

Splitting the Conditional Gradient Algorithm

Zev Woodstock* and Sebastian Pokutta

Zuse Institute Berlin and
Institute of Mathematics, Technische Universität Berlin,
Berlin, Germany

January 31, 2024

Abstract

We propose a novel generalization of the conditional gradient (CG / Frank-Wolfe) algorithm for minimizing a smooth function f under an intersection of compact convex sets, using a first-order oracle for ∇f and linear minimization oracles (LMOs) for the individual sets. Although this computational framework presents many advantages, there are only a small number of algorithms which require one LMO evaluation per set per iteration; furthermore, these algorithms require f to be convex. Our algorithm appears to be the first in this class which is proven to also converge in the nonconvex setting. Our approach combines a penalty method and a product-space relaxation. We show that one conditional gradient step is a sufficient subroutine for our penalty method to converge, and we provide several analytical results on the product-space relaxation's properties and connections to other problems in optimization. We prove that our average Frank-Wolfe gap converges at a rate of $\mathcal{O}(\ln t / \sqrt{t})$, – only a log factor worse than the vanilla CG algorithm with one set.

Keywords. Conditional gradient, splitting, nonconvex, Frank-Wolfe, projection free

MSC Classification. 46N10, 65K10, 90C25, 90C26, 90C30

1 Introduction

Given a smooth function f which maps from a real Hilbert space \mathcal{H} to \mathbb{R} and a finite collection of m nonempty compact convex subsets $(C_i)_{i \in I}$ of \mathcal{H} , we seek to solve the following:

$$\text{minimize } f(x) \quad \text{subject to } x \in \bigcap_{i \in I} C_i, \tag{1}$$

which has many applications in imaging, signal processing, and data science [1, 2, 3]. Classical projection-based algorithms can be used to solve (1) if given access to the operator $\text{Proj}_{\bigcap_{i \in I} C_i}$. However, in practice, computing a projection onto $\bigcap_{i \in I} C_i$ is either impossible or numerically costly, and utilizing the individual projection operators $(\text{Proj}_{C_i})_{i \in I}$ is more tractable. This issue has given rise to the advent of *splitting* algorithms, which seek to solve (1) by utilizing

*Corresponding author: woodstock@zib.de

operators associated with the individual sets – not their intersection. Projection-based splitting algorithms – which use the collection of operators $(\text{Proj}_{C_i})_{i \in I}$ instead of $\text{Proj}_{\bigcap_{i \in I} C_i}$ – have made previously-intractable problems of the form (1) solvable with simpler tools on a larger scale [1, 3, 4].

While splitting methods have successfully been applied to projection-based algorithms, relatively little has been done for the splitting of *conditional gradient* (CG / Frank-Wolfe) algorithms. Standard CG algorithms minimize a smooth function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ over one closed convex constraint set $C \subset \mathbb{R}^n$. While the iterates of this algorithm do not converge in general [5], at iteration $t \in \mathbb{N}$, the average Frank-Wolfe gap (which is closely related to showing Clarke stationarity [8]) converges at a rate of $\mathcal{O}(1/\sqrt{t})$, and the primal gap converges at a rate of $\mathcal{O}(1/t)$ when f is convex [6]. A key ingredient of these algorithms is the *linear minimization oracle*, LMO_C , which computes for a linear objective $c \in \mathbb{R}^n$ a point in $\text{Argmin}_{x \in C} \langle c, x \rangle$. Similarly to traditional projection-based methods, computing $\text{LMO}_{\bigcap_{i \in I} C_i}$ is often prohibitively costly, so an algorithm which relies on the individual operators $(\text{LMO}_{C_i})_{i \in I}$ would be more tractable.

In principle, if two sets C_1 and C_2 are polytopes, one could compute $\text{LMO}_{C_1 \cap C_2}$ by solving a linear program which incorporates the LP formulations of both C_1 and C_2 . However, since the number of inequalities in an LP formulation can scale exponentially with dimension [7, 8], LPs are usually only used to implement a polyhedral LMO if there are no alternatives. In reality, many polyhedra used in applications, e.g., the Birkhoff polytope and the ℓ_1 ball, have highly specialized algorithms for computing their LMO which are faster than using a linear program [9]. Hence, splitting algorithms which rely on evaluating the specialized algorithms for $(\text{LMO}_{C_i})_{i \in I}$ gain the favorable scalability of existing LMO implementations.

Conditional gradient methods have seen a resurgence in popularity since, particularly for high-dimensional settings, LMOs can be more computationally efficient than projections. For instance, a common constraint in matrix completion problems is the spectrahedron

$$S = \{x \in \mathbb{S}_+^n \mid \text{Trace}(x) = 1\}, \quad (2)$$

where \mathbb{S}_+^n is the set of positive semidefinite $n \times n$ matrices. Evaluating Proj_S requires a full eigen-decomposition, while computing LMO_S only requires determining a dominant eigenpair [10]. Clearly, there are high-dimensional settings where evaluating LMO_S is possible while Proj_S is too costly [9]. Thus, we are particularly motivated by high-dimensional problems in data science (e.g., cluster analysis, graph refinement, and matrix decomposition) with these LMO-advantaged constraints, e.g., the nuclear norm ball, the Birkhoff polytope of doubly stochastic matrices, and the ℓ_1 ball [8, 10, 11, 12, 13, 14, 15].

Inexact proximal splitting methods are a natural choice for solving (1) in our computational setting, since LMO-based subroutines can approximate a projection. In the convex case, this approach appears in [16, 17, 18, 19]. However, there is often no bound on the number of LMO calls required to meet the relative error tolerance required of the subroutine, e.g., in [19]. Methods which require increasingly-accurate approximations can drive the number of LMO calls in each subroutine to infinity [18], and even if a bound on the number of LMO calls exists, it often depends on the conditioning of the projection subproblem.

We are interested in algorithms with low iteration complexity, since they are more tractable on large-scale problems. It appears that, for this computational setting, the lowest iteration complexity currently requires one LMO per set per iteration [2, 12, 20, 21, 22, 23]. To the best of our knowledge, all algorithms in this class are restricted to the convex setting. The case when $f = 0$ is addressed by [20], the case when $(C_i)_{i \in I}$ have additional structure is addressed in [21], and a matrix recovery problem is addressed in [22]. The approaches in [2, 12, 23] essentially show that one CG step is a sufficient subroutine for an inexact augmented Lagrangian (AL) approach. These methods prove convergence of different optimality criteria at various rates, e.g., arbitrarily close to $\mathcal{O}(t^{-1/3})$ [2], $\mathcal{O}(1/\sqrt{t})$ [21, 23], and (under restrictions

on m or $(C_i)_{i \in I}$ $\mathcal{O}(1/t)$ [12, 20, 22]. All of these methods, similarly to many projection-based splitting algorithms, achieve approximate feasibility in the sense that a point in the intersection $\bigcap_{i \in I} C_i$ is only found asymptotically.

Our contributions are as follows. We propose a new algorithm in this class for solving (1) which requires one LMO per set per iteration. Our algorithm generalizes the vanilla CG algorithm in the sense that, when $m = 1$, both algorithms are identical. It appears that our algorithm is the first in this class possessing convergence guarantees for solving (1) in the setting when f is nonconvex. As is standard in the CG literature, we analyze convergence of the average of Frank-Wolfe gaps, and we prove a rate of $\mathcal{O}(\ln t/\sqrt{t})$ – only a log factor slower than the rate for nonconvex CG over a single constraint ($m = 1$) [6]. We also prove primal gap convergence for the convex case. Our theory deviates from the AL approach and shows convergence with direct CG analysis, without imposing additional structure on our problem. By recasting (1) in a product space, we derive a penalized relaxation which is tractable with the vanilla CG algorithm. At each step of our algorithm, we perform one vanilla CG step on our product space relaxation; then, we update the objective function via a penalty. We provide an analytical and geometric exploration about the properties of this subproblem as its penalty changes, as well as its relationships to (1) and related optimization problems. In particular, we show that for any sequence of penalty parameters which approach ∞ , our subproblems converge (in several ways) to the original problem.

Our method combines two classical tools from optimization: a product-space reformulation and a penalty method. Penalty methods with CG-based subroutines received some attention several decades ago [24, 25]. Our algorithm is related the Regularized Frank-Wolfe algorithm of [24], however their requirements do not apply in our setting.

In the remainder of this section, we introduce notation, background, and standing assumptions. In Section 2, we demonstrate our product space approach and we establish analytical results. In Section 3, we introduce our algorithm and prove it converges.

1.1 Notation, standing assumptions, and auxiliary results

Let \mathcal{H} be a real Hilbert space with inner product $\langle \cdot | \cdot \rangle$ and identity operator Id . A closed ball centered at $x \in \mathcal{H}$ of radius $\varepsilon > 0$ is denoted $B(x; \varepsilon)$. Let $m \in \mathbb{N}$, set $I = \{1, \dots, m\}$, and let $\{\omega_i\}_{i \in I} \subset]0, 1]$ satisfy $\sum_{i \in I} \omega_i = 1$ (e.g., $\omega_i \equiv 1/m$).

$$\mathcal{H} = \mathcal{H}^m \text{ is the real Hilbert space with inner product } \langle \cdot | \cdot \rangle_{\mathcal{H}} = \sum_{i \in I} \omega_i \langle \cdot | \cdot \rangle_{\mathcal{H}}. \quad (3)$$

We use bold to denote points \mathbf{x} in \mathcal{H} , and their subcomponents are $\mathbf{x} = (x^1, x^2, \dots, x^m) \in \mathcal{H}$. We call $\mathbf{D} = \{\mathbf{x} \in \mathcal{H} \mid x^1 = x^2 = \dots = x^m\}$ the *diagonal subspace* of \mathcal{H} . The block averaging operation and its adjoint are

$$A: \mathcal{H} \rightarrow \mathcal{H}: \mathbf{x} \mapsto \sum_{i \in I} \omega_i x^i \quad \text{and} \quad A^*: \mathcal{H} \rightarrow \mathcal{H}: x \mapsto (x, \dots, x). \quad (4)$$

The *projection* operator onto a closed convex set $C \subset \mathcal{H}$ is denoted $\text{Proj}_C: \mathcal{H} \rightarrow \mathcal{H}: x \mapsto \text{Argmin}_{c \in C} \|x - c\|$. The *distance* and *indicator* functions of the set \mathbf{D} are denoted

$$\text{dist}_{\mathbf{D}}: \mathcal{H} \rightarrow \mathbb{R}: \mathbf{x} \mapsto \inf_{\mathbf{z} \in \mathbf{D}} \|\mathbf{x} - \mathbf{z}\| \quad \text{and} \quad \iota_{\mathbf{D}}: \mathcal{H} \rightarrow [0, +\infty]: \mathbf{x} \mapsto \begin{cases} 0 & \text{if } \mathbf{x} \in \mathbf{D} \\ +\infty & \text{if } \mathbf{x} \notin \mathbf{D}. \end{cases} \quad (5)$$

Note $\|A\| \leq 1$ and the identities

$$A^*A = \text{Proj}_{\mathbf{D}}, \quad \frac{1}{2} \text{dist}_{\mathbf{D}}^2(\mathbf{x}) = \frac{1}{2} \sum_{i \in I} \omega_i \|A\mathbf{x} - x^i\|^2, \quad \text{and} \quad \nabla \frac{1}{2} \text{dist}_{\mathbf{D}}^2 = \text{Id} - \text{Proj}_{\mathbf{D}}. \quad (6)$$

Unless otherwise stated, let $(C_i)_{i \in I}$ be a collection of nonempty compact convex subsets of \mathcal{H} , let $L_f > 0$, and let $f: \mathcal{H} \rightarrow \mathbb{R}$ be a Gâteaux differentiable function which is L_f -smooth,

$$(\forall (x, y) \in \mathcal{H}^2) \quad f(y) - f(x) \leq \langle \nabla f(x) \mid y - x \rangle + \frac{L_f}{2} \|y - x\|^2 \quad (7)$$

and, when restricted to Section 3.1, also convex,

$$(\forall (x, y) \in \mathcal{H}^2) \quad \langle \nabla f(x) \mid y - x \rangle \leq f(y) - f(x). \quad (8)$$

Fact 1.1 Since $\nabla \text{dist}_D^2/2 = \text{Id} - \text{Proj}_D = \text{Proj}_{D^\perp}$ is a projection operator onto a nonempty closed convex set, it is 1-Lipschitz continuous and therefore $\text{dist}_D^2/2$ is 1-smooth.

For every $i \in I$ and every $x \in \mathcal{H}$, the operation LMO_i returns a point in $\text{Argmin}_{z \in C_i} \langle x \mid z \rangle$. The *Frank-Wolfe gap* (F-W gap) of f over a compact convex set $C \subset \mathcal{H}$ at $x \in \mathcal{H}$ is $G_{f,C}(x) := \sup_{v \in C} \langle \nabla f(x) \mid x - v \rangle = \langle \nabla f(x) \mid x - \text{LMO}_C(\nabla f(x)) \rangle$. Note that, for every $x \in \mathcal{H}$,

$$x \text{ is a stationary point of } \underset{x \in C}{\text{minimize}} f(x) \quad \Leftrightarrow \quad \begin{cases} x \in C \\ G_{f,C}(x) \leq 0. \end{cases} \quad (9)$$

Note that if $x \in C$, we always have $G_{f,C}(x) \geq 0$.

Lemma 1.2 Let f and h be real-valued functions on a nonempty set $C \subset \mathcal{H}$, let $\lambda, \Delta > 0$, and suppose that

$$x \in \underset{x \in C}{\text{Argmin}} f(x) + \lambda h(x) \quad \text{and} \quad z \in \underset{z \in C}{\text{Argmin}} f(z) + (\lambda + \Delta)h(z).$$

Then $f(x) \leq f(z)$ and $h(z) \leq h(x)$.

Proof. Since x and z are optimal solutions, we have $f(x) + \lambda h(x) \leq f(z) + \lambda h(z)$ and $f(z) + (\lambda + \Delta)h(z) \leq f(x) + (\lambda + \Delta)h(x)$, so in particular,

$$(\lambda + \Delta)(h(z) - h(x)) \leq f(x) - f(z) \leq \lambda(h(z) - h(x)). \quad (10)$$

Subtracting $\lambda(h(z) - h(x))$ from (10) implies that $h(z) - h(x) \leq 0$ which, in view of (10), yields $f(x) - f(z) \leq 0$. \square

We assume the ability to compute ∇f , $(\text{LMO}_i)_{i \in I}$, and basic linear algebra operations, e.g., those in (4). Let $f: \mathcal{H} \rightarrow]-\infty, +\infty]$. The *subdifferential* of f at $x \in \mathcal{H}$ is given by $\partial f(x) = \{u \in \mathcal{H} \mid (\forall y \in \mathcal{H}) \quad f(x) + \langle u \mid y - x \rangle \leq f(y)\}$. The *epigraph* of f is $\text{epi } f = \{(x, \eta) \in \mathcal{H} \times \mathbb{R} \mid f(x) \leq \eta\}$. The *graph* of an operator $M: \mathcal{H} \rightarrow 2^{\mathcal{H}}$ is $\text{gra } M = \{(x, u) \in \mathcal{H}^2 \mid u \in M(x)\}$. Some of our analytical results rely on the theory of convergence of sets and set-valued operators; for a broad review, see [26].

Definition 1.3 Let $(C_n)_{n \in \mathbb{N}}$ be a sequence of subsets of \mathbb{R}^n , and let $(f_n)_{n \in \mathbb{N}}$ be functions on \mathbb{R}^n . The *outer limit* of $(C_n)_{n \in \mathbb{N}}$ is $\limsup_{n \in \mathbb{N}} (C_n)_{n \in \mathbb{N}} = \{x \in \mathbb{R}^n \mid \limsup_{n \rightarrow +\infty} \text{dist}_{C_n}(x) = 0\}$; the *inner limit* of $(C_n)_{n \in \mathbb{N}}$ is $\liminf_{n \in \mathbb{N}} (C_n)_{n \in \mathbb{N}} = \{x \in \mathbb{R}^n \mid \liminf_{n \rightarrow +\infty} \text{dist}_{C_n}(x) = 0\}$ [26, Ex. 4.2]. If both limits exist and coincide, this set is the *limit* of $(C_n)_{n \in \mathbb{N}}$. The sequence $(f_n)_{n \in \mathbb{N}}$ *converges epigraphically* to a function f on \mathbb{R}^n if the sequence of epigraphs $(\text{epi } f_n)_{n \in \mathbb{N}}$ converge to $\text{epi } f$. The sequence $(\partial f_n)_{n \in \mathbb{N}}$ *converges graphically* to ∂f if $(\text{gra } \partial f_n)_{n \in \mathbb{N}}$ converges to $\text{gra } \partial f$.

2 Splitting constraints with a product space

This section outlines our algorithm and provides additional analysis relating our approach to similar problems in optimization.

2.1 Algorithm design

The vanilla conditional gradient algorithm solves

$$\underset{x \in C}{\text{minimize}} \quad f(x) \quad (11)$$

using LMO_C and gradients of f . However, one of the central hurdles in designing a tractable CG-based splitting algorithm is finding a way to enforce membership in the constraint $\bigcap_{i \in I} C_i$ without access to its projection or LMO. Our approach to solving this issue comes from the following construction on the product space \mathcal{H} (see Section 1.1 for notation).

Proposition 2.1 *Let $(C_i)_{i \in I}$ be a collection of nonempty subsets of \mathcal{H} , and let $D \subset \mathcal{H}$ denote the diagonal subspace. Then*

$$(\forall x \in \mathcal{H}) \quad (x, \dots, x) \in D \cap \bigtimes_{i \in I} C_i \Leftrightarrow x \in \bigcap_{i \in I} C_i \quad (12)$$

$$(\forall x \in \mathcal{H}) \quad x \in D \cap \bigtimes_{i \in I} C_i \Leftrightarrow (\exists x \in \mathcal{H}) \quad \begin{cases} x = (x, \dots, x) \\ x \in \bigcap_{i \in I} C_i, \end{cases} \quad (13)$$

Proof. Clear from construction.¹ □

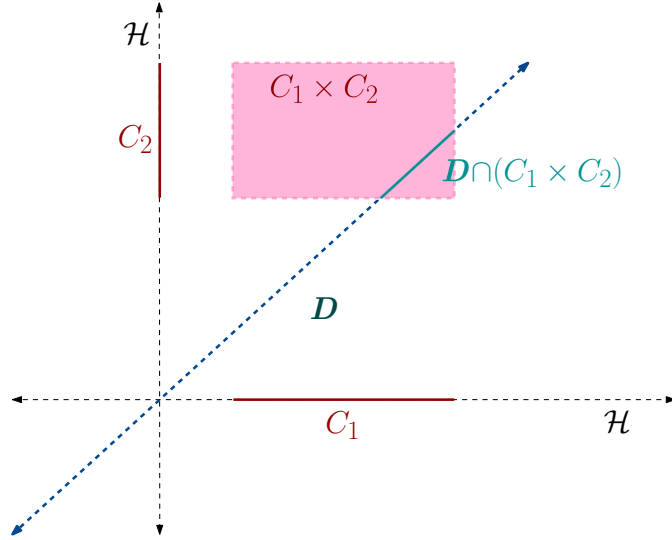


Fig. 1: Visualization of the product space for $\mathcal{H} = \mathbb{R}$ and $m = 2$. Our algorithm produces iterates x_t which are always inside the shaded constraint set, and their averages $A^* A x_t$ are always on the diagonal subspace D . The solid segment where $C_1 \times C_2$ and D intersect corresponds precisely to our split feasibility constraint via Proposition 2.1.

¹The type of construction in Proposition 2.1 goes back to the work of Pierra [27].

Proposition 2.1 provides a decomposition of the split feasibility constraint $\bigcap_{i \in I} C_i$ in terms of two simpler sets D and $\times_{i \in I} C_i$. This yields a product space reformulation of (1)

$$\underset{\mathbf{x} \in \times_{i \in I} C_i}{\text{minimize}} \quad f(A\mathbf{x}) + \iota_D(\mathbf{x}). \quad (14)$$

The constraints D and $\times_{i \in I} C_i$ are simpler in the sense that, even in our restricted computational setting, we can compute operators to enforce them. In particular, the projection onto D is computed by simply repeating the average of all components in every component $\text{Proj}_D \mathbf{x} = A^*(\sum_{i \in I} \omega_i \mathbf{x}^i)$. Critically, this operation is cheap, so one can actually evaluate the gradient $\nabla \text{dist}_D^2/2 = \text{Id} - \text{Proj}_D$ even though it involves a projection. The constraint $\times_{i \in I} C_i$ is readily processed using the following property.

Fact 2.2 Let $(C_i)_{i \in I}$ be a collection of nonempty compact convex subsets of \mathcal{H} . Then

$$(\forall \mathbf{x} \in \mathcal{H}) \quad \text{LMO}_{(\times_{i \in I} C_i)}(\mathbf{x}) = (\text{LMO}_{C_1} \mathbf{x}^1, \dots, \text{LMO}_{C_m} \mathbf{x}^m). \quad (15)$$

In particular, to evaluate an LMO for the product $\times_{i \in I} C_i$, it suffices to evaluate the individual operators $(\text{LMO}_i)_{i \in I}$ once.

With these ideas in mind, let us introduce the penalized function

$$F_\lambda: \mathcal{H} \rightarrow \mathbb{R}: \mathbf{x} \mapsto f(A\mathbf{x}) + \lambda \frac{1}{2} \text{dist}_D^2(\mathbf{x}), \quad (16)$$

which, for every $\lambda \geq 0$, is $(L_f + \lambda)$ -smooth (cf. Fact 1.1). We observe that for every penalty parameter $\lambda_t \geq 0$, even under our restricted computational setting, the following relaxation of (14) is still tractable with the vanilla CG algorithm

$$\underset{\mathbf{x} \in \times_{i \in I} C_i}{\text{minimize}} \quad F_{\lambda_t}(\mathbf{x}). \quad (17)$$

Indeed, vanilla CG requires the ability to compute the gradient of the objective function and the LMO of the constraint. Computing $\nabla F_\lambda = \nabla f + \lambda(\text{Id} - \text{Proj}_D)$ amounts to one evaluation of ∇f , computing one average, and some algebraic manipulations. By promoting membership of D via the objective function, we are left with the LMO-amenable constraint $\times_{i \in I} C_i$.

The core idea of our algorithm is, at each iteration $t \in \mathbb{N}$, to perform one Frank-Wolfe step to the relaxed subproblem (17). Then, between iterations, we update the objective function in (17) via λ_t to promote feasibility. Although (17) is a relaxation of the intractable problem (14), taking $\lambda_t \rightarrow \infty$ suffices to show convergence in F-W gap (and primal gap, in the convex case) to solutions of (14) and hence (1); this is substantiated in Sections 2.2.2 and 3. For every $\mathbf{x} \in \mathcal{H}$, the i th component of the gradient is given by $\nabla F_\lambda(\mathbf{x})^i = \nabla f(A\mathbf{x}) + \lambda(\mathbf{x}^i - A\mathbf{x})$. So, a CG step applied to (17) yields Algorithm 1. While Section 3 contains the precise schedules for Lines 3 and 4, the parameters behave like $(\lambda_t, \gamma_t) = (\mathcal{O}(\ln t), \mathcal{O}(1/\sqrt{t}))$.

CG-based algorithms possess the advantage that, at every iteration, the iterates are feasible (i.e., for (11), $\mathbf{x}_t \in C$). Our approach inherits this familiar property; however, since we solve a product space relaxation, $\mathbf{x}_t \in \times_{i \in I} C_i$ and hence, for every $i \in I$, the i th component of our sequence is feasible for the i th constraint, i.e., $(\mathbf{x}_t^i)_{t \in \mathbb{N}} \in C_i$. Importantly, this does not guarantee that any subcomponent \mathbf{x}_t^i resides in $\bigcap_{i \in I} C_i$, so they are not feasible for the splitting problem (1); feasibility in $\bigcap_{i \in I} C_i$ is acquired “in the limit”, by showing that $\mathbf{x}_t \in \times_{i \in I} C_i$ and $\text{dist}_D(\mathbf{x}_t) \rightarrow 0$ (proven in Section 3).

In practice, one needs a route to construct an approximate solution to (1) in \mathcal{H} from an iterate of Algorithm 1 in the product space \mathcal{H} . Instead of taking a component, we use the average computed in Line 10 as our approximate solution, since

$$(\forall \mathbf{x} \in \mathcal{H}) \quad \mathbf{x} \in D \cap \times_{i \in I} C_i \quad \Rightarrow \quad A\mathbf{x} \in \bigcap_{i \in I} C_i \quad (18)$$

Algorithm 1 Split conditional gradient (SCG) algorithm

Require: Smooth function f , weights $\{\omega_i\}_{i \in I} \subset]0, 1]$ such that $\sum_{i \in I} \omega_i = 1$, point $x_0 \in \times_{i \in I} C_i$

```

1:  $x_0 \leftarrow \sum_{i \in I} \omega_i x_0^i$ 
2: for  $t = 0, 1$  to  $\dots$  do
3:   Choose penalty parameter  $\lambda_t \in ]0, +\infty[$ 
4:   Choose step size  $\gamma_t \in ]0, 1]$ 
5:    $g_t \leftarrow \nabla f(x_t)$  # Store  $\nabla f(Ax_t)$  for CG step on (17)
6:   for  $i = 1$  to  $m$  do
7:      $v_t^i \leftarrow \text{LMO}_i(g_t + \lambda_t(x_t^i - x_t))$  # LMO applied to  $\nabla F_{\lambda_t}(x_t)^i$ 
8:      $x_{t+1}^i \leftarrow x_t^i + \gamma_t(v_t^i - x_t^i)$  # CG step in  $i$ th component
9:   end for
10:   $x_{t+1} \leftarrow \sum_{i \in I} \omega_i x_{t+1}^i$  # Approximate solution by averaging
11: end for
  
```

is a strict implication. Hence the condition $Ax \in \bigcap_{i \in I} C_i$ is easier to satisfy than $x \in D \cap \times_{i \in I} C_i$ (see also Sec. 2.2.1).

Remark 2.3 If we have only $m = 1$ set constraint, then $A = \text{Id}$, and $\mathcal{H} = \mathcal{H} = D$, so at every iteration $t \in \mathbb{N}$, $F_{\lambda_t} = f(x)$. Therefore, the classical CG algorithm is a special case of Algorithm 1.

2.2 Analysis

Here we gather analytical results pertaining to our algorithm, the geometry of our product-space construction, and how our relaxed problem relates to other classical problems in optimization. While these results are interesting in their own right, many are also used to show convergence in Section 3.

2.2.1 Geometry (and tractability) of penalty functions on the Cartesian product

As seen in Section 2.1, Algorithm 1 promotes split feasibility by, at every iteration $t \in \mathbb{N}$, requiring that $x_t \in \times_{i \in I} C_i$ and penalizing the distance from x_t to D . However, as seen in (18), $\text{dist}_D(x_t) = 0$ is a sufficient (but not necessary) condition to acquire a feasible average $Ax_t \in \bigcap_{i \in I} C_i$; see Fig. 2. In this section, we present a penalty function which precisely characterizes this condition. Via a simple geometric argument based on the projection theorem, we guarantee that although utilizing this penalty is not computationally tractable, it is nonetheless minimized when dist_D vanishes. These results also further substantiate the claim that $x \in D \cap \times_{i \in I} C_i$ is a stricter condition than $Ax \in \bigcap_{i \in I} C_i$, which is our motivation to use the average in Line 10 of Algorithm 1 as our approximate solution to (1).

Proposition 2.4 Let $(C_i)_{i \in I}$ be a collection of nonempty closed convex subsets of \mathcal{H} , let $D \subset \mathcal{H}$ denote the diagonal subspace, and set

$$d: \mathcal{H} \rightarrow]-\infty, +\infty] : x \mapsto \sum_{i \in I} \omega_i \text{dist}_{C_i}^2(Ax). \quad (19)$$

Then, for every $x \in \mathcal{H}$, the following are equivalent.

- (i) $d(x) = 0$.

(ii) $A\mathbf{x} \in \bigcap_{i \in I} C_i$.

(iii) $\text{Proj}_D(\mathbf{x}) \in \times_{i \in I} C_i$.

Proof. (i) \Rightarrow (ii): For every $i \in I$, $0 \leq \omega_i \text{dist}_{C_i}^2(A\mathbf{x}) \leq d(\mathbf{x}) = 0$. Since $\omega_i > 0$, it follows that $\text{dist}_{C_i}^2(A\mathbf{x}) = 0$ and hence $A\mathbf{x} \in C_i$.

(ii) \Rightarrow (iii): By applying A^* to the inclusion (ii), (6) implies that

$$\text{Proj}_D \mathbf{x} = A^* A \mathbf{x} \in A^* \bigcap_{i \in I} C_i = \left\{ (x, \dots, x) \in \mathcal{H} \mid x \in \bigcap_{i \in I} C_i \right\}. \quad (20)$$

So, by Proposition 2.1, $\text{Proj}_D \mathbf{x} \in \left\{ \mathbf{x} \in \mathcal{H} \mid \mathbf{x} \in D \cap \times_{i \in I} C_i \right\} \subset \times_{i \in I} C_i$.

(iii) \Rightarrow (i): We begin by observing that

$$\text{Proj}_{\times_{i \in I} C_i}(\text{Proj}_D \mathbf{x}) = \underset{\mathbf{c} \in \times_{i \in I} C_i}{\text{Argmin}} \|\mathbf{c} - A^* A \mathbf{x}\|_{\mathcal{H}}^2 = \underset{\mathbf{c} \in \times_{i \in I} C_i}{\text{Argmin}} \sum_{i \in I} \omega_i \|\mathbf{c}^i - A \mathbf{x}\|_{\mathcal{H}}^2 \quad (21)$$

is a separable problem whose solution is $(\text{Proj}_{C_1}(A\mathbf{x}), \dots, \text{Proj}_{C_m}(A\mathbf{x}))$. Therefore,

$$d(\mathbf{x}) = \sum_{i \in I} \omega_i \|A\mathbf{x} - \text{Proj}_{C_i}(A\mathbf{x})\|_{\mathcal{H}}^2 = \|\text{Proj}_D \mathbf{x} - \text{Proj}_{\times_{i \in I} C_i}(\text{Proj}_D \mathbf{x})\|_{\mathcal{H}}^2. \quad (22)$$

Since $\text{Proj}_D(\mathbf{x}) = \text{Proj}_{\times_{i \in I} C_i}(\text{Proj}_D \mathbf{x})$, we conclude $d(\mathbf{x}) = 0$. \square

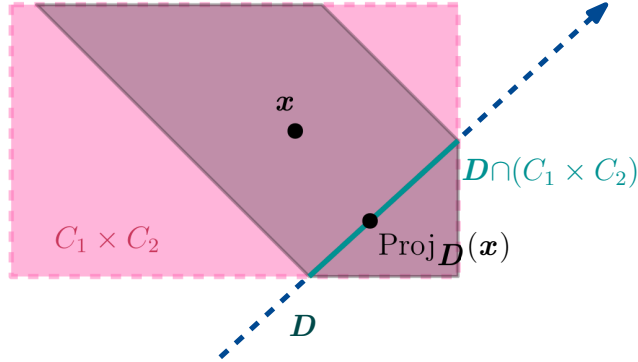


Fig. 2: Zoomed view of Fig. 1. The darker shaded area is the collection of points $\mathbf{x} \in \times_{i \in I} C_i$ for which $\text{Proj}_D(\mathbf{x})$ remains in $\times_{i \in I} C_i$. By Proposition 2.4, this is the set of points satisfying $A\mathbf{x} \in \bigcap_{i \in I} C_i$. This exemplifies that the implication (18) is strict.

Since we do not assume the ability to project onto the sets $(C_i)_{i \in I}$, evaluating $\nabla d = 2 \sum_{i \in I} \omega_i (\text{Id} - \text{Proj}_{C_i})$ is not possible. Therefore, replacing F_λ in (17) with the composite function $f(A\mathbf{x}) + \lambda d(\mathbf{x})$ is not tractable with a vanilla CG-based approach. However, d is closely related to our penalty function dist_D^2 via the following result.

Corollary 2.5 *In the setting of Proposition 2.4, let $\mathbf{x} \in \times_{i \in I} C_i$, set $\mathbf{y} = \text{Proj}_D \mathbf{x}$ and set $\mathbf{p} = \text{Proj}_{(\times_{i \in I} C_i)}(\mathbf{y})$. Then*

$$d(\mathbf{x}) = \text{dist}_D^2(\mathbf{x}) - \|\mathbf{x} - \mathbf{p}\|^2 + 2 \underbrace{\langle \mathbf{x} - \mathbf{p} \mid \mathbf{y} - \mathbf{p} \rangle}_{\leq 0}. \quad (23)$$

In consequence, $0 \leq d(\mathbf{x}) \leq \text{dist}_D^2(\mathbf{x})$.

Proof. Follows from Lemma 2.12 and Theorem 3.16 of [28]. \square

Since the iterates of Algorithm 1 always reside in $\times_{i \in I} C_i$, Corollary 2.5 reinforces our choice of Ax_t as our approximate solution of (1). Firstly, its implication that $\text{dist}_D^2(x) = 0 \Rightarrow d(x) = 0$ underlines the observation from (18) that $Ax_t \in \bigcap_{i \in I} C_i$ is easier to satisfy than $x \in D \cap \times_{i \in I} C_i$. Furthermore, by characterizing the gap between d and dist_D^2 , we see that there are plenty of points for which the inequality between d and dist_D^2 is strict, e.g., those $x \in \times_{i \in I} C_i$ for which $x \neq p$ (see also Fig. 2). Due to this strictness, $d(x_t)$ may vanish far *before* $\text{dist}_D^2(x_t)$ vanishes over the iterations of Algorithm 1. This is consistent with our preliminary numerical observations that $Ax_t \in \bigcap_{i \in I} C_i$ often occurs before $\text{dist}_D^2(x_t)$ vanishes.

Remark 2.6 The function $g: x \mapsto \text{dist}_{\bigcap_{i \in I} C_i}^2(Ax) = \|\text{Proj}_D x - \text{Proj}_{D \cap (C_1 \times C_2)} x\|^2$ is also a natural penalty to consider, although evaluating $\nabla g = 2A^*(\text{Id} - \text{Proj}_{\bigcap_{i \in I} C_i})(Ax)$ involves computing an intractable projection. While, for every $x \in \mathcal{H}$, g and d (see (19)) have the order

$$d(x) = \sum_{i \in I} \omega_i \inf_{c \in C_i} \|Ax - c\|^2 \leq \sum_{i \in I} \omega_i \inf_{c \in \bigcap_{i \in I} C_i} \|Ax - c\|^2 = g(x), \quad (24)$$

there is no general ordering between g and our penalty dist_D^2 for $\dim(\mathcal{H}) \geq 2$. However, they are related in the following geometric sense

$$g(x) + \text{dist}_D^2(x) = \sum_{i \in I} \omega_i \|x^i - P_{\bigcap_{i \in I} C_i}(Ax)\|^2 + \underbrace{2\langle A^* P_{\bigcap_{i \in I} C_i}(Ax) - A^* Ax \mid x - A^* Ax \rangle}_{=0}. \quad (25)$$

Since the lefthand and righthand vectors in the scalar product are in D and D^\perp respectively, g and dist_D^2 describe the squared magnitude of two orthogonal vectors.

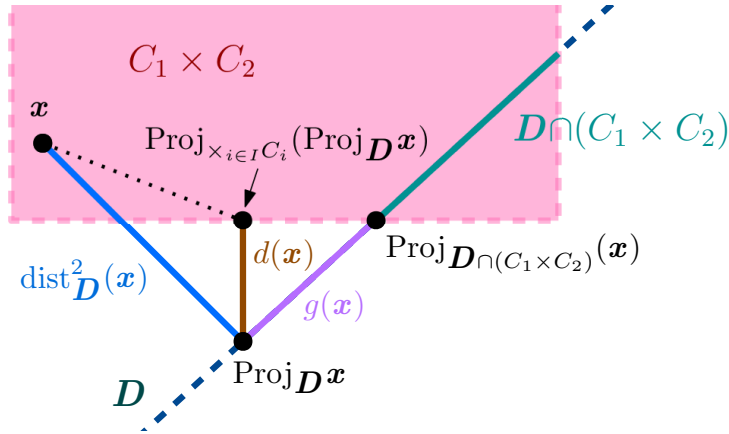


Fig. 3: Zoomed view of Fig. 1 displaying geometric relationships between dist_D^2 , d , and g (see (22) and Remark 2.6). This shows a feasible point $x \in \times_{i \in I} C_i$ and several vectors whose squared magnitude are equal to the given labels. Corollary 2.5 describes the relative magnitude of dist_D^2 and d , as well as the obtuse angle between x , $\text{Proj}_{\times_{i \in I} C_i}(\text{Proj}_D x)$, and $\text{Proj}_D x$. Remark 2.6 describes the orthogonality seen in the angle between x , $\text{Proj}_D x$, and $\text{Proj}_{D \cap \times_{i \in I} C_i}(x)$. We also see $\text{Proj}_{D \cap \times_{i \in I} C_i} x = \text{Proj}_{D \cap \times_{i \in I} C_i}(\text{Proj}_D x)$, which holds in general [28, Prop. 24.18].

2.2.2 Interpolating constraints: From the Minkowski sum to the intersection

This section presents an analysis of how our subproblem (17) changes with the parameter λ . In addition to their utility in Section 3 to prove that our sequence of relaxations (17) actually solves the correct problem (14), the results in this section show that (17) connects two classical problems in optimization.

From a certain perspective, (17) “interpolates” from the following problem (when $\lambda = 0$) over the Minkowski sum

$$\underset{x \in \sum_{i \in I} \omega_i C_i}{\text{minimize}} \quad f(x), \quad \text{where} \quad \sum_{i \in I} \omega_i C_i = \left\{ \sum_{i \in I} \omega_i c^i \mid (\forall i \in I) \ c^i \in C_i \right\}, \quad (26)$$

to the splitting problem (1) (when $\lambda \nearrow +\infty$). We shall make this latter observation precise via several notions of convergence in Proposition 2.12.

Remark 2.7 While this article is predominantly focused on (1), it is worth noting that, when $\lambda = 0$, the problems (17) and (26) coincide in the sense that, for every solution x^* of (17), Ax^* solves (26) (and for every solution $\sum_{i \in I} \omega_i x^i$ of (26), $(x^i)_{i \in I}$ solves (17)). Therefore, Fact 2.2 leads to a Frank-Wolfe approach to solving (26). The Minkowski sum constraint arises in Bayesian learning, placement problems, and robot motion planning [29, 30, 31, 32].

We begin with the following observations about how F_λ relates as λ varies.

Lemma 2.8 *Let $f: \mathcal{H} \rightarrow \mathbb{R}$, let $\lambda, \Delta \in \mathbb{R}$, let $D \subset \mathcal{H}$ be nonempty, and set $F_\lambda: x \mapsto f(Ax) + \lambda \text{dist}_D^2(x)/2$. Then,*

$$(\forall x \in \mathcal{H}) \quad F_\lambda(x) = F_{\lambda+\Delta}(x) - \Delta \frac{1}{2} \text{dist}_D^2(x). \quad (27)$$

In consequence, if $\Delta \geq 0$, then $F_\lambda(x) \leq F_{\lambda+\Delta}(x)$ and $\inf F_\lambda(\times_{i \in I} C_i) \leq \inf F_{\lambda+\Delta}(\times_{i \in I} C_i)$.

Proof. $F_\lambda(x) = f(Ax) + (\lambda + \Delta) \text{dist}_D^2(x)/2 - \Delta \text{dist}_D^2(x)/2 = F_{\lambda+\Delta}(x) - \Delta \text{dist}_D^2(x)/2$. \square

Next, we show that the optimal value of (17) is sandwiched between that of the splitting problem (1) and the Minkowski sum problem (26).

Proposition 2.9 *Let $f: \mathcal{H} \rightarrow \mathbb{R}$, let $\lambda \geq 0$, let $D \subset \mathcal{H}$ be nonempty, set $F_\lambda: x \mapsto f(Ax) + \lambda \text{dist}_D^2(x)/2$, and let $(C_i)_{i \in I}$ be a collection of nonempty compact convex subsets of \mathcal{H} such that $\bigcap_{i \in I} C_i \neq \emptyset$. Then*

$$\inf_{x \in \bigcap_{i \in I} C_i} f(x) \geq \inf_{x \in \times_{i \in I} C_i} F_\lambda \geq \inf_{x \in \sum_{i \in I} \omega_i C_i} f(x). \quad (28)$$

Proof. To show the first inequality, we note that for every $x \in D$, $\text{dist}_D^2(x) = 0$, so using the product space formulation (14) of (1),

$$\inf_{x \in \bigcap_{i \in I} C_i} f(x) = \inf_{x \in D \cap \times_{i \in I} C_i} f(Ax) + \frac{\lambda}{2} \text{dist}_D^2(x) \geq \inf_{x \in \times_{i \in I} C_i} F_\lambda(x). \quad (29)$$

The second inequality follows from the observation that (26) coincides with $\inf_{x \in \times_{i \in I} C_i} F_0(x)$, so by Lemma 2.8 we have $\inf_{x \in \times_{i \in I} C_i} F_\lambda(x) \geq \inf_{x \in \times_{i \in I} C_i} F_0(x)$. \square

It turns out that, for an increasing sequence of penalty parameters $(\lambda_n)_{n \in \mathbb{N}}$, the ordering of Proposition 2.9 is preserved if we only consider the optimal values of f (instead of F_λ). Intuitively, the order is reversed when we compare optimal values of the penalty dist_D^2 .

Corollary 2.10 Let $f: \mathcal{H} \rightarrow \mathbb{R}$, let $D \subset \mathcal{H}$ be nonempty, set $F_\lambda: x \mapsto f(Ax) + \lambda \text{dist}_D^2(x)/2$, and let $(C_i)_{i \in I}$ be a collection of nonempty compact convex subsets of \mathcal{H} such that $\bigcap_{i \in I} C_i \neq \emptyset$. Suppose that $(\lambda_t)_{t \in \mathbb{N}}$ is an increasing sequence of real numbers and, for every $t \in \mathbb{N}$, let x_t^* be a minimizer of F_{λ_t} over $\bigtimes_{i \in I} C_i$. Then

$$\inf_{x \in \bigcap_{i \in I} C_i} f(x) \geq f(Ax_{t+1}^*) \geq f(Ax_t^*) \geq \inf_{x \in \sum_{i \in I} \omega_i C_i} f(x). \quad (30)$$

If $z \in \text{Argmin}_{x \in \bigtimes_{i \in I} C_i} f(Ax)$ (i.e., Az solves (26)), then

$$0 \leq \text{dist}_D^2(x_{t+1}^*) \leq \text{dist}_D^2(x_t^*) \leq \text{dist}_D^2(z). \quad (31)$$

Proof. Follows from Lemma 1.2 and Proposition 2.9. \square

The following example demonstrates that the penalty sequence $(\lambda_t)_{t \in \mathbb{N}}$ may need to tend to $+\infty$ in order for the solutions of (17) and (1) to coincide.

Example 2.11 Set $\mathcal{H} = \mathbb{R}$, set $f = \|x\|^2/2$, let $z \geq 0$, set $C_1 = \{z\}$, and set $C_2 = [-z-1, z+1]$. Clearly, $z = \text{Argmin}_{x \in C_1 \cap C_2} f(x)$. However, it is straightforward to verify that, for every $\lambda \geq 0$, $x_\lambda^* = ((\lambda-1)z/(1+\lambda), z)$ is the unique minimizer of F_λ over $C_1 \times C_2$. Since $Ax_\lambda^* = \lambda z/(1+\lambda) \neq z$, the solutions of (17) and (1) (via (14)) do not coincide for finite λ ; taking $\lambda \rightarrow +\infty$ implies $Ax_\lambda^* \rightarrow z$.

The following result establishes three notions of convergence (see Definition 1.3) relating the problems (17) and (1) (via its equivalent product space formulation (14)). For this result, we rely on the fact that every constrained optimization problem can be described using a single objective function via the use of indicator functions.

Proposition 2.12 Let $f: \mathcal{H} \rightarrow \mathbb{R}$, let $(C_i)_{i \in I}$ be a collection of nonempty compact convex subsets of \mathcal{H} such that $\bigcap_{i \in I} C_i \neq \emptyset$, and let D denote the diagonal subspace of \mathcal{H} . Suppose that $(\lambda_t)_{t \in \mathbb{N}} \rightarrow +\infty$ and, for every $t \in \mathbb{N}$, set $f_t = f \circ A + \lambda_t \text{dist}_D^2/2 + \iota_{\bigtimes_{i \in I} C_i}$. Then the following hold.

- (i) f_t converges pointwise to $f \circ A + \iota_{D \cap \bigtimes_{i \in I} C_i}$.
- (ii) Suppose $\mathcal{H} = \mathbb{R}^n$. Then f_t converges epigraphically to $f \circ A + \iota_{D \cap \bigtimes_{i \in I} C_i}$.
- (iii) Suppose $\mathcal{H} = \mathbb{R}^n$ and f is convex. Then ∂f_n converges graphically to $\partial(f \circ A + \iota_{D \cap \bigtimes_{i \in I} C_i})$.

Proof. Since

$$\iota_{\bigtimes_{i \in I} C_i} + \iota_D = \iota_{D \cap \bigtimes_{i \in I} C_i}, \quad (32)$$

it suffices to show that $\lambda_t \text{dist}_D^2/2$ converges to ι_D under each notion of convergence.

(i): Let $x \in \mathcal{H}$. If $x \in D$, then for every $n \in \mathbb{N}$, $\lambda_t \text{dist}_D^2(x)/2 = 0 = \iota_D(x)$. On the other hand, if $x \notin D$, then $0 < \lambda_t \text{dist}_D^2(x)/2 \rightarrow +\infty = \iota_D(x)$.

(ii): Let $x \in \mathcal{H}$. By [26, Proposition 7.2], it suffices to show both of the following.

$$\text{For some sequence } (x_t)_{t \in \mathbb{N}} \text{ converging to } x, \quad \limsup_{t \in \mathbb{N}} \frac{\lambda_t}{2} \text{dist}_D^2(x_t) \leq \iota_D(x). \quad (33)$$

$$\text{For every sequence } (x_t)_{t \in \mathbb{N}} \text{ converging to } x, \quad \liminf_{t \in \mathbb{N}} \frac{\lambda_t}{2} \text{dist}_D^2(x_t) \geq \iota_D(x). \quad (34)$$

To realize (33), we consider the constant sequence $(x_t)_{t \in \mathbb{N}} \equiv x$. By (i),

$$\limsup_{t \in \mathbb{N}} \lambda_t \text{dist}_D^2(x_t)/2 = \lim_{t \in \mathbb{N}} \lambda_t \text{dist}_D^2(x)/2 = \iota_D(x), \quad (35)$$

so this is always satisfied with equality. To show (34), let $(x_t)_{t \in \mathbb{N}}$ be a sequence converging to x . If $x \in D$, then since $\text{dist}_D^2 \geq 0$ and $\iota_D(x) = 0$, (34) holds. Otherwise, if $x \notin D$, then there exists a radius $\varepsilon > 0$ such that $B(x; \varepsilon) \cap D = \emptyset$. Since dist_D^2 is continuous and only vanishes on D , we know $\eta := \inf_{y \in B(x; \varepsilon/2)} \text{dist}_D^2(y)/2 > 0$. Therefore, since $x_t \rightarrow x$, we have that, for some $N \in \mathbb{N}$, $n > N$ implies that $x_t \in B(x; \varepsilon/2)$, hence

$$\frac{\lambda_t}{2} \text{dist}_D^2(x_t) \geq \lambda_t \eta \rightarrow +\infty. \quad (36)$$

In particular, $\lim_{t \in \mathbb{N}} \lambda_t \text{dist}_D^2(x_t)/2 = +\infty = \iota_D(x)$ so we are done.

(iii): Follows from (ii) and Attouch's Theorem [26, Theorem 12.35]. \square

In general, the functions in Proposition 2.12 do not converge uniformly². In spite of this, it turns out that one can nonetheless commute the limit with an infimum, hence showing that the optimal values of our subproblems (17) converge to the optimal value of (1).

Proposition 2.13 *Let $f: \mathcal{H} \rightarrow \mathbb{R}$, let $(C_i)_{i \in I}$ be a collection of nonempty compact convex subsets of \mathcal{H} , let D denote the diagonal subspace of \mathcal{H} , and for every $\lambda \geq 0$, set $F_\lambda: x \mapsto f(Ax) + \lambda \text{dist}_D^2(x)/2$. Suppose that $(\lambda_n)_{n \in \mathbb{N}} \nearrow +\infty$. Then*

$$\lim_{t \rightarrow +\infty} \left(\inf_{x \in \times_{i \in I} C_i} F_{\lambda_t}(x) \right) \rightarrow \inf_{x \in \times_{i \in I} C_i} \left(\lim_{t \rightarrow \infty} F_{\lambda_t}(x) \right) = \inf_{x \in \bigcap_{i \in I} C_i} f(x). \quad (37)$$

Proof. First, we point out that the equality in (37) follows from Proposition 2.12 and the fact that the minimal values of (1) and (14) coincide. Let $\mu < \inf_{x \in \bigcap_{i \in I} C_i} f(x) = \inf_{x \in \times_{i \in I} C_i} f(Ax) + \iota_D(x)$. By Proposition 2.12, for every $x \in \times_{i \in I} C_i$, $\lim_{t \rightarrow \infty} F_{\lambda_t}(x) = f(Ax) + \iota_D(x) > \mu$. Since $\times_{i \in I} C_i$ is compact, for $t \in \mathbb{N}$ sufficiently large, $\inf_{x \in \times_{i \in I} C_i} F_{\lambda_t}(x) \geq \mu$, which implies (via Proposition 2.9 for the second inequality)

$$\mu \leq \lim_{t \rightarrow \infty} \left(\inf_{x \in \times_{i \in I} C_i} F_{\lambda_t}(x) \right) \leq \inf_{x \in \bigcap_{i \in I} C_i} f(x). \quad (38)$$

Taking $\mu \uparrow \inf_{x \in \bigcap_{i \in I} C_i} f(x)$ completes the result. \square

3 Convergence of Algorithm 1

We first prove that Algorithm 1 converges in function value when f is convex (Section 3.1). Then, we establish guarantees for stationarity in general (Section 3.2). We begin with an estimate which is used for both settings.

Lemma 3.1 *Let $(C_i)_{i \in I}$ be a finite collection of nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty]$, and let D denote the diagonal subspace of \mathcal{H} . Suppose that $\bigcap_{i \in I} C_i \neq \emptyset$. Then*

$$\left(\forall x, y \in \times_{i \in I} C_i \right) \quad \text{dist}_D^2(x) \leq \sum_{i \in I} \omega_i R_i^2 \quad \text{and} \quad \|x - y\|^2 \leq \sum_{i \in I} \omega_i R_i^2. \quad (39)$$

²Uniform convergence for extended-real valued functions is defined in [26].

Proof. Since, for every $i \in I$, $\text{Proj}_{\cap_{i \in I} C_i}(A\mathbf{x}) \in C_i$, (25) yields the upper bound

$$\text{dist}_D^2(\mathbf{x}) \leq \sum_{i \in I} \omega_i \|\mathbf{x}^i - \text{Proj}_{\cap_{i \in I} C_i}(A\mathbf{x})\|^2 \leq \sum_{i \in I} \omega_i R_i^2. \quad (40)$$

For the second bound, $\|\mathbf{x} - \mathbf{y}\|^2 = \sum_{i \in I} \omega_i \|\mathbf{x}^i - \mathbf{y}^i\|^2 \leq \sum_{i \in I} \omega_i R_i^2$. \square

3.1 Convex setting

Here we show that, if f is convex, Algorithm 1 achieves an $\mathcal{O}(\ln t / \sqrt{t})$ convergence rate in terms of the primal value gap of our subproblems (17). In tandem with Proposition 2.12, this establishes function value convergence. Unlike the Augmented Lagrangian approaches [2, 12, 23], our analysis does not require further assumptions concerning the relative interiors of $(C_i)_{i \in I}$, making it consistent with traditional Frank-Wolfe theory [8, Section 2.1].

Lemma 3.2 *Let f be convex and L_f -smooth, let D denote the diagonal subspace of \mathcal{H} , let $(C_i)_{i \in I}$ be a finite collection of nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty[$ such that $\cap_{i \in I} C_i \neq \emptyset$, and for every $\lambda \geq 0$, set $F_\lambda: \mathcal{H} \rightarrow]-\infty, +\infty]: \mathbf{x} \mapsto f(A\mathbf{x}) + \lambda \text{dist}_D^2(\mathbf{x})/2$, set $\mathbf{x}_t^* \in \text{Argmin}_{\mathbf{x} \in \cap_{i \in I} C_i} F_{\lambda_t}(\mathbf{x})$, and set $H_t = F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_t^*)$. Suppose that $(\lambda_t)_{t \in \mathbb{N}}$ is an increasing sequence. Then the iterates of Algorithm 1 satisfy*

$$H_{t+1} \leq (1 - \gamma_t)H_t + \frac{(\lambda_{t+1} - \lambda_t)}{2} \sum_{i \in I} \omega_i R_i^2 + \gamma_t^2 \frac{(\lambda_t + L_f)}{2} \sum_{i \in I} \omega_i R_i^2. \quad (41)$$

Proof. Let us begin by observing that F_{λ_t} is convex and $L_f + \lambda_t$ -smooth (cf. Fact 1.1). Since Algorithm 1 performs one step of the vanilla CG algorithm to (17), a standard CG argument [8] (relying on smoothness (7), Line 7 and Fact 2.2, then convexity (8)) shows

$$F_{\lambda_t}(\mathbf{x}_{t+1}) - F_{\lambda_t}(\mathbf{x}_t) \leq \gamma_t (F_{\lambda_t}(\mathbf{x}_t^*) - F_{\lambda_t}(\mathbf{x}_t)) + \gamma_t^2 \frac{L_f + \lambda_t}{2} \sum_{i \in I} \omega_i R_i^2. \quad (42)$$

Using Lemma 2.8, then adding $F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_t^*)$ to both sides of (42) reveals

$$H_{t+1} \leq F_{\lambda_{t+1}}(\mathbf{x}_{t+1}) - F_{\lambda_t}(\mathbf{x}_t^*) \quad (43)$$

$$= F_{\lambda_t}(\mathbf{x}_{t+1}) - F_{\lambda_t}(\mathbf{x}_t^*) + \frac{\lambda_{t+1} - \lambda_t}{2} \text{dist}_D^2(\mathbf{x}_{t+1}) \quad (44)$$

$$\leq (1 - \gamma_t)H_t + \frac{\lambda_{t+1} - \lambda_t}{2} \text{dist}_D^2(\mathbf{x}_{t+1}) + \gamma_t^2 \frac{L_f + \lambda_t}{2} \sum_{i \in I} \omega_i R_i^2. \quad (45)$$

Finally, Lemma 3.1 finishes the result. \square

Theorem 3.3 *In the setting of Lemma 3.2, for every $t \geq 0$ set $\gamma_t = 2/(\sqrt{t} + 2)$. Let $\lambda_0 > 0$ and for every $t \geq 1$ set $\lambda_{t+1} = \lambda_t + \lambda_0(\sqrt{t} + 2)^{-2}$. Then, for every $t \in \mathbb{N}$, the iterates of Algorithm 1 satisfy*

$$0 \leq H_t \leq 2 \sum_{i \in I} \omega_i R_i^2 \left(\frac{\lambda_0(2 \ln(\sqrt{t} + 2) + \frac{1}{4}) + L_f}{\sqrt{t} + 2} + \frac{4\lambda_0}{(\sqrt{t} + 2)^2} \right). \quad (46)$$

In particular, $F_{\lambda_t}(\mathbf{x}_t) \rightarrow \inf_{\mathbf{x} \in \cap_{i \in I} C_i} f(\mathbf{x})$ and $\text{dist}_D(\mathbf{x}_t) \rightarrow 0$. Furthermore, every accumulation point \mathbf{x}_∞ of $(\mathbf{x}_t)_{t \in \mathbb{N}}$ produces a solution $A\mathbf{x}_\infty \in \cap_{i \in I} C_i$ such that $f(A\mathbf{x}_\infty) = \inf_{\mathbf{x} \in \cap_{i \in I} C_i} f(\mathbf{x})$.

Proof. For notational convenience, set $R = \sum_{i \in I} \omega_i R_i^2$ and $\xi: \mathbb{R} \rightarrow \mathbb{R}: s \mapsto 2 \ln(\sqrt{s} + 2) + 4/(\sqrt{s} + 2)$. By calculus, for every $t \in \mathbb{N}$ such that $t \geq 1$, $\lambda_t - \lambda_0 \leq \lambda_0 \xi(t-1) - \lambda_0 \xi(1)$, so $\lambda_t \leq \lambda_0 \xi(t-1)$. We shall proceed by induction. The base case for $t = 0$ follows from (7), (9), and Lemma 3.1. Next, we suppose that (46) holds for $t \in \mathbb{N}$. Our inductive hypothesis, bound on λ_t , and (41) yield

$$H_{t+1} \leq (1 - \gamma_t) \left(2R \frac{\lambda_0 \xi(t) + L_f + \frac{\lambda_0}{4}}{\sqrt{t} + 2} \right) + \frac{\lambda_{t+1} - \lambda_t}{2} R + \gamma_t^2 \frac{(L_f + \lambda_0 \xi(t-1))R}{2} \quad (47)$$

$$= \frac{\sqrt{t}}{\sqrt{t} + 2} \left(2R \frac{\lambda_0 \xi(t) + L_f + \frac{\lambda_0}{4}}{\sqrt{t} + 2} \right) + 2R \left(\frac{L_f + \frac{\lambda_0}{4}}{(\sqrt{t} + 2)^2} + \frac{\lambda_0 \xi(t-1)}{(\sqrt{t} + 2)^2} \right) \quad (48)$$

$$\leq \frac{\sqrt{t} + 1}{(\sqrt{t} + 2)^2} (2R(L_f + \frac{\lambda_0}{4} + \lambda_0 \xi(t+1))) \quad (49)$$

$$\leq \frac{1}{\sqrt{t} + 1 + 2} \left(2R(L_f + \frac{\lambda_0}{4} + \lambda_0 \xi(t+1)) \right), \quad (50)$$

where (49) is because ξ is increasing and (50) is because $(\sqrt{t} + 1)(\sqrt{t} + 1 + 2) \leq (\sqrt{t} + 2)^2$. Having shown (46), we point out that Proposition 2.13 implies $\lim_{t \rightarrow \infty} F_{\lambda_t}(\mathbf{x}_t^*) = \inf_{x \in \bigcap_{i \in I} C_i} f(x)$. Hence $\lim_{t \rightarrow \infty} F_{\lambda_t}(\mathbf{x}_t)$ exists and, via (46), is equal to $\inf_{x \in \bigcap_{i \in I} C_i} f(x)$. Since $\lambda_t \rightarrow \infty$, it must be that $\text{dist}_D^2(\mathbf{x}_t) \rightarrow 0$. Therefore, every accumulation point $x_\infty \in \bigcap_{i \in I} C_i$ must also reside in D , so $A\mathbf{x}_\infty \in \bigcap_{i \in I} C_i$. Passing to a subsequence, since f is continuous we have

$$\inf_{x \in \bigcap_{i \in I} C_i} f(x) \leq f(A\mathbf{x}_\infty) = \lim_{k \rightarrow \infty} f(A\mathbf{x}_{t_k}) \leq \lim_{k \rightarrow \infty} F_{\lambda_{t_k}}(\mathbf{x}_{t_k}) = \inf_{x \in \bigcap_{i \in I} C_i} f(x). \quad (51)$$

□

Note that, although Theorem 3.3 shows convergence of the primal gaps of the subproblem (17), these gaps are never actually computed in practice, since \mathbf{x}_t^* is inaccessible. We also point out that, for the choice of $\lambda_0 = L_f$, our convergence rate becomes scale-invariant.

The convergence rate in Theorem 3.3 is atypical of CG algorithms with convex objective functions, because they usually have an $\mathcal{O}(1/t)$ convergence rate. This was achieved in the split-LMO setting under the condition $m = 2$ in [20, 22] and with a Slater-type condition in [12] by choosing stepsizes of magnitude $\gamma_t = \mathcal{O}(1/t)$. However, in order to achieve convergence in the proof of Theorem 3.3 with this larger stepsize, this would necessitate that $\lambda_{t+1} - \lambda_t \leq \mathcal{O}(1/t^2)$, i.e., $\lambda_t \not\rightarrow \infty$. Since Example 2.11 establishes that $\lambda_t \rightarrow \infty$ can be necessary (supported also by Proposition 2.12), we would no longer be able to show that the sequence of relaxed subproblems (17) converges to the original splitting problem (1). So, using a faster stepsize schedule would still yield a convergent algorithm, but it would not necessarily solve (1). We shall consider the topic of achieving a faster rate with extra assumptions in future work.

Remark 3.4 Without additional assumptions, Algorithm 1 does not guarantee iterate convergence of $(\mathbf{x}_t)_{t \in \mathbb{N}}$, which is consistent with other CG methods [5]. If, for instance, f is also μ -strongly convex, then Theorem 3.3 can be strengthened to provide convergence of the averages, because $A\mathbf{x}_t^*$ converges to the unique solution x^* of (1) and $0 \leq \mu \|A\mathbf{x}_t - A\mathbf{x}_t^*\|/2 \leq F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_t^*) \rightarrow 0$, so $A\mathbf{x}_t \rightarrow x^*$ as well.

3.2 Nonconvex setting

For CG methods which address (11) in the case when f is nonconvex, it is standard to show that the Frank-Wolfe gap at $x \in \mathcal{H}$, $G_{f,C}(x) := \sup_{v \in C} \langle \nabla f(x) | x - v \rangle$, converges to zero, because f is stationary at $x \in C$ whenever the F-W gap vanishes (9) [8]. Since F-W gaps are highly variable between iterations, convergence rates are typically derived for the average of F-W gaps. In this section, we consider the F-W gaps for our subproblems (17) which converge to (1) (in the sense of Proposition 2.12).

We begin by connecting the F-W gaps of our subproblems (17) to that of the original problem (1). In particular, for every $\lambda \geq 0$ the Frank-Wolfe gaps of our subproblems at $x \in \times_{i \in I} C_i$ provide an upper bound to both the penalty λdist_D^2 and the F-W gap of the original problem (1) at Ax . Interestingly, although $G_{F_\lambda, \times_{i \in I} C_i}(x_t) \geq 0$ is guaranteed, the F-W gap for the splitting problem (1), namely $G_{f, \cap_{i \in I} C_i}(Ax_t)$, may actually be negative since Ax_t is not guaranteed to reside in $\cap_{i \in I} C_i$ after a finite number of iterations.

Lemma 3.5 *Let f be smooth, set $\beta_f = \sup_{x \in \times_{i \in I} C_i} \|\nabla f(x)\|$, let $D \subset \mathcal{H}$ denote the diagonal subspace of \mathcal{H} , let $(C_i)_{i \in I}$ be a finite collection of nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty[$ such that $\cap_{i \in I} C_i \neq \emptyset$, and for every $\lambda \geq 0$, set $F_\lambda : \mathcal{H} \rightarrow]-\infty, +\infty] : x \mapsto f(Ax) + \lambda \text{dist}_D^2(x)/2$. Then, for every $x \in \mathcal{H}$,*

$$\sup_{v \in \times_{i \in I} C_i} \langle \nabla F_\lambda(x) | x - v \rangle \geq \sup_{v \in \cap_{i \in I} C_i} \langle \nabla f(Ax) | Ax - v \rangle + \lambda \text{dist}_D^2(x) \geq -\beta_f \sum_{i \in I} \omega_i R_i. \quad (52)$$

Proof. First, by infimizing over a subset of $\times_{i \in I} C_i$, we find

$$\inf_{v \in \times_{i \in I} C_i} \langle \nabla F_\lambda(x) | v - x \rangle = \inf_{v \in \times_{i \in I} C_i} \langle A^* \nabla f(Ax) + \lambda(x - A^* Ax) | v - x \rangle \quad (53)$$

$$\leq \inf_{v \in D \cap \times_{i \in I} C_i} \langle \nabla f(Ax) | Av - Ax \rangle + \lambda \langle x - A^* Ax | v - x \rangle. \quad (54)$$

Since $x - A^* Ax \in D^\perp$ and $A^* Ax \in D$, we have the following identity for every $v \in D$

$$\langle x - A^* Ax | v - x \rangle = \langle x - A^* Ax | -A^* Ax - (x - A^* Ax) \rangle = -\|x - A^* Ax\|^2. \quad (55)$$

So, using Proposition 2.1 for a change of variables, we set $p = \text{Proj}_{\cap_{i \in I} C_i}(Ax)$ to find that

$$\inf_{v \in \times_{i \in I} C_i} \langle \nabla F_\lambda(x) | v - x \rangle \leq \inf_{v \in \cap_{i \in I} C_i} \langle \nabla f(Ax) | v - Ax \rangle - \lambda \text{dist}_D^2(x) \quad (56)$$

$$\leq \langle \nabla f(Ax) | p - Ax \rangle - \lambda \text{dist}_D^2(x) \quad (57)$$

$$\leq \beta_f \text{dist}_{\cap_{i \in I} C_i}(Ax) - \lambda \text{dist}_D^2(x) \quad (58)$$

$$\leq \beta_f \sum_{i \in I} \omega_i R_i, \quad (59)$$

since $\text{dist}_{\cap_{i \in I} C_i}(Ax) = \|\sum_{i \in I} \omega_i(x^i - p)\| \leq \sum_{i \in I} \omega_i R_i$. Finally, negation yields (52). \square

With these results in-hand, we can now prove our main result.

Theorem 3.6 Let f be L_f -smooth, let $\mathbf{D} \subset \mathcal{H}$ denote the diagonal subspace of \mathcal{H} , let $(C_i)_{i \in I}$ be a finite collection of nonempty compact convex subsets of \mathcal{H} with diameters $\{R_i\}_{i \in I} \subset [0, +\infty[$ such that $\bigcap_{i \in I} C_i \neq \emptyset$, and for every $\lambda \geq 0$, set $F_\lambda: \mathcal{H} \rightarrow]-\infty, +\infty]: \mathbf{x} \mapsto f(A\mathbf{x}) + \lambda \text{dist}_{\mathbf{D}}^2(\mathbf{x})/2$. Set $\gamma_t = 1/\sqrt{t+1}$, let $\lambda_0 > 0$, and for every $t \geq 1$, set $\lambda_t = \lambda_0 \sum_{k=0}^{t-1} 1/(k+1)$. Then, for every $t \geq 1$, the iterates of Algorithm 1 satisfy³

$$0 \leq \frac{1}{t} \sum_{k=0}^{t-1} \sup_{\mathbf{v} \in \times_{i \in I} C_i} \langle \nabla F_{\lambda_k}(\mathbf{x}_k) \mid \mathbf{x}_k - \mathbf{v} \rangle \leq \mathcal{O} \left(\frac{\ln t}{\sqrt{t}} + \frac{1}{\sqrt{t}} \right). \quad (60)$$

In particular, there exists a subsequence $(t_k)_{k \in \mathbb{N}}$ such that $(\langle \nabla F_{\lambda_{t_k}}(\mathbf{x}_{t_k}) \mid \mathbf{x}_{t_k} - \mathbf{v}_{t_k} \rangle)_{k \in \mathbb{N}} \rightarrow 0$. Furthermore, every accumulation point \mathbf{x}_∞ of $(\mathbf{x}_{t_k})_{k \in \mathbb{N}}$ yields a stationary point $A\mathbf{x}_\infty \in \bigcap_{i \in I} C_i$ of the problem (1).

Proof. We begin by setting $\mathbf{v}_t = (\mathbf{v}_t^i)_{i \in I} \in \times_{i \in I} C_i$ and recalling that F_{λ_t} is $(L_f + \lambda_t)$ -smooth. For every $t \in \mathbb{N}$, let \mathbf{x}_t^* be a minimizer of F_{λ_t} over $\times_{i \in I} C_i$. For notational convenience, set $H_t = F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_t^*)$, $R = \sum_{i \in I} \omega_i R_i^2$, $R_A = \sum_{i \in I} \omega_i R_i$, and $B = \max\{\beta_p \sqrt{R}, R\}$. By the optimality of \mathbf{v}_t (Fact 2.2 and Line 7) and the smoothness inequality (7),

$$0 \leq \gamma_t \langle \nabla F_{\lambda_t}(\mathbf{x}_t) \mid \mathbf{x}_t - \mathbf{v}_t \rangle \leq F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_t}(\mathbf{x}_{t+1}) + \gamma_t^2 \frac{L_f + \lambda_t}{2} \|\mathbf{v}_t - \mathbf{x}_t\|^2. \quad (61)$$

So, using Lemma 2.8 and Lemma 3.1 twice,

$$0 \leq \langle \nabla F_{\lambda_t}(\mathbf{x}_t) \mid \mathbf{x}_t - \mathbf{v}_t \rangle \quad (62)$$

$$\leq \frac{F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_{t+1}}(\mathbf{x}_{t+1})}{\gamma_t} + \frac{\lambda_{t+1} - \lambda_t}{\gamma_t} \text{dist}_{\mathbf{D}}^2(\mathbf{x}_{t+1}) + \gamma_t \frac{L_f + \lambda_t}{2} R \quad (63)$$

$$\leq \frac{F_{\lambda_t}(\mathbf{x}_t) - F_{\lambda_{t+1}}(\mathbf{x}_{t+1})}{\gamma_t} + \frac{\lambda_{t+1} - \lambda_t}{\gamma_t} R + \gamma_t \frac{L_f + \lambda_t}{2} R. \quad (64)$$

Furthermore, since f and $\text{dist}_{\mathbf{D}}^2/2$ are smooth and $\sum_{i \in I} \omega_i C_i$ and $\times_{i \in I} C_i$ are compact, it follows that their gradients are bounded. Hence, f and $\text{dist}_{\mathbf{D}}^2/2$ are Lipschitz continuous on these sets, with constants $\beta_f := \sup_{c \in \sum_{i \in I} \omega_i C_i} \|\nabla f(c)\|$ and $\beta_p := \sup_{c \in \times_{i \in I} C_i} \|\nabla \text{dist}_{\mathbf{D}}^2(c)/2\|$ respectively. Therefore, we find that by Jensen's inequality and Lemma 3.1,

$$H_t \leq \beta_f \|A\mathbf{x}_t - A\mathbf{x}_t^*\| + \lambda_t \beta_p \|\mathbf{x}_t - \mathbf{x}_t^*\| \leq \beta_f \sum_{i \in I} \omega_i R_i + \lambda_t \beta_p \sqrt{\sum_{i \in I} \omega_i R_i^2} = \beta_f R_A + \lambda_t \beta_p \sqrt{R}. \quad (65)$$

³Precise constants are in (77).

By Lemma 2.8, we have $\gamma_t^{-1}(F_{\lambda_{t+1}}(\mathbf{x}_{t+1}^*) - F_{\lambda_t}(\mathbf{x}_t^*)) \geq 0$. Combining all of these facts, we find

$$0 \leq \sum_{k=0}^{t-1} \left\langle \nabla F_{\lambda_k}(\mathbf{x}_k) \mid \mathbf{x}_k - \mathbf{v}_k \right\rangle \quad (66)$$

$$\leq \sum_{k=0}^{t-1} \left(\frac{F_{\lambda_k}(\mathbf{x}_k) - F_{\lambda_{k+1}}(\mathbf{x}_{k+1})}{\gamma_k} + \frac{\lambda_{k+1} - \lambda_k}{\gamma_k} R + \gamma_k \frac{L_f + \lambda_k}{2} R \right) \quad (67)$$

$$\leq \sum_{k=0}^{t-1} \left(\frac{H_k - H_{k+1}}{\gamma_k} + \frac{\lambda_{k+1} - \lambda_k}{\gamma_k} R + \gamma_k \frac{L_f + \lambda_k}{2} R \right) \quad (68)$$

$$= \frac{H_0}{\gamma_0} - \frac{H_t}{\gamma_{t-1}} + \sum_{k=1}^{t-1} \left(\frac{1}{\gamma_k} - \frac{1}{\gamma_{k-1}} \right) H_k + \sum_{k=0}^{t-1} \left(\frac{\lambda_{k+1} - \lambda_k}{\gamma_k} R + \gamma_k \frac{L_f + \lambda_k}{2} R \right) \quad (69)$$

$$\leq \frac{\beta_f R_A + \lambda_0 \beta_p \sqrt{R}}{\gamma_0} + \sum_{k=1}^{t-1} \left(\frac{1}{\gamma_k} - \frac{1}{\gamma_{k-1}} \right) (\beta_f R_A + \lambda_k \beta_p \sqrt{R}) \quad (70)$$

$$+ \sum_{k=0}^{t-1} \left(\frac{\lambda_{k+1} - \lambda_k}{\gamma_k} R + \gamma_k \frac{L_f + \lambda_k}{2} R \right) \quad (71)$$

$$\leq \frac{\beta_f R_A + \lambda_0 B}{\gamma_0} + \sum_{k=1}^{t-1} \left(\frac{1}{\gamma_k} - \frac{1}{\gamma_{k-1}} \right) (\beta_f R_A + \lambda_k B) + \sum_{k=0}^{t-1} \frac{\lambda_{k+1} - \lambda_k}{\gamma_k} B$$

$$+ \sum_{k=0}^{t-1} \gamma_k \frac{L_f + \lambda_k}{2} R$$

$$= \frac{\beta_f R_A + \lambda_t B}{\gamma_{t-1}} + \sum_{k=0}^{t-1} \gamma_k \frac{L_f + \lambda_k}{2} R, \quad (72)$$

where we use Lemma 2.8 in (68), drop a negative term and use (65) in (70), and simplify in (72). Next, we note that $\sum_{k=0}^{t-1} \gamma_k \leq 2\sqrt{t}$ and $\lambda_t \leq \lambda_0(\ln(t+1) + 1)$, so

$$0 \leq \frac{1}{t} \sum_{k=0}^{t-1} \left\langle \nabla F_{\lambda_k}(\mathbf{x}_k) \mid \mathbf{x}_k - \mathbf{v}_k \right\rangle \quad (73)$$

$$\leq \frac{\beta_f R_A + \lambda_t B}{\sqrt{t}} + \frac{1}{t} \sum_{k=0}^{t-1} \gamma_k \frac{L_f + \lambda_k}{2} R \quad (74)$$

$$\leq \frac{\beta_f R_A + \lambda_t B}{\sqrt{t}} + \frac{1}{t} (L_f + \lambda_{t-1}) \sum_{k=0}^{t-1} \gamma_k \frac{1}{2} R \quad (75)$$

$$\leq \frac{\beta_f R_A + \lambda_0(\ln(t+1) + 1)B}{\sqrt{t}} + \frac{1}{\sqrt{t}} (L_f + \lambda_0(\ln(t) + 1)) R \quad (76)$$

$$\leq \frac{1}{\sqrt{t}} \left(\beta_f \sum_{i \in I} \omega_i R_i + (L_f + \lambda_0) \sum_{i \in I} \omega_i R_i^2 + \lambda_0 B \right) + \frac{\ln(t+1)}{\sqrt{t}} \lambda_0 \left(\sum_{i \in I} \omega_i R_i^2 + B \right), \quad (77)$$

which establishes (60). Since the Frank-Wolfe gaps $(\langle \nabla F_{\lambda_t}(\mathbf{x}_t) \mid \mathbf{x}_t - \mathbf{v}_t \rangle)_{t \in \mathbb{N}}$ are positive and the sequence of averages goes to zero, the existence of a subsequence $(t_k)_{k \in \mathbb{N}}$ such that $\langle \nabla F_{\lambda_{t_k}}(\mathbf{x}_{t_k}) \mid \mathbf{x}_{t_k} - \mathbf{v}_{t_k} \rangle \rightarrow 0$ follows. Lemma 3.5 implies that

$$\left(\sup_{v \in \bigcap_{i \in I} C_i} \langle \nabla f(A\mathbf{x}_{t_k}) \mid A\mathbf{x}_{t_k} - v \rangle + \lambda_{t_k} \text{dist}_D^2(\mathbf{x}_{t_k}) \right)_{k \in \mathbb{N}} \quad (78)$$

is bounded. So, since $\lambda_{t_k} \rightarrow \infty$, we must have $\text{dist}_D^2(\mathbf{x}_{t_k}) \rightarrow 0$. Therefore, for every accumulation point \mathbf{x}_∞ of $(\mathbf{x}_{t_k})_{k \in \mathbb{N}}$, $\mathbf{x}_\infty \in D \cap \bigcap_{i \in I} C_i$, so $A\mathbf{x}_\infty \in \bigcap_{i \in I} C_i$ and

$$0 \leq \sup_{v \in \bigcap_{i \in I} C_i} \langle \nabla f(A\mathbf{x}_\infty) \mid A\mathbf{x}_\infty - v \rangle. \quad (79)$$

Finally, we can bound the gap above using continuity and Lemma 3.5:

$$\sup_{v \in \bigcap_{i \in I} C_i} \langle \nabla f(A\mathbf{x}_\infty) \mid A\mathbf{x}_\infty - v \rangle \leq \limsup_{k \rightarrow \infty} \left(\sup_{v \in \bigcap_{i \in I} C_i} \langle \nabla f(A\mathbf{x}_{t_k}) \mid A\mathbf{x}_{t_k} - v \rangle \right) \quad (80)$$

$$\leq \limsup_{k \rightarrow \infty} \left(\langle \nabla F_{\lambda_{t_k}}(\mathbf{x}_{t_k}) \mid \mathbf{x}_{t_k} - v_{t_k} \rangle \right) \quad (81)$$

$$= 0. \quad (82)$$

Since $G_{f, \bigcap_{i \in I} C_i}(A\mathbf{x}_\infty) = 0$, we conclude from (9) that $A\mathbf{x}_\infty$ is a stationary point. \square

Remark 3.7 We emphasize that, for the cost of one extra inner product, the Frank-Wolfe gap $\langle \nabla F_{\lambda_t}(\mathbf{x}_t) \mid \mathbf{x}_t - v_t \rangle$ can be computed while Algorithm 1 is running. So, checking for stationarity in the subproblems (17) is tractable in practice. Also, similarly to the convex-case, the choice of $\lambda_0 = L_f$ makes our convergence rate in (77) scale-invariant.

4 Conclusion and Future Work

Theorem 3.6 appears to be the first convergence guarantee for solving (1) in the nonconvex split-LMO setting. Furthermore, our rate of convergence is only one log factor less than the rate of CG for one constraint ($m = 1$) [6]. While it is unclear if this log factor can be removed for the nonconvex setting, we believe that the analysis for the convex rate can be improved since typically the nonconvex average-FW-gap rate is quadratically slower than the convex primal gap rate [6]. This speed-up has been achieved in some settings with algorithms which require one LMO call per iteration [12, 22], but it appears that the question of whether or not $\mathcal{O}(1/t)$ convergence is possible in the split-LMO setting without additional assumptions remains open.

In addition to the question above, there are several interesting theoretical and numerical investigations to be performed. One topic is the use of alternative stepsizes and penalty parameter schedules. The proofs of Theorems 3.3 and 3.6 can easily be extended to a short-step selection for γ_t similar to [6] by minimizing the upper bound arising from (7). Another direction is investigating Algorithm 1 under additional assumptions on the objective or constraints. For instance, CG algorithms possess accelerated convergence rates when the objective function or constraints are strongly convex [33, 34]; extending this analysis to Algorithm 1 is also a topic of future interest. Many projection-based splitting methods have an advantage of being block-iterative, i.e., instead of requiring a computation for all constraints indexed by I (as is required in the for loop in Algorithm 1, Line 6) at every iteration t , only a subset $I_t \subset I$ of updates are performed. This can significantly reduce the computational load per iteration, and block-iterative projection methods enjoy convergence under very mild assumptions on the blocks $(I_t)_{t \in \mathbb{N}}$ [3, 4]. It is worth noting that the inner loop of Algorithm 1 can be parallelized, and a block-iterative capability would further improve the per-iteration cost. Several LMO-based block-iterative algorithms have been proposed for solving problems like the relaxation (17) [35, 36], but extending them to solve (1) remains to be done.

Acknowledgements

The work for this article has been supported by MODAL-Synlab, and took place on the Research Campus MODAL funded by the German Federal Ministry of Education and Research

(BMBF) (fund numbers 05M14ZAM, 05M20ZBM). This research was also supported by the DFG Cluster of Excellence MATH+ (EXC-2046/1, project ID 390685689) funded by the Deutsche Forschungsgemeinschaft (DFG).

We thank Kamiar Asgari, Gábor Braun, Mathieu Besançon, Ibrahim Ozaslan, Christophe Roux, Antonio Silveti-Falls, David Martínez-Rubio, and Elias Wirth for their valuable feedback and discussions.

References

- [1] Y. Censor, A. Cegielski, Projection methods: an annotated bibliography of books and reviews, *Optimization* 64 (11) (2015) 2343–2358. doi:[10.1080/02331934.2014.957701](https://doi.org/10.1080/02331934.2014.957701).
- [2] A. Silveti-Falls, C. Molinari, J. Fadili, Generalized conditional gradient with augmented Lagrangian for composite minimization, *SIAM J. Optim.* 30 (4) (2020) 2687–2725. doi:[10.1137/19M1240460](https://doi.org/10.1137/19M1240460).
- [3] P. L. Combettes, J.-C. Pesquet, *Proximal Splitting Methods in Signal Processing*, Springer New York, New York, NY, 2011. doi:[10.1007/978-1-4419-9569-8_10](https://doi.org/10.1007/978-1-4419-9569-8_10).
- [4] P. L. Combettes, Z. C. Woodstock, Reconstruction of functions from prescribed proximal points, *J. Approx. Theory* 268 (2021) 105606. doi:<https://doi.org/10.1016/j.jat.2021.105606>.
- [5] J. Bolte, C. W. Combettes, E. Pauwels, The iterates of the Frank-Wolfe algorithm may not converge, *Math. Oper. Res.* (to appear).
- [6] F. Pedregosa, G. Negiar, A. Askari, M. Jaggi, Linearly convergent Frank-Wolfe with backtracking line-search, in: *International conference on artificial intelligence and statistics*, PMLR, 2020, pp. 1–10.
- [7] T. Rothvoss, The matching polytope has exponential extension complexity, *J. ACM* 64 (6) (2017). doi:[10.1145/3127497](https://doi.org/10.1145/3127497).
- [8] G. Braun, A. Carderera, C. Combettes, H. Hassani, A. Karbasi, A. Mokhtari, S. Pokutta, Conditional gradient methods (2022). [arXiv:2211.14103](https://arxiv.org/abs/2211.14103).
- [9] C. W. Combettes, S. Pokutta, Complexity of linear minimization and projection on some sets, *Oper. Res. Lett.* 49 (4) (2021) 565–571.
- [10] D. Garber, A. Kaplan, S. Sabach, Improved complexities of conditional gradient-type methods with applications to robust matrix recovery problems, *Math. Program.* 186 (2021) 185–208.
- [11] T. Ding, D. Lim, R. Vidal, B. D. Haeffele, Understanding doubly stochastic clustering, in: K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, S. Sabato (Eds.), *Proceedings of the 39th International Conference on Machine Learning*, Vol. 162 of *Proceedings of Machine Learning Research*, PMLR, 2022, pp. 5153–5165.
- [12] G. Gidel, F. Pedregosa, S. Lacoste-Julien, Frank-Wolfe splitting via augmented Lagrangian method, in: A. Storkey, F. Perez-Cruz (Eds.), *Proceedings of the 21st International Conference on Artificial Intelligence and Statistics*, Vol. 84 of *Proceedings of Machine Learning Research*, PMLR, 2018, pp. 1456–1465.
- [13] E. Richard, P. Savalle, N. Vayatis, Estimation of simultaneously sparse and low rank matrices, in: *Proceedings of the 29th International Conference on Machine Learning*, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012, [icml.cc / Omnipress](http://icml.cc/), 2012.
- [14] Z. Yang, J. Corander, E. Oja, Low-rank doubly stochastic matrix decomposition for cluster analysis, *J. Mach. Learn. Res.* 17 (1) (2016) 6454–6478.
- [15] Z. Zhang, Z. Zhai, L. Li, Graph refinement via simultaneously low-rank and sparse approximation, *SIAM J. Sci. Comput.* 44 (3) (2022) A1525–A1553.
- [16] N. He, Z. Harchaoui, Semi-proximal mirror-prox for nonsmooth composite minimization, in: C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, R. Garnett (Eds.), *Advances in Neural Information Processing Systems*, Vol. 28, Curran Associates, Inc., 2015.
- [17] V. Kolmogorov, T. Pock, One-sided Frank-Wolfe algorithms for saddle problems, in: M. Meila,

- T. Zhang (Eds.), Proceedings of the 38th International Conference on Machine Learning, Vol. 139 of Proceedings of Machine Learning Research, PMLR, 2021, pp. 5665–5675.
- [18] Y.-F. Liu, X. Liu, S. Ma, On the nonergodic convergence rate of an inexact augmented Lagrangian framework for composite convex programming, *Math. Oper. Res.* 44 (2) (2019) 632–650.
 - [19] R. D. Millán, O. P. Ferreira, L. F. Prudente, Alternating conditional gradient method for convex feasibility problems 80 (1) (2021) 245–269. [doi:10.1080/10556788.2013.796683](https://doi.org/10.1080/10556788.2013.796683).
 - [20] G. Braun, S. Pokutta, R. Weismantel, Alternating linear minimization: Revisiting von Neumann’s alternating projections (2022). [arXiv:2212.02933](https://arxiv.org/abs/2212.02933).
 - [21] G. Lan, E. Romeijn, Z. Zhou, Conditional gradient methods for convex optimization with general affine and nonlinear constraints, *SIAM J. Optim.* 31 (3) (2021) 2307–2339.
 - [22] C. Mu, Y. Zhang, J. Wright, D. Goldfarb, Scalable robust matrix recovery: Frank-Wolfe meets proximal methods, *SIAM J. Sci. Comput.* 38 (5) (2016) A3291–A3317.
 - [23] A. Yurtsever, O. Fercoq, V. Cevher, A conditional-gradient-based augmented Lagrangian framework, in: K. Chaudhuri, R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning, Vol. 97 of Proceedings of Machine Learning Research, PMLR, 2019, pp. 7272–7281.
 - [24] A. Migdalas, A regularization of the Frank–Wolfe method and unification of certain nonlinear programming methods, *Math. Program.* 65 (1994) 331–345.
 - [25] I. Chrysosoverghi, A. Bacopoulos, B. Kokkinis, J. Coletsos, Mixed Frank–Wolfe penalty method with applications to nonconvex optimal control problems, *J. Optim. Theory Appl.* 94 (1997) 311–334.
 - [26] R. T. Rockafellar, R. J.-B. Wets, *Variational Analysis*, Vol. 317, Springer Science & Business Media, 2009.
 - [27] G. Pierra, Decomposition through formalization in a product space, *Math. Program.* 28 (1984) 96–115.
 - [28] H. H. Bauschke, P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, 2nd ed., Springer, 2017.
 - [29] K. Lange, J.-H. Won, J. Xu, Projection onto Minkowski sums with application to constrained learning, in: K. Chaudhuri, R. Salakhutdinov (Eds.), Proceedings of the 36th International Conference on Machine Learning, Vol. 97 of Proceedings of Machine Learning Research, PMLR, 2019, pp. 3642–3651.
 - [30] L. L. Duan, A. L. Young, A. Nishimura, D. B. Dunson, Bayesian constraint relaxation, *Biometrika* 107 (1) (2019) 191–204. [doi:10.1093/biomet/asz069](https://doi.org/10.1093/biomet/asz069).
 - [31] T. Bernholt, F. Eisenbrand, T. Hofmeister, Constrained Minkowski sums: A geometric framework for solving interval problems in computational biology efficiently, *Discrete Comput. Geom.* 42 (1) (2009) 22–36.
 - [32] T. Lozano-Pérez, M. A. Wesley, An algorithm for planning collision-free paths among polyhedral obstacles, *Commun. ACM* 22 (10) (1979) 560–570.
 - [33] D. Garber, E. Hazan, Faster rates for the Frank-Wolfe method over strongly-convex sets, in: F. Bach, D. Blei (Eds.), Proceedings of the 32nd International Conference on Machine Learning, Vol. 37 of Proceedings of Machine Learning Research, PMLR, Lille, France, 2015, pp. 541–549.
 - [34] E. Wirth, T. Kerdreux, S. Pokutta, Acceleration of Frank-Wolfe algorithms with open-loop step-sizes, in: F. Ruiz, J. Dy, J.-W. van de Meent (Eds.), Proceedings of The 26th International Conference on Artificial Intelligence and Statistics, Vol. 206 of Proceedings of Machine Learning Research, PMLR, 2023, pp. 77–100.
 - [35] A. Beck, E. Pauwels, S. Sabach, The cyclic block conditional gradient method for convex optimization problems, *SIAM J. Optim.* 25 (4) (2015) 2024–2049.
 - [36] I. Bomze, F. Rinaldi, D. Zeffiro, Projection free methods on product domains (2023). [arXiv:2302.04839](https://arxiv.org/abs/2302.04839).