



Published in final edited form as:

J Optim Theory Appl. 2018 July ; 178(1): 240–263. doi:10.1007/s10957-018-1287-4.

Adaptive Restart of the Optimized Gradient Method for Convex Optimization

Donghwan Kim,

Jeffrey A. Fessler

University of Michigan, Ann Arbor, MI

Abstract

First-order methods with momentum such as Nesterov’s fast gradient method are very useful for convex optimization problems, but can exhibit undesirable oscillations yielding slow convergence for some applications. An adaptive restarting scheme can improve the convergence rate of the fast gradient method, when the parameter of a strongly convex cost function is unknown or when the iterates of the algorithm enter a locally well-conditioned region. Recently, we introduced an optimized gradient method, a first-order algorithm that has an inexpensive per-iteration computational cost similar to that of the fast gradient method, yet has a worst-case cost function convergence bound that is twice smaller than that of the fast gradient method and that is optimal for large-dimensional smooth convex problems. Building upon the success of accelerating the fast gradient method using adaptive restart, this paper investigates similar heuristic acceleration of the optimized gradient method. We first derive new step coefficients of the optimized gradient method for a strongly convex quadratic problem with known function parameters, yielding a convergence rate that is faster than that of the analogous version of the fast gradient method. We then provide a heuristic analysis and numerical experiments that illustrate that adaptive restart can accelerate the convergence of the optimized gradient method. Numerical results also illustrate that adaptive restart is helpful for a proximal version of the optimized gradient method for nonsmooth composite convex functions.

Keywords

Convex optimization; First-order algorithms; Accelerated gradient method; Optimized gradient method; Restarting

Mathematics Subject Classification (2000)

80M50; 90C06; 90C25

1 Introduction

The computational expense of first-order methods depends only mildly on the problem dimension, so they are attractive for solving large-dimensional optimization problems [1].

In particular, Nesterov's fast gradient method (FGM) [2,3,4] is used widely because it has a worst-case cost function convergence bound that is optimal up to a constant for large-dimensional smooth convex problems [3]. In addition, for smooth and strongly convex problems where the strong convexity parameter is known, a version of FGM has a linear convergence rate [3] that improves upon that of a standard gradient method. However, without knowledge of the function parameters, conventional FGM does not guarantee a linear convergence rate.

When the strong convexity parameter is unknown, a simple adaptive restarting scheme [5] for FGM heuristically improves its convergence rate (see also [6,7] for theory and [1,8,9] for applications). In addition, adaptive restart is useful even when the function is only locally strongly convex near the minimizer [5]. First-order methods are known to be suitable when only moderate solution accuracy is required, and adaptive restart can help first-order methods achieve medium to high accuracy.

Recently we proposed the optimized gradient method (OGM) [10] (built upon [11]) that has efficient per-iteration computation similar to FGM yet that achieves the optimal worst-case convergence bound for decreasing a large-dimensional smooth convex function among all first-order methods with fixed or dynamic step sizes [12]. (See [13,14,15] for further analysis and extensions of OGM.) This paper examines OGM for strongly convex *quadratic* functions and develops an OGM variant that provides a linear convergence rate that is faster than that of FGM. The analysis reveals that, like FGM, OGM may exhibit undesirable oscillating behavior in some cases. Building on the quadratic analysis of FGM in [5], we propose an adaptive restart scheme [5] that heuristically accelerates the convergence rate of OGM when the function is strongly convex or even when it is only locally well-conditioned. This restart scheme circumvents the oscillating behavior. Numerical results illustrate that the proposed OGM with restart performs better than FGM with restart in [5].

Sec. 2 describes convex problem and reviews first-order algorithms for convex problems such as gradient method, FGM, and OGM. Sec. 3 studies OGM for strongly convex quadratic problems. Sec. 4 suggests an adaptive restart scheme for OGM using the quadratic analysis in Sec. 3. Sec. 5 illustrates the proposed adaptive version of OGM that we use for numerical experiments on various convex problems in Sec. 6, including nonsmooth composite convex functions, and Sec. 7 concludes.

2 Problem and Algorithms

2.1 Smooth and Strongly Convex Problem

We first consider the smooth and strongly convex minimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \tag{M}$$

that satisfies the following smooth and strongly convex conditions:

$$- f: \mathbb{R}^d \rightarrow \mathbb{R} \text{ has Lipschitz continuous gradient with Lipschitz constant } L > 0, \text{ i.e.,}$$

$$\| \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) \| \leq L \| \mathbf{x} - \mathbf{y} \|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d, \tag{1}$$

– f is strongly convex with strong convexity parameter $\mu > 0$, i.e.,

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{\mu}{2} \| \mathbf{x} - \mathbf{y} \|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d. \tag{2}$$

We let $\mathcal{F}_{\mu, L}(\mathbb{R}^d)$ denote the class of functions f that satisfy the above two conditions hereafter, and let \mathbf{x}^* denote the unique minimizer of f . We let $q := \frac{\mu}{L}$ denote the reciprocal of the condition number of a function $f \in \mathcal{F}_{\mu, L}(\mathbb{R}^d)$. We also let $\mathcal{F}_{0, L}(\mathbb{R}^d)$ denote the class of smooth convex functions f that satisfy the above two conditions with $\mu = 0$, and let \mathbf{x}^* denote a minimizer of f .

Some algorithms discussed in this paper require knowledge of both μ and L , but in many cases estimating μ is challenging compared to computing L .¹ Therefore, this paper focuses on the case where the parameter μ is unavailable while L is available. Even without knowing μ , the adaptive restart approach in [5] and the proposed approach in this paper both exhibit linear convergence rates.

We next review known accelerated first-order algorithms for solving (M).

2.2 Accelerated First-order Algorithms

This paper focuses on accelerated first-order algorithms of the form shown in Alg. 1. The fast gradient method (FGM) [2,3,4] (with $\gamma_k = 0$ in Alg. 1) accelerates the gradient method (GM) (with $\beta_k = \gamma_k = 0$) using the *momentum* term $\beta_k(\mathbf{y}_{k+1} - \mathbf{y}_k)$ with negligible additional computation. The optimized gradient method (OGM) [10,14] uses an over-relaxation term $\gamma_k(\mathbf{y}_{k+1} - \mathbf{x}_k) = -\gamma_k \alpha \nabla f(\mathbf{x}_k)$ for further acceleration.

Algorithm 1 Accelerated First-order Algorithms

```

1: Input:  $f \in \mathcal{F}_{\mu, L}(\mathbb{R}^d)$ ,  $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^d$ .
2: for  $k \geq 0$  do
3:    $\mathbf{y}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)$ 
4:    $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \beta_k(\mathbf{y}_{k+1} - \mathbf{y}_k) + \gamma_k(\mathbf{y}_{k+1} - \mathbf{x}_k)$ 

```

2.2.1 Fast Gradient Method (FGM)—For the function class $\mathcal{F}_{0, L}(\mathbb{R}^d)$, the following coefficients are the standard choice for FGM [4, 2]:

$$\alpha = \frac{1}{L}, \quad \beta_k = \frac{t_k - 1}{t_k + 1}, \quad \gamma_k = 0, \quad t_k = \begin{cases} 1, & k = 0, \\ \frac{1}{2} \left(1 + \sqrt{1 + 4t_{k-1}^2} \right), & \text{otherwise,} \end{cases} \tag{3}$$

where β_k (3) increases from 0 towards 1 as $k \rightarrow \infty$, and the resulting primary iterates $\{\mathbf{y}_k\}$ satisfy the following bound [4, Thm. 4.4]:

¹For some applications even estimating L is expensive, and one must employ a backtracking scheme [4] or similar approaches. We assume L is known throughout this paper. An estimate of μ could be found by a backtracking scheme as described in [16, Sec. 5.3].

$$f(\mathbf{y}_k) - f(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2t_k^2 - 1} \leq \frac{2L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{(k+1)^2}. \quad (4)$$

For the function class $\mathcal{F}_{\mu, L}(\mathbb{R}^d)$ with known $q > 0$, a typical choice for the coefficients of FGM [3, Eqn. (2.2.11)] is:

$$\alpha = \frac{1}{L}, \quad \beta_k = \frac{1 - \sqrt{q}}{1 + \sqrt{q}}, \quad \gamma_k = 0, \quad (5)$$

for which the primary iterates $\{\mathbf{y}_k\}$ of FGM satisfy the following linear convergence bound [3, Thm. 2.2.3]:

$$f(\mathbf{y}_k) - f(\mathbf{x}^*) \leq (1 - \sqrt{q})^k \frac{(1+q)L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2}. \quad (6)$$

Although FGM converges faster than GM [3], The convergence rate of FGM is not optimal for both function classes $\mathcal{F}_{0, L}(\mathbb{R}^d)$ and $\mathcal{F}_{\mu, L}(\mathbb{R}^d)$, and finding optimal algorithms for such classes is of interest. We next review the recently proposed OGM [10,14] (built upon [11]) that has an optimal worst-case cost function convergence rate for the function class $\mathcal{F}_{0, L}(\mathbb{R}^d)$ for large-scale problems [12].

2.2.2 Optimized Gradient Method (OGM)—For function class $\mathcal{F}_{0, L}(\mathbb{R}^d)$, the usual coefficients for OGM are [10]:

$$\alpha = \frac{1}{L}, \beta_k = \frac{\theta_k - 1}{\theta_{k+1}}, \gamma_k = \frac{\theta_k}{\theta_{k+1}}, \theta_k = \begin{cases} 1, & k = 0, \\ \frac{1}{2} \left(1 + \sqrt{1 + 4\theta_{k-1}^2} \right), & k = 1, \dots, N-1, \\ \frac{1}{2} \left(1 + \sqrt{1 + 8\theta_{k-1}^2} \right) & k = N, \end{cases} \quad (7)$$

for a given total number of iterations N . For these coefficients, the last secondary iterate \mathbf{x}_N of OGM satisfies the following bound for $\mathcal{F}_{0, L}(\mathbb{R}^d)$ [10, Thm. 2]:

$$f(\mathbf{x}_N) - f(\mathbf{x}^*) \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2\theta_N^2} \leq \frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{(N+1)(N+1+\sqrt{2})}, \quad (8)$$

which is twice smaller than the bound (4) of FGM and is optimal for first-order methods (with fixed or dynamic step sizes) for the function class $\mathcal{F}_{0, L}(\mathbb{R}^d)$ under the large-scale condition $d \gg N+1$ [12].

In addition, for the following coefficients [14] that are independent of N :

$$\alpha = \frac{1}{L}, \quad \beta_k = \frac{t_k - 1}{t_{k+1}}, \quad \gamma_k = \frac{t_k}{t_{k+1}}, \quad t_k = \begin{cases} 1, & k = 0, \\ \frac{1}{2} \left(1 + \sqrt{1 + 4t_{k-1}^2} \right), & k = 1, \dots, \end{cases} \quad (9)$$

the primary iterates $\{\mathbf{y}_k\}$ of OGM satisfy the following bound [14, Thm. 4.1]:

$$f(\mathbf{y}_k) - f(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{4t_k^2 - 1} \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{(k+1)^2} \quad (10)$$

for $\mathcal{F}_{0,L}(\mathbb{R}^d)$. Here, as $k \rightarrow \infty$, β_k and γ_k in (9) increase from 0 and $\frac{\sqrt{5}-1}{2} \approx 0.618$ to both 1, respectively. Note that the coefficients in (7) and (9) differ only at the last iteration N . Interestingly, OGM-type acceleration with the coefficients (9) was studied for accelerating the proximal point method long ago in [17, Appx.].

It is yet unknown whether some choice of OGM coefficients will yield a linear convergence rate for the general function class $\mathcal{F}_{\mu,L}(\mathbb{R}^d)$ that is faster than the rate (6) for FGM; this topic is left as an interesting future work.² Towards this direction, Sec. 3 studies OGM for strongly convex *quadratic* problems, improving upon FGM. Sec. 4 uses this quadratic analysis to analyze an adaptive restart scheme for OGM.

3 Analysis of OGM for Quadratic Functions

This section analyzes the behavior of OGM for minimizing a strongly convex quadratic function. We optimize the coefficients of OGM for such quadratic function, yielding a linear convergence rate that is faster than that of FGM. The quadratic analysis of OGM in this section is similar in spirit to the analyses of a heavy-ball method [19, Sec. 3.2] and FGM [20, Appx. A] [5, Sec. 4].

The resulting OGM requires the knowledge of q , and we show that using the coefficients (7) or (9) instead (without the knowledge of q) will cause the OGM iterates to oscillate when the momentum is larger than a critical value. This analysis stems from the dynamical system analysis of FGM in [5, Sec. 4].

3.1 Quadratic Analysis of OGM

This section considers minimizing a strongly convex quadratic function:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{p}^\top \mathbf{x} \in \mathcal{F}_{\mu,L}(\mathbb{R}^d) \quad (11)$$

where $\mathbf{Q} \in \mathbb{R}^{d \times d}$ is a symmetric positive definite matrix, $\mathbf{p} \in \mathbb{R}^d$ is a vector. Here, $\nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{p}$ is the gradient, and $\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{p}$ is the optimum. The smallest and the largest eigenvalues of \mathbf{Q} correspond to the parameters μ and L of the function respectively. For

²Very recently, [18] developed a new first-order method with known q that achieves a linear convergence rate $(1 - \sqrt{q})^2$ for the cost function decrease that is faster than the linear rate $(1 - \sqrt{q})$ in (6) for FGM.

simplicity, for the quadratic analysis we consider the version of OGM that has constant coefficients (α, β, γ) .

Defining the vectors $\xi_k := (\mathbf{x}_k^\top, \mathbf{x}_{k-1}^\top)^\top \in \mathbb{R}^{2d}$ and $\xi_* := (\mathbf{x}_*^\top, \mathbf{x}_*^\top)^\top \in \mathbb{R}^{2d}$, and extending the analysis for FGM in [20, Appx. A], OGM with constant coefficients (α, β, γ) has the following equivalent form for $k \geq 1$:

$$\xi_{k+1} - \xi_* = \mathbf{T}(\alpha, \beta, \gamma)(\xi_k - \xi_*), \tag{12}$$

where the system matrix $\mathbf{T}(\alpha, \beta, \gamma)$ of OGM is defined as

$$\mathbf{T}(\alpha, \beta, \gamma) := \begin{bmatrix} (1 + \beta)(\mathbf{I} - \alpha\mathbf{Q}) - \gamma\alpha\mathbf{Q} & -\beta(\mathbf{I} - \alpha\mathbf{Q}) \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{2d \times 2d} \tag{13}$$

for an identity matrix $\mathbf{I} \in \mathbb{R}^{d \times d}$. The sequence $\{\tilde{\xi}_k := (\mathbf{y}_k^\top, \mathbf{y}_{k-1}^\top)^\top\}_{k \geq 1}$ also satisfies the recursion (12), implying that (12) characterizes the behavior of both the primary sequence $\{\mathbf{y}_k\}$ and the secondary sequence $\{\mathbf{x}_k\}$ of OGM with constant coefficients.

The spectral radius $\rho(\mathbf{T}(\cdot))$ of matrix $\mathbf{T}(\cdot)$ determines the convergence rate of the algorithm. Specifically, for any $\epsilon > 0$, there exists $K \geq 0$ such that $[\rho(\mathbf{T})]^k \|\mathbf{T}^k\| \leq (\rho(\mathbf{T}) + \epsilon)^k$ for all $k \geq K$, establishing the following convergence bound:

$$\|\xi_{k+1} - \xi_*\| \leq (\rho(\mathbf{T}(\alpha, \beta, \gamma)) + \epsilon)^k \|\xi_1 - \xi_*\|. \tag{14}$$

We next analyze $\rho(\mathbf{T}(\alpha, \beta, \gamma))$ for OGM.

Considering the eigen-decomposition of \mathbf{Q} in $\mathbf{T}(\cdot)$ as in [20, Appx. A], the spectral radius of $\mathbf{T}(\cdot)$ is:

$$\rho(\mathbf{T}(\alpha, \beta, \gamma)) = \max_{\mu \leq \lambda \leq L} \rho(\mathbf{T}_\lambda(\alpha, \beta, \gamma)), \tag{15}$$

where for any eigenvalue λ of matrix \mathbf{Q} we define a matrix $\mathbf{T}_\lambda(\alpha, \beta, \gamma) \in \mathbb{R}^{2 \times 2}$ by plugging in λ and 1 instead of \mathbf{Q} and \mathbf{I} in $\mathbf{T}(\alpha, \beta, \gamma)$ respectively. Similar to the analysis of FGM in [20, Appx. A], the spectral radius of $\mathbf{T}_\lambda(\alpha, \beta, \gamma)$ for OGM is:

$$\rho(\mathbf{T}_\lambda(\alpha, \beta, \gamma)) = \max\{|r_1(\alpha, \beta, \gamma, \lambda)|, |r_2(\alpha, \beta, \gamma, \lambda)|\} \tag{16}$$

$$= \begin{cases} \frac{1}{2}(|(1 + \beta)(1 - \alpha\lambda) - \gamma\alpha\lambda| + \sqrt{\Delta(\alpha, \beta, \gamma, \lambda)}), & \Delta(\alpha, \beta, \gamma, \lambda) \geq 0, \\ \sqrt{\beta(1 - \alpha\lambda)}, & \text{otherwise,} \end{cases}$$

where $r_1(\alpha, \beta, \gamma, \lambda)$ and $r_2(\alpha, \beta, \gamma, \lambda)$ denote the roots of the characteristic polynomial of $\mathbf{T}_\lambda(\cdot)$:

$$r^2 - ((1 + \beta)(1 - \alpha\lambda) - \gamma\alpha\lambda)r + \beta(1 - \alpha\lambda), \quad (17)$$

and $(\alpha, \beta, \gamma, \lambda) = ((1 + \beta)(1 - \alpha\lambda) - \gamma\alpha\lambda)^2 - 4\beta(1 - \alpha\lambda)$ denotes the corresponding discriminant. For fixed (α, β, γ) , the spectral radius $\rho(\mathbf{T}_\lambda(\alpha, \beta, \gamma))$ in (16) is a continuous and quasi-convex³ function of λ ; thus its maximum over λ occurs at one of its boundary points $\lambda = \mu$ or $\lambda = L$.

The next section optimizes the coefficients (α, β, γ) of OGM to provide the fastest convergence rate, *i.e.*, the smallest spectral radius $\rho(\mathbf{T}(\cdot))$ in (15).

3.2 Optimizing OGM Coefficients

We would like to choose OGM coefficients that provide the fastest convergence for minimizing a strongly convex quadratic function *i.e.*, to solve

$$\arg \min_{\alpha, \beta, \gamma} \rho(\mathbf{T}(\alpha, \beta, \gamma)) = \arg \min_{\alpha, \beta, \gamma} \max \{ \rho(\mathbf{T}_\mu(\alpha, \beta, \gamma)), \rho(\mathbf{T}_L(\alpha, \beta, \gamma)) \}. \quad (18)$$

Note that it is yet unknown which form of first-order algorithm with fixed coefficients is optimal for decreasing a strongly convex quadratic function (*e.g.*, [18]), so our focus here is simply to optimize the coefficients within the OGM algorithm class.

Similar coefficient optimization was studied previously for GM and FGM, which is equivalent to optimizing (18) with additional constraints on (α, β, γ) . GM corresponds to the choice $\beta = \gamma = 0$, for which it is well known that optimizing (18) over α yields the optimal GM step size $\alpha = \frac{2}{\mu + L}$. Similarly, FGM with the standard choice (5) results from optimizing (18) over β for the choice⁴ $\alpha = \frac{1}{L}$ and $\gamma = 0$. Another version of FGM corresponds to the choice $\gamma = 0$, for which optimizing (18) over (α, β) yields coefficients $\alpha = \frac{4}{\mu + 3L}$, $\beta = \frac{\sqrt{3+q} - 2\sqrt{q}}{\sqrt{3+q} + 2\sqrt{q}}$, in [20, Prop. 1].

Although a general unconstrained solution to (18) would be an interesting future direction, here we focus on optimizing (18) over (β, γ) for the choice $\alpha = \frac{1}{L}$. This choice simplifies the problem (18) and is useful for analyzing an adaptive restart scheme for OGM in Sec. 4.

³It is straightforward to show that $\rho(\mathbf{T}_\lambda(\alpha, \beta, \gamma))$ in (16) is quasi-convex over λ . First, $\sqrt{\beta(1 - \alpha\lambda)}$ is quasi-convex over λ (for $(\alpha, \beta, \gamma, \lambda) < 0$). Second, the eigenvalue λ satisfying $(\alpha, \beta, \gamma, \lambda) = 0$ is in the region where the function $\frac{1}{2}(|(1 + \beta)(1 - \alpha\lambda) - \gamma\alpha\lambda| + \sqrt{\Delta(\alpha, \beta, \gamma, \lambda)})$ either monotonically increases or decreases, which overall makes the continuous function $\rho(\mathbf{T}_\lambda(\alpha, \beta, \gamma))$ quasi-convex over λ . This proof can be simply applied to other variables, *i.e.*, $\rho(\mathbf{T}_\lambda(\alpha, \beta, \gamma))$ is quasi-convex over either α, β or γ .

⁴For FGM with (5), the value of $\rho(\mathbf{T}_L(1/L, \beta, 0))$ is 0, and the function $\rho(\mathbf{T}_\mu(1/L, \beta, 0))$ is continuous and quasi-convex over β (see footnote 3). The minimum of $\rho(\mathbf{T}_\mu(1/L, \beta, 0))$ occurs at the point $\beta = \frac{1 - \sqrt{q}}{1 + \sqrt{q}}$ in (5) satisfying $(1/L, \beta, 0, \mu) = 0$, verifying the statement that FGM with (5) results from optimizing (18) over β given $\alpha = \frac{1}{L}$ and $\gamma = 0$.

3.3 Optimizing the Coefficients (β, γ) of OGM When $\alpha = 1/L$

When $\alpha = \frac{1}{L}$ and $\lambda = L$, the characteristic polynomial (17) becomes $r^2 + \gamma r = 0$. The roots are $r = 0$ and $r = -\gamma$, so $\rho(\mathbf{T}_L(1/L, \beta, \gamma)) = |\gamma|$. In addition, because $\rho(\mathbf{T}_\mu(1/L, \beta, \gamma))$ is continuous and quasi-convex over β (see footnote 3), it can be easily shown that the smaller value of β satisfying the following equation:

$$\begin{aligned} \Delta(1/L, \beta, \gamma, \mu) &= ((1 + \beta)(1 - q) - \gamma q)^2 - 4\beta(1 - q) \\ &= (1 - q)^2 \beta^2 - 2(1 - q)(1 + q + \gamma q)\beta + (1 - q)(1 - q - 2\gamma q) + q^2 \gamma^2 = 0 \end{aligned} \tag{19}$$

minimizes $\rho(\mathbf{T}_\mu(1/L, \beta, \gamma))$ for any given γ satisfying $|\gamma| \leq 1$. The optimal β is

$$\beta^*(\gamma) := (1 - \sqrt{q(1 + \gamma)})^2 / (1 - q), \tag{20}$$

which reduces to $\beta = \beta^*(0) = \frac{1 - \sqrt{q}}{1 + \sqrt{q}}$ in (5) for FGM (with $\gamma = 0$). Substituting (20) into (16) yields $\rho(\mathbf{T}_\mu(1/L, \beta^*(\gamma), \gamma)) = |1 - \sqrt{q(1 + \gamma)}|$, leading to the following simplification of (18) with $\alpha = \frac{1}{L}$ and $\beta = \beta^*(\gamma)$ from (20):

$$\gamma^* := \arg \min_{\gamma} \max \left\{ |1 - \sqrt{q(1 + \gamma)}|, |\gamma| \right\}. \tag{21}$$

The minimizer of (21) satisfies $1 - \sqrt{q(1 + \gamma)} = \pm \gamma$, and with simple algebra, we get the following solutions to (18) with $\alpha = \frac{1}{L}$ (and (21)):

$$\beta^* := \beta^*(\gamma^*) = \frac{(\gamma^*)^2}{1 - q} = \frac{(2 + q - \sqrt{q^2 + 8q})^2}{4(1 - q)}, \quad \gamma^* = \frac{2 + q - \sqrt{q^2 + 8q}}{2}, \tag{22}$$

for which the spectral radius is $\rho^* := \rho(\mathbf{T}(1/L, \beta^*, \gamma^*)) = 1 - \sqrt{q(1 + \gamma^*)} = \gamma^*$.

Table 1 compares the spectral radius of the new optimally tuned OGM to existing optimally tuned GM and FGM. Simple algebra shows that the spectral radius of OGM is smaller than those of FGM, *i.e.*, $\frac{2 + q - \sqrt{q^2 + 8q}}{2} \leq 1 - \frac{2\sqrt{q}}{\sqrt{3 + q}} \leq 1 - \sqrt{q}$. Therefore, OGM based on (22) achieves a worst-case convergence rate of $\|\xi_k - \xi^*\|$ that is faster than that of FGM for a strongly convex quadratic function.

To further understand the behavior of OGM for each eigen-mode, Fig. 1 plots $\rho(\mathbf{T}_\lambda(1/L, \beta, \gamma))$ for $\mu = \lambda = L$ for $q = 0.1$ as an example, where $(\beta^*, \gamma^*) = (0.4, 0.6)$. Fig. 1 first compares the OGM spectral radius values with optimally tuned coefficients $(\alpha, \beta, \gamma) = (1/L, \beta^*(\gamma^*), \gamma^*)$ from (22) to those of the optimally tuned $\beta^*(\gamma)$ in (20) for other choices of $\gamma = 0, 0.4, 0.8$. The optimal choice (β^*, γ^*) (upper red curve in Fig. 1) has worst-case spectral radius values at both the smallest and the largest eigenvalues, unlike other choices of γ (with

$\beta^*(\gamma)$ where either $\rho(\mathbf{T}_\mu(1/L, \beta, \gamma))$ or $\rho(\mathbf{T}_L(1/L, \beta, \gamma))$ are largest. The other choices thus have a spectral radius larger than that of the optimally tuned OGM.

Fig. 1 also illustrates spectral radius values for different choices of β for given $\gamma = \gamma^*$, showing that suboptimal β value will slow down convergence. OGM with $(\alpha, \beta, \gamma) = (1/L, 0, \gamma)$ is equivalent to GM with $\alpha = \frac{1}{L}(1 + \gamma)$, and Fig. 1 illustrates this choice for comparison. Interestingly, GM has some modes for mid-valued λ values that will converge faster than in the accelerated methods, but its overall convergence rate is worse. Apparently no one algorithm can have superior convergence rates for all modes.

Although using the optimized coefficients (β^*, γ^*) leads to OGM having the smallest possible overall spectral radius $\rho(\mathbf{T}(\cdot))$, the upper red and blue curves in Fig. 1 illustrate that this “tuned” OGM will have modes for large eigenvalues that converge slower than with OGM with $\gamma = 0$ (*i.e.*, FGM). This behavior may be undesirable when such modes dominate the overall convergence behavior. Interestingly, Sec. 3.4 describes that the convergence of the primary sequence $\{\mathbf{y}_k\}$ of OGM is not governed by such modes unlike the secondary sequence $\{\mathbf{x}_k\}$ of OGM. Fig. 1 reveals change points across λ meaning that there are different regimes; Sec. 3.4 elaborates on this behavior building upon a dynamical system analysis of FGM [5, Sec. 4].

3.4 Convergence Properties of OGM When $\alpha = 1/L$

[5, Sec. 4] analyzed a constant-step FGM as a linear dynamical system for minimizing a strongly convex quadratic function (11), and showed that there are three regimes of behavior for the system; low momentum, optimal momentum, and high momentum regimes. This section similarly analyzes OGM to better understand its convergence behavior when solving a strongly convex quadratic problem (11), complementing the previous section’s spectral radius analysis of OGM

We use the eigen-decomposition of $\mathbf{Q} = \mathbf{V}\mathbf{A}\mathbf{V}^\top$ with $\mathbf{A} = \text{diag}\{\lambda_i\}$, where the eigenvalues $\{\lambda_i\}$ are in an ascending order, *i.e.*, $\mu = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_d = L$. And for simplicity, we let $\mathbf{p} = \mathbf{0}$ without loss of generality, leading to $\mathbf{x}^* = \mathbf{0}$. By defining $\mathbf{w}_k := (w_{k,1}, \dots, w_{k,d})^\top = \mathbf{V}^\top \mathbf{y}_k \in \mathbb{R}^d$ and $\mathbf{v}_k := (v_{k,1}, \dots, v_{k,d})^\top = \mathbf{V}^\top \mathbf{x}_k \in \mathbb{R}^d$ as the mode coefficients of the primary and secondary sequences respectively and using (12), we have the following d independently evolving identical recurrence relations for the evolution of $w_{\cdot,i}$ and $v_{\cdot,i}$ of the constant-step OGM respectively:

$$w_{k+2,i} = ((1 + \beta)(1 - \lambda_i/L) - \gamma\lambda_i/L)w_{k+1,i} - \beta(1 - \lambda_i/L)w_{k,i}, \quad (23)$$

$$v_{k+2,i} = ((1 + \beta)(1 - \lambda_i/L) - \gamma\lambda_i/L)v_{k+1,i} - \beta(1 - \lambda_i/L)v_{k,i},$$

for $i = 1, \dots, d$, although the initial conditions differ as follows:

$$w_{1,i} = (1 - \lambda_i/L)w_{0,i}, \quad v_{1,i} = ((1 + \beta + \gamma)(1 - \lambda_i/L) - (\beta + \gamma))v_{0,i} \quad (24)$$

with $w_{0,i} = v_{0,i}$. The convergence behavior of the i th dynamical system of both $w_{i,j}$ and $v_{i,j}$ in (23) is determined by the characteristic polynomial (17) with $\alpha = \frac{1}{L}$ and $\lambda = \lambda_i$. Unlike the previous sections that studied only the worst-case convergence performance using the largest absolute value of the roots of the polynomial (17), we next discuss the convergence behavior of OGM more comprehensively using (17) with $\alpha = \frac{1}{L}$ and $\lambda = \lambda_i$ for the cases where 1) $\lambda_i = L$ and 2) $\lambda_i < L$.

1) $\lambda_i = L$: The characteristic polynomial (17) of the mode of $\lambda_i = L$ reduces to $r^2 + \gamma r = 0$ with two roots 0 and $-\gamma$ regardless of the choice of β . Thus we have monotone convergence for this (d th) mode of the dynamical system [21, Sec. 17.1]:

$$w_{k,d} = 0^k + c_d(-\gamma)^k, \quad v_{k,d} = 0^k + \hat{c}_d(-\gamma)^k, \tag{25}$$

where c_d and \hat{c}_d are constants depending on the initial conditions (24). Substituting $w_{1,d} = 0$ and $v_{1,d} = -(\beta + \gamma)v_{0,d}$ (24) into (23) yields $c_d = 0$ and $\hat{c}_d = v_{0,d}\left(1 + \frac{\beta}{\gamma}\right)$, illustrating that the primary sequence $\{w_{k,d}\}$ reaches its optimum after one iteration, whereas the secondary sequence $\{v_{k,d}\}$ has slow monotone convergence of the distance to the optimum, while exhibiting undesirable oscillation due to the term $(-\gamma)^k$, corresponding to overshooting over the optimum.

2) $\lambda_i < L$: In (22) we found the optimal overall β^* for OGM. One can alternatively explore what the best value of β would be for any given mode of the system for comparison. The polynomial (17) has repeated roots for the following β , corresponding to the smaller zero of the discriminant $(1/L, \beta, \gamma, \lambda_i)$ for given γ and λ_i :

$$\beta_i^*(\gamma) := \left(1 - \sqrt{(1 + \gamma)\lambda_i/L}\right)^2 / (1 - \lambda_i/L). \tag{26}$$

This root satisfies $\beta^* = \beta^*(\gamma^*) = \beta_1^*(\gamma^*)$ (22), because λ_1 is the smallest eigenvalue. Next we examine the convergence behavior of OGM in the following three regimes, similar to FGM in [5, Sec. 4.3]:⁵

- $\beta < \beta_i^*(\gamma)$: low momentum, over-damped,
- $\beta = \beta_i^*(\gamma)$: optimal momentum, critically damped,
- $\beta > \beta_i^*(\gamma)$: high momentum, under-damped.

If $\beta \leq \beta_i^*(\gamma)$, the polynomial (17) has two real roots, $r_{1,i}$ and $r_{2,i}$ where we omit $(1/L, \beta, \gamma, \lambda_i)$ in $r_{i,j} = r(1/L, \beta, \gamma, \lambda_i)$ for simplicity. Then, the system evolves as [21, Sec. 17.1]:

$$w_{k,i} = c_{1,i}r_{1,i}^k + c_{2,i}r_{2,i}^k, \quad v_{k,i} = \hat{c}_{1,i}r_{1,i}^k + \hat{c}_{2,i}r_{2,i}^k, \tag{27}$$

⁵For simplicity in the momentum analysis, we restricted the choice of β within $[0, 1]$, containing the β_k values of FGM in (3), (5) and OGM in (7), (9). This restriction simply discards the effect of a larger solution β of $(1/L, \beta, \gamma, \lambda_i) = 0$ in the analysis, which is larger than 1.

where constants $c_{1,i}$, $c_{2,i}$, $\hat{c}_{1,i}$ and $\hat{c}_{2,i}$ depend on the initial conditions (24). In particular, when $\beta = \beta_i^*(\gamma)$, we have the repeated root $r_i^*(\gamma) := 1 - \sqrt{(1 + \gamma)\lambda_i/L}$, corresponding to critical damping, yielding the fastest monotone convergence among (27) for any β s.t. $\beta \leq \beta_i^*(\gamma)$. This property is due to the quasi-convexity of $\rho(\mathbf{T}_\lambda(1/L, \beta, \gamma))$ over β . If $\beta < \beta_i^*(\gamma)$, the system is over-damped, which corresponds to the low momentum regime, where the system is dominated by the larger root that is greater than $r_i^*(\gamma)$, and thus has slow monotone convergence. However, depending on the initial conditions (24), the system may only be dominated by the smaller root, as noticed for the case $\lambda_i = L$ in (25). Also note that the mode of $\lambda_i = L$ is always in the low momentum regime regardless of the value of β .

If $\beta > \beta_i^*(\gamma)$, the system is under-damped, which corresponds to the high momentum regime. This means that the system evolves as [21, Sec. 17.1]:

$$w_{k,i} = c_i(\sqrt{\beta(1 - \lambda_i/L)})^k \cos(k\psi_i(\beta, \gamma) - \delta_i), \tag{28}$$

$$v_{k,i} = \hat{c}_i(\sqrt{\beta(1 - \lambda_i/L)})^k \cos(k\psi_i(\beta, \gamma) - \hat{\delta}_i),$$

where the frequency of the oscillation is given by

$$\psi_i(\beta, \gamma) = \cos^{-1}(((1 + \beta)(1 - \lambda_i/L) - \gamma\lambda_i/L)/(2\sqrt{\beta(1 - \lambda_i/L)})), \tag{29}$$

and c_i , δ_i , \hat{c}_i and $\hat{\delta}_i$ denote constants that depend on the initial conditions (24); in particular for $\beta \approx 1$, we have $\delta_i \approx 0$ and $\hat{\delta}_i \approx 0$ so we will ignore them.

We categorize the behavior of the i th mode of OGM for each λ_j based on the above momentum analysis. Regimes with two curves and one curve in Fig. 1 correspond to the low- and high-momentum regimes, respectively. In particular, for $\beta = \beta^*(\gamma)$ in Fig. 1, most λ_j values experience high momentum (and the optimal momentum for λ_j satisfying $\beta^*(\gamma) = \beta_i^*(\gamma)$, e.g., $\lambda_j = \mu$), whereas modes where $\lambda_j \approx L$ experience low momentum. The fast convergence of the primary sequence $\{w_{k,d}\}$ in (25) generalizes to the case $\lambda_j \approx L$, corresponding to the lower curves in Fig. 1. In addition, for β smaller than $\beta^*(\gamma)$ in Fig. 1, both $\lambda \approx \mu$ and $\lambda \approx L$ experience low momentum so increasing β improves the convergence rate.

Based on the quadratic analysis, we would like to use appropriately large β and γ coefficients, namely (β^*, γ^*) , to have fast monotone convergence (for the dominating modes). However, such values require knowing the function parameter $q = \mu/L$ that is usually unavailable in practice. Using OGM (with coefficients (7) and (9)) without knowing q will likely lead to oscillation due to the high momentum (or under-damping) for strongly convex functions. The next section describes restarting schemes inspired by [5] that we suggest to use with OGM to avoid such oscillation and thus heuristically accelerate the rate

of OGM for a strongly convex quadratic function and even for a convex function that is locally well-conditioned.

4 Restarting Schemes

Restarting an algorithm (*i.e.*, starting the algorithm again by using the current iterate as the new starting point) after a certain number of iterations or when some restarting condition is satisfied has been found useful, *e.g.*, for the conjugate gradient method [22,23], called “fixed restart” and “adaptive restart” respectively. The fixed restart approach was also studied for accelerated gradient schemes such as FGM in [24, Sec. 11.4] [16]. Recently adaptive restart of FGM was shown to provide dramatic practical acceleration without requiring knowledge of function parameters [5,6,7]. Building upon those ideas, this section reviews and applies restarting approaches for OGM. A quadratic analysis in [5] justified using a restarting condition for FGM; this section extends that analysis to OGM by studying an observable quantity of oscillation that serves as an indicator for restarting the momentum of OGM.

4.1 Fixed Restart

Restarting an algorithm every k iterations can yield a linear rate for decreasing a function in $\mathcal{F}_{\mu, L}(\mathbb{R}^d)$ [24, Sec. 11.4] [16, Sec. 5.1]. We examine this restart approach for OGM here. Let $\mathbf{x}_{j,i}$ denote the j th outer iteration and i th inner iteration of an OGM variant that is restarted every k (inner) iterations. Specifically, this OGM uses $\mathbf{x}_{j+1,0} = \mathbf{x}_{j,k}$ to initialize the next $(j+1)$ th outer iteration. Combining the OGM bound (8) and the strong convexity inequality (2) yields the following linear rate of cost function decrease for k inner iterations of OGM:

$$f(\mathbf{x}_{j,k}) - f(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}_{j,0} - \mathbf{x}^*\|^2}{k^2} \leq \frac{2L}{\mu k^2} (f(\mathbf{x}_{j,0}) - f(\mathbf{x}^*)). \quad (30)$$

This bound is smaller than the $4L/\mu k^2$ bound for FGM with fixed restart (using the FGM bound (4)). Here, an optimal restarting interval k that minimizes the bound (30) for a given total number of steps jk is $k_{\text{fixed}} := e\sqrt{2Lq}$.

There are two drawbacks of the fixed restart approach [5, Sec. 3.1]. First, computing the optimal interval k_{fixed} requires knowledge of q that is usually unavailable in practice. Second, using a global parameters q may be too conservative when the iterates enter locally well-conditioned region. Therefore, adaptive restarting [5] has been found useful in practice, which we review next and then apply to OGM. The above two drawbacks also apply to the algorithms in Sec. 2 that assume knowledge of the global parameter q .

4.2 Adaptive Restart

To circumvent the drawbacks of fixed restart, [5] proposes the following two adaptive restart schemes for FGM:

– Function scheme for restarting (FR): restart whenever

$$f(\mathbf{y}_{k+1}) > f(\mathbf{y}_k), \quad (31)$$

– Gradient scheme for restarting (GR): restart whenever

$$\langle -\nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle < 0. \quad (32)$$

These schemes heuristically improve convergence rates of FGM with coefficients (3) and both performed similarly [5,7]. Although the function scheme guarantees monotonic decreasing function values, the gradient scheme has two advantages over the function scheme [5]; the gradient scheme involves only arithmetic operations with already computed quantities, and it is numerically more stable.

These two schemes encourage algorithm restart whenever the iterates take a “bad” direction, *i.e.*, when the function value increases or the negative gradient and the momentum have an obtuse angle, respectively. However, a convergence proof that justifies their empirical acceleration is yet unknown, so [5] analyzes such restarting schemes for strongly convex quadratic functions. An alternative scheme in [7] that restarts whenever the magnitude of the momentum decreases, *i.e.*, $\|\mathbf{y}_{k+1} - \mathbf{y}_k\| < \|\mathbf{y}_k - \mathbf{y}_{k-1}\|$, has a theoretical convergence analysis for the function class $\mathcal{F}_{\mu, L}(\mathbb{R}^d)$. However, empirically both the function and gradient schemes performed better in [7]. Thus, this paper focuses on adapting practical restart schemes to OGM and extending the analysis in [5] to OGM. First we introduce a new additional adaptive scheme designed specifically for OGM.

4.3 Adaptive Decrease of γ for OGM

Sec. 3.4 described that the secondary sequence $\{\mathbf{x}_k\}$ of OGM might experience overshooting and thus slow convergence, unlike the primary sequence $\{\mathbf{y}_k\}$, when the iterates enter a region where the mode of the largest eigenvalue dominates. (Sec. 6.1.2 illustrates such an example.) From (25), the overshoot of \mathbf{x}_k has magnitude proportional to γ , yet a suitably large γ , such as γ^* (21), is essential for overall acceleration.

To avoid (or reduce) such overshooting, we suggest the following adaptive scheme:

– Gradient scheme for decreasing γ (GD γ): decrease γ whenever

$$\langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_{k-1}) \rangle < 0. \quad (33)$$

Because the primary sequence $\{\mathbf{y}_k\}$ of OGM is unlikely to overshoot, one could choose to simply use the primary sequence $\{\mathbf{y}_k\}$ as algorithm output instead of the secondary sequence $\{\mathbf{x}_k\}$. However, if one needs to use the secondary sequence of OGM (*e.g.*, Sec. 5.2), adaptive scheme (33) can help.

4.4 Observable OGM Quantities

This section revisits Sec. 3.4 that suggested that observing the evolution of the mode coefficients $\{w_{k,i}\}$ and $\{v_{k,i}\}$ can help identify the momentum regime. However, in practice that evolution is unobservable because the optimum \mathbf{x}^* is unknown, whereas Sec. 3.4 assumed $\mathbf{x}^* = 0$. Instead we can observe the evolution of the function values, which are related to the mode coefficients as follows:

$$f(\mathbf{y}_k) = \frac{1}{2} \sum_{i=1}^d \lambda_i w_{k,i}^2, \quad f(\mathbf{x}_k) = \frac{1}{2} \sum_{i=1}^d \lambda_i v_{k,i}^2, \quad (34)$$

and also the inner products of the gradient and momentum, *i.e.*,

$$\langle -\nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle = - \sum_{i=1}^a \lambda_i v_{k,i} (w_{k+1,i} - w_{k,i}), \quad (35)$$

$$\langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_{k-1}) \rangle = \sum_{i=1}^d \lambda_i^2 v_{k,i} v_{k-1,i}. \quad (36)$$

These quantities appear in the conditions for the adaptive schemes (31), (32), and (33).

One would like to increase β and γ as large as possible for acceleration up to β^* and γ^* (22). However, without knowing q (and β^* , γ^*), we could end up placing the majority of the modes in the high momentum regime, eventually leading to slow convergence with oscillation as described in Sec. 3.4. To avoid such oscillation, we hope to detect it using (34) and (35) and restart the algorithm. We also hope to detect the overshooting (25) of the modes of the large eigenvalues (in the low momentum regime) using (36) so that we can then decrease γ and avoid such overshooting.

We focus on the case where $\beta > \beta_1(\gamma)$ for given γ , when the most of the modes are in the high momentum regime. Because the maximum of $\rho(\mathbf{T}_\lambda(1/L, \beta, \gamma))$ occurs at the points $\lambda = \mu$ or $\lambda = L$, we expect that (34), (35), and (36) will be quickly dominated by the mode of the smallest or the largest values. Using (25) and (28) leads to the following approximations:

$$f(\mathbf{y}_k) \approx \frac{1}{2} \mu c_1^2 \beta^k (1 - \mu/L)^k \cos^2(k\psi_1), \quad (37)$$

$$f(\mathbf{x}_k) \approx \frac{1}{2} \mu \hat{c}_1^2 \beta^k (1 - \mu/L)^k \cos^2(k\psi_1) + \frac{1}{2} L \hat{c}_d^2 \gamma^{2k}$$

$$\langle -\nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle \approx -\mu c_1 \hat{c}_1 \beta^k (1 - \mu/L)^k \cos(k\psi_1) \\ \times (\sqrt{\beta(1 - \mu/L)} \cos((k+1)\psi_1) - \cos(k\psi_1)),$$

$$\langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_{k-1}) \rangle \approx \mu^2 \hat{c}_1^2 \beta^k - \frac{1}{2} (1 - \mu/L)^k - \frac{1}{2} \cos(k\psi_1) \cos((k-1)\psi_1) \\ - L^2 \hat{c}_d^2 \gamma^{2k-1},$$

where $\psi_1 = \psi_1(\beta, \gamma)$. It is likely that these expressions will be dominated by the mode of either the smallest or largest eigenvalues. We next analyze each case separately.

4.4.1 Case 1: the Mode of the Smallest Eigenvalue Dominates—When the mode of the smallest eigenvalue dominates, we further approximate (37) as

$$\begin{aligned}
 f(\mathbf{y}_k) &\approx \frac{1}{2} \mu \hat{c}_1^2 \beta^k (1 - \mu/L)^k \cos^2(k\psi_1), & f(\mathbf{x}_k) &\approx \frac{1}{2} \mu \hat{c}_1^2 \beta^k (1 - \mu/L)^k \cos^2(k\psi_1), \\
 \langle -\nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle & & & \\
 &\approx -\mu \hat{c}_1 \hat{c}_1 \beta^k (1 - \mu/L)^k \cos(k\psi_1) (\cos((k+1)\psi_1) - \cos(k\psi_1)) \\
 &= 2\mu \hat{c}_1 \hat{c}_1 \beta^k (1 - \mu/L)^k \cos(k\psi_1) \sin((k+1/2)\psi_1) \sin(\psi_1/2) \\
 &\approx 2\mu \hat{c}_1 \hat{c}_1 \sin(\psi_1/2) \beta^k (1 - \mu/L)^k \sin(2k\psi_1)
 \end{aligned} \tag{38}$$

using simple trigonometric identities and the approximations $\sqrt{\beta(1 - \mu/L)} \approx 1$ and $\sin(k\psi_1) \approx \sin((k+1/2)\psi_1)$. The values (38) exhibit oscillations at a frequency proportional to $\psi_1(\beta, \gamma)$ in (29). This oscillation can be detected by the conditions (31) and (32) and is useful in detecting the high momentum regime where a restart can help improve the convergence rate.

4.4.2 Case 2: the Mode of the Largest Eigenvalue Dominates—Unlike the primary sequence $\{\mathbf{y}_k\}$ of OGM, convergence of the secondary sequence $\{\mathbf{x}_k\}$ of OGM may be dominated by the mode of the largest eigenvalue in (25). By further approximating (37) for the case when the mode of the largest eigenvalue dominates, the function value $f(\mathbf{x}_k) \approx \frac{1}{2} L \hat{c}_d^2 \gamma^{2k}$ decreases slowly but monotonically, whereas $f(\mathbf{y}_k) \approx f(\mathbf{x}_*) = 0$ and $\langle -\nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle \approx 0$. Therefore, neither restart condition (31) or (32) can detect such non-oscillatory observable values, even though the secondary mode $\{w_{k,d}\}$ of the largest eigenvalue is oscillating (corresponding to overshooting over the optimum). However, the inner product of two sequential gradients $\langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_{k-1}) \rangle \approx -L^2 \hat{c}_d^2 \gamma^{2k-1}$, can detect the overshoot of the secondary sequence $\{\mathbf{x}_k\}$, suggesting that the algorithm should adapt by decreasing γ when condition (33) holds. Decreasing γ too much may slow down the overall convergence rate when the mode of the smallest eigenvalue dominates. Thus, we use (33) only when using the secondary sequence of OGM as algorithm output (e.g., Sec. 5.2).

5 Proposed Adaptive Schemes for OGM

5.1 Adaptive Scheme of OGM for Smooth and Strongly Convex Problems

Alg. 2 illustrates a new adaptive version of OGM that is used in our numerical experiments in Sec. 6. When a restart condition is satisfied in Alg. 2, we reset $t_k = 1$ to discard the previous momentum that has a bad direction. When the decreasing γ condition is satisfied in Alg. 2, we decrease σ to suppress undesirable overshoot of the secondary sequence $\{\mathbf{x}_k\}$. Although the analysis in Sec. 3 considered only strongly convex quadratic functions, the

numerical experiments in Sec. 6 illustrate that the adaptive scheme is also useful more generally for smooth convex functions in $\mathcal{F}_{0,L}(\mathbb{R}^d)$, as described in [5, Sec. 4.6].

5.2 Adaptive Scheme of the Proximal Version of OGM for Nonsmooth Composite Convex Problems

Modern applications often involve nonsmooth composite convex problems:

$$\arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + \phi(\mathbf{x}) \right\}, \quad (39)$$

where $f \in \mathcal{F}_{0,L}(\mathbb{R}^d)$ is a smooth convex function (typically not strongly convex) and $\phi \in \mathcal{F}_{0,\infty}(\mathbb{R}^d)$ is a convex function that is possibly nonsmooth and “proximal-friendly” [25], such as the ℓ_1 regularizer $\phi(\mathbf{x}) = \|\mathbf{x}\|_1$. Our numerical experiments in Sec. 6 show that a new adaptive version of a proximal variant of OGM can be useful for solving such problems.

Algorithm 2 OGM with restarting momentum and decreasing γ

```

1: Input:  $f \in \mathcal{F}_{\mu,L}(\mathbb{R}^d)$  or  $\mathcal{F}_{0,L}(\mathbb{R}^d)$ ,  $\mathbf{x}_{-1} = \mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^d$ ,  $t_0 = \sigma = 1$ ,  $\bar{\sigma} \in [0, 1]$ .
2: for  $k \geq 0$  do
3:    $\mathbf{y}_{k+1} = \mathbf{x}_k - \frac{1}{2} \nabla f(\mathbf{x}_k)$ 
4:   if  $f(\mathbf{y}_{k+1}) > f(\mathbf{y}_k)$  (or  $\langle -\nabla f(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle < 0$ ) then  $\triangleright$  Restart condition
5:      $t_k = 1$ ,  $\sigma \leftarrow 1$ 
6:   else if  $\langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_{k-1}) \rangle < 0$  then  $\triangleright$  Decreasing  $\gamma$  condition
7:      $\sigma \leftarrow \bar{\sigma} \sigma$ 
8:      $t_{k+1} = \frac{1}{2} \left( 1 + \sqrt{1 + 4t_k^2} \right)$ 
9:    $\mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \frac{t_k - 1}{t_{k+1}} (\mathbf{y}_{k+1} - \mathbf{y}_k) + \sigma \frac{t_k}{t_{k+1}} (\mathbf{y}_{k+1} - \mathbf{x}_k)$ 

```

To solve such problems using first-order information, [4] developed a fast proximal gradient method, popularized under the name fast iterative shrinkage-thresholding algorithm (FISTA), that directly extends FGM with coefficients (3) for solving (39) while preserving the $O(1/k^2)$ rate of FGM. Variants of FISTA with adaptive restart were studied in [5, Sec. 5.2].

Inspired by the fact that OGM converges faster than FGM, [15] studied a proximal variant⁶ of OGM (POGM) with coefficients (7). It is natural to pursue acceleration of POGM by using variations of any (or all) of the three adaptive schemes (31), (32), (33), as illustrated in Alg. 3 where the proximity operator is defined as $\text{prox}_h(\mathbf{z}) := \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|^2 + h(\mathbf{x}) \right\}$.

As a function restart condition for POGM, we use $F(\mathbf{x}_{k+1}) > F(\mathbf{x}_k)$ instead of $F(\mathbf{y}_{k+1}) > F(\mathbf{y}_k)$, because $F(\mathbf{y}_k)$ can be un-bounded (e.g., \mathbf{y}_k can be unfeasible for constrained problems).

For gradient conditions of POGM, we consider the composite gradient mapping $G(\mathbf{x}_k)$ in Alg. 3 that differs from the standard composite gradient mapping in [16]. We then use the gradient conditions $\langle -G(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle < 0$ and $\langle G(\mathbf{x}_k), G(\mathbf{x}_{k-1}) \rangle < 0$ for restarting POGM or decreasing γ of POGM respectively. Here POGM must output the secondary sequence $\{\mathbf{x}_k\}$ because the function value $F(\mathbf{y}_k)$ of the primary sequence may be unbounded. This situation was the motivation for (33) (or $\langle G(\mathbf{x}_k), G(\mathbf{x}_{k-1}) \rangle < 0$) and Sec. 4.3. When $\phi(\mathbf{x}) = 0$, Alg. 3

⁶Applying the proximity operator to the primary sequence $\{\mathbf{y}_k\}$ of OGM, similar to the extension of FGM to FISTA, leads to a poor worst-case convergence bound [15]. Therefore, [15] applied the proximity operator to the secondary sequence of OGM and showed numerically that this version has a convergence bound about twice smaller than that of FISTA.

reduces to an algorithm that is similar to Alg. 2, where only the location of the restart and decreasing γ conditions differs.

The worst-case bound for POGM in [15] requires choosing the number of iterations N in advance for computing $\theta_N(7)$, which seems incompatible with adaptive restarting of POGM. Like (7) in Alg. 1, the fact that $\theta_N(7)$ is larger than the standard t_N in (9) at the last iteration helps to dampen (by reducing the values of β and γ) the final update to guarantee fast convergence in the worst-case. (This property was studied for smooth convex minimization in [14].) We could perform one last update using θ_N after a restart condition is satisfied, but this step appears unnecessary because restarting already has the effect of slowing down the algorithm. Thus, we did not include any such extra update in Alg. 3 in our experiment in the next section.

Algorithm 3 POGM with restarting momentum and decreasing γ

```

1: Input:  $f \in \mathcal{F}_{0,L}(\mathbb{R}^d)$ ,  $\phi \in \mathcal{F}_{0,\infty}(\mathbb{R}^d)$ ,  $\mathbf{x}_{-1} = \mathbf{x}_0 = \mathbf{y}_0 = \mathbf{u}_0 = \mathbf{z}_0 \in \mathbb{R}^d$ ,
2:  $t_0 = \zeta_0 = \sigma = 1$ ,  $\bar{\sigma} \in [0, 1]$ .
3: for  $k \geq 0$  do
4:    $\mathbf{u}_{k+1} = \mathbf{x}_k - \frac{1}{L} \nabla f(\mathbf{x}_k)$ 
5:    $t_{k+1} = \frac{1}{2} \left( 1 + \sqrt{1 + 4t_k^2} \right)$ 
6:    $\mathbf{z}_{k+1} = \mathbf{u}_{k+1} + \frac{t_k - 1}{t_{k+1}} (\mathbf{u}_{k+1} - \mathbf{u}_k) + \sigma \frac{t_k}{t_{k+1}} (\mathbf{u}_{k+1} - \mathbf{x}_k) - \frac{t_k - 1}{t_{k+1}} \frac{1}{L \zeta_k} (\mathbf{x}_k - \mathbf{z}_k)$ 
7:    $\zeta_{k+1} = \frac{1}{L} \left( 1 + \frac{t_k - 1}{t_{k+1}} + \sigma \frac{t_k}{t_{k+1}} \right)$ 
8:    $\mathbf{x}_{k+1} = \text{prox}_{\zeta_{k+1} \phi}(\mathbf{z}_{k+1})$ 
9:    $G(\mathbf{x}_k) = \nabla f(\mathbf{x}_k) - \frac{1}{\zeta_{k+1}} (\mathbf{x}_{k+1} - \mathbf{z}_{k+1})$ 
10:   $\mathbf{y}_{k+1} = \mathbf{x}_k - \frac{1}{L} G(\mathbf{x}_k)$ 
11:  if  $F(\mathbf{x}_{k+1}) > F(\mathbf{x}_k)$  (or  $\langle -G(\mathbf{x}_k), \mathbf{y}_{k+1} - \mathbf{y}_k \rangle < 0$ ) then  $\triangleright$  Restart condition
12:     $t_{k+1} = 1$ ,  $\sigma \leftarrow 1$ 
13:  else if  $\langle G(\mathbf{x}_k), G(\mathbf{x}_{k-1}) \rangle < 0$  then  $\triangleright$  Decreasing  $\gamma$  condition
14:     $\sigma \leftarrow \bar{\sigma} \sigma$ 

```

6 Numerical Results

This section shows the results of applying OGM (and POGM) with adaptive schemes to various numerical examples including both strongly convex quadratic problems and non-strongly convex problems.⁷ The results illustrate that OGM (or POGM) with adaptive schemes converges faster than FGM (or FISTA) with adaptive restart. The plots show the decrease of $F(\mathbf{y}_k)$ of the primary sequence for FGM (FISTA) and OGM unless specified. For POGM, we use the secondary sequence $\{\mathbf{x}_k\}$ as an output and plot $F(\mathbf{x}_k)$, since $F(\mathbf{y}_k)$ can be unbounded.

6.1 Strongly Convex Quadratic Examples

This section considers two types of strongly convex quadratic examples, where the mode of either the smallest eigenvalue or the largest eigenvalue dominates, providing examples of the analysis in Sec. 4.4.1 and 4.4.2 respectively.

6.1.1 Case 1: the Mode of the Smallest Eigenvalue Dominates—Fig. 2 compares GM, FGM and OGM, with or without the knowledge of q , for minimizing a strongly convex quadratic function (11) in $d = 500$ dimensions with $q = 10^{-4}$, where we generated \mathbf{A} (for $\mathbf{Q} = \mathbf{A}^\top \mathbf{A}$) and \mathbf{p} randomly. In Fig. 2, ‘GM’ and ‘GM- q ’ denote GM with $\alpha = \frac{1}{L}$ and $\alpha = \frac{2}{\mu + L}$

⁷Software for the algorithms and for producing the figures in Sec. 6 is available at <https://gitlab.eecs.umich.edu/michigan-fast-optimization/ogm-adaptive-restart>.

respectively. ‘FGM’ and ‘FGM- q ’ denote step coefficients (3) and (5) respectively. ‘OGM’ and ‘OGM- q ’ denote step coefficients (9) and (22) respectively. As expected, knowing q accelerates convergence.

Fig. 2 also illustrates that adaptive restart helps FGM and OGM to nearly achieve the fast linear converge rate of their non-adaptive versions that know q . As expected, OGM converges faster than FGM for all cases. In Fig. 2, ‘FR’ and ‘GR’ stand for function restart (31) and gradient restart (32), respectively, and both behave nearly the same.

6.1.2 Case 2: the Mode of the Largest Eigenvalue Dominates—Consider the strongly convex quadratic function with $\mathbf{Q} = \begin{bmatrix} q & 0 \\ 0 & 1 \end{bmatrix}$, $q = 0.01$, $\mathbf{p} = \mathbf{0}$ and $\mathbf{x}^* = \mathbf{0}$. When starting the algorithm from the initial point $\{\mathbf{x}_0\} = (0.2, 1)$, the secondary sequence $\{\mathbf{x}_k\}$ of OGM-GR⁸ (or equivalently OGM-GR-GD $\gamma(\bar{\sigma} = 1.0)$) is dominated by the mode of largest eigenvalue in Fig. 3, illustrating the analysis of Sec. 4.4.2. Fig. 3 illustrates that the primary sequence of OGM-GR converges faster than that of FGM-GR, whereas the secondary sequence of OGM-GR initially converges even slower than GM. To deal with such slow convergence coming from the overshooting behavior of the mode of the largest eigenvalue of the secondary sequence of OGM, we employ the decreasing γ scheme in (33). Fig. 3 shows that using $\bar{\sigma} < 1$ in Alg. 2 leads to overall faster convergence of the secondary sequence $\{\mathbf{x}_k\}$ than the standard OGM-GR where $\bar{\sigma} = 1$. We leave optimizing the choice of $\bar{\sigma}$ or studying other strategies for decreasing γ as future work.

6.2 Non-strongly Convex Examples

This section applies adaptive OGM (or POGM) to three non-strongly convex numerical examples in [5, 7]. The numerical results show that adaptive OGM (or POGM) converges faster than FGM (or FISTA) with adaptive restart.

6.2.1 Log-Sum-Exp—The following function from [5] is smooth but non-strongly convex:

$$f(\mathbf{x}) = \eta \log \left(\sum_{i=1}^m \exp \left(\frac{1}{\eta} (\mathbf{a}_i^\top \mathbf{x} - b_i) \right) \right).$$

It approaches $\max_{i=1, \dots, m} (\mathbf{a}_i^\top \mathbf{x} - b_i)$ as $\eta \rightarrow 0$. Here, η controls the function smoothness $L = \frac{1}{\eta} \lambda_{\max}(\mathbf{A}^\top \mathbf{A})$ where $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_m]^\top \in \mathbb{R}^{m \times d}$. The region around the optimum is approximately quadratic since the function is smooth, and thus the adaptive restart can be useful without knowing the local condition number.

For $(m, d) = (100, 20)$, we randomly generated $\mathbf{a}_i \in \mathbb{R}^d$ and $b_i \in \mathbb{R}$ for $i = 1, \dots, m$, and investigated $\eta = 1, 10$. Fig. 4 shows that OGM with adaptive restart converges faster than FGM with the adaptive restart. The benefit of adaptive restart is dramatic here; apparently

⁸Fig. 3 only compares the results of the gradient restart (GR) scheme for simplicity, where the function restart (FR) behaves similarly.

FGM and OGM enter a locally well-conditioned region after about 100 – 200 iterations, where adaptive restart then provide a fast linear rate.

6.2.2 Sparse Linear Regression—Consider the following cost function used for sparse linear regression:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2, \quad \phi(\mathbf{x}) = \tau \|\mathbf{x}\|_1,$$

for $\mathbf{A} \in \mathbb{R}^{m \times d}$, where $L = \lambda_{\max}(\mathbf{A}^\top \mathbf{A})$ and the parameter τ balances between the measurement error and signal sparsity. The proximity operator becomes a soft-thresholding operator, e.g., $\text{prox}_{\zeta_{k+1}} \phi(\mathbf{x}) = \text{sgn}(\mathbf{x}) \max\{|\mathbf{x}| - \zeta_{k+1}, 0\}$. The minimization seeks a sparse solution \mathbf{x}^* , and often the cost function is strongly convex with respect to the non-zero elements of \mathbf{x}^* . Thus we expect to benefit from adaptive restarting.

For each choice of (m, d, s, τ) in Fig. 5, we generated an s -sparse true vector \mathbf{x}_{true} by taking the s largest entries of a randomly generated vector. We then simulated $\mathbf{b} = \mathbf{Ax}_{\text{true}} + \mathbf{e}$, where the entries of matrix \mathbf{A} and vector \mathbf{e} were sampled from a zero-mean normal distribution with variances 1 and 0.1 respectively. Fig. 5 illustrates that POGM with adaptive schemes provide acceleration over FISTA with adaptive restart. While Sec. 3.4 discussed the undesirable overshooting behavior that a secondary sequence of OGM (or POGM) may encounter, these examples rarely encountered such behavior. Therefore the choice of $\bar{\sigma}$ in the adaptive POGM was not significant in this experiment, unlike Sec. 6.1.2.

6.2.3 Constrained Quadratic Programming—Consider the following box-constrained quadratic program:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{p}^\top \mathbf{x}, \quad \phi(\mathbf{x}) = \begin{cases} 0, & \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}, \\ \infty, & \text{otherwise,} \end{cases}$$

where $L = \lambda_{\max}(\mathbf{Q})$. The algorithms ISTA (a proximal variant of GM), FISTA and POGM use the projection operator: $\text{prox}_{\frac{1}{L} \phi}(\mathbf{x}) = \text{prox}_{\zeta_{k+1}} \phi(\mathbf{x}) = \min\{\max\{\mathbf{x}, \mathbf{l}\}, \mathbf{u}\}$. Fig. 6 denotes each algorithm by a projected GM, a projected FGM, and a projected OGM respectively. Similar to Sec. 6.2.2, after the algorithm identifies the active constraints the problem typically becomes a strongly convex quadratic problem where we expect to benefit from adaptive restart.

Fig. 6 studies two examples with problem dimensions $d = 500, 1000$, where we randomly generate a positive definite matrix \mathbf{Q} having a condition number 10^7 (i.e., $q = 10^{-7}$), and a vector \mathbf{p} . Vectors \mathbf{l} and \mathbf{u} correspond to the interval constraints $-1 \leq x_j \leq 1$ for $\mathbf{x} = \{x_j\}$. The optimum \mathbf{x}^* had 47 and 81 active constraints out of 500 and 1000 respectively. In Fig. 6, the projected OGM with adaptive schemes converged faster than FGM with adaptive restart and other non-adaptive algorithms.

7 Conclusions

We introduced adaptive restarting schemes for the optimized gradient method (OGM) that heuristically exhibits a fast linear convergence rate when the function is strongly convex or even when the function is not globally strongly convex. The method resets the momentum when it makes a bad direction. We provided a heuristic dynamical system analysis to justify the practical acceleration of the adaptive scheme of OGM, by extending the existing analysis of FGM. On the way, we described new optimized constant step coefficients for OGM for strongly convex quadratic problems. Numerical results illustrate that the proposed adaptive approach practically accelerates the convergence rate of OGM, and in particular, performs faster than FGM with adaptive restart. An interesting open problem is to determine the worst-case bounds for OGM (and FGM) with adaptive restart.

Acknowledgements

This research was supported in part by NIH grant U01 EB018753.

References

1. Cevher V, Becker S, Schmidt M: Convex optimization for big data: scalable, randomized, and parallel algorithms for big data analytics. *IEEE Sig. Proc. Mag* 31(5), 32–43 (2014)
2. Nesterov Y: A method for unconstrained convex minimization problem with the rate of convergence $O(1/k^2)$. *Dokl. Akad. Nauk. USSR* 269(3), 543–7 (1983)
3. Nesterov Y: *Introductory lectures on convex optimization: A basic course*. Kluwer (2004)
4. Beck A, Teboulle M: A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci* 2(1), 183–202 (2009)
5. O’Donoghue B, Candès E: Adaptive restart for accelerated gradient schemes. *Found. Comp. Math* 15(3), 715–32 (2015)
6. Giselsson P, Boyd S: Monotonicity and restart in fast gradient methods. In: *Proc. Conf. Decision and Control*, pp. 5058–63 (2014)
7. Su W, Boyd S, Candès EJ: A differential equation for modeling Nesterov’s accelerated gradient method: theory and insights. *J. Mach. Learning Res* 17(153), 1–43 (2016)
8. Muckley MJ, Noll DC, Fessler JA: Fast parallel MR image reconstruction via B1-based, adaptive restart, iterative soft thresholding algorithms (BARISTA). *IEEE Trans. Med. Imag* 34(2), 578–88 (2015)
9. Monteiro RDC, Ortiz C, Svaiter BF: An adaptive accelerated first-order method for convex optimization. *Comput. Optim. Appl* 64(1), 31–73 (2016)
10. Kim D, Fessler JA: Optimized first-order methods for smooth convex minimization. *Mathematical Programming* 159(1), 81–107 (2016) [PubMed: 27765996]
11. Drori Y, Teboulle M: Performance of first-order methods for smooth convex minimization: A novel approach. *Mathematical Programming* 145(1–2), 451–82 (2014)
12. Drori Y: The exact information-based complexity of smooth convex minimization. *J. Complexity* 39, 1–16 (2017)
13. Kim D, Fessler JA: Generalizing the optimized gradient method for smooth convex minimization (2016). Arxiv 1607.06764
14. Kim D, Fessler JA: On the convergence analysis of the optimized gradient method. *J. Optim. Theory Appl* 172(1), 187–205 (2017) [PubMed: 28461707]
15. Taylor AB, Hendrickx JM, Glineur François.: Exact worst-case performance of first-order algorithms for composite convex optimization (2015). Arxiv 1512.07516
16. Nesterov Y: Gradient methods for minimizing composite functions. *Mathematical Programming* 140(1), 125–61 (2013)

17. Güler O: New proximal point algorithms for convex minimization. *SIAM J. Optim* 2(4), 649–64 (1992)
18. Van Scoy B, Freeman RA, Lynch KM: The fastest known globally convergent first-order method for the minimization of strongly convex functions (2017). URL http://www.optimization-online.org/DB_HTML/2017/03/5908.html
19. Polyak BT: Introduction to optimization. Optimization Software Inc, New York (1987)
20. Lessard L, Recht B, Packard A: Analysis and design of optimization algorithms via integral quadratic constraints. *SIAM J. Optim* 26(1), 57–95 (2016)
21. Chiang A: Fundamental methods of mathematical economics. McGraw-Hill, New York (1984)
22. Powell MJD: Restart procedures for the conjugate gradient method. *Mathematical Programming* 12(1), 241–54 (1977)
23. Nocedal J, Wright SJ: Numerical optimization. Springer, New York (2006). DOI 10.1007/978-0-387-40065-5. 2nd edition.
24. Nemirovski A: Efficient methods in convex programming (1994). URL http://www2.isye.gatech.edu/~nemirovs/Lect_EMCO.pdf. Lecture notes
25. Combettes PL, Pesquet JC: Proximal splitting methods in signal processing (2011). DOI 10.1007/978-1-4419-9569-8_10. Fixed-Point Algorithms for Inverse Problems in Science and Engineering. Springer, Optimization and Its Applications, pp 185–212

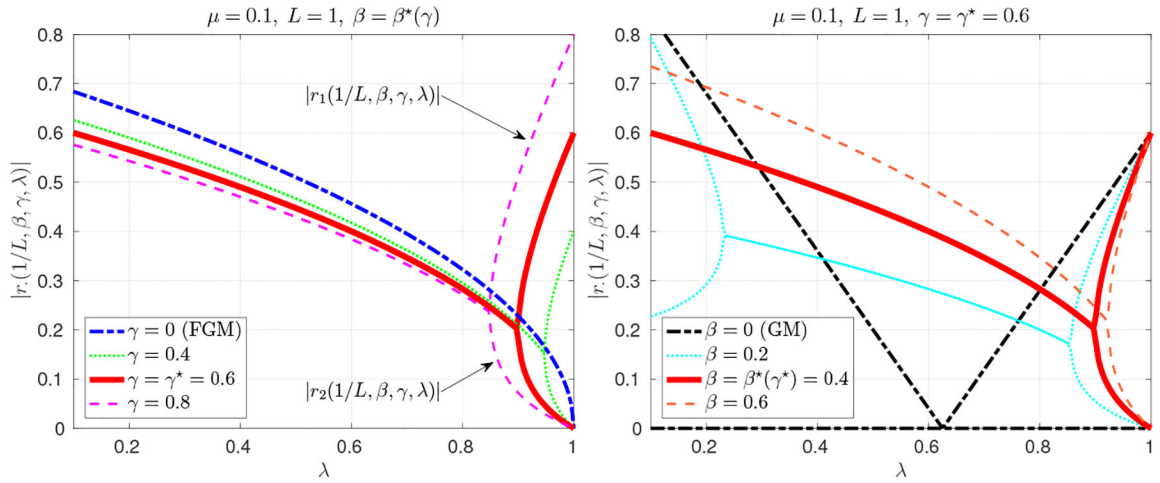


Fig. 1. Plots of $|r_1(1/L, \beta, \gamma, \lambda)|$ and $|r_2(1/L, \beta, \gamma, \lambda)|$ over $\mu \ \lambda \ L$ for various (Left) γ values for given $\beta = \beta^*(\gamma)$, and (Right) β values for given $\gamma = \gamma^*$, for a strongly convex quadratic problem with $q = 0.1$, where $(\beta^*, \gamma^*) = (0.4, 0.6)$. Note that the maximum of $|r_1(1/L, \beta, \gamma, \lambda)|$ and $|r_2(1/L, \beta, \gamma, \lambda)|$, i.e. the upper curve in the plot, corresponds to the value of $\rho(T_\lambda(1/L, \beta, \gamma))$ in (16).

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

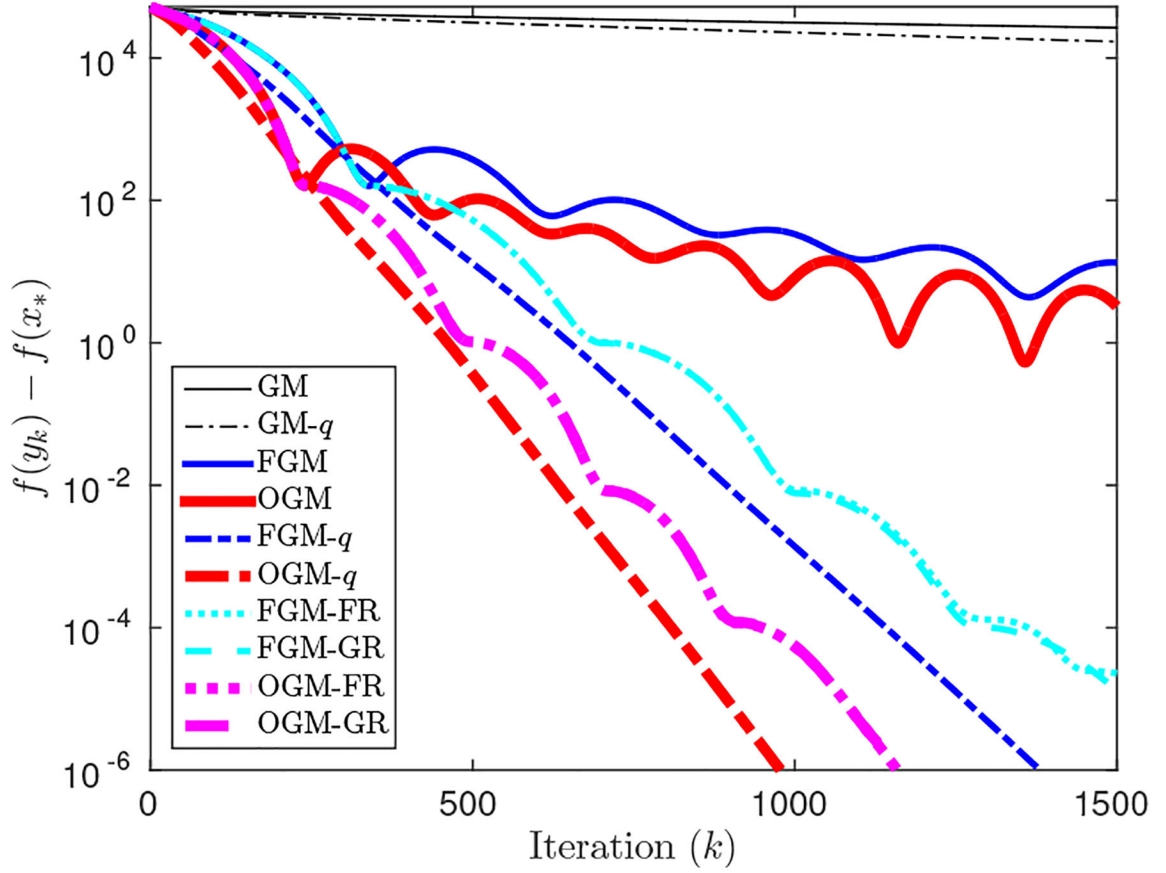


Fig. 2. Minimizing a strongly convex quadratic function - Case 1: the mode of the smallest eigenvalue dominates. (FGM-FR and FGM-GR are almost indistinguishable, as are OGM-FR and OGM-GR.)

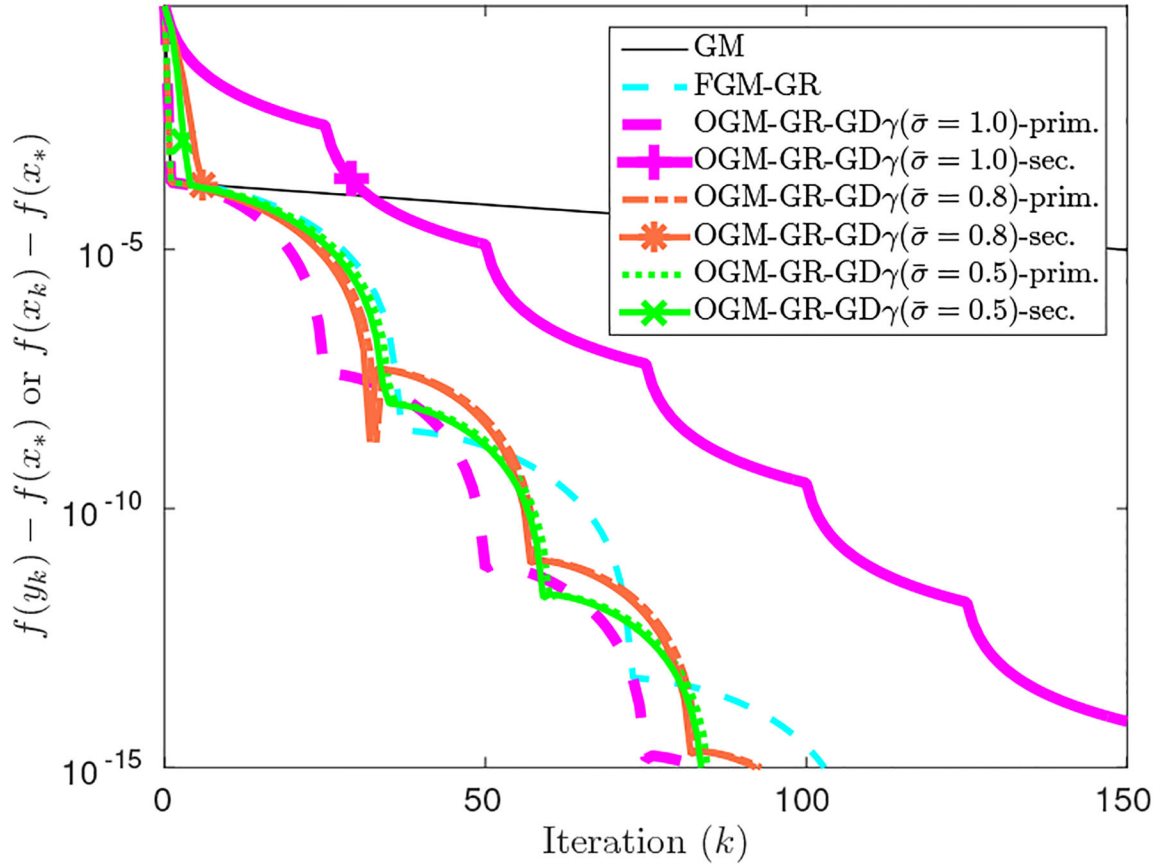


Fig. 3. Minimizing a strongly convex quadratic function - Case 2: the mode of the largest eigenvalue dominates for the secondary sequence $\{x_k\}$ of OGM. Using $\text{GD}\gamma$ (33) with $\bar{\sigma} < 1$ accelerates convergence of the secondary sequence of OGM-GR, where both the primary and secondary sequences behave similarly after first few iterations, unlike $\bar{\sigma} = 1$.

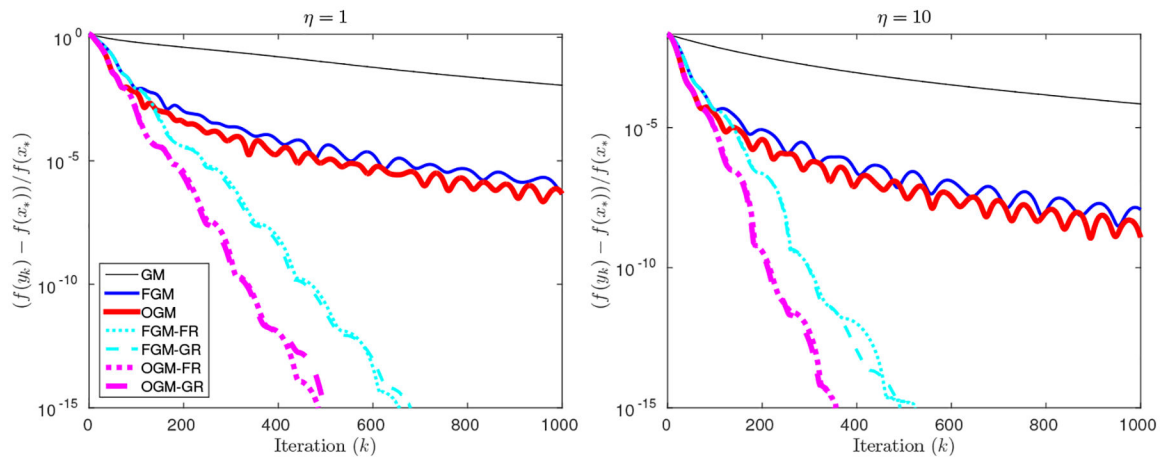


Fig. 4. Minimizing a smooth but non-strongly convex Log-Sum-Exp function.

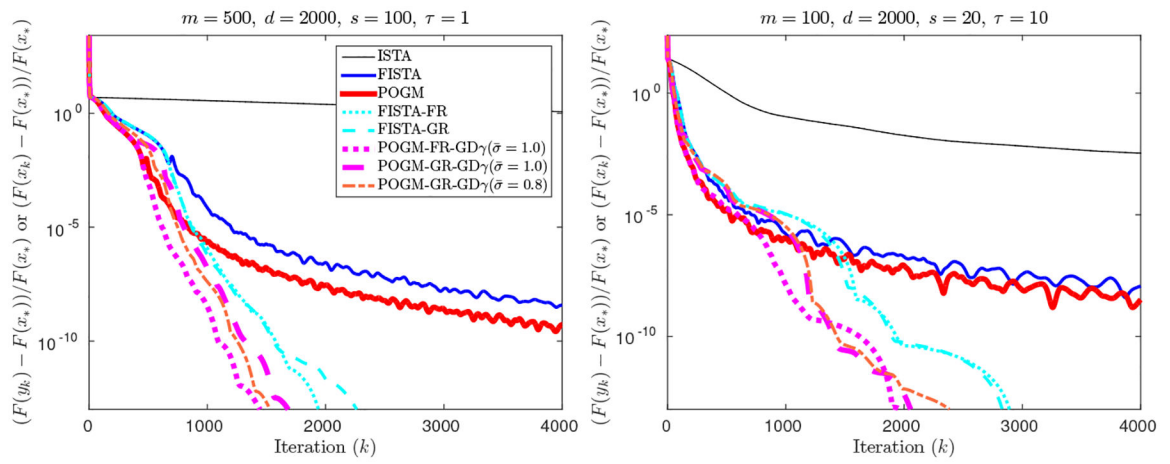


Fig. 5. Solving a sparse linear regression problem. (ISTA is a proximal variant of GM.)

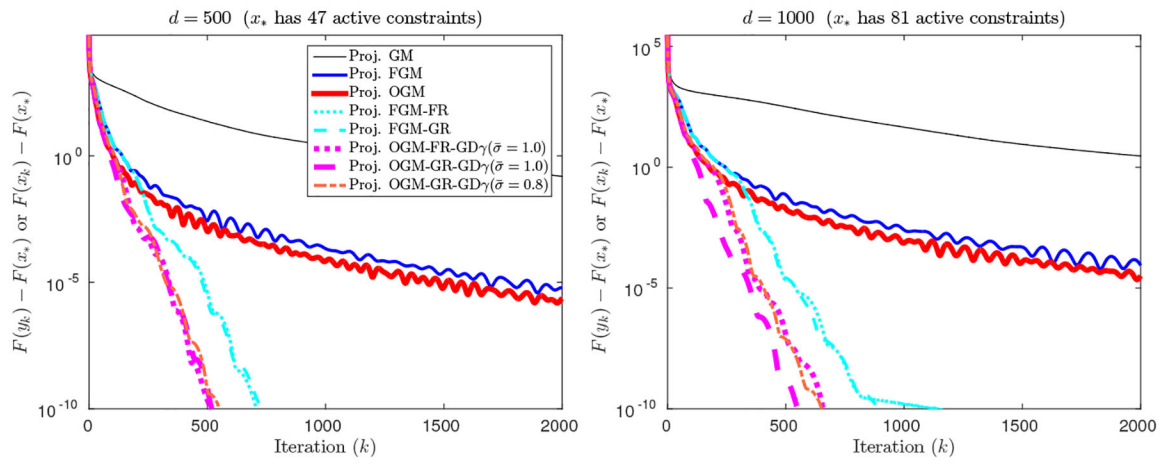


Fig. 6. Solving a box-constrained quadratic programming problem.

Table 1

Optimally tuned coefficients (α, β, γ) of GM, FGM, and OGM, and their spectral radius $\rho(T(\alpha, \beta, \gamma))$ (15). These optimal coefficients result from solving (18) with the shaded coefficients fixed.

Algorithm	α	β	γ	$\rho(T(\alpha, \beta, \gamma))$
GM	$\frac{2}{\mu + L}$	0	0	$\frac{1 - q}{1 + q}$
FGM	$\frac{1}{L}$	$\frac{1 - \sqrt{q}}{1 + \sqrt{q}}$	0	$1 - \sqrt{q}$
	$\frac{4}{\mu + 3L}$	$\frac{\sqrt{3 + q} - 2\sqrt{q}}{\sqrt{3 + q} + 2\sqrt{q}}$	0	$1 - \frac{2\sqrt{q}}{\sqrt{3 + q}}$
OGM	$\frac{1}{L}$	$\frac{(2 + q - \sqrt{q^2 + 8q})^2}{4(1 - q)}$	$\frac{2 + q - \sqrt{q^2 + 8q}}{2}$	$\frac{2 + q - \sqrt{q^2 + 8q}}{2}$

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript