{S \sim \text{Vol}_{\tau+1}(\mathbf{B})} [\mathbf{I}_S (\mathbf{B}_{S \times S})^{-1} \mathbf{I}_S^\top], \quad (17)$$

where $S \subseteq \{1, \dots, n\}$ is a random $(\tau + 1)$ -element subset of coordinates, $\mathbf{I}_S \in \mathbb{R}^{n \times (\tau+1)}$ is the matrix obtained from the identity matrix by retaining only the columns with indices from S , $\mathbf{B}_{S \times S} \stackrel{\text{def}}{=} \mathbf{I}_S^\top \mathbf{B} \mathbf{I}_S \in \mathbb{R}^{(\tau+1) \times (\tau+1)}$, and $\text{Vol}_{\tau+1}(\mathbf{B})$ is the volume sampling distribution prescribing to pick S with probability $\propto \det(\mathbf{B}_{S \times S})$.

The idea of applying volume sampling in Optimization was first proposed in (Rodomanov & Kropotov, 2020) for accelerating coordinate descent methods. It was shown that using this particular nonuniform sampling of coordinates leads to a provable acceleration by a factor whose magnitude depends on gaps in the spectrum of the curvature matrix.

Thus, we can interpret our basic Gradient Method (Algorithm 1) with a fixed Symmetric Polynomial Preconditioner \mathbf{P}_τ as a deterministic counterpart of the randomized coordinate descent method from (Rodomanov & Kropotov, 2020) with $(\tau + 1)$ -element volume sampling of coordinates. Correspondingly, both methods have very similar convergence properties and theoretical efficiency estimates.

Nevertheless, this work offers several significant advantages

over (Rodomanov & Kropotov, 2020). First, in addition to the basic method, we have an accelerated one (Algorithm 2), while the accelerated version of coordinate descent with volume sampling is an open question. Second, volume sampling is an expensive operation which is difficult to carry out already when $\tau = 2$. In contrast, the corresponding preconditioner \mathbf{P}_2 for our gradient methods is still computationally efficient (see Section 4.1). Finally, as we will show next, the basic Gradient Method can be improved to *automatically* choose the best possible polynomial preconditioner of degree τ (including the one we have been discussing in this section), and the resulting algorithm can easily handle much bigger values of τ .

5. Krylov Subspace Preconditioning

Our new symmetric polynomial preconditioners, introduced in the previous section, can be viewed as a certain family of polynomials that we apply to our curvature matrix \mathbf{B} . Thus, for a fixed degree $\tau > 0$, we use the matrix $\mathbf{P} = p_\tau(\mathbf{B})$ as a preconditioner, where p_τ is a specifically constructed polynomial of degree τ such that $\mathbf{P} \succ 0$.

A natural question is *how optimal is this choice of a polynomial?* Indeed, the problem of polynomial approximation has a long and rich history with an affirmative answer provided by the classical Chebyshev polynomials (Mason & Handscomb, 2002) for the uniform approximation bound. We present a new adaptive algorithm that automatically achieves the *best* polynomial preconditioning. Then, we study what are the complexity guarantees that we can get with this optimal approach. In this section, we focus on the non-composite case only, i.e. the problem of unconstrained minimization of a smooth function: $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$.

5.1. Gradient Method with Krylov Preconditioning

We denote $\mathbf{P}_a \stackrel{\text{def}}{=} a_0 \mathbf{I} + a_1 \mathbf{B} + \dots + a_\tau \mathbf{B}^\tau$, where vector $\mathbf{a} = (a_0, \dots, a_\tau) \in \mathbb{R}^{\tau+1}$ is a parameter. In each iteration, it is found by solving the linear system:

$$\mathbf{a} = \mathbf{A}_\tau^{-1} \mathbf{g}_\tau \in \mathbb{R}^{\tau+1}, \quad (18)$$

where $\mathbf{A}_\tau = \mathbf{A}_\tau(\mathbf{x}) \in \mathbb{R}^{(\tau+1) \times (\tau+1)}$ is the Gram matrix with the following structure ($0 \leq i, j \leq \tau$):

$$[\mathbf{A}_\tau(\mathbf{x})]^{(i,j)} \stackrel{\text{def}}{=} L \cdot \langle \nabla f(\mathbf{x}), \mathbf{B}^{i+j+1} \nabla f(\mathbf{x}) \rangle, \quad (19)$$

and $\mathbf{g}_\tau = \mathbf{g}_\tau(\mathbf{x}) \in \mathbb{R}^{\tau+1}$ is defined by ($0 \leq i \leq \tau$):

$$[\mathbf{g}_\tau(\mathbf{x})]^{(i)} \stackrel{\text{def}}{=} \langle \nabla f(\mathbf{x}), \mathbf{B}^i \nabla f(\mathbf{x}) \rangle. \quad (20)$$

Note that this operation is exactly the projection of the direction $\frac{1}{L} \mathbf{B}^{-1} \nabla f(\mathbf{x}_k)$ onto the *Krylov subspace* (3):

$$\mathbf{x}_{k+1} - \mathbf{x}_k := \underset{\mathbf{h} \in \mathcal{K}_{\mathbf{x}_k}^\tau}{\text{argmin}} \|\mathbf{h} + \frac{1}{L} \mathbf{B}^{-1} \nabla f(\mathbf{x}_k)\|_{\mathbf{B}}^2.$$

Fortunately, for computing this projection we indeed do not need to invert the curvature matrix \mathbf{B} , but to solve only a small linear system (18) of size $\tau + 1$. We are ready to formulate our new adaptive method.

Algorithm 3 Gradient Method with Krylov Preconditioning

Initialization: $\mathbf{x}_0 \in \mathbb{R}^n$, $\tau \geq 0$, $L > 0$.
for $k = 0, 1, \dots$ **do**
 Form matrix $\mathbf{A}_\tau(\mathbf{x}_k)$ and vector $\mathbf{g}_\tau(\mathbf{x}_k)$ by (19), (20).
 Compute $\mathbf{a}_k = \mathbf{A}_\tau(\mathbf{x}_k)^{-1} \mathbf{g}_\tau(\mathbf{x}_k) \in \mathbb{R}^{\tau+1}$.
 Set $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{P}_{\mathbf{a}_k} \nabla f(\mathbf{x}_k)$.
end for

We prove the following optimality result.

Theorem 5.1. *Let $\mathbf{P} \succ 0$ be any preconditioner that is a polynomial of degree τ of the curvature matrix:*

$$\mathbf{P} = p_\tau(\mathbf{B}), \quad p_\tau \in \mathbb{R}[s], \quad \deg(p_\tau) = \tau.$$

Then, for the iteration of Algorithm 3 we have the global rates (8),(10) with the condition number that is attributed to \mathbf{P} (7): $\frac{\beta}{\alpha} = \frac{\beta(\mathbf{P})}{\alpha(\mathbf{P})}$.

Hence, our method *automatically* chooses the best possible preconditioning matrix from the polynomial class. Let us understand what are the bounds for $\frac{\beta}{\alpha}$ that we can achieve in this case.

5.2. Bounds for the Condition Number

Let us assume that the top $\tau > 0$ eigenvalues of \mathbf{B} are all *separated*. Then, we can easily cut them off with the following simple construction. Define

$$q_\tau(s) \stackrel{\text{def}}{=} \left(1 - \frac{s}{\lambda_1}\right) \left(1 - \frac{s}{\lambda_2}\right) \cdots \left(1 - \frac{s}{\lambda_\tau}\right). \quad (21)$$

Proposition 5.2. *For any τ , taking $\mathbf{P} = p_\tau(\mathbf{B})$, where $p_\tau(s) := \frac{1+q_\tau(s) \cdot (\alpha s - 1)}{s}$ with q_τ defined by (21) and $\alpha = \frac{2}{\lambda_{\tau+1} + \lambda_n}$, the condition number is bounded by $\frac{\beta}{\alpha} \leq \frac{\lambda_{\tau+1}}{\lambda_n}$.*

The worst case instance for the cutting strategy is when all the eigenvalues except one share the same value. A better approach in such a situation would be to find a bound from the *uniform* polynomial approximation for the whole interval $[\lambda_n, \lambda_1]$, which is achieved with the Chebyshev polynomials (Nemirovski, 1995).

Proposition 5.3. *For a fixed $0 < \epsilon < 1$, let $\tau := \lfloor \sqrt{\frac{\lambda_1}{\lambda_n} \ln \frac{8}{\epsilon}} \rfloor$. Then, taking $\mathbf{P} = p_\tau(\mathbf{B})$, where $p_\tau(s) := \frac{1-Q_\tau(s)}{s}$ with Q_τ is a normalized Chebyshev polynomial⁷ of the first kind of degree $\tau + 1$, the condition number is bounded by $\frac{\beta}{\alpha} \leq 1 + \epsilon$.*

⁷See Appendix B.9 for the precise definition.

5.3. Discussion

We see that in the case of unconstrained smooth minimization, it is possible to achieve the guarantee of the *best polynomial* of a fixed degree τ , by computing a certain projection onto the corresponding Krylov subspace. Namely, we can achieve $\frac{\beta}{\alpha} \leq \frac{\lambda_{\tau+1}}{\lambda_n}$ (Proposition 5.2), which cuts off the top τ eigenvalues of the spectrum completely, if they are separated from the others. At the same time, by using the Chebyshev polynomials, we can contract a part of the spectrum *uniformly*, with an appropriate degree τ (Proposition 5.3). It remains to be an open question where we can incorporate adaptive Krylov preconditioning into the Fast Gradient Method, which would give us a further improvement of the condition number.

6. Experiments

Huber Loss. Let us present an illustrative experiment, with the regression model (Example 2.4) with the Huber loss function:

$$\phi(t) := \begin{cases} \frac{t^2}{2\mu}, & \text{if } |t| \leq \mu, \\ |t| - \frac{\mu}{2}, & \text{otherwise,} \end{cases}$$

where $\mu := 0.1$ is a parameter. The data is generated with a fixed distribution of eigenvalues: $\lambda_1 > \lambda_2 > \lambda_3 = \dots = \lambda_n = 1$. Thus, we have two gaps between the leading eigenvalues. We use the Gradient Method (Algorithm 1), with the adaptive search to fit the parameter M . The results are shown in Figure 4. Using the preconditioner \mathbf{P}_1 helps the method to deal with the large gap between λ_1 and λ_2 , while \mathbf{P}_2 makes the method to be insensitive to the gap between λ_1 and λ_3 , as predicted by our theory.

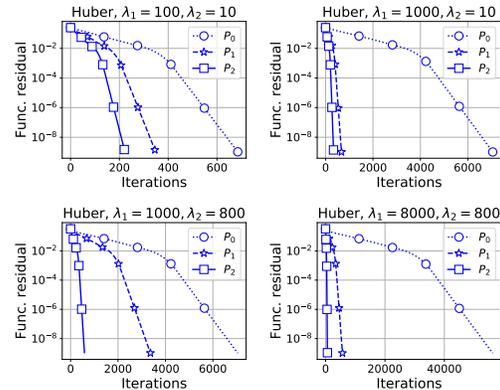


Figure 4: Minimizing the Huber loss.

Logistic Regression. We examine the training of logistic regression on real data. In Figure 5, we see that the best convergence is achieved by the Fast Gradient Method (FGM, Algorithm 2) with \mathbf{P}_2 . Using Symmetric Polynomial Preconditioning makes the methods to converge much

better (*two times faster* for GM using P_2 instead of $P_0 \equiv I$, and about *1.5 times faster* for FGM). Among the versions of GM, the most encouraging performance belongs to the Krylov preconditioning, which is consistent with the theory.

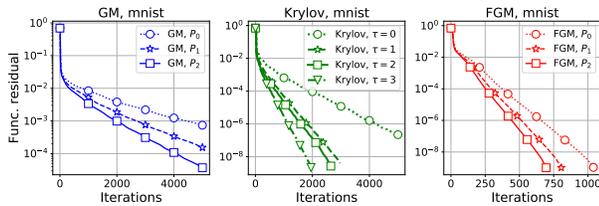


Figure 5: Training logistic regression.

Acknowledgements

The work of the first author was supported by the Swiss State Secretariat for Education, Research and Innovation (SERI) under contract number 22.00133. The work of the second author received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No. 788368).

References

- Bishop, C. M. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- Bullins, B. Highly smooth minimization of non-smooth problems. In *Conference on Learning Theory*, pp. 988–1030. PMLR, 2020.
- 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.
- Deshpande, A., Rademacher, L., Vempala, S. S., and Wang, G. Matrix approximation and projective clustering via volume sampling. *Theory of Computing*, 2(12):225–247, 2006.
- Dubois, P. F., Greenbaum, A., and Rodrigue, G. H. Approximating the inverse of a matrix for use in iterative algorithms on vector processors. *Computing*, 22(3):257–268, 1979.
- Dummit, D. S. and Foote, R. M. *Abstract algebra*, volume 3. Wiley Hoboken, 2004.
- 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.
- Gupta, S. D., Freund, R. M., Sun, X. A., and Taylor, A. Nonlinear conjugate gradient methods: worst-case convergence rates via computer-assisted analyses. *arXiv preprint arXiv:2301.01530*, 2023.
- Hardy, G. H., Littlewood, J. E., and Pólya, G. *Inequalities*. Cambridge University Press, second edition, 1952.
- 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, 1952.
- Johnson, O. G., Micchelli, C. A., and Paul, G. Polynomial preconditioners for conjugate gradient calculations. *SIAM Journal on Numerical Analysis*, 20(2):362–376, 1983.
- Kalman, D. A matrix proof of Newton’s identities. *Mathematics Magazine*, 73(4):313–315, 2000.
- Kovalev, D., Richtárik, P., Gorbunov, E., and Gasanov, E. Stochastic spectral and conjugate descent methods. *Advances in Neural Information Processing Systems*, 31, 2018.
- Liu, Q., Morgan, R. B., and Wilcox, W. Polynomial preconditioned gmres and gmres-dr. *SIAM Journal on Scientific Computing*, 37(5):S407–S428, 2015.
- Loe, J. A. and Morgan, R. B. Toward efficient polynomial preconditioning for gmres. *Numerical Linear Algebra with Applications*, 29(4):e2427, 2022.
- Mardal, K.-A. and Winther, R. Preconditioning discretizations of systems of partial differential equations. *Numerical Linear Algebra with Applications*, 18(1):1–40, 2011.
- Mason, J. C. and Handscomb, D. C. *Chebyshev polynomials*. Chapman and Hall/CRC, 2002.
- Nemirovski, A. Information-based complexity of convex programming. *Lecture Notes*, 834, 1995.
- Nemirovski, A. and Yudin, D. Problem complexity and method efficiency in optimization. 1983.
- Nesterov, Y. A method for solving the convex programming problem with convergence rate $O(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pp. 543–547, 1983.
- Nesterov, Y. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.

- Nesterov, Y. Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161, 2013.
- Nesterov, Y. *Lectures on convex optimization*, volume 137. Springer, 2018.
- Nocedal, J. and Wright, S. *Numerical optimization*. Springer Science & Business Media, 2006.
- Polyak, B. T. *Introduction to optimization*. Optimization Software, 1987.
- Qu, Z., Richtárik, P., Takác, M., and Fercoq, O. Sdna: stochastic dual newton ascent for empirical risk minimization. In *International Conference on Machine Learning*, pp. 1823–1832. PMLR, 2016.
- Rodomanov, A. and Kropotov, D. A randomized coordinate descent method with volume sampling. *SIAM Journal on Optimization*, 30(3):1878–1904, 2020.
- Saad, Y. Practical use of polynomial preconditionings for the conjugate gradient method. *SIAM Journal on Scientific and Statistical Computing*, 6(4):865–881, 1985.
- 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.
- 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.
- Van Gijzen, M. A polynomial preconditioner for the gmres algorithm. *Journal of Computational and Applied Mathematics*, 59(1):91–107, 1995.
- Vishnoi, N. K. $Lx=b$. *Foundations and Trends® in Theoretical Computer Science*, 8(1–2):1–141, 2013.

Supplementary Material

A. Extra Experiments

Logistic Regression. Let us present experimental results for our preconditioning strategies, for the training of Logistic Regression with several real datasets. We investigate both the number of iterations and the number of matrix-vector products (the most difficult operation) required to reach a certain accuracy level in the functional residual. The results are shown in Figure 6.

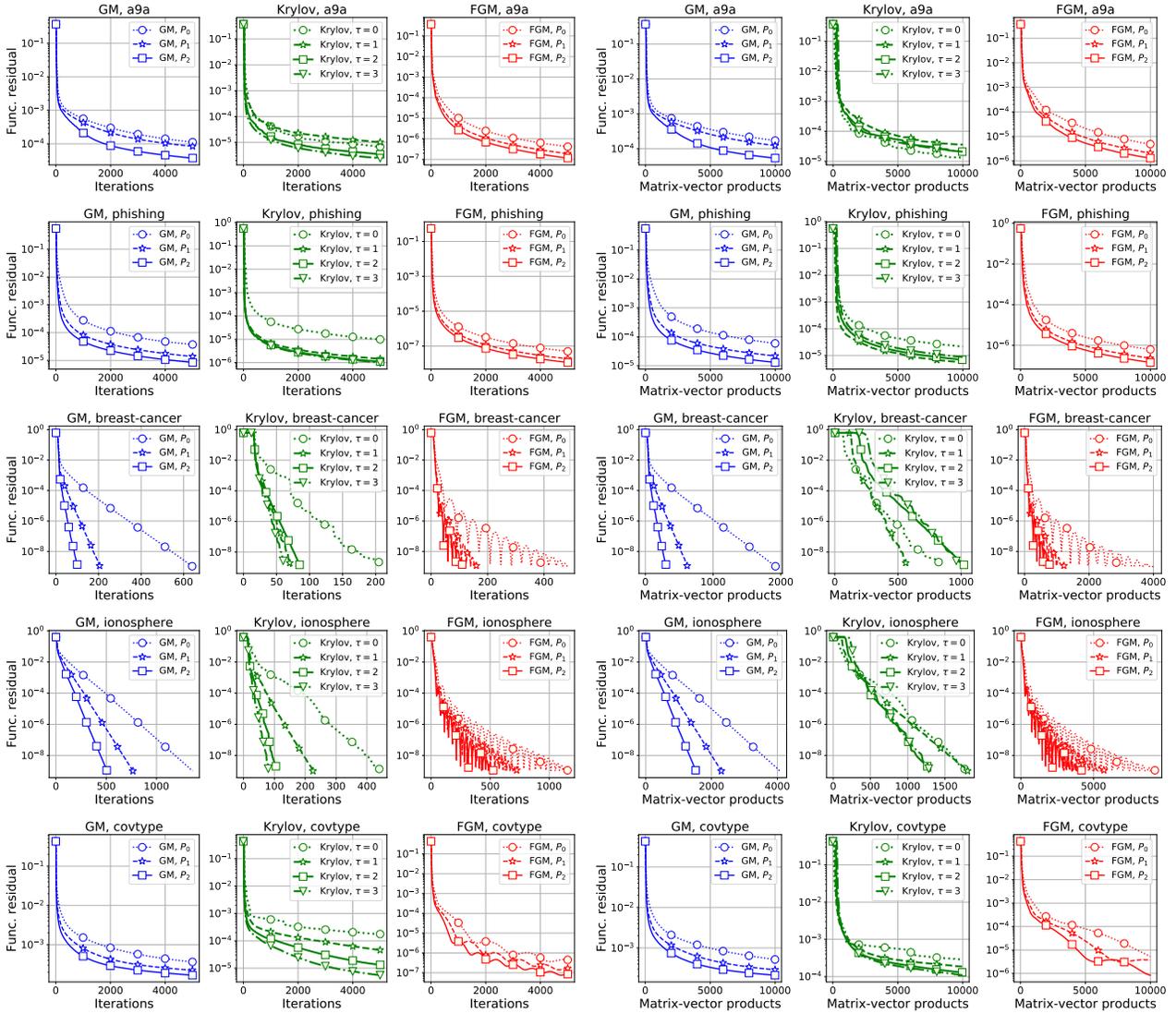


Figure 6: Training logistic regression with Algorithm 1 (GM) and Algorithm 2 (FGM) employing Symmetric Polynomial Preconditioning (11); and with Algorithm 3 (Krylov).

We see that using Symmetric Polynomial Preconditioning (P_1 and P_2) significantly accelerates both the Gradient Method (GM) and the Fast Gradient Method (GM), without extra arithmetic efforts during each iteration. Using the Krylov

preconditioning is more costly, while it equips GM with the best possible iteration rates.

Typical Distributions of the Data Spectrum. Let us provide the plots with the distributions of the leading eigenvalues (in the logarithmic scale) of the curvature matrix B for a selection of typical machine learning datasets⁸.

We see (Figure 7) that it is quite common to have several top eigenvalues well separated from others. In these cases, our new Symmetric Polynomial Preconditioners provides the gradient methods with the best acceleration.

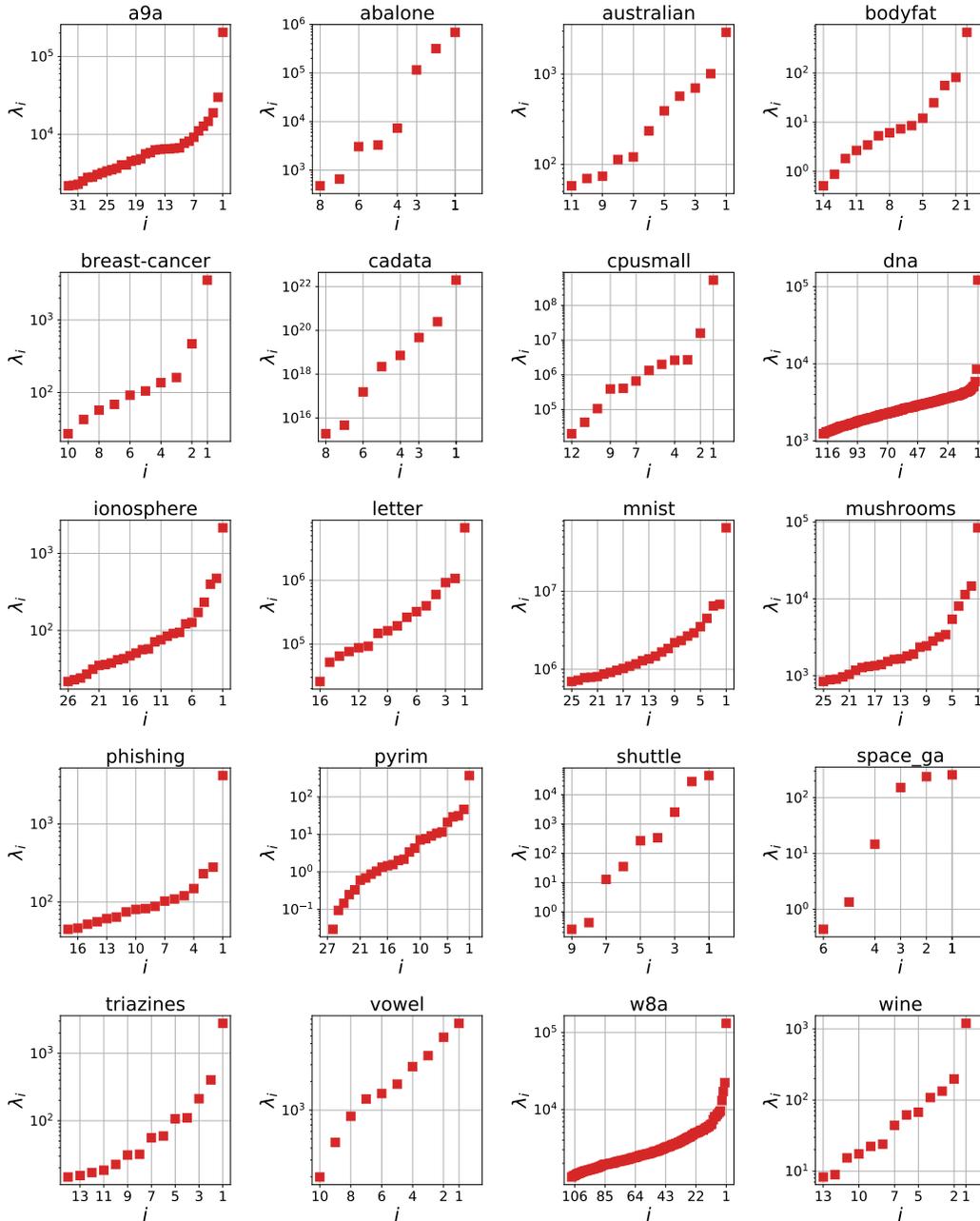


Figure 7: Leading eigenvalues (in the logarithmic scale) of the curvature matrix for several real datasets.

⁸<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

Comparison with BFGS and L-BFGS. In the following experiment, we compare the performance of the Gradient Method (GM) and Fast Gradient Method (FGM) equipped with our Symmetric Polynomial Preconditioning, and GM with Krylov preconditioning, with the classical BFGS and L-BFGS optimization schemes (Nocedal & Wright, 2006).

The results are presented in Figure 8. We show both the number of matrix-vector products (the most expensive operation) and the total computational time⁹ required to reach a given accuracy in terms of the functional residual.

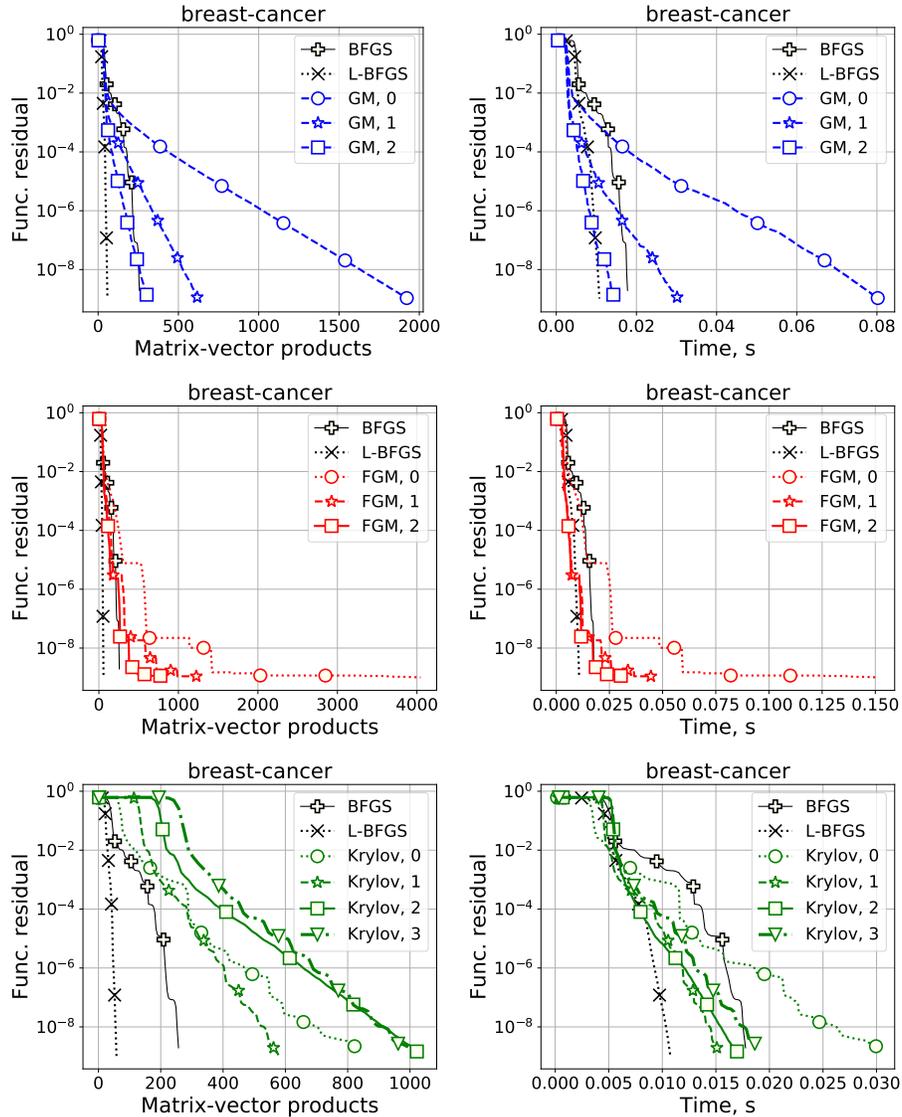


Figure 8: Comparison of our methods with Quasi-Newton methods: training Logistic Regression

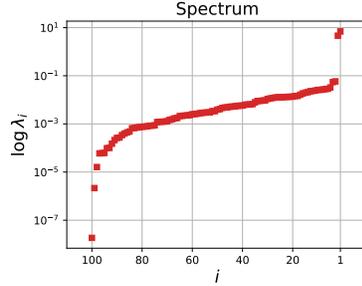
We see that the performance of the Gradient Method and Fast Gradient Method with Symmetric Polynomial preconditioner of order $\tau = 2$ (GM, 2 and FGM, 2 correspondingly) is comparable to that one of the BFGS and L-BFGS methods.

⁹Clock time was evaluated using the machine with Intel Core i5 CPU, 1.6GHz; 8 GB RAM. All methods were implemented in Python.

In the next experiment, we consider the problem of minimizing the Soft Maximum objective (log-sum-exp):

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \mu \ln \left(\sum_{i=1}^m \exp \left(\frac{\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i}{\mu} \right) \right) \approx \max_{1 \leq i \leq m} [\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i], \quad (22)$$

where $\mu > 0$ is a sufficiently small number. The problems of this type are important in applications with minimax strategies for matrix games and for training ℓ_∞ -regression (Nesterov, 2005; Bullins, 2020). The vectors $\mathbf{a}_1, \dots, \mathbf{a}_m \in \mathbb{R}^n$ and $\mathbf{b} \in \mathbb{R}^m$ form our data, with $n = 100, m = 200$. That the structure of this objective satisfies Example 2.3 and the corresponding curvature matrix has the following distribution of the spectrum (in the double-logarithmic scale):



Note that the function (22) does not have a finite-sum structure, thus it is impossible to apply to this problem stochastic optimization methods. We use the value $\mu = 0.005$ for the smoothing parameter. The results are shown in Figure 9.

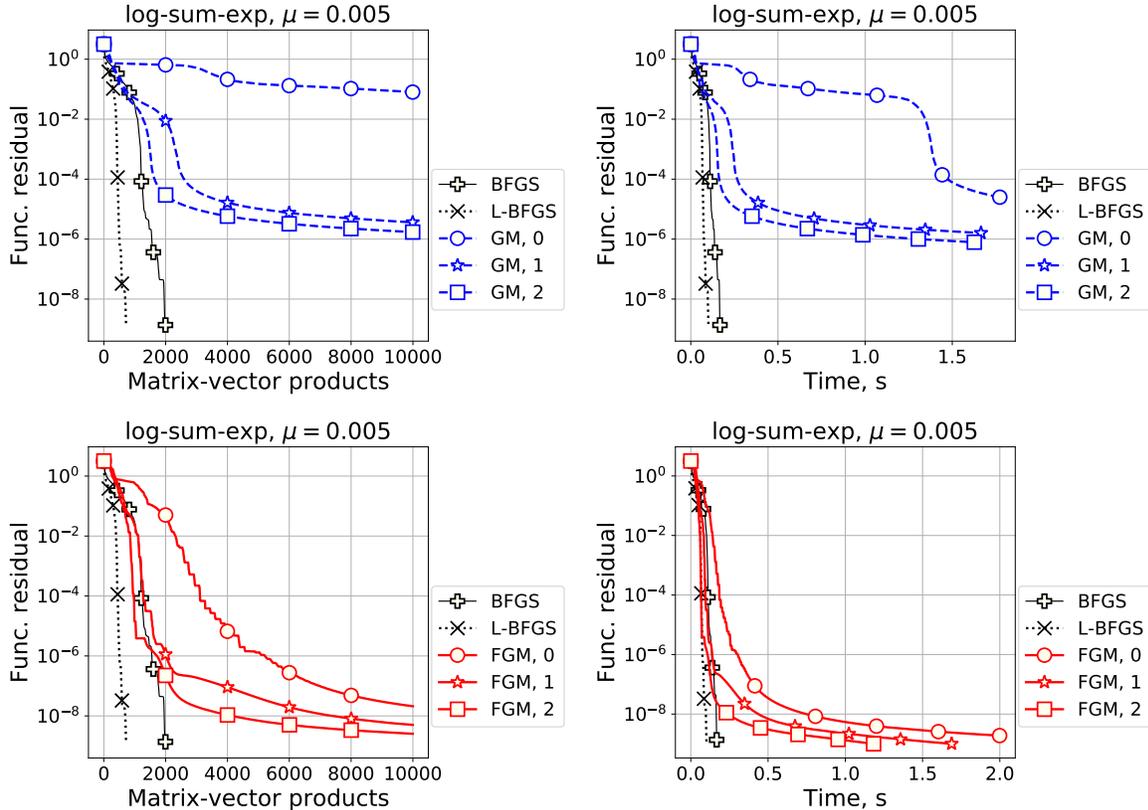


Figure 9: Comparison of our methods with Quasi-Newton methods: training Soft Maximum (log-sum-exp objective)

We see that using our Symmetric Polynomial Preconditioning significantly helps the Gradient Method (GM) and the Fast Gradient Method (FGM). The performance of FGM with preconditioner of order $\tau = 2$ is comparable to that one of the BFGS algorithm, both in terms of the matrix-vector products and total computational time.

B. Proofs

B.1. Proof of Theorem 3.1

Let us consider one iteration of the method, for some $k \geq 0$. By definition, $\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{y} \in \operatorname{dom} \psi} \{\Omega_k(\mathbf{y})\}$, where

$$\Omega_k(\mathbf{y}) \stackrel{\text{def}}{=} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y} - \mathbf{x}_k \rangle + \frac{M}{2} \|\mathbf{y} - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2 + \psi(\mathbf{y})$$

is strongly convex with respect to \mathbf{P}^{-1} norm with parameter $M := \beta L$. Thus, we have, for any $\mathbf{y} \in \operatorname{dom} \psi$:

$$\begin{aligned} \frac{M}{2} \|\mathbf{y} - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2 + F(\mathbf{y}) &\geq \Omega_k(\mathbf{y}) \geq \Omega_k(\mathbf{x}_{k+1}) + \frac{M}{2} \|\mathbf{y} - \mathbf{x}_{k+1}\|_{\mathbf{P}^{-1}}^2 \\ &\geq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{B}}^2 + \psi(\mathbf{x}_{k+1}) + \frac{M}{2} \|\mathbf{y} - \mathbf{x}_{k+1}\|_{\mathbf{P}^{-1}}^2 \\ &\stackrel{(6)}{\geq} F(\mathbf{x}_{k+1}) + \frac{M}{2} \|\mathbf{y} - \mathbf{x}_{k+1}\|_{\mathbf{P}^{-1}}^2. \end{aligned} \quad (23)$$

Hence, substituting $\mathbf{y} := \mathbf{x}^*$ (solution to the problem), we establish the boundness for all iterates:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_{\mathbf{P}^{-1}} \leq \|\mathbf{x}_k - \mathbf{x}^*\|_{\mathbf{P}^{-1}}. \quad (24)$$

Further, let us take $\mathbf{y} := \gamma_k \mathbf{x}^* + (1 - \gamma_k) \mathbf{x}_k$, for some $\gamma_k \in [0, 1]$. We obtain

$$\begin{aligned} F(\mathbf{x}_{k+1}) &\stackrel{(23)}{\leq} F(\gamma_k \mathbf{x}^* + (1 - \gamma_k) \mathbf{x}_k) + \frac{\gamma_k^2 M}{2} \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2 \\ &\leq \gamma_k F^* + (1 - \gamma_k) F(\mathbf{x}_k) + \frac{\gamma_k^2 M}{2} \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2. \end{aligned} \quad (25)$$

Now, setting $A_k \stackrel{\text{def}}{=} k \cdot (k + 1)$, $a_{k+1} \stackrel{\text{def}}{=} A_{k+1} - A_k = 2(k + 1)$, and $\gamma_k := \frac{a_{k+1}}{A_{k+1}} = \frac{2}{k+2}$, we obtain

$$\begin{aligned} A_{k+1} (F(\mathbf{x}_{k+1}) - F^*) &\stackrel{(25)}{\leq} A_k (F(\mathbf{x}_k) - F^*) + \frac{a_{k+1}^2}{A_{k+1}} \cdot \frac{M}{2} \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2 \\ &\stackrel{(24)}{\leq} A_k (F(\mathbf{x}_k) - F^*) + \frac{a_{k+1}^2}{A_{k+1}} \cdot \frac{M}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_{\mathbf{P}^{-1}}^2 \\ &\stackrel{(7)}{\leq} A_k (F(\mathbf{x}_k) - F^*) + \frac{a_{k+1}^2}{A_{k+1}} \cdot \frac{\beta}{\alpha} \cdot \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_{\mathbf{B}}^2. \end{aligned} \quad (26)$$

Telescoping this bound for the first k iterations, we get

$$F(\mathbf{x}_k) - F^* \stackrel{(26)}{\leq} \frac{\beta}{\alpha} \cdot \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_0\|_{\mathbf{B}}^2 \cdot \frac{1}{A_k} \sum_{i=1}^k a_i^2 = \mathcal{O}\left(\frac{\beta}{\alpha} \cdot \frac{L}{2k} \|\mathbf{x}^* - \mathbf{x}_0\|_{\mathbf{B}}^2\right).$$

To prove the linear rate for the strongly convex case, we continue as follows

$$F(\mathbf{x}_{k+1}) \stackrel{(25),(6)}{\leq} \gamma_k F^* + (1 - \gamma_k) F(\mathbf{x}_k) + \gamma_k^2 \cdot \frac{\beta L}{\alpha \mu} \cdot (F(\mathbf{x}_k) - F^*).$$

Choosing $\gamma_k := \frac{\alpha \mu}{2\beta L} < 1$, we get the exponential rate

$$F(\mathbf{x}_{k+1}) - F^* \leq \left(1 - \frac{\alpha \mu}{4\beta L}\right) (F(\mathbf{x}_k) - F^*),$$

which completes the proof. \square

B.2. Proof of Theorem 3.2

Let $\mathbf{x} \in \operatorname{dom} \psi$ and $k \geq 0$ be arbitrary. From (6), (7), and the fact that $\rho = \alpha \mu$, it follows that

$$F(\mathbf{x}) = f(\mathbf{x}) + \psi(\mathbf{x}) \geq \ell_k(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{x} - \mathbf{y}_k\|_{\mathbf{P}^{-1}}^2, \quad \ell_k(\mathbf{x}) \stackrel{\text{def}}{=} f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x} - \mathbf{y}_k \rangle + \psi(\mathbf{x}).$$

Hence,

$$\begin{aligned}
 A_k F(\mathbf{x}_k) + a_{k+1} F(\mathbf{x}) + \frac{1 + \rho A_k}{2} \|\mathbf{x} - \mathbf{v}_k\|_{\mathcal{P}^{-1}}^2 \\
 &\geq A_k \ell_k(\mathbf{x}_k) + a_{k+1} \ell_k(\mathbf{x}) + \frac{1 + \rho A_k}{2} \|\mathbf{x} - \mathbf{v}_k\|_{\mathcal{P}^{-1}}^2 + \frac{\rho a_{k+1}}{2} \|\mathbf{x} - \mathbf{y}_k\|_{\mathcal{P}^{-1}}^2 \\
 &\geq A_k \ell_k(\mathbf{x}_k) + a_{k+1} \ell_k(\mathbf{x}) + \frac{1 + \rho A_{k+1}}{2} \|\mathbf{x} - \hat{\mathbf{v}}_k\|_{\mathcal{P}^{-1}}^2 \stackrel{\text{def}}{=} \zeta_k(\mathbf{x}),
 \end{aligned} \tag{27}$$

where the final inequality follows from the convexity of the squared norm and the fact that, according to our definitions,

$$\frac{(1 + \rho A_k) \mathbf{v}_k + \rho a_{k+1} \mathbf{y}_k}{1 + \rho A_{k+1}} = (1 - \omega_k) \mathbf{v}_k + \omega_k \mathbf{y}_k = (1 - \omega_k) \mathbf{v}_k + \omega_k [(1 - \theta_k) \mathbf{x}_k + \theta_k \hat{\mathbf{v}}_k] = \hat{\mathbf{v}}_k.$$

Note that ζ_k is a $(1 + \rho A_{k+1})$ -strongly convex function w.r.t. $\|\cdot\|_{\mathcal{P}^{-1}}$, and \mathbf{v}_{k+1} is precisely its minimizer. Therefore,

$$\zeta_k(\mathbf{x}) \geq \zeta_k(\mathbf{v}_{k+1}) + \frac{1 + \rho A_{k+1}}{2} \|\mathbf{x} - \mathbf{v}_{k+1}\|_{\mathcal{P}^{-1}}^2. \tag{28}$$

Since ℓ_k is a convex function, we have, by our definition of \mathbf{x}_{k+1} ,

$$A_k \ell_k(\mathbf{x}_k) + a_{k+1} \ell_k(\mathbf{v}_{k+1}) \geq A_{k+1} \ell_k(\mathbf{x}_{k+1}).$$

On the other hand, by the definition of \mathbf{x}_{k+1} and \mathbf{y}_k ,

$$\mathbf{x}_{k+1} - \mathbf{y}_k = \theta_k (\mathbf{v}_{k+1} - \hat{\mathbf{v}}_k) = \frac{a_{k+1}}{A_{k+1}} (\mathbf{v}_{k+1} - \hat{\mathbf{v}}_k).$$

Therefore,

$$\begin{aligned}
 \zeta_k(\mathbf{v}_{k+1}) &= A_k \ell_k(\mathbf{x}_k) + a_{k+1} \ell_k(\mathbf{v}_{k+1}) + \frac{1 + \rho A_{k+1}}{2} \|\mathbf{v}_{k+1} - \hat{\mathbf{v}}_k\|_{\mathcal{P}^{-1}}^2 \\
 &\geq A_{k+1} \left[\ell_k(\mathbf{x}_{k+1}) + \frac{A_{k+1} (1 + \rho A_{k+1})}{2 a_{k+1}^2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|_{\mathcal{P}^{-1}}^2 \right].
 \end{aligned}$$

In view of our choice of a_{k+1} , we have the following identity:

$$\frac{M a_{k+1}^2}{A_{k+1}} = 1 + \rho A_{k+1}. \tag{29}$$

Combining this with the fact that $M = \beta L$ and using (7) and (6), we get

$$\begin{aligned}
 \zeta_k(\mathbf{v}_{k+1}) &\geq A_{k+1} \left[\ell_k(\mathbf{x}_{k+1}) + \frac{M}{2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|_{\mathcal{P}^{-1}}^2 \right] \geq A_{k+1} \left[\ell_k(\mathbf{x}_{k+1}) + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|_B^2 \right] \\
 &= A_{k+1} \left[f(\mathbf{y}_k) + \langle \nabla f(\mathbf{y}_k), \mathbf{x}_{k+1} - \mathbf{y}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{y}_k\|_B^2 + \psi(\mathbf{x}_{k+1}) \right] \geq A_{k+1} F(\mathbf{x}_{k+1}).
 \end{aligned}$$

Substituting the above bound into (28), and that one into (28), we thus obtain

$$A_k F(\mathbf{x}_k) + a_{k+1} F(\mathbf{x}) + \frac{1 + \rho A_k}{2} \|\mathbf{x} - \mathbf{v}_k\|_{\mathcal{P}^{-1}}^2 \geq A_{k+1} F(\mathbf{x}_{k+1}) + \frac{1 + \rho A_{k+1}}{2} \|\mathbf{x} - \mathbf{v}_{k+1}\|_{\mathcal{P}^{-1}}^2.$$

This inequality is valid for any $k \geq 0$.

Fixing an arbitrary $k \geq 1$ and summing up the previous inequalities for all indices $k' = 0, \dots, k-1$, we get

$$A_k F(\mathbf{x}_k) \leq A_k F(\mathbf{x}) + \frac{1 + \rho A_0}{2} \|\mathbf{x} - \mathbf{v}_0\|_{\mathcal{P}^{-1}}^2 = A_k F(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_0\|_{\mathcal{P}^{-1}}^2.$$

Substituting further $\mathbf{x} = \mathbf{x}^*$ (an optimal solution) and using (7), gives us the following convergence rate estimate:

$$F(\mathbf{x}_k) - F^* \leq \frac{\|\mathbf{x}^* - \mathbf{x}_0\|_{\mathcal{P}^{-1}}^2}{2A_k} \leq \frac{\|\mathbf{x}^* - \mathbf{x}_0\|_B^2}{2\alpha A_k}. \tag{30}$$

To complete the proof, it remains to use standard lower bounds on A_k (see (Nesterov, 2018)). Specifically, dropping the second term from the right-hand side of (29) and rearranging, we obtain, for any $k \geq 0$,

$$\sqrt{\frac{A_{k+1}}{M}} \leq a_{k+1} = A_{k+1} - A_k = (\sqrt{A_{k+1}} - \sqrt{A_k})(\sqrt{A_{k+1}} + \sqrt{A_k}) \leq 2(\sqrt{A_{k+1}} - \sqrt{A_k})\sqrt{A_{k+1}}.$$

Cancelling $\sqrt{A_{k+1}}$ on both sides and using the fact that $A_0 = 0$, we obtain, for any $k \geq 1$,

$$\sqrt{A_k} \geq \frac{k}{2\sqrt{M}}.$$

Squaring both sides, substituting the resulting inequality into (30) and replacing $M = \beta L$, we get (10).

When $\mu > 0$, we can drop the first term from the right-hand side of (29). This gives us

$$a_{k+1}^2 \geq \frac{\rho}{M} A_{k+1}^2.$$

Hence, for any $k \geq 0$,

$$A_{k+1} - A_k = a_{k+1} \geq q A_{k+1}, \quad q \stackrel{\text{def}}{=} \sqrt{\frac{\rho}{M}} \leq 1,$$

or, equivalently,

$$A_{k+1} \geq \frac{A_k}{1 - q}.$$

Consequently, for any $k \geq 1$,

$$A_k \geq \frac{A_1}{(1 - q)^{k-1}} \geq \frac{1}{M(1 - q)^{k-1}},$$

where the final inequality is due to (29) combined with the fact that $A_0 = 0$. Substituting this inequality into (30) and replacing $M = \beta L$, $\rho = \alpha\mu$, we get the second bound from Theorem 3.2. \square

B.3. Proof of Lemma 4.1

Let us denote by $u_k(\mathbf{a})$ the k -th power sum of the variables:

$$u_k(\mathbf{a}) \stackrel{\text{def}}{=} \sum_{i=1}^{n-1} a_i^k, \quad \forall \mathbf{a} \in \mathbb{R}^{n-1}.$$

Then, the classical Newton-Girard identities (see, e.g. (Kalman, 2000)) state the following relation between the elementary symmetric polynomials:

$$\sigma_\tau(\mathbf{a}) \equiv \frac{1}{\tau} \sum_{i=1}^{\tau} (-1)^{i-1} \sigma_{\tau-i}(\mathbf{a}) \cdot u_i(\mathbf{a}). \quad (31)$$

Note that for the matrix $U_\tau \stackrel{\text{def}}{=} \text{tr}(\mathbf{B}^\tau) \mathbf{I} - \mathbf{B}^\tau$, the following spectral decomposition holds:

$$\begin{aligned} U_\tau &= \mathbf{Q} \text{Diag} \left(\sum_{i=1}^n \lambda_i^\tau - \lambda_1^\tau, \sum_{i=1}^n \lambda_i^\tau - \lambda_2^\tau, \dots, \sum_{i=1}^n \lambda_i^\tau - \lambda_n^\tau \right) \mathbf{Q}^\top \\ &= \mathbf{Q} \text{Diag} \left(u_\tau(\boldsymbol{\lambda}_{-1}), u_\tau(\boldsymbol{\lambda}_{-2}), \dots, u_\tau(\boldsymbol{\lambda}_{-n}) \right) \mathbf{Q}^\top. \end{aligned} \quad (32)$$

Now, the identity that we need to prove is

$$\mathbf{P}_\tau = \mathbf{Q} \text{Diag} \left(\sigma_\tau(\boldsymbol{\lambda}_{-1}), \sigma_\tau(\boldsymbol{\lambda}_{-2}), \dots, \sigma_\tau(\boldsymbol{\lambda}_{-n}) \right) \mathbf{Q}^\top. \quad (33)$$

We justify (33) by induction. By definition, $P_0 \stackrel{\text{def}}{=} I$ and $\sigma_0(\mathbf{a}) \equiv 1$, therefore (33) holds for $\tau = 0$, which is our base. Let us fix $\tau \geq 1$ and assume that (33) is true for all smaller indices. Then,

$$\begin{aligned} P_\tau &\stackrel{\text{def}}{=} \frac{1}{\tau} \sum_{i=1}^{\tau} (-1)^{i-1} P_{\tau-i} U_i \\ &\stackrel{(33),(32)}{=} Q \text{Diag} \left(\sum_{i=1}^{\tau} (-1)^{i-1} \sigma_{\tau-i}(\lambda_{-1}) \cdot u_i(\lambda_{-1}), \dots, \sum_{i=1}^{\tau} (-1)^{i-1} \sigma_{\tau-i}(\lambda_{-n}) \cdot u_i(\lambda_{-n}) \right) Q^\top \\ &\stackrel{(31)}{=} Q \text{Diag} \left(\sigma_\tau(\lambda_{-1}), \dots, \sigma_\tau(\lambda_{-n}) \right) Q^\top. \end{aligned}$$

Hence, (33) is proven for all $0 \leq \tau \leq n-1$. \square

B.4. Proof of Theorem 4.2

By Lemma 4.1, we have the following representation of our preconditioner:

$$P_\tau = Q \text{Diag}(\sigma_\tau(\lambda_{-1}), \sigma_\tau(\lambda_{-2}), \dots, \sigma_\tau(\lambda_{-n})) Q^\top.$$

Further, for the spectrum of the matrix

$$B^{1/2} P_\tau B^{1/2} = Q \text{Diag}(\lambda_1 \cdot \sigma_\tau(\lambda_{-1}), \lambda_2 \cdot \sigma_\tau(\lambda_{-2}), \dots, \lambda_n \cdot \sigma_\tau(\lambda_{-n})) Q^\top,$$

it holds, according to Lemma D.10, that

$$\lambda_1 \cdot \sigma_\tau(\lambda_{-1}) \geq \lambda_2 \cdot \sigma_\tau(\lambda_{-2}) \geq \dots \geq \lambda_n \cdot \sigma_\tau(\lambda_{-n}). \quad (34)$$

Consequently,

$$\lambda_n \cdot \sigma_\tau(\lambda_{-n}) I \leq B^{1/2} P_\tau B^{1/2} \leq \lambda_1 \cdot \sigma_\tau(\lambda_{-1}) I,$$

which proves the required bound. \square

B.5. Proof of Theorem 4.5

Let $B = Q \text{Diag}(\lambda) Q^\top$ be a spectral decomposition of B , where $\lambda = \lambda(B)$ and $Q \in \mathbb{R}^{n \times n}$ is an orthogonal matrix. Formula (3.5) in (Rodomanov & Kropotov, 2020) states that

$$\mathbb{E}_{S \sim \text{Vol}_{\tau+1}(B)} [I_S (B_{S \times S})^{-1} I_S^\top] = \frac{1}{\sigma_{\tau+1}(\lambda)} Q \text{Diag}(\sigma_\tau(\lambda_{-1}), \dots, \sigma_\tau(\lambda_{-n})) Q^\top.$$

(Their σ_τ is $\sigma_{\tau+1}$ in our notation.) But, according to Lemma 4.1,

$$Q \text{Diag}(\sigma_\tau(\lambda_{-1}), \dots, \sigma_\tau(\lambda_{-n})) Q^\top = P_\tau,$$

and the claim follows. \square

B.6. Proof of Lemma 4.3

Clearly, when $\tau = 0$, we have $\xi_0(\lambda) \equiv 1$. For $\tau = n-1$, inequalities (16) are in fact identities, and $\xi_{n-1}(\lambda) \stackrel{(15)}{=} \frac{\lambda_n}{\lambda_1}$.

To prove that $\xi_\tau(\lambda)$ is decreasing in τ , we need to justify that, for any $0 \leq \tau \leq n-2$, we have

$$\xi_\tau(\lambda) \stackrel{\text{def}}{=} \frac{\sigma_\tau(\lambda_{-1})}{\sigma_\tau(\lambda_{-n})} \geq \xi_{\tau+1}(\lambda) \stackrel{\text{def}}{=} \frac{\sigma_{\tau+1}(\lambda_{-1})}{\sigma_{\tau+1}(\lambda_{-n})}.$$

This follows from Lemma D.11 (recall that, by our assumptions, $\lambda_1 \geq \dots \geq \lambda_n > 0$).

To prove the limit, let us divide the right hand side of

$$\xi_\tau(\lambda) = \frac{\sigma_\tau(\lambda_{-1})}{\sigma_\tau(\lambda_{-n})} \leq \frac{\sum_{2 \leq i_1 < \dots < i_\tau \leq n-1} \lambda_{i_1} \dots \lambda_{i_\tau}}{\lambda_1 \dots \lambda_\tau}$$

by the biggest element from the sum, which is $\lambda_2 \dots \lambda_{\tau+1}$. Thus, we get

$$\xi_\tau(\lambda) \leq \frac{1+E}{(\lambda_1/\lambda_{\tau+1})},$$

where E is a finite sum of numbers that are smaller than 1. Hence, $\xi_\tau(\lambda) \rightarrow 0$ when $\frac{\lambda_1}{\lambda_{\tau+1}} \rightarrow \infty$. \square

B.7. Proof of Lemma 4.4

We can assume that $\tau \geq 1$ since otherwise the inequality is trivial ($\xi_0(\boldsymbol{\lambda}) \leq 1$).

Using the definition of $\xi_\tau(\boldsymbol{\lambda})$ and applying Lemma D.12 (recalling that $\lambda_1 \geq \dots \geq \lambda_n > 0$ by our assumption), we obtain

$$\xi_\tau(\boldsymbol{\lambda}) = \frac{\sigma_\tau(\boldsymbol{\lambda}_{-1})}{\sigma_\tau(\boldsymbol{\lambda}_{-n})} \leq \frac{\lambda_n + s_{\tau-1}}{\lambda_1 + s_{\tau-1}},$$

where $s_{\tau-1}$ is the sum of all but the $\tau - 1$ largest elements of $\boldsymbol{\lambda}_{-\{1,n\}} = (\lambda_2, \dots, \lambda_{n-1})$:

$$s_{\tau-1} = \sum_{i=\tau+1}^{n-1} \lambda_i.$$

Substituting this into the previous display, we get

$$\xi_\tau(\boldsymbol{\lambda}) \leq \frac{\lambda_n + \sum_{i=\tau+1}^{n-1} \lambda_i}{\lambda_1 + \sum_{i=\tau+1}^{n-1} \lambda_i} = \frac{\sum_{i=\tau+1}^n \lambda_i}{\lambda_1 + \sum_{i=\tau+1}^{n-1} \lambda_i},$$

which is exactly the desired inequality.

B.8. Proof of Proposition 5.2

The problem of finding the best polynomial preconditioner can be reformulated as minimizing the norm of a symmetric matrix over the set of (positive) polynomials of a fixed degree $\tau \geq 0$:

$$\min_{p_\tau \in \mathcal{P}_\tau} \left\{ \gamma(p_\tau) \stackrel{\text{def}}{=} \| \mathbf{B} p_\tau(\mathbf{B}) - \mathbf{I} \| \right\},$$

where $\mathcal{P}_\tau \stackrel{\text{def}}{=} \{p_\tau \in \mathbb{R}[s] : \deg(p_\tau) = \tau, p_\tau(\mathbf{B}) \succ 0\}$. Here we use the spectral norm to measure the size of a symmetric matrix, and the objective can be rewritten as

$$\gamma(p_\tau) = \max_{s \in \text{Spec}(\mathbf{B})} |s p_\tau(s) - 1|, \quad (35)$$

where $\text{Spec}(\mathbf{B})$ is the discrete set of eigenvalues of the curvature matrix. For any value of $\gamma := \gamma(p_\tau)$, our original approximation guarantee (7) clearly satisfied with $\beta = 1 + \gamma$ and $\alpha = 1 - \gamma$, and the condition number becomes¹⁰

$$\frac{\beta}{\alpha} = \frac{1+\gamma}{1-\gamma}. \quad (36)$$

Now, we take

$$q_\tau(s) := \left(1 - \frac{s}{\lambda_1}\right) \left(1 - \frac{s}{\lambda_2}\right) \cdots \left(1 - \frac{s}{\lambda_\tau}\right). \quad (37)$$

First, note that $q_\tau(0) = 1$ and thus the polynomial $1 + q_\tau(s) \cdot (\alpha s - 1)$ is divisible by s . Hence the degree of the polynomial

$$p_\tau(s) := \frac{1 + q_\tau(s) \cdot (\alpha s - 1)}{s}$$

is exactly τ . Then, we obtain

$$\begin{aligned} \gamma &= \max_{s \in \text{Spec}(\mathbf{B})} |s p_\tau(s) - 1| = \max_{s \in \text{Spec}(\mathbf{B})} |q_\tau(s) \cdot (\alpha s - 1)| \\ &\leq \max_{s \in \{\lambda_{\tau+1}, \dots, \lambda_n\}} |\alpha s - 1| = \frac{\lambda_{\tau+1} - \lambda_n}{\lambda_{\tau+1} + \lambda_n}, \end{aligned} \quad (38)$$

and the optimal value is $\alpha = \frac{2}{\lambda_{\tau+1} + \lambda_n}$, where we put formally $\lambda_{n+1} \equiv \lambda_n$. It remains to substitute this bound into

$$\frac{\beta}{\alpha} = \frac{1+\gamma}{1-\gamma},$$

which is monotone in γ . □

¹⁰We are interested in $\gamma < 1$, since $\gamma = 1$ trivially holds for zero polynomial.

B.9. Proof of Proposition 5.3

Let us use an upper bound on γ from (35), which is the *uniform* polynomial approximation for the whole interval $[\lambda_n, \lambda_1]$:

$$\gamma(p_\tau) \leq \max_{s \in [\lambda_n, \lambda_1]} |sp_\tau(s) - 1|. \quad (39)$$

Then, we use

$$Q_\tau(s) \stackrel{\text{def}}{=} T_{\tau+1}\left(\frac{\lambda_1 + \lambda_n - 2s}{\lambda_1 - \lambda_n}\right) \cdot T_{\tau+1}\left(\frac{\lambda_1 + \lambda_n}{\lambda_1 - \lambda_n}\right)^{-1}, \quad (40)$$

where $T_{\tau+1}(\cdot)$ is the standard Chebyshev polynomial of the first kind of degree $\tau + 1$. Namely, we can define them recursively:

$$T_0(x) \stackrel{\text{def}}{=} 1, \quad T_1(x) \stackrel{\text{def}}{=} x, \quad T_{k+1}(x) \stackrel{\text{def}}{=} 2x \cdot T_k(x) - T_{k-1}(x), \quad k \geq 1.$$

Note that $Q_\tau(0) = 1$, thus the polynomial $1 - Q_\tau(s)$ is divisible by s . Then, we take

$$p_\tau(s) := \frac{1 - Q_\tau(s)}{s},$$

which is the polynomial of degree τ . This choice ensures that

$$\begin{aligned} \gamma &\stackrel{(39),(40)}{\leq} \max_{x \in [-1, 1]} |T_{\tau+1}(x)| \cdot T_{\tau+1}\left(\frac{\lambda_1 + \lambda_n}{\lambda_1 - \lambda_n}\right)^{-1} \\ &= T_{\tau+1}\left(\frac{\lambda_1 + \lambda_n}{\lambda_1 - \lambda_n}\right)^{-1} \leq 2 \left(\frac{\sqrt{\lambda_1} - \sqrt{\lambda_n}}{\sqrt{\lambda_1} + \sqrt{\lambda_n}}\right)^{\tau+1}, \end{aligned}$$

where the last inequality is the classical bound for the Chebyshev polynomials (see, e.g. Section 16.4 in (Vishnoi, 2013)). Thus, the condition number

$$\frac{\beta}{\alpha} = \frac{1 + \gamma}{1 - \gamma}$$

decreases exponentially with τ . □

B.10. Proof of Theorem 5.1

Let us fix an arbitrary $\mathbf{P} = \mathbf{P}^\top \succ 0$ such that $\mathbf{P} = p_\tau(\mathbf{B})$, for some polynomial $p_\tau \in \mathbb{R}[s]$ and $\deg(p_\tau) = \tau$. We take $\beta := \beta(\mathbf{P})$ and $\alpha := \alpha(\mathbf{P})$ (from the definition (7)) and denote

$$\bar{\mathbf{P}} := \frac{1}{\beta L} \mathbf{P}.$$

Let us consider an arbitrary iteration $k \geq 0$, and denote the following step

$$\mathbf{T} := \mathbf{x}_k - \bar{\mathbf{P}} \nabla f(\mathbf{x}_k).$$

Recall also that

$$\mathbf{x}_{k+1} := \mathbf{x}_k - \mathbf{P}_{\mathbf{a}_k} \nabla f(\mathbf{x}_k).$$

By the optimality of \mathbf{a}_k as the projection of $\mathbf{B}^{-1} \nabla f(\mathbf{x}_k)$ onto the Krylov subspace, we have:

$$\begin{aligned} \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{B}}^2 &\equiv \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k + \frac{1}{L} \mathbf{B}^{-1} \nabla f(\mathbf{x}_k)\|_{\mathbf{B}}^2 - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_{\mathbf{B}^{-1}}^2 \\ &\leq \frac{L}{2} \|\mathbf{T} - \mathbf{x}_k + \frac{1}{L} \mathbf{B}^{-1} \nabla f(\mathbf{x}_k)\|_{\mathbf{B}}^2 - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_{\mathbf{B}^{-1}}^2 \equiv \langle \nabla f(\mathbf{x}_k), \mathbf{T} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{T} - \mathbf{x}_k\|_{\mathbf{B}}^2. \end{aligned} \quad (41)$$

Hence, we obtain, for any $\mathbf{y} \in \mathbb{R}^n$:

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\stackrel{(6)}{\leq} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{B}}^2 \\ &\stackrel{(41)}{\leq} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{T} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{T} - \mathbf{x}_k\|_{\mathbf{B}}^2 \\ &\leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{T} - \mathbf{x}_k \rangle + \frac{\beta L}{2} \|\mathbf{T} - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2 \\ &\leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{y} - \mathbf{x}_k \rangle + \frac{\beta L}{2} \|\mathbf{y} - \mathbf{x}_k\|_{\mathbf{P}^{-1}}^2. \end{aligned}$$

where we used that \mathbf{T} is the minimizer of the last upper bound in \mathbf{y} . Thus, using convexity, we get, for any $\mathbf{y} \in \mathbb{R}^n$:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}) + \frac{\beta L}{2} \|\mathbf{y} - \mathbf{x}_k\|_{\mathcal{P}^{-1}}^2 \leq f(\mathbf{y}) + \frac{\beta}{\alpha} \cdot \frac{L}{2} \|\mathbf{y} - \mathbf{x}_k\|_{\mathcal{B}}^2. \quad (42)$$

In particular, substituting $\mathbf{y} := \mathbf{x}_k$ we justify that the method is *monotone*: $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k), \forall k \geq 0$. Therefore,

$$\mathbf{x}_k \in \mathcal{F}_0 \stackrel{\text{def}}{=} \left\{ \mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \leq f(\mathbf{x}_0), \right\}$$

and we assume that the initial level set \mathcal{F}_0 is *bounded*, denoting

$$D_0 \stackrel{\text{def}}{=} \sup_{\mathbf{x} \in \mathcal{F}_0} \|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{B}} < +\infty.$$

Substituting $\mathbf{y} := \gamma_k \mathbf{x}^* + (1 - \gamma_k) \mathbf{x}_k, \gamma_k \in [0, 1]$ into (42), we obtain

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq \gamma_k f^* + (1 - \gamma_k) f(\mathbf{x}_k) + \gamma_k^2 \frac{\beta}{\alpha} \cdot \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathcal{B}}^2 \\ &\leq \gamma_k f^* + (1 - \gamma_k) f(\mathbf{x}_k) + \gamma_k^2 \frac{\beta}{\alpha} \cdot \frac{L}{2} D_0^2. \end{aligned} \quad (43)$$

Substituting $\gamma_k := \frac{2}{k+1}$ and using the standard technique (see the proof of Theorem 3.1), we establish the global rate for the convex case:

$$f(\mathbf{x}_k) - f^* \leq \mathcal{O}\left(\frac{\beta}{\alpha} \cdot \frac{L D_0^2}{k}\right).$$

For strongly convex functions ($\mu > 0$), we continue as

$$f(\mathbf{x}_{k+1}) \stackrel{(43),(6)}{\leq} \gamma_k f^* + (1 - \gamma_k) f(\mathbf{x}_k) + \gamma_k^2 \frac{\beta L}{\alpha \mu} \cdot (f(\mathbf{x}_k) - f^*),$$

and choosing $\gamma_k := \frac{\alpha \mu}{2\beta L}$ we establish the exponential rate. □

C. Adaptive Search

In this section, we briefly present adaptive versions of Algorithms 1 and 2 which do not require the knowledge of the constant $M = \beta L$ and can automatically “tune” it in iterations yet preserving the original worst-case efficiency estimates. This is achieved by using a standard “backtracking line search” which can be found, e.g., in (Nesterov, 2013).

In what follows, for any $\mathbf{x}, \mathbf{y} \in \text{dom } \psi, M > 0$ and $\mathbf{P} = \mathbf{P}^\top \succ 0$, we define the following predicate:

$$\text{QuadGrowth}_{M, \mathbf{P}}(\mathbf{x}, \mathbf{y}): \quad f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{M}{2} \|\mathbf{y} - \mathbf{x}\|_{\mathcal{P}^{-1}}^2.$$

According to our assumptions (6) and (7), we know that this predicate is surely satisfied for any pair of points once $M \geq \beta L$.

The adaptive version of Algorithm 1 is presented in Algorithm 4. This method starts with a certain initial guess \tilde{M}_0 for the constant βL and then, at every iteration, repeatedly increases the current guess in two times until the predicate becomes satisfied. This process is guaranteed to terminate (when M_k becomes bigger or equal to βL , or even sooner). After that, we accept the new point \mathbf{x}_{k+1} and choose a new “optimistic” guess of the constant M for the next iteration by halving the value of M_k that we have accepted at the current iteration.

We assume that the preconditioner \mathbf{P} is sufficiently simple so that we can efficiently check the predicate $\text{QuadGrowth}_{M_k, \mathbf{P}}(\mathbf{x}_k, \mathbf{x}_{k+1})$. For example, if $\psi = 0$, then $\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{1}{M_k} \mathbf{P} \nabla f(\mathbf{x}_k)$ and $M_k \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathcal{P}^{-1}}^2 = \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}_{k+1} \rangle$ can be efficiently computed.

For Algorithm 4, we can prove exactly the same rates as in Theorem 3.1 (up to absolute constants) provided that

$$\tilde{M}_0 \leq \beta L. \quad (44)$$

The proof is essentially the same as in Appendix B.1 with only two minor differences: 1) inequality (23) is now guaranteed by our predicate; 2) instead of using $M = \beta L$ in (26), we should use the bound $M_k \leq 2\beta L$ which follows from (44) and

Algorithm 4 Adaptive Preconditioned GM

Input: $\mathbf{x}_0 \in \text{dom } \psi$, $\mathbf{P} = \mathbf{P}^\top \succ 0$, $\tilde{M}_0 > 0$.

for $k = 0, 1, \dots$ **do**

Find smallest integer $i_k \geq 0$ such that

$$\mathbf{x}_{k+1} = \text{GradStep}_{M_k, \mathbf{P}}(\mathbf{x}_k, \nabla f(\mathbf{x}_k)), \quad M_k = 2^{i_k} \tilde{M}_k$$

satisfies the predicate $\text{QuadGrowth}_{M_k, \mathbf{P}}(\mathbf{x}_k, \mathbf{x}_{k+1})$.

Set $\tilde{M}_{k+1} = M_k/2$.

end for

the fact that any value of $M \geq \beta L$ is always acceptable in the line search. Using a classical argument from (Nesterov, 2013), it is not difficult to show that, on average, Algorithm 4 makes only ~ 2 gradient steps at each iteration.

In contrast to an upper estimate of the constant βL , an initial guess satisfying (44) can be easily generated. One simple recipe is to make a trial step $\mathbf{x}'_1 = \text{GradStep}_{M'_0, \mathbf{P}}(\mathbf{x}_0, \nabla f(\mathbf{x}_0))$ for an *arbitrarily chosen* $M'_0 > 0$ and then compute

$$\tilde{M}_0 = \frac{f(\mathbf{x}'_1) - f(\mathbf{x}_0) - \langle \nabla f(\mathbf{x}_0), \mathbf{x}'_1 - \mathbf{x}_0 \rangle}{\frac{1}{2} \|\mathbf{x}'_1 - \mathbf{x}_0\|_{\mathbf{P}^{-1}}^2}.$$

Alternatively, we can find a suitable \tilde{M}_0 by choosing an arbitrary $M'_0 > 0$ and then repeatedly halving it until the predicate $\text{QuadGrowth}(\mathbf{x}_0, \mathbf{x}'_1(M))$ stops being satisfied for $\mathbf{x}'_1(M) = \text{GradStep}_{M, \mathbf{P}}(\mathbf{x}_0, \nabla f(\mathbf{x}_0))$. This auxiliary procedure either terminates in a logarithmic number of steps, in which case we get a suitable \tilde{M}_0 , or, otherwise, we quickly find an approximate solution of our problem.

Similar technique can be applied for the Fast Gradient Method. Specifically, let us introduce an auxiliary procedure shown in Algorithm 5 for computing one iteration of Algorithm 2 for a given value of M . Then, the adaptive FGM method can be constructed as shown in Algorithm 6. As in the basic method, we can show that the rates from Theorem 3.2 still remain valid (up to absolute constants) for Algorithm 6, provided that \tilde{M}_0 satisfies (44). For generating the initial guess \tilde{M}_0 , we can use exactly the same techniques as before.

Algorithm 5 $(\mathbf{x}_+, \mathbf{v}_+, A_+; \mathbf{y}) = \text{FastGradStep}_{M, \rho, \mathbf{P}}(\mathbf{x}, \mathbf{v}, A)$

Require: $M > 0$; $\rho \geq 0$; $\mathbf{P} = \mathbf{P}^\top \succ 0$; $\mathbf{x}, \mathbf{v} \in \text{dom } \psi$; $A > 0$.

Find a_+ from eq. $\frac{Ma_+^2}{A+a_+} = 1 + \rho(A + a_+)$.

Set $A_+ = A + a_+$, $H = \frac{1+\rho A_+}{a_+}$, $\theta = \frac{a_+}{A_+}$, $\omega = \frac{\rho}{H}$, $\gamma = \frac{\omega(1-\theta)}{1-\omega\theta}$.

Set $\hat{\mathbf{v}} = (1 - \gamma)\mathbf{v} + \gamma\mathbf{x}$, $\mathbf{y} = (1 - \theta)\mathbf{x} + \theta\hat{\mathbf{v}}$.

Compute $\mathbf{v}_+ = \text{GradStep}_{H, \mathbf{P}}(\hat{\mathbf{v}}, \nabla f(\mathbf{y}))$.

Set $\mathbf{x}_+ = (1 - \theta)\mathbf{x} + \theta\mathbf{v}_+$.

Return $(\mathbf{x}_+, \mathbf{v}_+, A_+; \mathbf{y})$.

Algorithm 6 Adaptive Preconditioned FGM

Input: $\mathbf{x}_0 \in \text{dom } \psi$, $\mathbf{P} = \mathbf{P}^\top \succ 0$, $\tilde{M}_0 > 0$.

Set $\mathbf{v}_0 = \mathbf{x}_0$, $A_0 = 0$.

for $k = 0, 1, \dots$ **do**

Find smallest integer $i_k \geq 0$ such that

$$(\mathbf{x}_{k+1}, \mathbf{v}_{k+1}, A_{k+1}; \mathbf{y}_k) = \text{FastGradStep}_{M_k, \mathbf{P}}(\mathbf{x}_k, \mathbf{v}_k, A_k), \quad M_k = 2^{i_k} \tilde{M}_k$$

satisfies the predicate $\text{QuadGrowth}_{M_k, \mathbf{P}}(\mathbf{y}_k, \mathbf{x}_{k+1})$.

Set $\tilde{M}_{k+1} = M_k/2$.

end for

D. Auxiliary Results

D.1. Elementary Symmetric Polynomials

The *elementary symmetric polynomial* in variables $\mathbf{x} \in \mathbb{R}^n$ of degree k (integer, $1 \leq k \leq n$) is defined as

$$\sigma_k(\mathbf{x}) := \sum_{1 \leq i_1 < \dots < i_k \leq n} x_{i_1} \dots x_{i_k}.$$

It will be convenient to extend this definition to arbitrary integer degrees k and also to the case $n = 0$ (which corresponds to the empty vector \mathbf{x}). For this, we additionally define, for any $\mathbf{x} \in \mathbb{R}^n$ with $n \geq 0$,

$$\sigma_k(\mathbf{x}) := \begin{cases} 1, & \text{if } k = 0, \\ 0, & \text{if } k < 0 \text{ or } k > n. \end{cases}$$

Thus, $\sigma_k(\mathbf{x})$ is defined for any $\mathbf{x} \in \mathbb{R}^n$ with $n \geq 0$ and any integer k .

The following three properties are obvious from the definition.

Observation D.1 (symmetry). *For any $\mathbf{x} \in \mathbb{R}^n$ with $n \geq 0$, any integer k , and any permutation $\boldsymbol{\pi} := (\pi_1, \dots, \pi_n)$ of indices $\{1, \dots, n\}$, we have $\sigma_k(\mathbf{x}) = \sigma_k(\mathbf{x}_{\boldsymbol{\pi}})$, where $\mathbf{x}_{\boldsymbol{\pi}} := (x_{\pi_1}, \dots, x_{\pi_n}) \in \mathbb{R}^n$ is the vector obtained from \mathbf{x} by rearranging¹¹ its components according to $\boldsymbol{\pi}$.*

Observation D.2. *For any $\mathbf{x} \in \mathbb{R}_+^n$ with $n \geq 0$ and any integer k , we have $\sigma_k(\mathbf{x}) \geq 0$.*

Observation D.3. *For any $\mathbf{x} \in \mathbb{R}_+^n$ with at least $1 \leq k \leq n$ strictly positive elements, we have $\sigma_k(\mathbf{x}) > 0$.*

Let us now establish a number of other properties that will be useful in our analysis.

In what follows, given a vector \mathbf{x} with $n \geq 1$ elements (indexed by $1, \dots, n$), and an index $1 \leq i \leq n$, we use the notation \mathbf{x}_{-i} to denote the $(n-1)$ -dimensional vector obtained from \mathbf{x} by removing its i th element. More generally, for a set of indices $I \subseteq \{1, \dots, n\}$, we denote by \mathbf{x}_{-I} the $(n-|I|)$ -dimensional vector obtained from \mathbf{x} by removing the elements with indices from I .

Also, for a vector \mathbf{x} with $n \geq 0$ elements, we use the notation $\mathbf{x} \geq 0$ to express the fact that each element of \mathbf{x} is nonnegative. For an empty vector \mathbf{x} (i.e., when $n = 0$), this inequality is always assumed to be satisfied (vacuously).

We start with a simple but very useful recursive decomposition.

Lemma D.4. *For any $\mathbf{x} \in \mathbb{R}^n$, any index $1 \leq i \leq n$, and any integer k , we have*

$$\sigma_k(\mathbf{x}) = x_i \sigma_{k-1}(\mathbf{x}_{-i}) + \sigma_k(\mathbf{x}_{-i}).$$

Proof. In view of Observation D.1, it suffices to prove the identity only for $i = 1$.

If $k < 0$ or $k > n$, then $\sigma_k(\mathbf{x}) = \sigma_{k-1}(\mathbf{x}_{-1}) = \sigma_k(\mathbf{x}_{-1}) = 0$, and we get the identity $0 = 0$ which is indeed valid.

If $k = 0$, then $\sigma_{k-1}(\mathbf{x}_{-1}) = 0$, while $\sigma_k(\mathbf{x}) = \sigma_k(\mathbf{x}_{-1}) = 1$, and we get the identity $1 = 1$ which is also valid.

If $k = 1$, then

$$\sigma_k(\mathbf{x}) = \sum_{j=1}^n x_j = x_1 + \sum_{j=2}^n x_j = x_1 + \sigma_1(\mathbf{x}_{-1}),$$

so the claim is valid since $\sigma_0(\mathbf{x}_{-1}) = 1$ by definition.

Finally, in the general case when $2 \leq k \leq n$, we have

$$\begin{aligned} \sigma_k(\mathbf{x}) &= \sum_{1 \leq i_1 < \dots < i_k \leq n} x_{i_1} \dots x_{i_k} = x_1 \sum_{2 \leq i_2 < \dots < i_k \leq n} x_{i_2} \dots x_{i_k} + \sum_{2 \leq i_1 < \dots < i_k \leq n} x_{i_1} \dots x_{i_k} \\ &= x_1 \sigma_{k-1}(\mathbf{x}_{-1}) + \sigma_k(\mathbf{x}_{-1}). \end{aligned}$$

□

¹¹By convention, the empty vector gets rearranged into the empty vector.

Next, we consider several inequalities between elementary symmetric polynomials of different degrees.

Lemma D.5 (Weak Newton's inequality). *For any $\mathbf{x} \in \mathbb{R}_+^n$ with $n \geq 0$ and any integer k , we have*

$$\sigma_k^2(\mathbf{x}) \geq \sigma_{k-1}(\mathbf{x})\sigma_{k+1}(\mathbf{x}).$$

Remark D.6. In the only nontrivial case $1 \leq k \leq n-1$, the above is a weaker version of the classical *Newton's inequality*:

$$\hat{\sigma}_k^2(\mathbf{x}) \geq \hat{\sigma}_{k-1}(\mathbf{x})\hat{\sigma}_{k+1}(\mathbf{x}),$$

where, for any $1 \leq k \leq n$, $\hat{\sigma}_k(\mathbf{x}) := \frac{\sigma_k(\mathbf{x})}{\binom{n}{k}}$ is the normalized elementary symmetric polynomial in variables \mathbf{x} , and $\binom{n}{k}$ is the binomial coefficient. For more information about the two inequalities, see, e.g., Section 2.22 in (Hardy et al., 1952). We present the proof of Lemma D.5 below for the reader's convenience.

Proof. Note that the claim is obvious if either $k < 1$ or $k \geq n$ since then either $\sigma_{k-1}(\mathbf{x}) = 0$ or $\sigma_{k+1}(\mathbf{x}) = 0$, and we get the trivial inequality $\sigma_k^2(\mathbf{x}) \geq 0$.

Let us prove the claim by induction on n .

For $n = 1$, the claim is obvious since, in this case, all values of k satisfy either $k < 1$ or $k \geq n$.

Assume that $n \geq 2$, and that we have already proved the claim for $n' = n - 1$. Let us prove it for $n' = n$.

According to the observation made at the beginning, we can assume that $1 \leq k \leq n - 1$ since otherwise the claim is obvious. Since both sides of the claimed inequality are continuous in \mathbf{x} (as certain polynomials), we can further assume w.l.o.g. that all the elements of \mathbf{x} are strictly positive.

According to Lemma D.4, we can decompose

$$\sigma_k(\mathbf{x}) = x_1\sigma_{k-1} + \sigma_k, \quad \sigma_{k-1}(\mathbf{x}) = x_1\sigma_{k-2} + \sigma_{k-1}, \quad \sigma_{k+1}(\mathbf{x}) = x_1\sigma_k + \sigma_{k+1},$$

where $\sigma_{k'} := \sigma_{k'}(\mathbf{x}_{-1}) \geq 0$ for any $k' \in \{k-2, k-1, k, k+1\}$ (see Observation D.2). Note that \mathbf{x}_{-1} has exactly $n-1 \geq 1$ elements, all of which are strictly positive (since we have assumed that all the elements of \mathbf{x} are so). In particular, we can assume that $\sigma_k > 0$ (see Observation D.3 and recall that $1 \leq k \leq n-1$ by our assumption).

The inequality we need to prove is then

$$(x_1\sigma_{k-1} + \sigma_k)^2 \geq (x_1\sigma_{k-2} + \sigma_{k-1})(x_1\sigma_k + \sigma_{k+1}).$$

After the expansion of both sides, the above inequality becomes

$$x_1^2\sigma_{k-1}^2 + 2x_1\sigma_{k-1}\sigma_k + \sigma_k^2 \geq x_1^2\sigma_{k-2}\sigma_k + x_1(\sigma_{k-2}\sigma_{k+1} + \sigma_{k-1}\sigma_k) + \sigma_{k-1}\sigma_{k+1}.$$

Making cancellations and rearranging, we come to the following inequality we need to prove:

$$x_1^2(\sigma_{k-1}^2 - \sigma_{k-2}\sigma_k) + x_1(\sigma_{k-1}\sigma_k - \sigma_{k-2}\sigma_{k+1}) + (\sigma_k^2 - \sigma_{k-1}\sigma_{k+1}) \geq 0.$$

Since $x_1 \geq 0$, it suffices to prove the following three inequalities:

$$\sigma_{k-1}^2 \geq \sigma_{k-2}\sigma_k, \quad \sigma_{k-1}\sigma_k \geq \sigma_{k-2}\sigma_{k+1}, \quad \sigma_k^2 \geq \sigma_{k-1}\sigma_{k+1}.$$

But this is simple: the first and the third ones are valid in view of our inductive assumption, while the second one follows from the other two (indeed, $\sigma_{k-1}\sigma_k^2 \geq \sigma_{k-1}^2\sigma_{k+1} \geq \sigma_{k-2}\sigma_k\sigma_{k+1}$, and it remains to cancel σ_k which is assumed to be positive). \square

Lemma D.7. *For any $\mathbf{x} \in \mathbb{R}^n$, any index $1 \leq i \leq n$, such that $\mathbf{x}_{-i} \geq 0$, and any integer k , we have*

$$\sigma_k(\mathbf{x})\sigma_k(\mathbf{x}_{-i}) \geq \sigma_{k+1}(\mathbf{x})\sigma_{k-1}(\mathbf{x}_{-i}).$$

Proof. Let us denote for brevity $\sigma_{k'} := \sigma_{k'}(\mathbf{x}_{-i})$ for any $k' \in \{k-1, k, k+1\}$. Then, according to Lemma D.4,

$$\sigma_k(\mathbf{x}) = x_i \sigma_{k-1} + \sigma_k, \quad \sigma_{k+1}(\mathbf{x}) = x_i \sigma_k + \sigma_{k+1}.$$

The inequality we need to justify is then

$$(x_i \sigma_{k-1} + \sigma_k) \sigma_k \geq (x_i \sigma_k + \sigma_{k+1}) \sigma_{k-1},$$

or, equivalently,

$$\sigma_k^2 \geq \sigma_{k-1} \sigma_{k+1}.$$

But this is indeed true according to Lemma D.5 (and our assumption that $\mathbf{x}_{-i} \geq 0$). □

Lemma D.8. For any $\mathbf{x} \in \mathbb{R}_+^n$ with $n \geq 0$ and any integer $0 \leq k \leq n$, we have

$$\sigma_{k+1}(\mathbf{x}) \leq \bar{s}_k^\uparrow(\mathbf{x}) \sigma_k(\mathbf{x}),$$

where $\bar{s}_k^\uparrow(\mathbf{x})$ is the sum of all but k largest elements of \mathbf{x} (i.e., $\bar{s}_k^\uparrow(\mathbf{x}) = \sum_{i=k+1}^n x_{[i]}$, where $x_{[1]} \geq \dots \geq x_{[n]}$ are the components of \mathbf{x} sorted in nonincreasing order, and $\bar{s}_k^\uparrow(\mathbf{x}) := 0$ for the empty vector \mathbf{x}).

Proof. We can assume that $n \geq 1$ since otherwise $k = n = 0$, the vector \mathbf{x} is empty, $\sigma_{k+1}(\mathbf{x}) = \bar{s}_k^\uparrow(\mathbf{x}) = 0$, and we get a trivial inequality $0 \leq 0$.

Further, in view of Observation D.1, we can assume w.l.o.g. that $x_1 \geq \dots \geq x_n$. Then, we need to prove that

$$\sigma_{k+1}(\mathbf{x}) \leq \left(\sum_{i=k+1}^n x_i \right) \sigma_k(\mathbf{x}).$$

Since both sides of the above inequality are continuous in \mathbf{x} (as certain polynomials), we can assume w.l.o.g. that all the components of \mathbf{x} are strictly positive.

Applying repeatedly Lemma D.7, we obtain

$$\frac{\sigma_{k+1}(\mathbf{x})}{\sigma_k(\mathbf{x})} \leq \frac{\sigma_k(\mathbf{x}_{-1})}{\sigma_{k-1}(\mathbf{x}_{-1})} \leq \frac{\sigma_{k-1}(\mathbf{x}_{-\{1,2\}})}{\sigma_{k-2}(\mathbf{x}_{-\{1,2\}})} \leq \dots \leq \frac{\sigma_1(\mathbf{x}_{-\{1,\dots,k\}})}{\sigma_0(\mathbf{x}_{-\{1,\dots,k\}})} = \sum_{i=k+1}^n x_i,$$

where the final identity follows from the definitions of σ_1 and σ_0 . (Note that each denominator in the above display is strictly positive, see Observation D.3 and recall that, by our assumptions, \mathbf{x} has strictly positive components and $0 \leq k \leq n$.) This is exactly the desired inequality. □

Lemma D.9. For any $\mathbf{x} \in \mathbb{R}^n$, any indices $1 \leq i, j \leq n$, any integer k , and any $a, b \in \mathbb{R}$, we have the implication

$$x_i \geq x_j \quad \text{and} \quad a \sigma_k(\mathbf{x}_{-\{i,j\}}) \& b \sigma_{k-1}(\mathbf{x}_{-\{i,j\}}) \quad \implies \quad (ax_i + b) \sigma_k(\mathbf{x}_{-i}) \& (ax_j + b) \sigma_k(\mathbf{x}_{-j}),$$

where “&” is either “ \leq ” or “ \geq ”. Furthermore, if $x_i > x_j$, then the reverse implication is also true.

Proof. We can assume that $i \neq j$ since otherwise the claim is trivial. In particular, we can assume that $n \geq 2$.

According to Lemma D.4, we can decompose

$$\sigma_k(\mathbf{x}_{-i}) = x_j \sigma_{k-1} + \sigma_k, \quad \sigma_k(\mathbf{x}_{-j}) = x_i \sigma_{k-1} + \sigma_k,$$

where $\sigma_{k'} := \sigma_{k'}(\mathbf{x}_{-\{i,j\}})$ for any $k' \in \{k-1, k\}$. The inequality after the implication sign is then

$$(ax_i + b)(x_j \sigma_{k-1} + \sigma_k) \& (ax_j + b)(x_i \sigma_{k-1} + \sigma_k).$$

After the expansion of both sides, this inequality reads

$$x_i x_j a \sigma_{k-1} + x_i a \sigma_k + x_j b \sigma_{k-1} + b \sigma_k \& x_i x_j a \sigma_{k-1} + x_i b \sigma_{k-1} + x_j a \sigma_k + b \sigma_k.$$

Making cancellations and rearranging, we obtain the following equivalent inequality:

$$(x_i - x_j)a\sigma_k \& (x_i - x_j)b\sigma_{k-1}.$$

This inequality is obviously true if $x_i \geq x_j$ and $a\sigma_k \& b\sigma_{k-1}$. On the other hand, if $x_i > x_j$, we can cancel $(x_i - x_j)$ on both sides and conclude that $a\sigma_k \& b\sigma_{k-1}$. \square

Lemma D.10. For any $\mathbf{x} \in \mathbb{R}^n$, any indices $1 \leq i, j \leq n$, such that $\mathbf{x}_{-\{i,j\}} \geq 0$, and any integer k , we have the implication

$$x_i \geq x_j \quad \implies \quad x_i\sigma_k(\mathbf{x}_{-i}) \geq x_j\sigma_k(\mathbf{x}_{-j}).$$

Proof. Follows from Lemma D.9 (applied to $a = 1$ and $b = 0$) since $\sigma_k(\mathbf{x}_{-\{i,j\}}) \geq 0$ in view of our assumption that $\mathbf{x}_{-\{i,j\}} \geq 0$ (see Observation D.2). \square

Lemma D.11. For any $\mathbf{x} \in \mathbb{R}^n$, any indices $1 \leq i, j \leq n$, such that $\mathbf{x}_{-\{i,j\}} \geq 0$, and any integer k , we have the implication

$$x_i \geq x_j \quad \implies \quad \sigma_k(\mathbf{x}_{-i})\sigma_{k+1}(\mathbf{x}_{-j}) \geq \sigma_{k+1}(\mathbf{x}_{-i})\sigma_k(\mathbf{x}_{-j}).$$

Proof. For brevity, let us denote $\sigma_{k'} := \sigma_{k'}(\mathbf{x}_{-\{i,j\}})$ for any $k' \in \{k-1, k, k+1\}$. According to Lemma D.5 and our assumption that $\mathbf{x}_{-\{i,j\}} \geq 0$, we have $\sigma_k^2 \geq \sigma_{k-1}\sigma_{k+1}$. Since we also assume that $x_i \geq x_j$, we can therefore apply Lemma D.9 with $a = \sigma_k$ and $b = \sigma_{k+1}$ to get

$$(x_i\sigma_k + \sigma_{k+1})\sigma_k(\mathbf{x}_{-i}) \geq (x_j\sigma_k + \sigma_{k+1})\sigma_k(\mathbf{x}_{-j}).$$

This is exactly the desired inequality since, according to Lemma D.4,

$$\sigma_{k+1}(\mathbf{x}_{-i}) = x_j\sigma_k + \sigma_{k+1}, \quad \sigma_{k+1}(\mathbf{x}_{-j}) = x_i\sigma_k + \sigma_{k+1}. \quad \square$$

Lemma D.12. For any $\mathbf{x} \in \mathbb{R}^n$, any indices $1 \leq i, j \leq n$, with $\mathbf{x}_{-\{i,j\}} \geq 0$, and any integer¹² $1 \leq k \leq \dim(\mathbf{x}_{-\{i,j\}}) + 1$, the following implication holds:

$$x_i \geq x_j \quad \implies \quad \sigma_k(\mathbf{x}_{-i})(x_i + \bar{s}_{k-1}^\uparrow(\mathbf{x}_{-\{i,j\}})) \leq \sigma_k(\mathbf{x}_{-j})(x_j + \bar{s}_{k-1}^\uparrow(\mathbf{x}_{-\{i,j\}})),$$

where $\bar{s}_{k-1}^\uparrow(\mathbf{x}_{-\{i,j\}})$ is the sum of all but $k-1$ largest elements of the vector $\mathbf{x}_{-\{i,j\}}$.

Proof. Follows from Lemma D.9 applied to $a = 1$ and $b = \bar{s}_{k-1}^\uparrow(\mathbf{x}_{-\{i,j\}})$ since $\sigma_k(\mathbf{x}_{-\{i,j\}}) \leq b\sigma_{k-1}(\mathbf{x}_{-\{i,j\}})$ according to Lemma D.8 (applied to $k' = k-1$ and $\mathbf{x}' = \mathbf{x}_{-\{i,j\}}$; note that $\dim(\mathbf{x}') = n-2 \geq 0$ and $0 \leq k' \leq \dim(\mathbf{x}')$ by our assumption). \square

¹²Here, $\dim(\mathbf{x})$ denotes the number of elements in the vector \mathbf{x} .