

2016/52

Exact Worst-Case Performance of First-Order Methods for Composite Convex Optimization

ADRIEN B. TAYLOR, JULIEN M. HENDRICKX AND
FRANÇOIS GLINEUR



50 YEARS OF
CORE
DISCUSSION PAPERS

CORE

Voie du Roman Pays 34, L1.03.01

Tel (32 10) 47 43 04

Fax (32 10) 47 43 01

Email: immaq-library@uclouvain.be

<http://www.uclouvain.be/en-44508.html>

EXACT WORST-CASE PERFORMANCE OF FIRST-ORDER METHODS FOR COMPOSITE CONVEX OPTIMIZATION*

ADRIEN B. TAYLOR[†], JULIEN M. HENDRICKX[†], AND FRANÇOIS GLINEUR[†]

Abstract. We provide a framework for computing the exact worst-case performance of any algorithm belonging to a broad class of oracle-based first-order methods for composite convex optimization, including those performing explicit, projected, proximal, conditional and inexact (sub)gradient steps. We simultaneously obtain tight worst-case guarantees and explicit instances of optimization problems on which the algorithm reaches this worst-case. We achieve this by reducing the computation of the worst-case to solving a convex semidefinite program, generalizing previous works on performance estimation by Drori and Teboulle [13] and the authors [43].

We use these developments to obtain a tighter analysis of the proximal point algorithm and of several variants of fast proximal gradient, conditional gradient, subgradient and alternating projection methods. In particular, we present a new analytical worst-case guarantee for the proximal point algorithm that is twice better than previously known, and improve the standard worst-case guarantee for the conditional gradient method by more than a factor of two.

We also show how the optimized gradient method proposed by Kim and Fessler in [22] can be extended by incorporating a projection or a proximal operator, which leads to an algorithm that converges in the worst-case twice as fast as the standard accelerated proximal gradient method [2].

Key words. convex optimization, composite convex optimization, first-order methods, worst-case analysis, performance estimation, semidefinite programming, convex interpolation

AMS subject classifications. 90C25, 90C30, 90C60, 68Q25, 90C22

1. Introduction. Consider the composite convex minimization problem

$$(CM) \quad \min_{x \in \mathbb{E}} \left\{ F(x) \equiv \sum_{k=1}^n F^{(k)}(x) \right\},$$

where \mathbb{E} is a finite-dimensional real vector space and each functional component $F^{(k)} : \mathbb{E} \rightarrow \mathbb{R} \cup \{\infty\}$ is a convex function belonging to some class $\mathcal{F}_k(\mathbb{E})$ — e.g., smooth or non-smooth, strongly convex or not, indicator functions, etc. — for which some operations are assumed to be available in closed-form (e.g. computing a gradient, projecting on the domain, computing a proximal step, etc.).

We are interested in the composite optimization problem (CM) because it naturally allows representing and exploiting a lot of the structure in many problems, which can play a major role in our ability to efficiently solve them (see [32] among others). In addition, the class of composite convex optimization problems arises very commonly in practice, as it contains for example constrained, ℓ_1 - and ℓ_2 -regularized convex optimization problems.

We focus on black-box oracle-based algorithms that use first-order information to approximately solve (CM), and in particular on obtaining exact and global worst-case guarantees on their performances. That is, for a given algorithm, we simultaneously seek to obtain worst-case guarantees — for example on objective function accuracy — and an instance of (CM) on which the algorithm behaves as such. In this work, we

*This paper presents research results of the Belgian Network DYSCO (Dynamical Systems, Control, and Optimization), funded by the Interuniversity Attraction Poles Programme, initiated by the Belgian State, Science Policy Office, and of the Concerted Research Action (ARC) programme supported by the Federation Wallonia-Brussels (contract ARC 14/19-060). The scientific responsibility rests with its authors. Adrien Taylor is a FRIA fellow.

[†]Université catholique de Louvain, ICTEAM Institute/CORE, B-1348 Louvain-la-Neuve, Belgium
E-mail: adrien.taylor@uclouvain.be, julien.hendrickx@uclouvain.be, francois.glineur@uclouvain.be

treat the case of fixed-step linear first-order methods, which includes among others fixed-step projected, proximal, conditional and inexact (sub)gradient methods.

This work builds on the recent idea of performance estimation, first developed by Drori and Teboulle in [13] and followed-up by Kim and Fessler [22] and the authors [43]. The approach was initially tailored for obtaining upper bounds on the worst-case behavior of fixed-step gradient methods for unconstrained minimization of a single smooth convex objective function. Motivated by subsequent results (see among others [21, 22]) we extend the framework of performance estimation to the composite case involving a much broader class of algorithms and function classes (see Section 1.4 for more details about previous works).

Our performance estimation framework relies on formulating the worst-case computation problem as a tractable semidefinite program (SDP), which can be tackled with standard solvers [24, 26, 41]. It enjoys the following attractive features:

- any primal feasible solution to this SDP leads to a lower bound on the worst-case performance of the method under consideration, by exhibiting a particular instance of (CM),
- any dual feasible solutions to this SDP corresponds to an upper bound on the worst-case performance of the method under consideration, that can be converted into an explicit proof based on a combination of valid inequalities.

1.1. Notations. In this paper, we work in a finite-dimensional real vector space \mathbb{E} and the corresponding dual space \mathbb{E}^* consisting of all linear functions on \mathbb{E} , and denote their dimension by $d = \dim \mathbb{E} = \dim \mathbb{E}^*$. We consider a dual pairing¹ between those spaces, denoted by $\langle \cdot, \cdot \rangle : \mathbb{E}^* \times \mathbb{E} \rightarrow \mathbb{R}$. We also consider a self-adjoint positive definite² linear operator $B : \mathbb{E} \rightarrow \mathbb{E}^*$ for $\langle \cdot, \cdot \rangle$, which allows defining the following primal and dual norms:

$$\|x\|_{\mathbb{E}}^2 = \langle Bx, x \rangle, \quad \forall x \in \mathbb{E}, \quad \|s\|_{\mathbb{E}^*}^2 = \langle s, B^{-1}s \rangle, \quad \forall s \in \mathbb{E}^*.$$

We denote $\langle x, y \rangle_{\mathbb{E}} = \langle Bx, y \rangle$ for $x, y \in \mathbb{E}$ and $\langle x, y \rangle_{\mathbb{E}^*} = \langle x, B^{-1}y \rangle$ for $x, y \in \mathbb{E}^*$. The usual case is simply $\mathbb{E} = \mathbb{E}^* = \mathbb{R}^d$ with $\langle x, y \rangle = x^\top y$ the standard Euclidean inner product and B the identity operator, for which we also have $\|x\|_{\mathbb{E}}^2 = \|x\|_{\mathbb{E}^*}^2 = \langle x, x \rangle$.

For a convex function $f : \mathbb{E} \rightarrow \mathbb{R} \cup \{\infty\}$, we denote by $f^* : \mathbb{E}^* \rightarrow \mathbb{R} \cup \{\infty\}$ its Legendre-Fenchel conjugate

$$f^*(y) = \sup_{x \in \mathbb{E}} \langle y, x \rangle - f(x),$$

by $\partial f(x)$ the subdifferential of f at x (set of all subgradients of f at x), and by $\tilde{\nabla} f(x)$ a particular subgradient of f at x . Similarly, the gradient of a differentiable function f at x is denoted by $\nabla f(x)$.

For notational convenience we denote by $K = \{1, \dots, n\}$ the set of indices corresponding to the different components $F^{(k)}$ in the objective function of (CM). We also denote by $\mathcal{F}_K(\mathbb{E})$ the set of functions of the form (CM) with components $F^{(k)} \in \mathcal{F}_k(\mathbb{E}) \forall k \in K$ — that is, $F \in \mathcal{F}_K(\mathbb{E})$.

Finally, we use the standard notation e_i for the unit vector having a single 1 as its i^{th} component.

¹The dual pairing is a real bilinear map $\langle \cdot, \cdot \rangle : \mathbb{E}^* \times \mathbb{E} \rightarrow \mathbb{R}$ satisfying (i) $\forall x \in \mathbb{E} \setminus \{0\}, \exists s \in \mathbb{E}^*$ such that $\langle s, x \rangle \neq 0$, and (ii) $\forall s \in \mathbb{E}^* \setminus \{0\}, \exists x \in \mathbb{E}$ such that $\langle s, x \rangle \neq 0$.

²That is, a linear operator B satisfying (i) $\langle Bx, y \rangle = \langle By, x \rangle \forall x, y \in \mathbb{E}$ (self-adjoint), and (ii) $\langle Bx, x \rangle > 0 \forall x \in \mathbb{E} \setminus \{0\}$ (positive definite).

1.2. Performance estimation problems. In [43], we introduced a formal definition for the performance estimation problem in the case of a black-box first-order methods for unconstrained minimization of a single convex function F . We now generalize the performance estimation framework for handling multiple components in the objective function.

First, we formalize black-box methods using the concept of *black-box oracles*. That means that, methods are only allowed to access the different components of the objective function by calling some routines, or oracles, returning some information about them at a given point. In particular, we focus on the standard first-order oracle for $F^{(k)}$: $\mathcal{O}_{F^{(k)}}(x) = (F^{(k)}(x), \tilde{\nabla}F^{(k)}(x))$ in the sequel, where $\tilde{\nabla}F^{(k)}(x) \in \partial F^{(k)}(x)$ is a subgradient of $F^{(k)}$ at x . The general formalism of the approach is nevertheless also valid for other standard oracles, as for examples zeroth-order or second-order ones — that is, $\mathcal{O}_{F^{(k)}}(x) = (F^{(k)}(x))$ or $\mathcal{O}_{F^{(k)}}(x) = (F^{(k)}(x), \nabla F^{(k)}(x), \nabla^2 F^{(k)}(x))$. However, as we will see, our ability to solve the corresponding performance estimation problems in an exact way is currently limited to first-order oracles.

Second, we consider a sequence of $N + 1$ iterates $\{x_i\}_{0 \leq i \leq N} \subset \mathbb{E}$, corresponding to a method that performs N steps from an initial iterate x_0 . For each of those iterates we consider the set of calls to the oracle for each functional component³ $\mathcal{O}_{F^{(k)}}: \{\mathcal{O}_{F^{(k)}}(x_i)\}_{0 \leq i \leq N}$.

Third, we consider a method \mathcal{M} whose iterates can be computed by combining past and current oracle information about F . This means that after the method has performed $i - 1$ steps, the next iterate x_i should be computable as a solution to an equation of the form:

$$\text{(EQ}_i\text{)} \quad \text{EQUATION}(x_0, \{\mathcal{O}_{F^{(k)}}(x_0)\}_{k \in K}, x_1, \{\mathcal{O}_{F^{(k)}}(x_1)\}_{k \in K}, \dots, x_i, \{\mathcal{O}_{F^{(k)}}(x_i)\}_{k \in K}).$$

Note that the only unknown in this equation is x_i , and that it thus provides an implicit definition for the next step. We will see later that this assumption on \mathcal{M} includes a large number of existing methods for composite optimization.

Finally, we consider a real-valued performance criterion \mathcal{P} for evaluating the efficiency of the method. In the sequel, we assume without loss of generality that the lower the value of \mathcal{P} , the better the corresponding method. Examples of such performance criteria include objective function accuracy $F(x_N) - F(x_*)$ (where x_* is any optimal solution of (CM)), and distance to an optimal solution $\|x_N - x_*\|_{\mathbb{E}}^2$.

In our framework, this performance criterion is generally allowed to depend on information returned by the oracles $\mathcal{O}_{F^{(k)}}$ at all the iterates $\{x_i\}_{0 \leq i \leq N}$, but also at an extra point $x_* \in \mathbb{E}$ assumed to be an optimal solution to problem (CM). Also, we allow \mathcal{P} to depend on the iterates themselves. For notational convenience we introduce an index set for all iterates (including optimal solution) $I = \{0, 1, \dots, N, *\}$.

The worst-case performance of method \mathcal{M} on (CM) is then the optimal value of the following optimization problem, with both functions $\{F^{(k)}\}_{k \in K}$ and iterates

³That is, we chose to associate a call to each oracle to every iterate. This is mostly for notational convenience and does not induce any loss of generality. Indeed, a method can always choose not to use the information returned by one of the oracles at some iterations.

$\{x_i\}_{i \in I}$ as variables, which we call a performance estimation problem (PEP).

$$\begin{aligned}
 \text{(PEP)} \quad & \sup_{\{F^{(k)}\}_{k \in K}, \{x_i\}_{i \in I}} \mathcal{P}(\{\mathcal{O}_{F^{(k)}}(x_i)\}_{i \in I, k \in K}, \{x_i\}_{i \in I}) \\
 & \text{subject to } F^{(k)} \in \mathcal{F}_k(\mathbb{E}) \text{ for all } k \in K, \\
 & x_0 \text{ satisfies some initialization condition,} \\
 & x_i \text{ is computed by } \mathcal{M} \text{ according to (EQ}_i\text{) for all } 1 \leq i \leq N, \\
 & x_* \text{ a minimizer of } F(x).
 \end{aligned}$$

That is, a solution to (PEP) corresponds to an instance of problem (CM) on which method \mathcal{M} behaves as badly as possible with respect to the performance criterion \mathcal{P} . The initialization condition on x_0 is required as most methods exhibit unbounded worst-case performance without it. In the sequel we will mostly restrict ourselves to the classical approach which consists in bounding the initial distance to an optimal solution with a constant R , i.e., assume $\|x_0 - x_*\|_{\mathbb{E}} \leq R$.

Note that (PEP) is inherently an infinite-dimensional optimization problem, as functions $F^{(k)}$ appear as variables. However, a crucial observation is that, due to the black-box assumption on the objective components, this problem can be cast completely equivalently in a finite-dimensional fashion. Indeed, introducing the *outputs* of the oracle calls as variables, namely $O_i^{(k)} = \mathcal{O}_{F^{(k)}}(x_i)$ for all iterates $i \in I$ and oracles $k \in K$, we observe that steps of method \mathcal{M} can be still be computed using only information contained in variables $O_i^{(k)}$, so that we can reformulate (PEP) as

$$\begin{aligned}
 \text{(PEP2)} \quad & \sup_{\{O_i^{(k)}\}_{i \in I, k \in K}, \{x_i\}_{i \in I}} \mathcal{P}\left(\left\{O_i^{(k)}\right\}_{i \in I, k \in K}, \{x_i\}_{i \in I}\right), \\
 & \text{subject to } \exists F^{(k)} \in \mathcal{F}_k(\mathbb{E}) \text{ satisfying } \mathcal{O}_{F^{(k)}}(x_i) = O_i^{(k)} \text{ for all } i \in I, k \in K, \\
 & x_0 \text{ satisfies some initialization condition,} \\
 & x_i \text{ is computed by } \mathcal{M} \text{ according to (EQ}_i\text{) for all } 1 \leq i \leq N, \\
 & x_* \text{ a minimizer of } F(x).
 \end{aligned}$$

Note the central role played by the interpolation conditions $\mathcal{O}_{F^{(k)}}(x_i) = O_i^{(k)}$, which enforce the existence of functions $F^{(k)}$ compatible with the output of the oracles. In the next subsection we describe situations for which this formulation is tractable.

1.3. First-order methods and first-order convex interpolation. In the remainder of this work, we restrict ourselves to first-order oracles and methods. We now investigate the concept of (first-order) convex interpolability, in order to make existence constraints from (PEP2) tractable — more precise requirements are detailed in Section 2. From the assumptions, the existence constraint for function $F^{(k)}$

$$\exists F^{(k)} \in \mathcal{F}_k(\mathbb{E}) \text{ satisfying } \mathcal{O}_{F^{(k)}}(x_i) = O_i^{(k)} \text{ for all } i \in I,$$

found in (PEP2) may be expressed in terms of first-order information only. Considering oracles returning first-order information $\mathcal{O}_{F^{(k)}}(x) = (F^{(k)}(x), \bar{\nabla} F^{(k)}(x))$, we denote their output at point x_i by $\mathcal{O}_{F^{(k)}}(x_i) = O_i^{(k)} = (f_i^{(k)}, g_i^{(k)})$. The above existence constraint can be rephrased into the following set of interpolation conditions

$$\text{(INT)} \quad \exists F^{(k)} \in \mathcal{F}_k(\mathbb{E}) \text{ satisfying } F^{(k)}(x_i) = f_i^{(k)} \text{ and } g_i^{(k)} \in \partial F^{(k)}(x_i),$$

which leads us to introduce the following general definition.

DEFINITION 1 ($\mathcal{F}(\mathbb{E})$ -interpolation). *Let I be an index set and $\mathcal{F}(\mathbb{E})$ a class of convex functions, and consider the set of triples $S = \{(x_i, g_i, f_i)\}_{i \in I}$ where $x_i \in \mathbb{E}$, $g_i \in \mathbb{E}^*$ and $f_i \in \mathbb{R}$ for all $i \in I$. The set S is $\mathcal{F}(\mathbb{E})$ -interpolable if and only if there exists a function $F \in \mathcal{F}(\mathbb{E})$ such that both $g_i \in \partial F(x_i)$ and $F(x_i) = f_i$ hold for all $i \in I$.*

The notion of $\mathcal{F}(\mathbb{E})$ -interpolation can be considered for any class of convex functions. It allows us to formulate our performance estimation problem in its final form

$$\begin{aligned}
 \text{(f-PEP)} \quad & \sup_{\{(f_i^{(k)}, g_i^{(k)})\}_{i \in I, k \in K}, \{x_i\}_{i \in I}} \mathcal{P} \left(\left\{ (f_i^{(k)}, g_i^{(k)}) \right\}_{i \in I, k \in K}, \{x_i\}_{i \in I} \right), \\
 \text{subject to} \quad & \left\{ (x_i, g_i^{(k)}, f_i^{(k)}) \right\}_{i \in I} \text{ is } \mathcal{F}_k\text{-interpolable for all } k \in K, \\
 & x_0 \text{ satisfies some initialization condition,} \\
 & x_i \text{ is computed by } \mathcal{M} \text{ according to (EQ}_i\text{) for all } 1 \leq i \leq N, \\
 & x_* \text{ a minimizer of } F(x).
 \end{aligned}$$

We conclude that identifying explicit conditions for convex interpolability by a given class of functions will be the key to eliminate the infinite-dimensional functional variables from (PEP) and transform it into a tractable estimation problem.

First-order convex interpolation was originally developed in [43] for classes of (possibly) L -smooth and (possibly) μ -strongly convex functions. In Section 3, we extend these results to classes of functions involving simultaneously strong convexity, smoothness, gradient boundedness and domain boundedness (for different norms). Those extensions also allow to consider interpolation by indicator or support functions, which may among others be used for problems involving constraints.

Also, note that the notion of first-order interpolability can be adapted for non-convex functions as well. Replacing the concept of subdifferentiability by standard differentiability can be used to study the convergence of first-order algorithms in the cases where some $F^{(k)}$ are not convex (see Section 3.4).

1.4. Prior work. The concept of performance estimation showed itself very promising in the pioneer work of Drori and Teboulle [13], and later in the work of Kim and Fessler [22]. In their work [13], Drori and Teboulle proposed a convex relaxation to obtain numerical upper bounds on the worst-case behaviour of fixed-step first-order algorithms minimizing a single smooth convex function over \mathbb{R}^d , which turned out to be tight in surprisingly many situations⁴. They also proposed a way to numerically optimize the step size parameters of a fixed-step algorithm by minimizing an upper bound on its worst-case. Their approach is based on semidefinite relaxations of (PEP), and was taken further by Kim and Fessler [22], who derived analytically the optimized gradient method previously identified numerically by Drori and Teboulle.

The performance estimation approach on the same smooth unconstrained minimization is further studied in [43], where convex interpolation allows the derivation of an exact convex reformulation of the problem, leading to tight worst-case estimates. The obtained semidefinite formulation also forms the basis for this work.

Another recent and closely related approach for studying performances of first-order methods consists in viewing optimization algorithms as dynamical systems, and

⁴An extension to provide upper bounds for the fixed-step projected gradient method is also provided in Drori's PhD thesis [11].

to use the related stability theory in order to numerically analyse them. This idea is proposed by Lessard et al. in [23], and is attractive because it requires solving a single semidefinite program to obtain a bound that is valid for all subsequent iterations. This technique is particularly efficient for problems involving strong convexity, for which tight linear convergence rates are often recovered. However, as they aim at finding global rates of convergence, they are naturally more conservative than the general performance estimation approach.

For more details on the general topic of convergence analysis of first-order methods, we refer to the seminal books of Yudin and Nemirovski [27], Polyak [36], Nesterov [29] and the more recent book of Bertsekas [4]. Concerning the development of accelerated methods, we specifically refer to the original work of Nesterov [28, 29], and to the later extensions to minimize smoothed convex functions [30] and composite functions [2, 32].

1.5. Paper organization and main contributions. This work is divided into three main parts. First, Section 2 is concerned with putting in place the performance estimation framework for large classes of first-order algorithms, objective functions, performance criteria and initialization conditions. The main idea of this section is to require every element of the performance estimation problem (PEP) to be *linearly Gram-representable* (defined in Section 2.2). This section contains multiple examples of standard settings for which the methodology applies — including among others those covering (sub)gradient methods (along with their projected and proximal counterparts) and conditional gradient methods.

Section 3 focuses on providing convex interpolation conditions for different classes of convex functions commonly arising in practice. Those classes include convex functions, possibly with strong convexity, smoothness, bounded domain and bounded (sub)gradient requirements. The subclasses of indicator and support functions are also explicitly handled. Those classes of functions can all be used directly in the performance estimation framework of Section 2, since their corresponding interpolation conditions are linearly Gram-representable. We end this section with an extension of the convex interpolation results to cope with smooth non-convex functions in a linear Gram-representable way.

In Section 4, we finally apply our approach to several concrete first-order algorithms. We obtain improvements on the analysis of several well-known methods, either analytically or numerically, including the proximal point algorithm and the conditional gradient method. We also use those results to provide an extension of the optimized gradient method proposed by Kim and Fessler [22] that incorporates a projection or a proximal operator to tackle constrained and composite problems.

2. Performance estimation framework for first-order algorithms. We start this section by formulating (f-PEP) in terms of a Gram matrix. This leads to a tractable convex formulation for (f-PEP) — once appropriate assumptions are made on the classes of objective function components, methods, performance criteria and initialization conditions. Those assumptions are motivated by practical applications, which we also provide in the following lines. The main point underlying those assumptions is to ensure that every element of the performance estimation problem can be formulated in a linear way in terms of the entries of a Gram matrix and the function values at the iterates.

2.1. Gram representations. Let us consider $N + 1$ iterates x_0, \dots, x_N and an optimal solution x_* , and the set of corresponding oracle outputs $\left\{ (f_i^{(k)}, g_i^{(k)}) \right\}_{i \in I, k \in K}$.

The accumulated information after those $N + 1$ calls can be gathered into an $d \times (n + 1)(N + 2)$ matrix⁵ P_N (using a slight abuse of notations) and a vector F_N of length $n(N + 2)$:

$$(1) \quad P_N = [Bx_0 \ \dots \ Bx_N \mid Bx_* \mid g_0^{(1)} \ \dots \ g_0^{(n)} \mid \dots \mid g_N^{(1)} \ \dots \ g_N^{(n)} \mid g_*^{(1)} \ \dots \ g_*^{(n)}],$$

$$(2) \quad F_N = [f_0^{(1)} \ \dots \ f_0^{(n)} \mid \dots \mid f_N^{(1)} \ \dots \ f_N^{(n)} \mid f_*^{(1)} \ \dots \ f_*^{(n)}].$$

We also denote by $B^{-1}P_N$ the matrix

$$B^{-1}P_N = [x_0 \ \dots \ x_N \mid x_* \mid B^{-1}g_0^{(1)} \ \dots \ B^{-1}g_*^{(n)}].$$

In order to formulate (PEP) in a tractable way for first-order methods, we use a Gram matrix. That is, we define the symmetric $(n + 1)(N + 2) \times (n + 1)(N + 2)$ Gram matrix $G_N \in \mathbb{S}^{(n+1)(N+2)}$, using the following construction :

$$G_N = \begin{pmatrix} \langle x_0, x_0 \rangle_{\mathbb{E}} & \dots & \langle x_0, x_N \rangle_{\mathbb{E}} & \langle x_0, x_* \rangle_{\mathbb{E}} & \langle g_0^{(1)}, x_0 \rangle & \dots & \langle g_*^{(n)}, x_0 \rangle \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \langle x_N, x_0 \rangle_{\mathbb{E}} & \dots & \langle x_N, x_N \rangle_{\mathbb{E}} & \langle x_N, x_* \rangle_{\mathbb{E}} & \langle g_0^{(1)}, x_N \rangle & \dots & \langle g_*^{(n)}, x_N \rangle \\ \langle x_*, x_0 \rangle_{\mathbb{E}} & \dots & \langle x_*, x_N \rangle_{\mathbb{E}} & \langle x_*, x_* \rangle_{\mathbb{E}} & \langle g_0^{(1)}, x_* \rangle & \dots & \langle g_*^{(n)}, x_* \rangle \\ \langle g_0^{(1)}, x_0 \rangle & \dots & \langle g_0^{(1)}, x_N \rangle & \langle g_0^{(1)}, x_* \rangle & \langle g_0^{(1)}, g_0^{(1)} \rangle_{\mathbb{E}^*} & \dots & \langle g_0^{(1)}, g_*^{(n)} \rangle_{\mathbb{E}^*} \\ \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \langle g_*^{(n)}, x_0 \rangle & \dots & \langle g_*^{(n)}, x_N \rangle & \langle g_*^{(n)}, x_* \rangle & \langle g_*^{(n)}, g_0^{(1)} \rangle_{\mathbb{E}^*} & \dots & \langle g_*^{(n)}, g_*^{(n)} \rangle_{\mathbb{E}^*} \end{pmatrix} \succeq 0.$$

This can be written more compactly as $[G_N]_{ij} = \langle P_N e_i, B^{-1}P_N e_j \rangle = \langle P_N e_i, P_N e_j \rangle_{\mathbb{E}^*}$, where $P_N e_k$ corresponds to the k th column of P_N . Also, note that the size of this matrix does not depend on the dimension d of the spaces we are working with.

REMARK 2. *Note that Gram matrix G_N is positive semidefinite for any matrix P_N (of the form (1)). The number of linearly independent columns of P_N is equal to the rank of G_N . Hence this rank is upper bounded by the dimension d of the ambient space of the iterates. It is possible to recover a matrix P_N of the form⁶ (1) from any Gram matrix $G_N \succeq 0$ satisfying $\text{Rank } G_N \leq d$.*

Our goal for the next subsections is to show that in a lot of situations, the performance estimation problem (f-PEP) can be expressed exactly as a semidefinite program in the F_N and G_N variables:

$$(\text{SDP-PEP}) \quad \sup_{G_N \succeq 0, F_N} c^\top F_N + \text{Tr}(CG_N) \quad \text{s.t.} \quad a_i + b_i^\top F_N + \text{Tr}(D_i G_N) \leq 0 \quad \forall i \in S$$

with S some index set related to the constraints, and elements a_i, b_i, c, D_i and C of appropriate dimensions for writing the constraints and objective function linearly in terms of the Gram matrix G_N and of the objective function values F_N .

2.2. Tractable formulation of the performance estimation problem.

In this section, we present our main result, stating that computing the exact worst-case performance of a method on a class of functions is tractable and can, in many cases, be formulated as (SDP-PEP). We start with the concept of Gram-representability for the different ingredients of the performance estimation problem.

⁵We remind the reader that $B : \mathbb{E} \rightarrow \mathbb{E}^*$ is a positive definite operator which is chosen as the identity operator in standard situations (see Section 1.1).

⁶In the case $\mathbb{E} = \mathbb{E}^* = \mathbb{R}^d$ with the usual inner product $\langle x, y \rangle = x^\top y$ and B the identity operator, this can be done using the standard Cholesky factorization. In the general cases the exact same idea can be used, using the chosen inner product $\langle \cdot, \cdot \rangle_{\mathbb{E}^*}$ in the process.

DEFINITION 3. *A class of functions is Gram-representable (resp. linearly Gram-representable) if and only if its interpolation conditions (INT) can be formulated using a finite number of convex (resp. linear) constraints involving only the matrix G_N and the corresponding function values F_N .*

The functional classes of smooth strongly convex functions, smooth convex functions with bounded (sub)gradients, and strongly convex functions with bounded domain are linearly Gram-representable. In addition, the particular subclasses of support and indicator convex functions share this same advantageous property. The details and proofs of these results are postponed to Section 3.

DEFINITION 4. *A performance measure is Gram-representable (resp. linearly Gram-representable) if and only if it can be expressed as a concave (resp. linear) function involving only the matrix G_N and the corresponding function values F_N .*

The class of linearly Gram-representable performance criteria contains a large variety of choices, including most standard measures we are aware of. For example, it is easy to check that standard optimality criteria in function values $F(x_N) - F(x_*)$, in residual subgradient norm $\|\tilde{\nabla}F(x_N)\|_{\mathbb{E}^*}^2$, distance to optimality $\|x_N - x_*\|_{\mathbb{E}}^2$, and distance to feasibility $\|x_N - \Pi_Q(x_N)\|_{\mathbb{E}}^2$ can be handled.

On the other hand, multiple examples of non-linear Gram-representable performance criteria can also be handled with no difficulty. This includes performance measures involving the best values among all iterates, for example $\min_{0 \leq i \leq N} F(x_i) - F(x_*)$, or the best residual gradient norm among the iterates $\min_{0 \leq i \leq N} \|\nabla F(x_i)\|_{\mathbb{E}^*}^2$ (see also [43, Sect. 4.3]).

DEFINITION 5. *An initialization condition is Gram-representable (resp. linearly Gram-representable) if and only if it can be expressed using a finite number of convex (resp. linear) constraints involving only the matrix G_N and the corresponding function values F_N .*

Standard examples of valid initial conditions include the classical bounds on the initial distance to optimality $\|x_0 - x_*\|_{\mathbb{E}}^2 \leq R^2$, on the initial function value $F_0 - F_* \leq R$, and on initial gradient value $\|\nabla F(x_0)\|_{\mathbb{E}^*}^2 \leq R^2$.

DEFINITION 6. *A first-order method is Gram-representable (resp. linearly Gram-representable) if and only if the computation of its iterates, implicitly defined by an equation of type (EQ_i), can be expressed using a finite number of convex (resp. linear) constraints involving only the matrix G_N and the corresponding function values F_N .*

We refer to the next section for examples of linearly Gram-representable methods.

We can now state our main results concerning Gram-representable situations. In the sequel, we recall that we use the notation $\mathcal{F}_K(\mathbb{E})$ to denote the set of functions of the form (CM) with components $F^{(k)} \in \mathcal{F}_k(\mathbb{E}) \forall k \in K$ — i.e., $F \in \mathcal{F}_K(\mathbb{E})$.

PROPOSITION 7. *Consider a class of composite objective functions $\mathcal{F}_K(\mathbb{E})$ with n components, a first-order method \mathcal{M} , a performance measure \mathcal{P} and an initial condition \mathcal{I} which are all Gram-representable.*

Computing the worst-case for criterion \mathcal{P} of method \mathcal{M} after N iterations on objective functions in class $\mathcal{F}_K(\mathbb{E})$ with initial condition \mathcal{I} can be formulated as a convex program when dimension of the space \mathbb{E} satisfies $d \geq (n+1)(N+2)$. Otherwise, it can be formulated as a convex program plus an additional non-convex rank constraint $\text{Rank } G_N \leq d$.

If in addition $\mathcal{F}_K(\mathbb{E})$, \mathcal{M} , \mathcal{P} and \mathcal{I} are linearly Gram-representable, then the

corresponding optimization problem is a semidefinite program of the form (SDP-PEP), whose variables are $F_N \in \mathbb{R}^{n(N+2)}$ and $G_N \in \mathbb{S}^{(n+1)(N+2)}$.

Proof. It directly follows from Remark 2 and from the definitions of (linear) Gram-representability for the class of functions, first-order methods, performance measures and initialization conditions: any solution to the corresponding optimization problem can be transformed into a particular instance of (CM) and vice versa. \square

REMARK 8. The optimal value of (PEP) increases with dimension d . When (PEP) with Gram-representable elements attains a finite optimal value, Proposition 7 implies the existence of a function with dimension at most $(n+1)(N+2)$ that achieves the worst-case value.

REMARK 9. The assumption $d \geq (n+1)(N+2)$ is referred to as the large-scale assumption in the sequel. In terms of performance estimation problems, this assumption allows to discard the non-convex rank constraint and lead to a tractable semidefinite programming problem, which can be solved to global optimality efficiently (see e.g., [45]). Without that assumption, our performance estimation problem is a nonconvex rank-constrained semidefinite program, equivalent to a quadratic programming problem that is NP-hard in general (e.g., it has max-cut [16] and other non-convex quadratic programs [34, 40] as particular cases). Approaches to handle rank constraints exist (e.g., via augmented Lagrangian techniques [6], via manifold optimization [20] or via Newton-like methods [33]), but in general only guarantee convergence to stationary points. This is not useful for in the case of (SDP-PEP), as this only provides lower bounds on the worst-case performance.

REMARK 10. Under the large-scale assumption, we obtain dimension-free guarantees (i.e. valid for any dimension, and tight as soon as $d \geq (n+1)(N+2)$), as is commonly found in the literature about first-order methods. In addition, we note that the dimension bound $(n+1)(N+2)$ is in fact (very) conservative for most standard algorithms — that is, the bound in the large-scale assumption can typically be significantly reduced, see Corollary 14 in the sequel.

REMARK 11. The worst-case results provided by the SDP from Proposition 7 provide a tight worst-case achievable for any operator B and any dual pairing $\langle \cdot, \cdot \rangle$.

REMARK 12. The necessary and sufficient condition for x_* to be optimal for F is linearly Gram-representable. Indeed, it corresponds to requiring $\tilde{\nabla}F(x_*) = 0$, i.e.

$$\sum_{k \in K} \tilde{\nabla}F^{(k)}(x_*) = \sum_{k \in K} g_*^{(k)} = 0 \Leftrightarrow \left\| \sum_{k \in K} g_*^{(k)} \right\|_{\mathbb{E}^*}^2 = \left\langle \sum_{k \in K} g_*^{(k)}, \sum_{k \in K} g_*^{(k)} \right\rangle_{\mathbb{E}^*} = 0,$$

where the last condition is linear in the entries of G_N .

2.3. Linearly Gram-representable first-order methods. This class of first-order methods contains as particular cases what we call in the following the class of *fixed-step linear first-order methods* (FSLFOM), whose iterations are defined by a linear equation (with known constant coefficients) involving the iterates and the corresponding (sub)gradients.

DEFINITION 13. A *fixed-step linear first-order method* (FSLFOM) is a method

which computes iterate x_{i+1} as the solution of ⁷

$$\text{(FSLFOM)} \quad t_{i+1,i+1} B x_{i+1} + \sum_{k \in K} h_{i+1,i+1}^{(k)} g_{i+1}^{(k)} = \sum_{j=0}^i t_{i+1,j} B x_j + \sum_{j=0}^i \sum_{k \in K} h_{i+1,j}^{(k)} g_j^{(k)},$$

where all coefficients $h_{i+1,j}^{(k)}, t_{i+1,j} \in \mathbb{R}$ are fixed beforehand.

Note the class of FSLFOM is exactly the class of methods whose iterations can be written in the form (using first-order optimality conditions, and the convexity of $F^{(k)}$):

$$x_{i+1} = \operatorname{argmin}_{x \in \mathbb{E}} \left\{ \sum_{k \in K} h_{i+1,i+1}^{(k)} F^{(k)}(x) + \frac{t_{i+1,i+1}}{2} \|x\|_{\mathbb{E}}^2 - \left\langle \sum_{j=1}^i t_{i+1,j} B x_j + \sum_{j=0}^i \sum_{k \in K} h_{i+1,j}^{(k)} \nabla F^{(k)}(x_j), x \right\rangle \right\};$$

which in some sense represents the most general method allowed in our framework. Those iterations can also be written by linearly combining the columns of the matrix P_N containing all the harvested first-order information about the problem:

$$0 = P_N \underline{\alpha}_k,$$

with $\underline{\alpha}_k \in \mathbb{R}^{(n+1)(N+2)}$ a vector containing appropriate coefficients. Therefore, we note that any FSLFOM is linearly Gram-representable using the following formulation:

$$(3) \quad 0 = P_N \underline{\alpha}_k \Leftrightarrow 0 = \|P_N \underline{\alpha}_k\|_{\mathbb{E}^*}^2 = \langle P_N \underline{\alpha}_k, B^{-1} P_N \underline{\alpha}_k \rangle,$$

which is clearly linear in terms of the Gram matrix G_N . Note that this can also be extended to cope with the more general class of linearly Gram-representable first-order methods⁸:

$$(4) \quad c_k^{(\text{low})\top} F_N + b_k^{(\text{low})} \leq \underline{\alpha}_k^\top G \underline{\alpha}_k \leq c_k^{(\text{up})\top} F_N + b_k^{(\text{up})},$$

where $c_k^{(\text{low})}, b_k^{(\text{low})}$ and $c_k^{(\text{up})}, b_k^{(\text{up})}$ are some fixed parameters. Those can for example be used in order to require a sufficient decrease condition, or an inexact version of (FSLFOM):

$$\text{(Inexact FSLFOM)} \quad \underline{\alpha}_k^\top G \underline{\alpha}_k \leq \varepsilon_k,$$

with $\varepsilon_k \geq 0$ some accuracy parameter for the computation of (FSLFOM).

Examples of FSLFOM. Before going into the details of the performance estimation problems for our class of linear fixed-step methods and over the different classes of convex functions, let us give several examples of methods fitting into the model provided by (FSLFOM) and (Inexact FSLFOM).

- Fixed-step subgradient and gradient algorithms: fixed-step subgradient methods for minimizing a convex function F are naturally described as $x_i = x_{i-1} - \alpha_i B^{-1} g_{i-1}$ with α_i some step size, and $g_{i-1} \in \partial F(x_{i-1})$. The method is clearly in the class of FSLFOM and its linear Gram matrix representation can be obtained using formulation (3).

⁷Note that the iteration is written as an equality on \mathbb{E} , but it is possible and totally equivalent to write it on \mathbb{E}^* using the operator B^{-1} — B is invertible by assumption.

⁸This formulation is just provided as an illustration to show that more general methods than (FSLFOM) can still be considered.

- Proximal methods and proximal gradient methods: fixed-step proximal gradient methods for minimizing $F^{(1)} + F^{(2)}$ is usually described as doing an explicit (sub)gradient step on $F^{(1)}$ followed by a minimization step on $F^{(2)}$:

$$\begin{aligned} x_i &= \text{prox}_{\alpha_i F^{(2)}} \left(x_{i-1} - \alpha_i B^{-1} \tilde{\nabla} F^{(1)}(x_{i-1}) \right) \\ &= \underset{x \in \mathbb{E}}{\text{argmin}} \left\{ \alpha_i F^{(2)}(x) + \frac{1}{2} \left\| x_{i-1} - \alpha_i B^{-1} \tilde{\nabla} F^{(1)}(x_{i-1}) - x \right\|_{\mathbb{E}}^2 \right\}. \end{aligned}$$

Optimality conditions on this last term allow writing each iterations as

$$Bx_i + \alpha_i \tilde{\nabla} F^{(2)}(x_i) = Bx_{i-1} - \alpha_i \tilde{\nabla} F^{(1)}(x_{i-1}),$$

with some $\tilde{\nabla} F^{(2)}(x_i) \in \partial F^{(2)}(x_i)$. This method is clearly a FSLFOM and therefore fits in the framework. Also, note that projected gradient methods are obtained using the same technique, but on the particular class of convex indicator functions, whereas proximal point algorithms correspond to the case where $F^{(1)} = 0$.

- Conditional gradient methods do also fit into the model provided by Equation (FSLFOM). Indeed, the iterations take the following form:

$$\begin{aligned} y_i &= \underset{z \in \mathbb{E}}{\text{argmin}} \left\{ \left\langle z - z_i, \tilde{\nabla} F^{(1)}(z_i) \right\rangle + F^{(2)}(z) \right\}, \\ z_{i+1} &= (1 - \lambda_i) z_i + \lambda_i y_i, \end{aligned}$$

with $\lambda_i \in [0, 1]$ chosen beforehand. Now, by imposing y_i using first-order necessary and sufficient optimality conditions on the intermediate optimization problem, we obtain

$$\tilde{\nabla} F^{(1)}(z_i) = -\tilde{\nabla} F^{(2)}(y_i).$$

Note that for conditional gradient-type methods, $F^{(2)}$ is usually chosen as the indicator function of some closed convex set Q . This algorithm can also clearly be written as a FSLFOM by artificially denoting for $i = 0, 1, \dots$ the iterates $x_{2i} = z_i$ and $x_{2i+1} = y_i$.

- Inexact (sub)gradient methods for a convex function $F^{(1)}$, with $x_{i+1} = x_i - \alpha_i B^{-1}(\tilde{\nabla} F^{(1)}(x_i) + \varepsilon_i)$ and $\|\varepsilon_i\|_{\mathbb{E}^*} \leq \epsilon_i$ for some $\epsilon_i \geq 0$ the tolerance on the (sub)gradient computation. This can be written in the inexact FSLFOM format:

$$\left\| \alpha_i (x_{i+1} - x_i) + \tilde{\nabla} F^{(1)}(x_i) \right\|_{\mathbb{E}^*}^2 \leq \epsilon_i^2.$$

Also, note that other noise models, as for example the one proposed by d'Aspremont [8] can also easily be used in the framework. On the other hand, the inexact (δ, L) -oracles developed by Devolder et al. [10] do not seem to easily fit into the approach⁹.

Note that a broad class of methods can be modelled using those operations, just by requiring the functions on which it is applied to belong to certain classes. As an example, alternate projection-type algorithms are special cases of proximal methods, applied on the class of convex indicator functions. Therefore, they can be represented in the FSLFOM format.

⁹This is due to the fact no necessary and sufficient interpolation conditions for this noise model were found — that is, standard conditions are only necessary to guarantee interpolability. Using necessary conditions that are not sufficient still allows obtaining upper bounds on the worst-case behavior, but those may not be tight.

2.4. Simplified performance estimation problems. Note that for standard algorithms such as the previous examples of FSLFOM, the SDP resulting from Proposition 7 can typically be further simplified, leading to a reduction in its size.

COROLLARY 14. *Consider a class of functions $\mathcal{F}_K(\mathbb{E})$, a performance measure \mathcal{P} and an initialization condition \mathcal{I} which are linearly Gram-representable, and a FSLFOM \mathcal{M} whose iterations are all linearly independent¹⁰.*

In addition, assume there are p points $(g_i^{(k)}, f_i^{(k)})$ such that neither $g_i^{(k)}$ nor $f_i^{(k)}$ are used in the performance measure \mathcal{P} , the initial condition \mathcal{I} and the method \mathcal{M} . Then, the performance estimation problem can be written as a convex SDP using variables $F_N \in \mathbb{R}^{n(N+2)-p}$ and $G_N \in \mathbb{S}^{(n+1)(N+2)-N-p-1}$, with the possible additional rank constraint $\text{rank } G_N \leq d$.

To see why this corollary holds, note that the variables in the simplified SDP correspond to the function values and the Gram matrix from which the p unnecessary points were removed, and from which N other variables were substituted using the N iteration constraints (FSLFOM)¹¹.

Under the assumptions of Corollary 14, the large-scale assumption becomes $d \geq (n+1)(N+2) - N - p - 1$. In the cases where only the output from a single oracle is used at each iteration, we have that $p = (n-1)(N+1)$, which leads to $d \geq N + n + 2$.

Furthermore, for standard performance measures (e.g. $F_N - F_*$, $\|x_N - x_*\|_{\mathbb{E}}^2$, $\|\tilde{\nabla}F(x_N)\|_{\mathbb{E}^*}^2$), one arbitrary point x_i may be fixed to 0 because solutions to the SDP are invariant with respect to translations. This would result in the large-scale assumption $d \geq N + n + 1$. For $n = 1$, we recover the standard $d \geq N + 2$ appearing in the case of a single component in the objective function [43].

The original SDP from Proposition 7 may be challenging to solve in practice, because of its potentially large size on the one hand, and because it may lack an interior on the other hand. We observe that the simplified performance estimation problem described above typically improves the situation for both issues, reducing the size of the problem and solving in a lot of cases the issue of a lack of interior points.

3. Convex interpolation. In this section, we study the convex interpolation problems for different standard classes of convex functions. The underlying motivation is to obtain discrete characterization of convex functions commonly arising in the context of convex optimization via first-order methods. More specifically, the classes of convex functions of interest for this section are all linearly Gram-representable (see Definition 3). Therefore, those classes can be used with tightness guarantees within the performance estimation framework.

Main technical tools from the next sections are borrowed from convex analysis, we refer to the seminal references [1, 18, 38, 39] for details.

3.1. Functional characteristics. In this section, we start by considering several desirable properties of proper closed convex functions. As previously underlined, the choice of examining those characteristics comes from their common appearance in the context of first-order convex optimization. Consider a proper, closed and convex function f , the main characteristics of interest for us are the following.

¹⁰That is, the vectors α_k used to characterize the iterations are linearly independent — this is very reasonable, as every method using new information at each iteration satisfies this. Remark that it does not imply that the points x_i themselves are linearly independent.

¹¹The additional -1 term appearing in the dimension of the Gram matrix comes from the fact that one of the $g_*^{(k)}$ may also be discarded, by substituting it using the optimality condition of x_* .

- (a) **Smoothness**: there exists some $L \in \mathbb{R}^{++} \cup \{\infty\}$ such that $\frac{1}{L}\|g_1 - g_2\|_{\mathbb{E}^*} \leq \|x_1 - x_2\|_{\mathbb{E}}$ holds for all pairs $x_1, x_2 \in \mathbb{E}$ and corresponding subgradients $g_1, g_2 \in \mathbb{E}^*$ (i.e. such that $g_1 \in \partial f(x_1)$ and $g_2 \in \partial f(x_2)$).
- (b) **Strong convexity**: there exists some $\mu \in \mathbb{R}^+$ such that the function $f(x) - \frac{\mu}{2}\|x\|_{\mathbb{E}}^2$ is convex.
- (c) **Gradient boundedness**: there exists some $M \in \mathbb{R}^+ \cup \{\infty\}$ such that $\|g\|_{\mathbb{E}^*} \leq M$ holds for all subgradient $g \in \mathbb{E}^*$ (i.e., such that $\exists x : g \in \partial f(x)$).
- (d) **Domain boundedness**: there exists some $D \in \mathbb{R}^+ \cup \{\infty\}$ such that $\forall x \in \{x \in \mathbb{E} : f(x) < \infty\}$ we have $\|x\|_{\mathbb{E}} \leq D$.

Alternatively, domain and gradient boundedness can be specified in terms of diameters instead of radii.

- (c') **Gradient boundedness**: there exists some $M \in \mathbb{R}^+ \cup \{\infty\}$ such that $\|g_1 - g_2\|_{\mathbb{E}^*} \leq M$ holds for all subgradients $g_1, g_2 \in \mathbb{E}^*$ (i.e., such that $\exists x_1, x_2 : g_1 \in \partial f(x_1)$ and $g_2 \in \partial f(x_2)$).
- (d') **Domain boundedness**: there exists some $D \in \mathbb{R}^+ \cup \{\infty\}$ such that $\forall x_1, x_2 \in \{x \in \mathbb{E} : f(x) < \infty\}$ we have $\|x_1 - x_2\|_{\mathbb{E}} \leq D$.

As some characteristics are incompatible (e.g., gradient boundedness is incompatible with strong convexity, domain boundedness is incompatible with smoothness), we define three classes of functions.

DEFINITION 15. *Let $f : \mathbb{E} \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper, closed and convex function.*

We say that:

- $f \in \mathcal{F}_{\mu,L}(\mathbb{E})$ (*L -smooth μ -strongly convex functions*) if it satisfies conditions (a) and (b) with $\mu < L$.
- $f \in \mathcal{C}_{M,L}(\mathbb{E})$ (*L -smooth M -Lipschitz convex functions*) if it satisfies conditions (a) and (c); alternatively, $f \in \mathcal{C}'_{M,L}(\mathbb{E})$ if it satisfies (a) and (c').
- $f \in \mathcal{S}_{D,\mu}(\mathbb{E})$ (*D -bounded μ -strongly convex functions*) if it satisfies conditions (b) and (d); alternatively $f \in \mathcal{S}'_{D,\mu}(\mathbb{E})$ if it satisfies (b) and (d').

Note that boundedness and smoothness constants are allowed to take the value ∞ , in order to embed the unbounded (domain or gradient) and non-smooth functions as well. We handle those by using the conventions $1/\infty = 0$ and $\infty - c = \infty$ for $c \in \mathbb{R}$.

A basic building block for the smooth convex interpolation conditions proposed in [43] comes from the Fenchel-Legendre conjugation. In particular, the duality correspondence $f \in \mathcal{F}_{0,\infty}(\mathbb{E}) : f \in \mathcal{F}_{\mu,L}(\mathbb{E}) \Leftrightarrow f^* \in \mathcal{F}_{1/L,1/\mu}(\mathbb{E}^*)$ was intensively used to require smoothness of the convex interpolant. In the following, we additionally use the duality $f \in \mathcal{F}_{0,\infty}(\mathbb{E}) : f \in \mathcal{C}_{M,L}(\mathbb{E}) \Leftrightarrow f^* \in \mathcal{S}_{M,1/L}(\mathbb{E}^*)$ (and the variant $f \in \mathcal{F}_{0,\infty}(\mathbb{E}) : f \in \mathcal{C}'_{M,L}(\mathbb{E}) \Leftrightarrow f^* \in \mathcal{S}'_{M,1/L}(\mathbb{E}^*)$) in order to include boundedness properties in the convex interpolating functions, along with smoothness.

THEOREM 16. *Consider a function $f \in \mathcal{F}_{0,\infty}(\mathbb{E})$. We have $f \in \mathcal{C}_{M,\infty}(\mathbb{E})$ (resp. $f \in \mathcal{C}'_{M,\infty}(\mathbb{E})$) if and only if $f^* \in \mathcal{S}_{M,0}(\mathbb{E}^*)$ (resp. $f^* \in \mathcal{S}'_{M,0}(\mathbb{E}^*)$).*

This theorem follows quite naturally using the equivalence: $\forall f \in \mathcal{F}_{0,\infty}(\mathbb{E}) : g \in \partial f(x) \Leftrightarrow x \in \partial f^*(g) \Leftrightarrow f(x) + f^*(g) = \langle g, x \rangle$.

3.2. Interpolation conditions. In this section, we provide interpolation conditions for the three classes of functions we previously introduced. We particularize those classes to convex indicator and support functions in the next section. We start by recalling the following known interpolation result [43, Theorem 6].

THEOREM 17. *The set $\{(x_i, g_i, f_i)\}_{i \in I}$ is $\mathcal{F}_{\mu,L}$ -interpolable if and only if the fol-*

lowing set of conditions holds for every pair of indices $i \in I$ and $j \in I$

$$f_i - f_j - \langle g_j, x_i - x_j \rangle \geq \frac{1}{2(1 - \mu/L)} \left(\frac{1}{L} \|g_i - g_j\|_{\mathbb{E}^*}^2 + \mu \|x_i - x_j\|_{\mathbb{E}}^2 - 2\frac{\mu}{L} \langle g_j - g_i, x_j - x_i \rangle \right).$$

In particular, the simpler closed, convex and proper interpolation conditions (i.e. $\mathcal{F}_{0,\infty}(\mathbb{E})$ interpolation) are $\forall i, j \in I$

$$(5) \quad f_i - f_j - \langle g_j, x_i - x_j \rangle \geq 0,$$

which will serve to develop our next interpolation conditions. We start with $\mathcal{S}_{D,\mu}(\mathbb{E})$ -interpolability, which will then serve to obtain $\mathcal{C}_{M,L}(\mathbb{E})$ -interpolation conditions using conjugation.

REMARK 18. Note that Theorem 17 is formally proven in [43] for the case of the standard inner product $\langle x, y \rangle = x^\top y$ (and therefore also only for $\|\cdot\|_2^2$). However, exactly the same approach can directly be used to obtain the desired result for general inner products on \mathbb{E} and self-adjoint positive definite linear operators B , and the corresponding induced primal and dual Euclidean norms.

REMARK 19. By assuming $\mu < L$, we ignore the classes $\mathcal{F}_{L,L}(\mathbb{E})$, for $L \geq 0$. That is, the classes of functions containing only functions of the form $f(x) = \frac{L}{2} \|x\|_{\mathbb{E}}^2 + \langle b, x \rangle + c$ for some $b \in \mathbb{E}^*$ and $c \in \mathbb{R}$. It is straightforward to adapt all the interpolation conditions for handling this case.

THEOREM 20. The set $\{(x_i, g_i, f_i)\}_{i \in I}$ is $\mathcal{S}_{D,\mu}$ (D -bounded, μ -strongly convex) (resp. $\mathcal{S}'_{D,\mu}$) interpolable if and only if the following set of conditions holds for every pair of indices $i \in I$ and $j \in I$

$$\begin{aligned} f_i - f_j - \langle g_j, x_i - x_j \rangle &\geq \frac{\mu}{2} \|x_i - x_j\|_{\mathbb{E}}^2, \\ \|x_j\|_{\mathbb{E}} &\leq D \quad (\text{resp. } \|x_j - x_i\|_{\mathbb{E}} \leq D). \end{aligned}$$

Proof. (Necessity) Every function $f \in \mathcal{S}_{D,\mu}(\mathbb{E})$ (resp. $f \in \mathcal{S}'_{D,\mu}(\mathbb{E})$) satisfies the conditions.

(Sufficiency) Consider the following construction:

$$f(x) = \begin{cases} \max_i \left\{ f_i + \langle g_i, x - x_i \rangle + \frac{\mu}{2} \|x - x_i\|_{\mathbb{E}}^2 \right\} & \text{if } x \in \text{conv}(\{x_i\}_{i \in I}) \\ \infty & \text{else,} \end{cases}$$

One can note that f is μ -strongly convex (convex domain, and maximum of μ -strongly convex functions), and that it indeed interpolates the set $\{(x_i, g_i, f_i)\}_{i \in I}$:

$$\begin{aligned} f(x_j) &= \max_i \left\{ f_i + \langle g_i, x_j - x_i \rangle + \frac{\mu}{2} \|x_j - x_i\|_{\mathbb{E}}^2 \right\}, \\ &\leq f_j, \end{aligned}$$

using interpolation conditions. By noting that the maximum is bigger than taking individually the component j , we also have $\max_i \left\{ f_i + \langle g_i, x_j - x_i \rangle + \frac{\mu}{2} \|x_j - x_i\|_{\mathbb{E}}^2 \right\} \geq$

f_j , which allows to conclude that $f(x_j) = f_j$. To obtain that $g_j \in \partial f(x_j)$, let us write

$$\begin{aligned} f(x) &= \max_i \left\{ f_i + \langle g_i, x - x_i \rangle + \frac{\mu}{2} \|x - x_i\|_{\mathbb{E}}^2 \right\}, \\ &\geq \max_i \{ f_i + \langle g_i, x - x_i \rangle \}, \\ &\geq f_j + \langle g_j, x - x_j \rangle. \end{aligned}$$

Finally, note that $\text{conv}(\{x_i\}_{i \in I}) \subseteq B_{\mathbb{E}}(0, D)$, with $B_{\mathbb{E}}(0, D)$ the ball of norm $\|\cdot\|_{\mathbb{E}}$ centered at the origin and with radius D . Indeed, choose $z = \sum_{i \in I} \lambda_i x_i$ with $\lambda_i \geq 0$ and $\sum_{i \in I} \lambda_i = 1$, we have $\|z\|_{\mathbb{E}} \leq \sum_{i \in I} \lambda_i \|x_i\|_{\mathbb{E}} \leq D$, and f has a bounded domain of radius D . Hence $\{(x_i, g_i, f_i)\}_{i \in I}$ is $\mathcal{S}_{D, \mu}$ -interpolable, which concludes the proof for the $\mathcal{S}_{D, \mu}$ part.

To obtain the same result for $\mathcal{S}'_{D, \mu}$, note that $\forall y, z \in \text{conv}(\{x_i\}_i)$, we can write $y = \sum_i \lambda_i x_i$ and $z = \sum_i \gamma_i x_i$ with $\lambda_i, \gamma_i \geq 0$ and $\sum_i \lambda_i = \sum_i \gamma_i = 1$. Hence, $\|y - z\|_{\mathbb{E}} \leq \sum_i \lambda_i \sum_j \gamma_j \|x_i - x_j\|_{\mathbb{E}} \leq D$. \square

This interpolation result can directly be used for developing interpolation conditions for the class of convex functions with bounded gradient, using the conjugate duality between smoothness and strong convexity on the one hand, and gradient and domain boundedness on the other hand.

THEOREM 21. *The set $\{(x_i, g_i, f_i)\}_{i \in I}$ is $\mathcal{C}_{M, L}$ (L -smooth with M -bounded subgradients) (resp $\mathcal{C}'_{M, L}$) interpolable if and only if the following set of conditions holds for every pair of indices $i \in I$ and $j \in I$*

$$\begin{aligned} (6) \quad & f_i - f_j - \langle g_j, x_i - x_j \rangle \geq \frac{1}{2L} \|g_i - g_j\|_{\mathbb{E}^*}^2, \\ (7) \quad & \|g_j\|_{\mathbb{E}^*} \leq M \quad (\text{resp. } \|g_j - g_i\|_{\mathbb{E}^*} \leq M). \end{aligned}$$

Proof. Note that a function $f \in \mathcal{C}_{M, L}(\mathbb{E})$ (resp. $f \in \mathcal{C}'_{M, L}(\mathbb{E})$) interpolates the set $\{(x_i, g_i, f_i)\}_{i \in I}$ if and only if there exists a corresponding conjugate function $f^* \in \mathcal{S}_{M, 1/L}(\mathbb{E}^*)$ (resp. $f^* \in \mathcal{S}'_{M, 1/L}(\mathbb{E}^*)$) interpolating the set $\{(g_i, x_i, \langle x_i, g_i \rangle - f_i)\}_{i \in I} = \{(\tilde{x}_i, \tilde{g}_i, \tilde{f}_i)\}_{i \in I}$ (see Section 3.1). Using interpolation conditions from Theorem 20, such a f^* exists if and only if

$$\begin{aligned} \tilde{f}_i - \tilde{f}_j - \langle \tilde{x}_i - \tilde{x}_j, \tilde{g}_j \rangle &\geq \frac{1}{2L} \|\tilde{x}_i - \tilde{x}_j\|_{\mathbb{E}^*}^2, \\ \|\tilde{x}_j\|_{\mathbb{E}^*} &\leq M \quad (\text{resp. } \|\tilde{x}_j - \tilde{x}_i\|_{\mathbb{E}^*} \leq M), \end{aligned}$$

which are respectively equivalent to conditions (6) and (7). \square

3.3. Indicator and support functions. Constraints and regularization are so recurrent in optimization that we dedicate the next lines specifically to interpolation procedures tailored to them. We remind the reader that we call a proper and closed convex function $i : \mathbb{E} \rightarrow \{0, \infty\}$ a D -bounded indicator function (which we denote by $f \in \mathcal{I}_D(\mathbb{E})$ — resp. $f \in \mathcal{I}'_D(\mathbb{E})$) if there exists a radius (resp. a diameter) $0 \leq D \leq \infty$ such that $\forall x \in \{x : i(x) = 0\}$ (resp. $\forall x_1, x_2 \in \{x : i(x) = 0\}$) we have $\|x\|_{\mathbb{E}} \leq D$ (resp. $\|x_1 - x_2\|_{\mathbb{E}} \leq D$).

Indicator functions. Let us now consider the special case of interpolating indicator functions. Basically, this problem is a particular case of the $\mathcal{S}_{D,\mu}$ (or $\mathcal{S}'_{D,\mu}$)-interpolation problem with $\mu = 0$ — however, note that indicator function interpolation is not straightforward from $\mathcal{S}'_{D,\mu}$ -interpolation, as for example requiring the corresponding interpolation constraints in addition to $f_i = 0$ would not a priori guarantee that the interpolated function from Theorem 20 would indeed satisfy $f(x) = 0$ on $\text{dom } f$. This class of functions is particularly interesting for projections in the context of performance estimation.

THEOREM 22. *The set $\{(x_i, g_i, f_i)\}_{i \in I}$ is \mathcal{I}_D (D -bounded indicator) (resp. \mathcal{I}'_D) interpolable if and only if the following inequalities hold $\forall i, j \in I$:*

$$(8) \quad \begin{aligned} f_i &= 0, \\ \langle g_j, x_i - x_j \rangle &\leq 0, \\ \|x_i\|_{\mathbb{E}} &\leq D \quad (\text{resp. } \|x_j - x_i\|_{\mathbb{E}} \leq D). \end{aligned}$$

Proof. (Necessity) Any function $f \in \mathcal{I}_D(\mathbb{E})$ (resp. $f \in \mathcal{I}'_D(\mathbb{E})$) satisfies those conditions.

(Sufficiency) Let us construct a convex set whose indicator function interpolate for the set $\{(x_i, g_i, 0)\}$. That is, we construct a closed convex set Q containing all x_i 's, and such that $\forall x \in Q$ we have $\langle g_i, x - x_i \rangle \leq 0$ and such that $\|x\|_{\mathbb{E}} \leq D \forall x \in Q$ (resp. $\|x - y\|_{\mathbb{E}} \leq D \forall x, y \in Q$). We start with the simpler case $D = \infty$, by considering the polyhedral set

$$Q = \{x \in \mathbb{E} \mid \langle a_j, x \rangle \leq b_j \forall j \in I\},$$

with $a_j = g_j$ and $b_j = \langle g_j, x_j \rangle$. The construction guarantees that $x_i \in Q$. Indeed, by Condition (8) we have $\langle g_j, x_i \rangle \leq \langle g_j, x_j \rangle$, which is equivalent to $\langle a_j, x_i \rangle \leq b_j$ using the definitions of a_j and b_j , and therefore guarantees that $x_i \in Q$. In order to add the boundedness requirement, we modify the set Q in the following way: $\tilde{Q} = Q \cap \text{conv}(\{x_i\}_{i \in I})$. This new set is still convex (intersection of two convex sets), it also trivially still satisfies $x_i \in \tilde{Q}$ (which are by construction both contained in Q and $\text{conv}(\{x_i\}_i)$) and $\langle g_i, x - x_i \rangle \leq 0 \forall x \in \tilde{Q}$ (since $\tilde{Q} \subseteq Q$). In addition, \tilde{Q} has a radius bounded above by D , because D is an upper bound on the radius (resp. diameter) of $\text{conv}(\{x_i\}_i)$. It is therefore clear that the indicator function $I_{\tilde{Q}} \in \mathcal{I}_D(\mathbb{E})$ (resp. $\mathcal{I}'_D(\mathbb{E})$) interpolates $\{(x_i, g_i, 0)\}_{i \in I}$ as the convex hull has a radius (resp. diameter) D (see proof of Theorem 20). \square

Support functions. Support functions — we denote the set of support function with a M -Lipschitz condition by $\mathcal{I}_M^*(\mathbb{E})$ (resp. $\mathcal{I}'_M^*(\mathbb{E})$) — are standardly defined as convex conjugate of indicator functions, and are also very common in applications (all norms are support functions, think e.g. of the l_1 norm being the support function of the unit ball for $\|\cdot\|_{\infty}$). Therefore, interpolation conditions for indicator function gives us the equivalent result for support functions for free. Indeed, the support function for the closed convex set $Q \subseteq \mathbb{E}$ is defined as

$$\sigma_Q(s) = \sup_{x \in Q} \langle s, x \rangle = \sup_{x \in \mathbb{E}} \langle s, x \rangle - I_Q(x).$$

Hence, requiring a set $S = \{(x_i, g_i, f_i)\}_{i \in I}$ to be \mathcal{I}_M^* (resp. \mathcal{I}'_M^*)-interpolable is equivalent to require the set $\tilde{S} = \{(g_i, x_i, \langle x_i, g_i \rangle - f_i)\}_{i \in I}$ to be \mathcal{I}_M (or \mathcal{I}'_M)-interpolable.

COROLLARY 23. *The set $\{(x_i, g_i, f_i)\}_{i \in I}$ is \mathcal{I}_M^* (support with M -bounded subgradients) (resp. \mathcal{I}_M^*)-interpolable if and only if the following inequalities hold $\forall i, j \in I$:*

$$\begin{aligned} \langle g_i, x_i \rangle - f_i &= 0, \\ \langle g_i - g_j, x_j \rangle &\leq 0, \\ \|g_i\|_{\mathbb{E}^*} &\leq M, \quad (\text{resp. } \|g_i - g_j\|_{\mathbb{E}^*} \leq M). \end{aligned}$$

3.4. Smooth non-convex interpolation. In this short section, we show how to extend convex interpolation to smooth non-linear interpolation. The underlying conditions are also linearly Gram-representable, and also allow obtaining tight versions of (f-PEP) for non-convex programming.

DEFINITION 24. *Let $L \in \mathbb{R}^+ \cup \{\infty\}$, a differentiable function $f : \mathbb{E} \rightarrow \mathbb{R} \cup \{\infty\}$ is L -smooth (notation $f \in \mathcal{F}_{-L,L}(\mathbb{E})$) if it satisfies the following condition $\forall x, y \in \mathbb{E}$:*

$$|f(x) + \langle \nabla f(x), y - x \rangle - f(y)| \leq \frac{L}{2} \|x - y\|_{\mathbb{E}}^2.$$

The following lemma allows obtaining simple interpolation conditions from the smooth convex case.

LEMMA 25. *Let $L \in \mathbb{R}^+$, and $f : \mathbb{E} \rightarrow \mathbb{R} \cup \{\infty\}$. We have $f \in \mathcal{F}_{-L,L}(\mathbb{E}) \Leftrightarrow f + \frac{L}{2} \|x\|_{\mathbb{E}}^2 \in \mathcal{F}_{0,2L}(\mathbb{E})$.*

Proof. Let $f : \mathbb{E} \rightarrow \mathbb{R} \cup \{\infty\}$ and define $h(x) = f(x) + \frac{L}{2} \|x\|_{\mathbb{E}}^2$. We have that $\nabla h(x) = \nabla f(x) + LBx$, and $\forall x, y \in \mathbb{E}$:

$$\begin{aligned} f(x) + \langle \nabla f(x), y - x \rangle - f(y) &\leq \frac{L}{2} \|x - y\|_{\mathbb{E}}^2 \Leftrightarrow h(y) \geq h(x) + \langle \nabla h(x), y - x \rangle, \\ -f(x) - \langle \nabla f(x), y - x \rangle + f(y) &\leq \frac{L}{2} \|x - y\|_{\mathbb{E}}^2 \Leftrightarrow h(y) \leq h(x) + \langle \nabla h(x), y - x \rangle \\ &\quad + L \|x - y\|_{\mathbb{E}}^2, \end{aligned}$$

where the equivalences are obtained by expressing f and ∇f in terms of h and ∇h (or reciprocally), which proves our statement. \square

From Lemma 25, it is now straightforward to establish the desired interpolation conditions.

THEOREM 26. *Let $L \in \mathbb{R}^{++}$, the set $\{(x_i, g_i, f_i)\}_{i \in I}$ is $\mathcal{F}_{-L,L}$ (L -smooth) - interpolable if and only if the following inequality holds $\forall i, j \in I$:*

$$f_i \geq f_j - \frac{L}{4} \|x_i - x_j\|_{\mathbb{E}}^2 + \frac{1}{2} \langle g_i + g_j, x_j - x_i \rangle + \frac{1}{4L} \|g_i - g_j\|_{\mathbb{E}^*}^2.$$

Proof. As L is positive and finite, it follows from the equivalence of $\mathcal{F}_{-L,L}$ -interpolability of $\{(x_i, g_i, f_i)\}_{i \in I}$ and the $\mathcal{F}_{0,2L}$ -interpolability of $\{(x_i, g_i + LBx_i, f_i + \frac{L}{2} \|x_i\|_{\mathbb{E}}^2)\}_{i \in I}$. \square

4. Algorithm analysis. In this section, we analytically and numerically study different algorithms for solving variants of (CM), and compare our results with stan-

standard guarantees from the literature¹². We begin with an analytical study of a proximal point algorithm (Section 4.1). This is followed by a comparison between several standard variants of fast proximal gradient methods (Section 4.2) using the PEP approach. On the way, we propose an extension to the optimized gradient method (OGM) proposed by Kim and Fessler [22]. Finally, we conclude by applying the approach on a conditional gradient (Section 4.4) and on two alternate projections schemes (Section 4.5). Those choices illustrate the applicability of the approach for studying a large variety of methods and performance measures.

4.1. A proximal point algorithm. Consider a simple model with only one convex (possibly non-smooth) term in the objective function,

$$\min_{x \in \mathbb{E}} F(x),$$

with $F \in \mathcal{F}_{0,\infty}(\mathbb{E})$. In this first example, we assume that the proximal operation is easy to compute for F :

$$x_{k+1} = \text{prox}_{\alpha_{k+1}F}(x_k) = \underset{x \in \mathbb{E}}{\text{argmin}} \left\{ \alpha_{k+1}F(x) + \frac{1}{2}\|x_k - x\|_{\mathbb{E}}^2 \right\}.$$

That is, the iterations can be written in the form of an implicit method $x_{k+1} = x_k - \alpha_{k+1}B^{-1}g_{k+1}$, for some $g_{k+1} \in \partial F(x_{k+1})$ and are therefore a particular case of (FSLFOM).

For recent overviews and motivations concerning proximal algorithms, we refer the reader to the work of Combettes and Pesquet¹³ [7], and to the review works of Bertsekas [3] and Parikh and Boyd [35]. For historical point of view on those methods, we refer to the pioneer works of Moreau [25], Rockafellar [37] and the analysis of Guler [17].

Proximal Point Algorithm (PPA)

Input: $F \in \mathcal{F}_{0,\infty}(\mathbb{E})$, $x_0 \in \mathbb{E}$. Parameters: $\{\alpha_k\}_k$ with $\alpha_k \geq 0$.

For $k = 1 : N$

$$x_k = \text{prox}_{\alpha_k F}(x_{k-1})$$

4.1.1. Convergence of PPA in function and gradient values. The standard convergence result for the proximal point algorithm is provided by Guler in [17, Theorem 2.1] :

$$F(x_N) - F_* \leq \frac{R^2}{2 \sum_{k=1}^N \alpha_k},$$

for any initial condition x_0 satisfying $\|x_0 - x_*\|_{\mathbb{E}} \leq R$. We improve this bound by a factor 2 using the PEP approach.

THEOREM 27. *Let $\{\alpha_k\}_k$ be a sequence of positive step sizes and x_0 some initial iterate satisfying $\|x_0 - x_*\|_{\mathbb{E}} \leq R$ for some optimal point x_* . Any sequence $\{x_k\}_k$*

¹²Note that most of the literature results are presented for B being the identity operator (and hence $\mathbb{E} = \mathbb{E}^*$). We will nevertheless compare our slightly more general results with the standard bounds from the literature (thus even when they are officially valid only for B being the identity) — we recall that our results are valid for general self-adjoint positive definite linear operator $B : \mathbb{E} \rightarrow \mathbb{E}^*$ (see Remark 11).

¹³This work among others features a large list of known proximal operators.

generated by the proximal point algorithm with step sizes $\{\alpha_k\}_k$ on a function $F \in \mathcal{F}_{0,\infty}(\mathbb{E})$ satisfies

$$F(x_N) - F_* \leq \frac{R^2}{4 \sum_{k=1}^N \alpha_k}.$$

In addition, this bound is tight, and it is attained on the l_1 -shaped one-dimensional function ($\dim \mathbb{E} = \dim \mathbb{E}^* = 1$) $F(x) = \frac{\sqrt{BR}|x|}{2 \sum_{k=1}^N \alpha_k} = \frac{R\|x\|_{\mathbb{E}}}{2 \sum_{k=1}^N \alpha_k}$ with $Bx_0^2 = R^2$.

Proof. The proof relies on finding a primal and dual form of (f-PEP) for the proximal point algorithm (details can be found in Appendix A). A dual solution allows to obtain the upper bound part. \square

Considering another convergence measure, the exact same idea allows to obtain strong numerical evidence for the following conjecture.

CONJECTURE 28. *Let $\{\alpha_k\}_k$ be a sequence of positive step sizes and x_0 some initial iterate satisfying $\|x_0 - x_*\|_{\mathbb{E}} \leq R$ for some optimal point x_* . For any sequence $\{x_k\}_k$ generated by the proximal point algorithm with step sizes $\{\alpha_k\}_k$ on a function $F \in \mathcal{F}_{0,\infty}(\mathbb{E})$, there exists a subgradient $g_N \in \partial F(x_N)$ such that*

$$\|g_N\|_{\mathbb{E}^*} \leq \frac{R}{\sum_{k=1}^N \alpha_k}.$$

In particular, the choice $g_N = \frac{Bx_{N-1} - Bx_N}{\alpha_N}$ is a subgradient satisfying the inequality.

Observe that this bound cannot be improved, as it is attained on the (one-dimensional) l_1 -shaped function $F(x) = \frac{\sqrt{BR}|x|}{\sum_{k=1}^N \alpha_k}$ with $Bx_0^2 = R^2$. The particular choice of subgradient suggested in the theorem corresponds to the subgradient appearing in the proximal operation when written as an implicit subgradient step.

This sort of convergence results in terms of residual (sub)gradient norm is particularly interesting when considering dual methods. In that case, the dual residual gradient norm corresponds to the primal distance to feasibility (see e.g., [9]).

4.2. Fast gradient methods. In this section, we consider the two-terms composite objective function

$$(9) \quad \min_{x \in \mathbb{E}} \left\{ F(x) \equiv F^{(1)}(x) + F^{(2)}(x) \right\},$$

with $F^{(1)} \in \mathcal{F}_{0,L}(\mathbb{E})$ (smooth convex function) and $F \in \mathcal{F}_{0,\infty}(\mathbb{E})$ (non-smooth convex function). We assume that gradients are easy to compute for $F^{(1)}$, and that the proximal operation is easy to compute for $F^{(2)}$:

$$\text{prox}_{\alpha F^{(2)}}(x) = \operatorname{argmin}_{y \in \mathbb{E}} \left\{ \alpha F(y) + \frac{1}{2} \|x - y\|_{\mathbb{E}}^2 \right\}.$$

In order to approximatively solve (9), it is common to use different variants of fast proximal gradient methods (FPGM). We numerically investigate the worst-case guarantees of two variants using different step sizes policies, and propose new variants with slightly better worst-case behaviors. Also, we illustrate differences in the worst-case performances obtained in the cases where $F^{(2)} = 0$ (unconstrained smooth convex minimization), $F^{(2)} \in \mathcal{I}_{\infty}(\mathbb{E})$ (constrained smooth convex minimization) or $F^{(2)} \in \mathcal{F}_{0,\infty}(\mathbb{E})$ (regularized smooth convex minimization).

In the following, we call FPGM1 the standard fast proximal gradient method (FISTA [2]), FPGM2 a variant with slightly better guarantees, and POGM a proximal version of the optimized gradient method [22]. FPGM2 and POGM illustrate how performance estimation problems can be used in the development of new optimization algorithms ; their study in this paper remains however entirely numerical.

4.2.1. Standard Fast Proximal Gradient Methods (FPMG1). The first variants of accelerated proximal methods we are considering use a standard proximal step after an explicit gradient step for generating the so-called *primary sequence* $\{y_k\}_k$.

Fast Proximal Gradient Method (FPGM1)

Input: $F^{(1)} \in \mathcal{F}_{0,L}(\mathbb{E})$, $F^{(2)} \in \mathcal{F}_{0,\infty}(\mathbb{E})$ $x_0 \in \mathbb{E}$, $y_0 = x_0$.

For $k = 1 : N$

$$y_k = \text{prox}_{F^{(2)}/L} \left(x_{k-1} - \frac{1}{L} B^{-1} \nabla F^{(1)}(x_{k-1}) \right)$$

$$x_k = y_k + \alpha_k (y_k - y_{k-1})$$

In this algorithm, we refer to α_k as *inertial parameters*. We use two standard variants: $\alpha_k^{(a)} = \frac{k-1}{k+2}$ — among others proposed in [42, 44] — and $\alpha_k^{(b)} = \frac{\theta_{k-1}-1}{\theta_k}$ (with $\theta_k = \frac{1+\sqrt{4\theta_{k-1}^2+1}}{2}$ and $\theta_0 = 1$) — see [2, 28, 44]. For both variants, the standard convergence result is (see e.g., [2, 42]):

$$(10) \quad F(y_N) - F_* \leq \frac{2LR^2}{(N+1)^2},$$

for any initial iterate x_0 such that $\|x_0 - x_*\|_{\mathbb{E}}^2 \leq R^2$. We numerically compare those two variants of FPGM1 using (f-PEP) on Fig. 1. After 100 iterations, both inertial parameter policies perform about the same way (parameters $\alpha_k^{(b)}$ performs only about 2% better than $\alpha_k^{(a)}$ in terms of worst-case performances). We also observe that the behavior of both variants of FPGM1 is well captured by the standard guarantee (10).

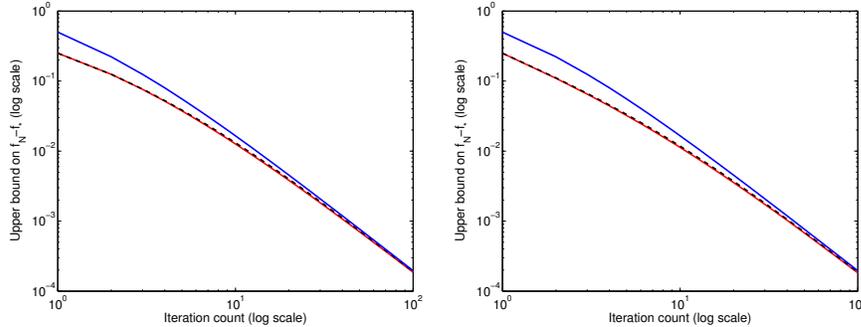


FIG. 1. Comparison of the worst-case convergence speed of the different variants of FPGM1 (left) and FPGM2 (right) for $N \in \{1, \dots, 100\}$, $L = 1$ and $R = 1$. The curves respectively corresponds to the different inertial coefficient, namely $\alpha_k^{(a)}$ (dashed, black) and $\alpha_k^{(b)}$ (red), and the standard guarantee (10) (blue).

4.2.2. New Fast Proximal Gradient Methods (FPGM2). Secondary sequences $\{x_k\}$ are usually converging slightly faster than primary sequences $\{y_k\}$ in

the unconstrained case ($F^{(2)} = 0$), as observed in [22, 43]. However, some issues may arise with the secondary sequences of FPGM1 when applied to constrained or proximal problems: iterates may in some cases become infeasible, or the objective may become unbounded (see Table 1 below). We therefore propose new variants of fast proximal gradient methods that do not suffer from these drawback, called FPGM2 (with two different step size policies). Part of the underlying motivation behind FPGM2 is also the ability to generalize it later to the optimized gradient method.

REMARK 29. *The design of FPGM2 is based on two ideas: on the one hand, it should be equivalent to the standard fast gradient method in the case of smooth unconstrained convex minimization, and on the other hand, it should not move after two consecutive iterates have reached the same optimal point for (9) (i.e., $x_{k-1} = x_k = x_*$ implies $x_{k+1} = x_*$).*

Fast Proximal Gradient Method 2 (FPGM2)

Input: $F^{(1)} \in \mathcal{F}_{0,L}(\mathbb{E})$, $F^{(2)} \in \mathcal{F}_{0,\infty}(\mathbb{E})$ $x_0 \in \mathbb{E}$, $z_0 = y_0 = x_0$.

For $k = 1 : N$

$$\begin{aligned} y_k &= x_{k-1} - \frac{1}{L} B^{-1} \nabla F^{(1)}(x_{k-1}) \\ z_k &= y_k + \alpha_k (y_k - y_{k-1}) + \frac{\alpha_k}{L \gamma_{k-1}} (z_{k-1} - x_{k-1}) \\ x_k &= \text{prox}_{\gamma_k F^{(2)}}(z_k) \end{aligned}$$

In this algorithm, we use the coefficients $\gamma_k = \frac{\alpha_k + 1}{L}$. Note that we introduced two intermediate sequences: on the one hand sequence $\{\gamma_k\}_k$, corresponding to the step sizes to be taken by the proximal steps, and on the other hand sequence $\{z_k\}_k$, which allows keeping track of the subgradient used in the proximal steps (note that $\frac{1}{\gamma_k}(z_k - x_k)$ corresponds to the subgradient used in the proximal step from z_k to x_k). Even if FPGM2 may look more intricate than the classical FPGM1, it is in fact simpler, as it involves only one sequence on which both implicit (proximal) and explicit (gradient) steps are being taken. Indeed, explicit steps are taken using gradient values of $F^{(1)}$ at x_k , and subgradients used in the proximal steps are subgradients of $F^{(2)}$ also at x_k . This can be seen by rewriting the iterations of FPGM2 using the secondary sequence $\{x_k\}_k$ only

$$\begin{aligned} x_{k+1} &= x_k + \alpha_{k+1}(x_k - x_{k-1}) \\ &+ \frac{\alpha_{k+1}}{L} B^{-1} \nabla F^{(1)}(x_{k-1}) - \frac{1}{L} B^{-1} \nabla F^{(1)}(x_k) - \frac{\alpha_{k+1}}{L} B^{-1} \nabla F^{(1)}(x_k) \\ &+ \frac{\alpha_{k+1}}{L} B^{-1} \tilde{\nabla} F^{(2)}(x_k) - \frac{1}{L} B^{-1} \tilde{\nabla} F^{(2)}(x_{k+1}) - \frac{\alpha_{k+1}}{L} B^{-1} \tilde{\nabla} F^{(2)}(x_{k+1}), \end{aligned}$$

with $\tilde{\nabla} F^{(2)}(x_k)$ the subgradient of $F^{(2)}$ used in the proximal operation generating x_k .

Comparing the different variants of FPGM2 on Fig. 1 leads to the same conclusion as for FPGM1: inertial parameters $\alpha^{(b)}$ perform slightly better than $\alpha^{(a)}$.

In Table 1, we report the different worst-case performances guarantees obtained numerically for FPGM1 (for both sequences) and FPGM2 (for the better secondary sequence only). We consider three situations: $F^{(2)} = 0$ (unconstrained smooth convex minimization), $F^{(2)} \in \mathcal{I}_\infty(\mathbb{E})$ (constrained smooth convex minimization with projected methods) and $F^{(2)} \in \mathcal{F}_{0,\infty}(\mathbb{E})$ (regularized smooth convex minimization with proximal methods).

Type	$F(y_N) - F_*$ (FPGM1)	$F(x_N) - F_*$ (FPGM1)	$F(x_N) - F_*$ (FPGM2)
Unconstrained ($F^{(2)} = 0$)	$\frac{LR^2}{2} \frac{4}{N^2+5N+6}$	$\frac{LR^2}{2} \frac{4}{N^2+7N+4}$	$\frac{LR^2}{2} \frac{4}{N^2+7N+4}$
Constrained ($F^{(2)} \in \mathcal{I}_\infty$)	$\frac{LR^2}{2} \frac{4}{N^2+5N+2}$	Infeasible	$\frac{LR^2}{2} \frac{4}{N^2+7N}$
Non-smooth ($F^{(2)} \in \mathcal{F}_{0,\infty}$)	$\frac{LR^2}{2} \frac{4}{N^2+5N+2}$	Unbounded	$\frac{LR^2}{2} \frac{4}{N^2+7N}$

TABLE 1

Worst-case obtained for FPGM1 and FPGM2 with inertial coefficient $\alpha_k = \frac{k-1}{k+2}$ and $N \geq 1$.

Convergence results reported in the table correspond to properly identified functions (i.e. they are rigorous lower bounds). After solving the corresponding PEPs numerically (for $L = R = 1$ and $1 \leq N \leq 100$), we conjecture them to be equal to the exact worst-case guarantees.

We observe that the worst-case guarantees for FPGM2 are slightly better than for FPGM1. Guarantees for the unconstrained case are slightly better than those for the constrained and proximal cases, which are equal. Note that the secondary sequence of FPGM1 is not guaranteed to be feasible in the constrained case, and that the corresponding objective value may be unbounded in the proximal case (for any $N \geq 1$).

The worst-case functions identified numerically for the unconstrained case are Huber-shaped functions [43]. In the constrained case, we identified one-dimensional linear optimization problems of the form $\min_{x \geq 0} cx$ as worst-cases, where c is a constant defined by

$$c = \frac{\sqrt{BR}}{2 \sum_{i=0}^{N-1} [\alpha_N]_i}.$$

Finally, for the proximal case, our worst-case has function $F^{(1)}(x) = cx$ with the same c as above, and function $F^{(2)}(x)$ may be chosen equal to zero for $x \geq 0$ and to sx for $x < 0$, for any negative value of the slope $s < 0$.

4.3. A proximal optimized gradient method. In this section, we consider again the regularized smooth convex minimization problem (9). In particular, we are concerned with the possibility of obtaining optimized methods for handling this sort of problems (i.e., methods whose worst-case performances are minimized).

The idea is to extend the optimized gradient method (OGM) developed by Kim and Fessler in [22], which is originally tailored for smooth unconstrained minimization ($F^{(2)} = 0$). In the unconstrained smooth minimization setting, this first-order method was recently shown in [12] to have the best achievable worst-case guarantee for the criterion $F_N - F_*$.

The method we propose has been obtained by combining the ideas obtained from the original OGM [22] and the non-standard placement of the proximal operator used for speeding up the convergence of fast proximal gradient methods (FPGM2). It was designed using the same two principles as FPGM2 (see Remark 29): on the one hand, it is equivalent to OGM when applied to smooth unconstrained convex minimization problems, and on the other hand, it remains at an optimal point when it reaches one.

Proximal Optimized Gradient Method (POGM)

Input: $F^{(1)} \in \mathcal{F}_{0,L}(\mathbb{E})$, $F^{(2)} \in \mathcal{F}_{0,\infty}(\mathbb{E})$, $x_0 \in \mathbb{E}$, $y_0 = x_0$, $\theta_0 = 1$.

For $k = 1 : N$

$$y_k = x_{k-1} - \frac{1}{L} B^{-1} \nabla F^{(1)}(x_{k-1})$$

$$z_k = y_k + \frac{\theta_{k-1} - 1}{\theta_k} (y_k - y_{k-1}) + \frac{\theta_{k-1}}{\theta_k} (y_k - x_{k-1}) + \frac{\theta_{k-1} - 1}{L\gamma_{k-1}\theta_k} (z_{k-1} - x_{k-1})$$

$$x_k = \text{prox}_{\gamma_k F^{(2)}}(z_k)$$

In this algorithm, we use the sequence $\gamma_k = \frac{1}{L} \frac{2\theta_{k-1} + \theta_k - 1}{\theta_k}$ and the inertial coefficients proposed in [22]:

$$\theta_k = \begin{cases} \frac{1 + \sqrt{4\theta_{k-1}^2 + 1}}{2}, & i \leq N - 1 \\ \frac{1 + \sqrt{8\theta_{k-1}^2 + 1}}{2}, & i = N \end{cases}$$

Simply trying to generalize OGM using the standard proximal step on the primary sequence $\{y_i\}$ (as for FPGM1) does not lead to a converging algorithm. We have numerical evidence, i.e. worst-case functions showing that this candidate algorithm does not see its worst-case improving after each iteration: in other words its worst-case rate is not converging to zero). Therefore we have to introduce the same idea used in FPGM2 concerning the place of the proximal operator.

We compare POGM to FPGM1 and FPGM2 with inertia $\alpha_k^{(b)}$ on Fig. 2. We obtain worst-case performances about twice better for POGM when compared to both FPGM1 and FPGM2. Of course, POGM suffers from the drawback of requiring the knowledge of the number of iterations in advance (this is because the rule to compute the last coefficient θ_N differs from the rule to compute all the previous ones). This practical disadvantage is not easily solved: if the last θ_N is updated with the same rule as all the previous θ_k , performance is degraded by a non-negligible factor, rendering it even slower than FPGM (note that this is already the case for smooth unconstrained minimization [21]).

4.4. A conditional gradient method. Consider the constrained smooth convex optimization problem

$$\min_{x \in Q} F(x),$$

with $F \in \mathcal{F}_{0,L}(\mathbb{E})$ and $Q \subset \mathbb{E}$ a bounded and closed convex set. In that setting, there exists different ways for treating the constraint set Q . In the previous section, we proposed to use fast gradient methods, which require the ability of projecting onto the closed convex set Q . In this section, we rather consider the standard conditional gradient method (also sometimes referred to as the Frank-Wolfe method), which originates from [15]. This algorithm has the advantage of not requiring to perform projections onto Q , but rather to perform linear optimization on this set (which is typically easier when Q is a polyhedral set).

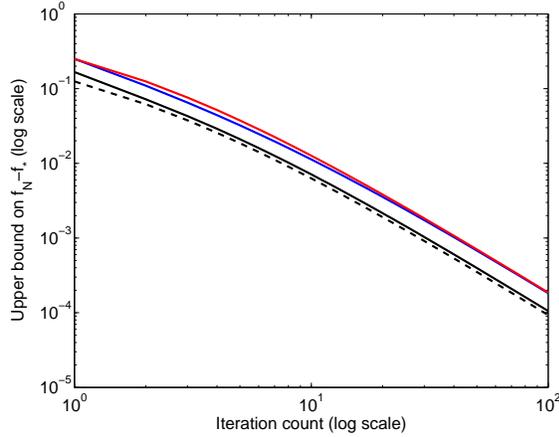


FIG. 2. Comparison between the worst-case performances of FPGM1 (with inertia $\alpha_k^{(b)}$) (red), FPGM2 (with inertia $\alpha_k^{(b)}$) (blue), POGM (black) and OGM (dashed, black) for $N \in \{1, \dots, 100\}$, $L = 1$ and $R = 1$. The worst-case performances of POGM are about twice better than the worst-case performance of FPGM between 1 and 100 iterations. Also, we observe that OGM [22] (equivalent to POGM with $F^{(2)} = 0$) behaves approximately 12% better than POGM in the worst-case.

Conditional Gradient Method (CGM)

Input: $F \in \mathcal{F}_{0,L}(\mathbb{E})$, closed convex $Q \subset \mathbb{E}$ with $\|x - y\|_{\mathbb{E}} \leq D \forall x, y \in Q$, $x_0 \in Q$.

For $k = 1 : N$

$$y_k = \operatorname{argmin}_{y \in Q} \{\langle \nabla F(x_{k-1}), y - x_{k-1} \rangle\}$$

$$\lambda_k = \frac{2}{1+k}$$

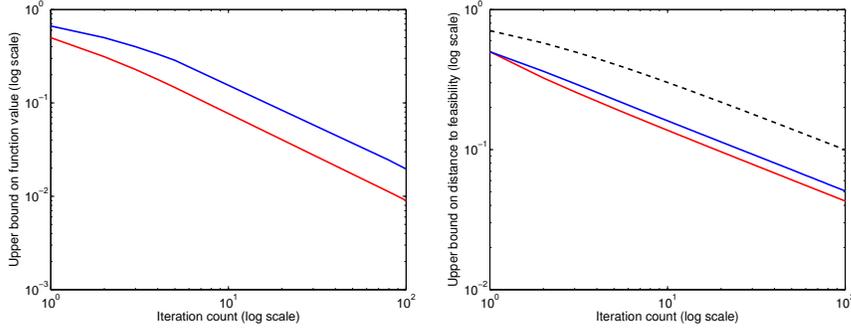
$$x_k = (1 - \lambda_k)x_{k-1} + \lambda_k y_k$$

The standard global convergence guarantee for this method (see e.g., [19, Theorem 1]) is

$$(11) \quad F(x_N) - F_* \leq \frac{2LD^2}{N+2},$$

which we compare with the exact bound provided by PEP on Fig. 3(a). As illustrated in Section 2.3, this algorithm fits into the (FSLFOM) format. The numerical guarantees we obtained by solving the performance estimation problem are between two and three times better than the standard guarantee, depending on the number of iterations.

4.5. Alternate projection and Dykstra methods. In this section, we numerically investigate the difference between the worst-case behaviors of the standard alternate projection method (APM) for finding a point in the intersection of two convex sets, and the Dykstra [5] method (DAPM) for finding the closest point in the intersection of two convex sets. APM is a particular instance of subgradient-type



(a) Worst-case performance of CGM (red) and its theoretical guarantee (11) (blue) for $N \in \{1, \dots, 100\}$, $L = 1$ and $D = 1$.
 (b) Worst-case performance of APM (red), DAPM (blue) and lower bound $\frac{MR}{\sqrt{N+1}}$ for subgradient methods (dashed, black) for $N \in \{1, \dots, 100\}$ and $R = 1$ ($M = 1$ by definition of the objective function (12)).

FIG. 3. Numerical analysis of a conditional gradient method (left) and of two variants of alternate projections algorithms (right).

descent¹⁴ applied to the problem

$$(12) \quad \min_{x \in \mathbb{E}} \{f(x) = \max_i \|x - \Pi_{Q_i}(x)\|_{\mathbb{E}}\},$$

whose objective function is convex and non-smooth (with Lipschitz constant $M = 1$). Therefore, its expected global convergence rate is $\mathcal{O}(\frac{1}{\sqrt{N}})$ (see [14, Theorem A.1]). We compare below the convergence of both APM and DAPM with the standard lower bound for subgradient schemes $\frac{MR}{\sqrt{N+1}}$ as a reference.

<p>Alternate Projection Method (APM) Input: $x_0 \in \mathbb{E}$, convex sets $Q_1, Q_2 \subseteq \mathbb{E}$, $\ x_0 - x_*\ _{\mathbb{E}} \leq R$, for some $x_* \in Q_1 \cap Q_2$. For $k = 1 : N$</p> $x_k = \Pi_{Q_2}(\Pi_{Q_1}(x_{k-1}))$
<p>Dykstra Alternate Projection Method (DAPM) Input: $x_0 \in \mathbb{E}$, convex sets $Q_1, Q_2 \subseteq \mathbb{E}$, $\ x_0 - x_*\ _{\mathbb{E}} \leq R$, for some $x_* \in Q_1 \cap Q_2$. Initialize $p_0 = q_0 = 0$. For $k = 0 : N - 1$</p> $y_k = \Pi_{Q_1}(x_k + p_k)$ $p_{k+1} = x_k + p_k - y_k$ $x_{k+1} = \Pi_{Q_2}(y_k + q_k)$ $q_{k+1} = y_k + q_k - x_{k+1}$

¹⁴It can be shown that $\frac{x - \Pi_{Q_k}(x)}{\|x - \Pi_{Q_k}(x)\|}$ is a subgradient of the function $f(x)$ (at x such that $f(x) = \|x - \Pi_{Q_k}(x)\|$). Therefore, in the case of two sets Q_1, Q_2 , and assuming that x is feasible for one of the two sets (say, Q_1), a projection on the other one corresponds to a subgradient step on f with step size $\|x - \Pi_{Q_2}(x)\|$. Hence, APM is an instance of subgradient method for $k > 1$ (when x_k is feasible for one of the two sets).

The optimality measure used is $\min_{x \in Q_1} \|x - x_N\|_{\mathbb{E}} = \|\Pi_{Q_1}(x_N) - x_N\|_{\mathbb{E}}$ (note that $x_N \in Q_2$). We do not give further details the corresponding performance estimation problem here, as it is very similar to the previous sections. The results for APM and DAPM are shown on Fig. 3(b), where the (expected) convergence in $\mathcal{O}(\frac{1}{\sqrt{N}})$ is clearly obtained. Interestingly, DAPM converges slightly slower than APM (more precisely, DAPM has a worst-case about 18% higher than APM), which is therefore more advisable in terms of worst-case performances for finding a point in the intersection of two convex sets when no additional structure is assumed. In addition, note that both APM and DAPM have a worst-case which is about twice better than the standard lower bound for explicit non-smooth schemes.

5. Conclusion. In this work, we presented a performance estimation approach for analyzing first-order algorithms for composite optimization problems. The results of [43] were largely extended to handle both larger classes of (composite) objective functions and larger classes of first-order algorithm (also in a more general setting for handling pairs of conjugate norms). Our contribution is essentially threefold: first, we developed specific interpolation conditions for different classes of convex and non-convex functions; then, we exploited those interpolation conditions to formulate the exact worst-case problem for fixed-step linear first-order methods and finally we applied that methodology to provide tight analyses for different first-order methods. Among others, we presented a new analytical guarantee for the proximal point algorithm that is twice better than previously known, and improved the standard worst-case guarantee for the conditional gradient method by more than a factor of two. On the way, we also proposed an extension of the optimized gradient method proposed by Kim and Fessler [22] that incorporates a projection or a proximal operator.

As further research, we believe this methodology should be applied to refine analyses of methods fitting in the context of fixed-step linear first-order methods, and possibly extended to handle dynamic step size rules. To this end, a possibility is to explore convex relaxations of the resulting possibly non-convex performance estimations problems. As an example, we believe it would be a real asset to be able to analyze algorithms involving line-search procedures such as backtracking or Armijo-Wolfe procedures. Moreover, it seems to us that the performance estimation approach could be used to refine the analyses of randomized coordinate descent-type algorithms [31]. Performance estimation problems also opened the door for looking towards optimized methods, as proposed by Kim and Fessler [22] for unconstrained smooth convex minimization. Such an optimized method in the proximal or conditional settings would be of great interest.

Finally, algorithmic analyses using performance estimation problems is intrinsically limited by our ability to solve semidefinite problems, both numerically (when the number of iterations is large) or analytically (to obtain results valid for any number of iterations). Therefore, any idea leading to (convex) programs that are easier to solve while maintaining reasonable guarantees would be very advantageous.

Software. MATLAB implementations of the performance estimation approach for different variants of gradient methods are available online. They can be downloaded from <http://perso.uclouvain.be/adrien.taylor>.

REFERENCES

- [1] H. H. BAUSCHKE AND P. COMBETTES, *Convex analysis and monotone operator theory in Hilbert spaces*, Springer, 2011.

- [2] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 183–202.
- [3] D. P. BERTSEKAS, *Incremental gradient, subgradient, and proximal methods for convex optimization: A survey*, Optimization for Machine Learning, (2010), pp. 1–38.
- [4] D. P. BERTSEKAS, *Convex Optimization Algorithms*, Athena Scientific, 2015.
- [5] J. P. BOYLE AND R. L. DYKSTRA, *A method for finding projections onto the intersection of convex sets in hilbert spaces*, in Advances in order restricted statistical inference, Springer, 1986, pp. 28–47.
- [6] S. BURER AND R. MONTEIRO, *A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization*, Mathematical Programming, 95 (2003), pp. 329–357.
- [7] P. L. COMBETTES AND J.-C. PESQUET, *Proximal splitting methods in signal processing*, in Fixed-point algorithms for inverse problems in science and engineering, Springer, 2011, pp. 185–212.
- [8] A. D’ASPREMONT, *Smooth optimization with approximate gradient*, SIAM Journal on Optimization, 19 (2008), pp. 1171–1183.
- [9] O. DEVOLDER, F. GLINEUR, AND Y. NESTEROV, *Double smoothing technique for large-scale linearly constrained convex optimization*, SIAM Journal on Optimization, 22 (2012), pp. 702–727.
- [10] O. DEVOLDER, F. GLINEUR, AND Y. NESTEROV, *First-order methods of smooth convex optimization with inexact oracle*, Mathematical Programming, 146 (2014), pp. 37–75.
- [11] Y. DRORI, *Contributions to the Complexity Analysis of Optimization Algorithms*, PhD thesis, Tel-Aviv University, 2014.
- [12] Y. DRORI, *The exact information-based complexity of smooth convex minimization*, arXiv preprint arXiv:1606.01424, (2016).
- [13] Y. DRORI AND M. TEOULLE, *Performance of first-order methods for smooth convex minimization: a novel approach*, Mathematical Programming, 145 (2014), pp. 451–482.
- [14] Y. DRORI AND M. TEOULLE, *An optimal variant of kelley’s cutting-plane method*, Mathematical Programming, (2016).
- [15] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, Naval research logistics quarterly, 3 (1956), pp. 95–110.
- [16] M. GOEMANS AND D. WILLIAMSON, *Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming*, Journal of the ACM (JACM), 42 (1995), pp. 1115–1145.
- [17] O. GÜLER, *On the convergence of the proximal point algorithm for convex minimization*, SIAM Journal on Control and Optimization, 29 (1991), pp. 403–419.
- [18] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms*, Springer Verlag, Heidelberg, 1996. Two volumes - 2nd printing.
- [19] M. JAGGI, *Revisiting Frank-Wolfe: Projection-free sparse convex optimization*, in Proceedings of the 30th International Conference on Machine Learning (ICML-13), 2013, pp. 427–435.
- [20] M. JOURNÉE, F. BACH, P.-A. ABSIL, AND R. SEPULCHRE, *Low-rank optimization on the cone of positive semidefinite matrices*, SIAM Journal on Optimization, 20 (2010), pp. 2327–2351.
- [21] D. KIM AND J. A. FESSLER, *On the convergence analysis of the optimized gradient methods*, Journal of Optimization Theory and Applications, (2016).
- [22] D. KIM AND J. A. FESSLER, *Optimized first-order methods for smooth convex minimization*, Mathematical Programming, (2016).
- [23] L. LESSARD, B. RECHT, AND A. PACKARD, *Analysis and design of optimization algorithms via integral quadratic constraints*, SIAM Journal on Optimization, 26 (2016), pp. 57–95.
- [24] J. LÖFBERG, *YALMIP: A toolbox for modeling and optimization in MATLAB*, in Proceedings of the CACSD Conference, 2004.
- [25] J.-J. MOREAU, *Proximité et dualité dans un espace hilbertien*, Bulletin de la Société mathématique de France, 93 (1965), pp. 273–299.
- [26] A. MOSEK, *The MOSEK optimization software*, Online at <http://www.mosek.com>, 54 (2010).
- [27] A. NEMIROVSKY AND D. YUDIN, *Problem complexity and method efficiency in optimization.*, Wiley-Interscience, New York, (1983).
- [28] Y. NESTEROV, *A method of solving a convex programming problem with convergence rate $O(1/k^2)$* , Soviet Mathematics Doklady, 27 (1983), pp. 372–376.
- [29] Y. NESTEROV, *Introductory lectures on convex optimization: a basic course*, Applied optimization, Kluwer Academic Publ., 2004.
- [30] Y. NESTEROV, *Smooth minimization of non-smooth functions*, Mathematical programming, 103 (2005), pp. 127–152.
- [31] Y. NESTEROV, *Efficiency of coordinate descent methods on huge-scale optimization problems*, SIAM Journal on Optimization, 22 (2012), pp. 341–362.

- [32] Y. NESTEROV, *Gradient methods for minimizing composite functions*, Mathematical Programming, 140 (2013), pp. 125–161.
- [33] R. ORSI, U. HELMKKE, AND J. B. MOORE, *A newton-like method for solving rank constrained linear matrix inequalities*, Automatica, 42 (2006), pp. 1875–1882.
- [34] P. M. PARDALOS AND S. A. VAVASIS, *Quadratic programming with one negative eigenvalue is np-hard*, Journal of Global Optimization, 1 (1991), pp. 15–22.
- [35] N. PARIKH AND S. BOYD, *Proximal algorithms*, Foundations and Trends in optimization, 1 (2013), pp. 123–231.
- [36] B. T. POLYAK, *Introduction to optimization*, Optimization Software New York, 1987.
- [37] R. T. ROCKAFELLAR, *Monotone operators and the proximal point algorithm*, SIAM journal on control and optimization, 14 (1976), pp. 877–898.
- [38] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, 1996.
- [39] R. T. ROCKAFELLAR AND R.-B. WETS, *Variational Analysis*, Springer, 1998.
- [40] S. SAHNI, *Computationally related problems*, SIAM Journal on Computing, 3 (1974), pp. 262–279.
- [41] J. F. STURM, *Using SeDuMi 1.02, a MATLAB toolbox for optimization over symmetric cones*, Optimization Methods and Software, 11–12 (1999), pp. 625–653.
- [42] W. SU, S. BOYD, AND E. CANDÈS, *A differential equation for modeling nesterovs accelerated gradient method: Theory and insights*, in Advances in Neural Information Processing Systems, 2014, pp. 2510–2518.
- [43] A. B. TAYLOR, J. M. HENDRICKX, AND F. GLINEUR, *Smooth strongly convex interpolation and exact worst-case performance of first-order methods*, Mathematical Programming, (2016).
- [44] P. TSENG, *On accelerated proximal gradient methods for convex-concave optimization*, Submitted to SIAM Journal on Optimization, (2008).
- [45] L. VANDENBERGHE AND S. BOYD, *Semidefinite programming*, SIAM Review, 38 (1994), pp. 49–95.

Appendix A. Proof of Theorem 27. We start by proving the lower bound, and then we prove the matching upper bound.

Lower bound. Let us show that applying PPA to the one-dimensional function $F(x) = \frac{\sqrt{BR}|x|}{2 \sum_{k=1}^N \alpha_k}$ with $x_0 = -\frac{R}{\sqrt{B}}$ allows to achieve:

$$F(x_N) - F(x_*) = \frac{R^2}{4 \sum_{k=1}^N \alpha_k},$$

which shows that the bound from Theorem 27 is tight.

First, note that for $x \neq 0$, we have $\nabla F(x) = \text{sign}(x) \frac{\sqrt{BR}}{2 \sum_{k=1}^N \alpha_k}$. Hence,

$$x_N = x_0 + B^{-1} \sum_{k=1}^N \alpha_k \frac{\sqrt{BR}}{2 \sum_{k=1}^N \alpha_k} = -\frac{R}{2\sqrt{B}}.$$

Therefore, by noting that $x_* = 0$ and $F(x_*) = 0$, we have the desired result.

Upper bound. In order to express the corresponding PEP in the simplest form, we heavily rely on some straightforward simplifications of (SDP-PEP) (see Corollary 14 and Remark 2.4). Let us denote by P_N the matrix containing the information harvested after N iterations (we use the notation g_i for subgradients $g_i \in \partial F(x_i)$): $P_N = [g_1 \ g_2 \ \dots \ g_N \ Bx_0]$, and by G_N its corresponding Gram matrix (see Section 2.1). Also, we introduce the step size matrices $\underline{\alpha}_k$ for expressing all intermediate iterates x_i 's in terms of x_0 and the subgradient g_i 's, that is: $x_k = P_N \underline{\alpha}_k$ ($k = 0, \dots, N$).

This results in the following explicit expressions for $\underline{\alpha}_k$: $\underline{\alpha}_k = e_{N+1} - \sum_{i=1}^k \alpha_i e_i$, along with $\underline{\alpha}_0 = e_{N+1}$ and $\underline{\alpha}_* = 0$ (i.e., we assume without loss of generality $x_* = 0$), where we use the standard notation e_i for the unit vector having a single 1 as its i^{th} component — we also denote $e_* = 0$. In order to perform the worst-case analysis for PPA, we now formulate the performance estimation problem (**f-PEP**) as the following

SDP (simplified version of (SDP-PEP) where the x_k 's ($k = 1, \dots, N$) have been substituted using the form of the algorithm $x_k = x_{k-1} - \alpha_k B^{-1} g_k$):

(PPA-PEP)

$$\begin{aligned} \max_{G_N \in \mathbb{S}^{N+1}, f_1, \dots, f_N, f_* \in \mathbb{R}^N} f_N - f_*, \text{ s.t. } f_j - f_i + \text{Tr}(A_{ij} G_N) \leq 0, \quad i, j \in \{1, \dots, N, *\} \\ \|x_0 - x_*\|_{\mathbb{E}}^2 \leq R^2, \\ G_N \succeq 0, \end{aligned}$$

with $2A_{ij} = e_j(\underline{\alpha}_i - \underline{\alpha}_j)^\top + (\underline{\alpha}_i - \underline{\alpha}_j)e_j^\top$, the matrices coming from the non-smooth convex interpolation inequalities (see Condition (5)). In order to obtain an analytical upper bound for PPA, we consider the Lagrangian dual to (PPA-PEP), which is given by the following:

$$\begin{aligned} \text{(PPA-dPEP)} \quad \min_{\lambda_{ij} \geq 0, \tau \geq 0} \tau R^2 \text{ s.t. } e_N - \sum_i \sum_{j \neq i} (\lambda_{ij} - \lambda_{ji}) e_j = 0, \\ \sum_i \sum_{j \neq i} \lambda_{ij} A_{ij} + \tau \underline{\alpha}_0 \underline{\alpha}_0^\top \succeq 0 \end{aligned}$$

(where the constraint corresponding to f_* can be discarded since it is clear that letting $f_* = 0$ does not change the optimal solution of (PPA-PEP)). Note that the set of equality constraints can be assimilated to a set of *flow* constraints on a complete directed graph. That is, considering a graph where the optimum and each iterate correspond to nodes, each $0 \leq \lambda_{ij} \leq 1$ corresponds to the flow on the edge going from node j to node i (we choose this direction by convention). This flow constraint imposes that the outgoing flow equals the ingoing flow for every node, except at iterate N where the outgoing flow should be equal to 1 and at the optimum, where the incoming flow should be equal to 1. We show that the following choice is a feasible point of the dual (PPA-dPEP).

$$\begin{aligned} \lambda_{i,i+1} &= \frac{\sum_{k=1}^i \alpha_k}{2 \sum_{k=1}^N \alpha_k - \sum_{k=1}^i \alpha_k}, & i \in \{1, \dots, N-1\} \\ \lambda_{*,i} &= \frac{2\alpha_i \sum_{k=1}^N \alpha_k}{\left(2 \sum_{k=1}^N \alpha_k - \sum_{k=1}^i \alpha_k\right) \left(2 \sum_{k=1}^N \alpha_k - \sum_{k=1}^{i-1} \alpha_k\right)}, & i \in \{1, \dots, N\} \\ \tau &= \frac{1}{4 \sum_{k=1}^N \alpha_k}, \end{aligned}$$

and $\lambda_{ij} = 0$ otherwise. First, we clearly have $\lambda_{ij} \geq 0$ and some basic computations allow to verify that the equality constraints from (PPA-dPEP) are satisfied:

$$\lambda_{*,1} - \lambda_{1,2} = 0, \quad \lambda_{*,i} + \lambda_{i-1,i} - \lambda_{i,i+1} = 0 \quad (i \in \{2, \dots, N-1\}), \quad \lambda_{*,N} + \lambda_{N-1,N} = 1.$$

It remains to show that the corresponding dual matrix S is positive semidefinite.

$$\begin{aligned} 2S &= \sum_{i=1}^{N-1} 2\alpha_{i+1} \lambda_{i,i+1} e_{i+1} e_{i+1}^\top + 2\tau e_{N+1} e_{N+1}^\top \\ &\quad + \sum_{i=1}^N \lambda_{*,i} \left[e_i \left(-e_{N+1} + \sum_{k=1}^i \alpha_k e_k \right)^\top + \left(-e_{N+1} + \sum_{k=1}^i \alpha_k e_k \right) e_i^\top \right]. \end{aligned}$$

In order to reduce the number of indices to be used, we note shortly $\lambda_i = \lambda_{i,i+1}$ and $\mu_i = \lambda_{*,i}$. Then, using the equality constraints, we arrive at the following dual matrix:

$$2S = \begin{pmatrix} 2\alpha_1\lambda_1 & \alpha_1\mu_2 & \alpha_1\mu_3 & \dots & \alpha_1\mu_{N-1} & \alpha_1\mu_N & -\mu_1 \\ \alpha_1\mu_2 & 2\alpha_2\lambda_2 & \alpha_2\mu_3 & \dots & \alpha_2\mu_{N-1} & \alpha_2\mu_N & -\mu_2 \\ \alpha_1\mu_3 & \alpha_2\mu_3 & 2\alpha_3\lambda_3 & \dots & \alpha_3\mu_{N-1} & \alpha_3\mu_N & -\mu_3 \\ \vdots & & \ddots & \ddots & & \vdots & \vdots \\ \alpha_1\mu_{N-1} & \alpha_2\mu_{N-1} & \alpha_3\mu_{N-1} & \dots & 2\alpha_{N-1}\lambda_{N-1} & \alpha_{N-1}\mu_N & -\mu_{N-1} \\ \alpha_1\mu_N & \alpha_2\mu_N & \alpha_3\mu_N & \dots & \alpha_{N-1}\mu_N & 2\alpha_N & -\mu_N \\ -\mu_1 & -\mu_2 & -\mu_3 & \dots & -\mu_{N-1} & -\mu_N & 2\tau \end{pmatrix}.$$

In order to prove $S \succeq 0$, we first use a Schur complement and then show that the resulting matrix is diagonally dominant with positive diagonal elements. After the Schur complement, we obtain the matrix \tilde{S} :

$$\tilde{S} = \begin{pmatrix} 2\alpha_1\lambda_1 & \alpha_1\mu_2 & \alpha_1\mu_3 & \dots & \alpha_1\mu_{N-1} & \alpha_1\mu_N \\ \alpha_1\mu_2 & 2\alpha_2\lambda_2 & \alpha_2\mu_3 & \dots & \alpha_2\mu_{N-1} & \alpha_2\mu_N \\ \alpha_1\mu_3 & \alpha_2\mu_3 & 2\alpha_3\lambda_3 & \dots & \alpha_3\mu_{N-1} & \alpha_3\mu_N \\ \vdots & & \ddots & \ddots & & \vdots \\ \alpha_1\mu_{N-1} & \alpha_2\mu_{N-1} & \alpha_3\mu_{N-1} & \dots & 2\alpha_{N-1}\lambda_{N-1} & \alpha_{N-1}\mu_N \\ \alpha_1\mu_N & \alpha_2\mu_N & \alpha_3\mu_N & \dots & \alpha_{N-1}\mu_N & 2\alpha_N \end{pmatrix} - \frac{1}{2\tau} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_N \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_N \end{pmatrix}^\top.$$

The first step to show the diagonally dominant character of \tilde{S} is to note that every non-diagonal element of \tilde{S} is non-positive: $\alpha_j\mu_i - \frac{\mu_i\mu_j}{2\tau} \leq 0$, $\forall i \neq j$. Indeed, this is equivalent to write this in the following form ($\mu_i > 0$):

$$\alpha_j - \frac{\mu_j}{2\tau} = \alpha_j \left(\frac{\left(2\sum_{k=1}^N \alpha_k - \sum_{k=1}^i \alpha_k\right) \left(2\sum_{k=1}^N \alpha_k - \sum_{k=1}^{i-1} \alpha_k\right) - \left(2\sum_{k=1}^N \alpha_k\right)^2}{\left(2\sum_{k=1}^N \alpha_k - \sum_{k=1}^i \alpha_k\right) \left(2\sum_{k=1}^N \alpha_k - \sum_{k=1}^{i-1} \alpha_k\right)} \right) \leq 0,$$

since $\alpha_k \geq 0$ by assumption. This allows to discard the absolute values in the diagonally dominance criteria. Then, using the equality constraints, we obtain an expression for the sum of all non-diagonal elements of line i of \tilde{S} :

$$\begin{aligned} \mu_i \sum_{j=1}^{i-1} \alpha_j + \alpha_i \sum_{j=i+1}^N \mu_j - \frac{\mu_i}{2\tau} \sum_{j \neq i} \mu_j \\ = \begin{cases} \mu_i \sum_{j=1}^{i-1} \alpha_j + \alpha_i(1 - \lambda_i) - \frac{1}{2\tau} \mu_i(1 - \mu_i), & \text{if } i < N \\ \mu_N \sum_{j=1}^{N-1} \alpha_j - \frac{1}{2\tau} \mu_N(1 - \mu_N) & \text{if } i = N \end{cases} \end{aligned}$$

Using the values of μ_i , λ_i and τ along with elementary computations allow to verify $\forall i \in \{1, \dots, N\}$:

$$\begin{cases} -(\mu_i \sum_{j=1}^{i-1} \alpha_j + \alpha_i(1 - \lambda_i) - \frac{1}{2\tau} \mu_i(1 - \mu_i)) & = 2\alpha_i\lambda_i - \frac{\mu_i^2}{2\tau} & \text{if } i = 1, \dots, N-1, \\ -(\mu_i \sum_{j=1}^{i-1} \alpha_j - \frac{1}{2\tau} \mu_i(1 - \mu_i)) & = 2\alpha_i - \frac{\mu_i^2}{2\tau} & \text{if } i = N, \end{cases}$$

which implies diagonal dominance of \tilde{S} (even more: the sum of the elements of each line equals 0).