

Frank-Wolfe Methods with an Unbounded Feasible Region and Applications to Structured Learning

Haoyue Wang* Haihao Lu† Rahul Mazumder‡

Abstract

The Frank-Wolfe (FW) method is a popular algorithm for solving large-scale convex optimization problems appearing in structured statistical learning. However, the traditional Frank-Wolfe method can only be applied when the feasible region is bounded, which limits its applicability in practice. Motivated by two applications in statistical learning, the ℓ_1 trend filtering problem and matrix optimization problems with generalized nuclear norm constraints, we study a family of convex optimization problems where the unbounded feasible region is the direct sum of an unbounded linear subspace and a bounded constraint set. We propose two new Frank-Wolfe methods: unbounded Frank-Wolfe method (uFW) and unbounded Away-Step Frank-Wolfe method (uAFW), for solving a family of convex optimization problems with this class of unbounded feasible regions. We show that under proper regularity conditions, the unbounded Frank-Wolfe method has a $O(1/k)$ sublinear convergence rate, and unbounded Away-Step Frank-Wolfe method has a linear convergence rate, matching the best-known results for the Frank-Wolfe method when the feasible region is bounded. Furthermore, computational experiments indicate that our proposed methods appear to outperform alternative solvers.

1 Introduction

The Frank-Wolfe (FW) method [17], also known as the conditional gradient method [11], is a well-studied first-order algorithm for smooth convex optimization with a bounded feasible region. Compared to other first-order methods, such as projected gradient methods and proximal type methods [39], where a projection operation onto the feasible set is required at every iteration, the Frank-Wolfe method avoids projection by minimizing a linear objective over the feasible set. Solving this problem is often computationally more attractive than a projection step in several large-scale problems arising in machine learning. For instance, when the constraint set S is polyhedral, the linear subproblem is given by a linear program, which may be a computationally friendlier alternative compared to solving a convex quadratic program in the projection step. When the constraint set S is a nuclear norm ball, the solution to the linear subproblem can be obtained by computing the leading singular vector/value pair, while the projection onto S may require computing several leading singular vectors/values. Due to its computational efficiency and projection-free nature, the Frank-Wolfe method has emerged as a popular choice for solving large-scale convex optimization problems arising in machine learning applications [28].

In many real-world applications however, the feasible region of the optimization problem may be unbounded, which limits the applicability of Frank-Wolfe methods. We present two motivating examples from statistical learning/high-dimensional statistics. In the generalized lasso problem [50], the problem of interest is to solve

$$\min_{x \in \mathbb{R}^n} \|b - Ax\|_2^2 \quad \text{s.t.} \quad \|Hx\|_1 \leq \delta, \quad (1)$$

*MIT Operations Research Center (email: haoyue@mit.edu).

†Booth School of Business, University of Chicago (email: haihao.lu@chicagobooth.edu).

‡MIT Sloan School of Management, Operations Research Center and MIT Center for Statistics (email: rahulmaz@mit.edu). This research was partially supported by awards from the Office of Naval Research ONR-N000141812298 (Young Investigator Award), the National Science Foundation (NSF-IIS-1718258), MIT-IBM Watson AI Lab to Rahul Mazumder.

where $x \in \mathbb{R}^n$ are the model coefficients or signal of interest (decision variable), $A \in \mathbb{R}^{N \times n}$ is the model matrix, $b \in \mathbb{R}^N$ is the response, and $H \in \mathbb{R}^{m \times n}$ is a general (usually non-square or singular) matrix imposing additional structure on the unknown signal x . A special case is the ℓ_1 trend filtering problem [49, 51, 30], where the matrix H is the r -th order discrete derivative matrix with $r \geq 1$. In generalized nuclear norm regularization problems [16], we are interested in the following problem

$$\min_{X \in \mathbb{R}^{m \times n}} f(X) \quad \text{s.t.} \quad \|PXQ\|_* \leq \delta \quad (2)$$

where, the objective f is a smooth convex function of the matrix variable $X \in \mathbb{R}^{m \times n}$, i.e., the gradient ∇f is Lipschitz continuous; $\|\cdot\|_*$ denotes the nuclear norm (i.e. the sum of singular values) of a matrix, and P, Q are general (usually singular) matrices. In example (2), the feasible region is unbounded when matrices P and Q are not full rank. Thus, the traditional Frank-Wolfe method is no longer applicable for these problems, even though the Frank-Wolfe method is known to work well for optimization problems involving the ℓ_1 norm ball or nuclear norm ball constraints [27, 19].

Notice that in Problems (1) and (2), the feasible regions can be expressed as the direct sum of a linear subspace T and a bounded set S . For example, the constraint set $\{x \mid \|Hx\|_1 \leq \delta\}$ in (1) can be formulated as $T \oplus S := \{x \mid x = s + t, t \in T, s \in S\}$, where $T = \ker(H)$ is a linear subspace, and $S = \{x \in \mathbb{R}^n \mid \|Hx\|_1 \leq \delta, x \in (\ker(H))^\perp\}$ is a bounded set in \mathbb{R}^n . The constraint $\{X \mid \|PXQ\|_* \leq \delta\}$ in (2) can also be expressed as $T \oplus S$ with a linear subspace T and a bounded set S (see Section 4.2 for details). Thus motivated, in this paper, we study the following smooth constrained convex optimization problem over an unbounded feasible region $T \oplus S$:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad x \in T \oplus S \quad (3)$$

where f is a smooth convex function on \mathbb{R}^n , T is a linear subspace of \mathbb{R}^n , S is a bounded convex set ($S \subset T^\perp$), and \oplus is the direct sum of two orthogonal spaces.

The major contribution of this paper is to generalize the traditional Frank-Wolfe method to solve (3) with computational guarantees. In Section 2, we present two new algorithms designed for (3): the unbounded Frank-Wolfe Method (uFW) and the unbounded Away-Step Frank-Wolfe Method (uAFW). We introduce new curvature constants generalizing those arising in the analysis of Frank-Wolfe methods over a bounded feasible region. Our key idea is to alternate between a Frank-Wolfe step along the bounded set S where we solve a linear subproblem; and a gradient descent step along the unbounded linear subspace T which admits a simple projection operation. In Section 3, we show that under suitable conditions, uFW converges to an optimal solution of (3) with a sublinear-rate $O(1/k)$, and uAFW converges to an optimal solution of (3) with a linear rate. In Section 4, we discuss how to apply our proposed algorithms to Problems (1) and (2), focusing on how to make the computations efficient (by exploiting problem structure). Section 5 presents numerical experiments suggesting that our proposal outperforms alternatives by a significant margin.

1.1 Related literature

As mentioned earlier, the Frank-Wolfe method has been extensively studied in the optimization community. The original Frank-Wolfe method, dating back to Frank and Wolfe [17], was designed to solve a smooth convex optimization problem over a polytope. The method was then extended to more general settings with a convex bounded feasible region, and was shown to attain an $O(1/k)$ sublinear rate of convergence [14, 13, 44]. Indeed, the $O(1/k)$ sublinear rate matches the lower complexity bound for solving a generic constrained convex optimization problem [32]. Recently, due to its projection-free nature, there has been renewed interest in the Frank-Wolfe method for solving large-scale optimization problems arising in structured statistical learning. Compared to other first-order methods (e.g., proximal gradient), which usually require us to compute a projection onto the constraint set (at every iteration), the Frank-Wolfe method solves a linear subproblem, which can be simpler than the projection problem in many applications [28, 26, 27, 10].

Recent works have explored different properties of the FW method and its variants. For example, [28] studies the invariance of the Frank-Wolfe method under linear transformations, and introduces a curvature

constant that improves the complexity analysis; [34] proposes gradient sliding schemes which achieve the lower complexity bound on both the number of linear optimization oracle calls as well as gradient computations; [5] proves support identification in finite time of two variants of the FW method on the standard simplex, and later work [6] develops complexity bounds for support identification. See for example, [33, Chapter 7] for a nice overview on the recent developments of the Frank-Wolfe method.

Recently, [22] proposes an interesting Frank-Wolfe type method to address problems with unbounded constraint sets, but their approach differs from what we propose. [22] use accelerated gradient descent as their base algorithm and use FW to approximately solve the projection step. This is a two-loop algorithm, which behaves similar to accelerated gradient descent. In contrast, our method is a single-loop algorithm, and resembles a Frank-Wolfe method. To deal with the unbounded constraint in each subproblem, [22] performs a Frank-Wolfe step on the intersection of $S \oplus T$ and an Euclidean ball centered at the initial point. This paper reports complexities of $O(1/\sqrt{\epsilon})$ for the number of gradient evaluations and $O(1/\epsilon)$ for the number of linear oracles in order to find a solution with primal optimality gap ϵ . Although their complexity bound on the number of gradient evaluations is better, each gradient evaluation is accompanied by solving many linear oracles. While the number of linear oracles is the same rate as ours, the linear oracles are different (due to the presence of an additional ball constraint) and can be harder to solve than the ones arising in our FW procedure. In addition, for the setting with a strongly convex objective function and polyhedral constraints, we prove linear convergence results—similar results are not discussed in [22]. Our numerical experiments (cf. Section 5) appear to suggest that our proposal is computationally more attractive compared to the proposal of [22].

The Away-step Frank-Wolfe method (AFW) is a variant of FW first proposed by Wolfe in 1970s [53], that attempts to ameliorate the so-called zigzagging behavior of the vanilla FW method. Later, [23] proposes a modified AFW that has a linear convergence guarantee for strongly convex objective function and under a strict complementarity condition. Recently, [20, 31] propose new variants of AFW which are linearly convergent for strongly convex objective function on a polytope. The version of AFW in [31] has been further investigated by [4, 43, 42] with a different analysis technique. In this paper, we generalize the version of AFW in [31] to the structured unbounded setting (3) and prove a linear convergence rate when the objective function is strongly convex and the constraint set is a polyhedron.

The efficacy of the Frank-Wolfe method depends upon how easily one can solve the linear optimization oracle. A recent survey [28] lists several constraint sets arising in modern applications that admit efficient linear optimization oracles. Examples include the ℓ_1 -norm ball [10] and the nuclear norm ball [27, 19], among others. Here we show that Frank-Wolfe type methods can also be applied to problems (1) and (2), where the constraint sets are transformed/unbounded versions of the ℓ_1 -norm ball and nuclear norm ball, respectively.

A special case of (1) is the ℓ_1 trend filtering problem [30]. Statistical and computational aspects of this problem have been studied in [30, 51, 50, 52]. Currently, popular algorithms for the ℓ_1 trend filtering problem include interior point methods [30], ADMM-based algorithms [52, 45], and parametric quadratic programming [50]. The generalized nuclear norm regularized matrix estimation problem (2), has important applications in computer vision [3], collaborative filtering [48], and matrix completion with side information and missing data [16, 15, 9]. To our knowledge, Frank-Wolfe type methods are yet to be explored for problems (1) and (2)—bridging this gap is the focus of our paper.

Notation: We let \mathbb{R}^n denote the n -dimensional Euclidean space, and $\mathbb{R}^{m \times n}$ the space of $m \times n$ real matrices. For $x, y \in \mathbb{R}^n$, $\langle x, y \rangle = \sum_{i=1}^n x_i y_i$ denotes the Euclidean inner product of x and y . Let $\mathbf{1}_n$ denote the vector in \mathbb{R}^n with all coordinates being 1. Let $\mathbf{0}_n$ denote the vector in \mathbb{R}^n with all coordinates being 0, and $\mathbf{0}_{m \times n}$ be the matrix in $\mathbb{R}^{m \times n}$ with all elements being 0. Let e_i be the unit vector in \mathbb{R}^n with the i -th coordinate being 1 and all other coordinates being 0. Let I_n be the identity matrix in $\mathbb{R}^{n \times n}$. For $A, B \in \mathbb{R}^{m \times n}$, $\langle A, B \rangle := \sum_{i=1}^m \sum_{j=1}^n A_{i,j} B_{i,j}$ defines the inner product of A and B , and $A \otimes B$ denotes their Kronecker product. For a function f , $\text{dom}(f)$ denotes the domain of f . For a given linear subspace $T \subset \mathbb{R}^n$, we let $\dim(T)$ denote the dimension of T , and T^\perp its orthogonal subspace. Furthermore, \mathcal{P}_T and \mathcal{P}_T^\perp denote the

projection operators onto the linear subspaces T and its orthogonal subspace T^\perp (respectively). For two subsets $T, S \subset \mathbb{R}^n$, $T \oplus S$ denotes the direct sum of the two subsets, i.e., $T \oplus S = \{x + y \mid x \in T, y \in S\}$. A function f is said to be L -smooth over a set Q if it is differentiable and $\|\nabla f(u) - \nabla f(v)\|_2 \leq L\|u - v\|_2$ for all $u, v \in Q$.

2 Unbounded Frank-Wolfe Algorithms

In this section, we present two new Frank-Wolfe algorithms (with and without away-step) designed to solve (3) with an unbounded feasible region. As mentioned above, we alternate between performing a Frank-Wolfe step on the bounded set S and a gradient descent step along the subspace T . Since the gradient descent direction can also be viewed as an extreme ray to the solution of the Frank-Wolfe linear subproblem along the linear subspace T , we call this the ‘‘Unbounded (Away Step) Frank-Wolfe Method’’.

Algorithm 1 presents our first algorithm, the unbounded Frank-Wolfe method, for solving (3). In each iteration, we first perform a gradient descent step along the unbounded subspace T . To do so, we first compute the negative projected gradient onto T , namely $-\mathcal{P}_T \nabla f(x^k)$, and move the current solution towards this direction with step-size η . While the projection onto a bounded set (e.g., S) can be expensive, the projection onto a linear subspace is usually much cheaper, and the solution can be computed in closed-form. Moreover, for ℓ_1 trend filtering, the subspace T is in fairly low dimension, making the projection step even simpler. Section 4 shows that a solution s^k of the linear subproblem within the bounded set S , can be solved efficiently for these two motivating examples by making use of problem structure. Finally, we perform the Frank-Wolfe step within set S using step-size α_k by moving along the direction $s^k - \mathcal{P}_T^\perp x^k$, i.e., the direction between the resulting extreme point and the current solution projected onto set S (note that $\mathcal{P}_T^\perp x^k = \mathcal{P}_T^\perp y^k$ as $\mathcal{P}_T^\perp \mathcal{P}_T = 0$). Different step-size rules have been explored in the literature of the Frank-Wolfe method [18]. Here, we study two step-size rules:

- Line search step-size rule:

$$\alpha_k \in \underset{\alpha \in [0,1]}{\operatorname{argmin}} f(y^k + \alpha(s^k - \mathcal{P}_T^\perp x^k)). \quad (4)$$

- Simple step-size rule:

$$\alpha_k = \frac{2}{k+2} \text{ if } f(y^k + \frac{2}{k+2}(s^k - \mathcal{P}_T^\perp x^k)) \leq f(x^0); \alpha_k = 0 \text{ otherwise.} \quad (5)$$

The simple step-size rule avoids line search. Indeed, $\alpha_k = 2/(k+2)$ is the standard step-size in the Frank-Wolfe literature [18]; and here, for technical reasons, we also want to make sure that x^k stays in a level set of f , i.e., $f(x^k) \leq f(x^0)$ for $k \geq 1$. In practice, this condition is always satisfied after the first several iterations.

Algorithm 1 Unbounded Frank-Wolfe Method (uFW)

Initialize. Initialize with $x^0 \in T \oplus S$, gradient descent step-size η and Frank-Wolfe step-size sequence $\alpha_k \in [0, 1]$ for all k .

Perform the following steps for iterations $k = 0, 1, 2, \dots$

- (1) **Gradient descent step** (along subspace T): $y^k = x^k - \eta \mathcal{P}_T \nabla f(x^k)$.
 - (2) **Solve linear oracle:** $s^k = \arg \min_{s \in S} \langle \nabla f(y^k), s \rangle$.
 - (3) **Frank-Wolfe step** (along subset S): $x^{k+1} = y^k + \alpha_k (s^k - \mathcal{P}_T^\perp x^k)$.
-

Unlike gradient descent, the traditional Frank-Wolfe method does not have linear convergence even if the objective function is strongly convex. An intuitive explanation is that the solutions to the linear subproblems in the Frank-Wolfe algorithm may alternate between two extreme points of the constraint set, the iterate solutions zigzag, slowing down the convergence of the algorithm (see [31] for a nice explanation). Away-step Frank-Wolfe method allows moving in a direction opposite to the maximal solution of the linear model

evaluated at the current solution. The iterate solution can land on a certain face of the constraint set, thereby avoiding the aforementioned zigzagging phenomenon. Moreover, with away-steps, the Frank-Wolfe method enjoys a linear convergence rate when the objective function is further strongly convex and the constraint set is a polytope [31, 43, 42]. In addition to faster convergence, the away-step Frank-Wolfe method can usually lead to a sparser solution [5, 6], which is helpful in many applications where sparsity and interpretability are desirable properties of a solution [19].

Algorithm 2 adapts the away-step Frank-Wolfe method to solve (3) with the unbounded feasible region $T \oplus S$. The major difference of Algorithm 2 with Algorithm 1 is that the former allows performing an away step for the update along S . Note that $\mathcal{P}_T^\perp x^k$ can be represented as a convex combination of at most k vertices $V(x^k) \subset \text{vertices}(S)$, which can be written as $\mathcal{P}_T^\perp x^k = \sum_{v \in V(x^k)} \lambda_v(x^k) v$, with $\lambda_v(x^k) > 0$ for $v \in V(x^k)$ and $\sum_{v \in V(x^k)} \lambda_v(x^k) = 1$. In Algorithm 2, we keep track of the set of vertices $V(x^k)$ and the weight vector $\lambda(x^k)$, and choose between the Frank-Wolfe step and away step accordingly. After a step is taken in the direction T^\perp , we update the set of vertices $V(x^{k+1})$ and the weights $\lambda(x^{k+1})$. Unlike Algorithm 1, we perform a line search to find the Frank-Wolfe step-size¹ (the simple step-size rule does not apply here). Note the cardinality of the vertex set $V(x^k)$ in the maximization problem in Step (2) of Algorithm 2 is usually small, so the linear oracle for computing v^k can be much simpler than the computation of s^k .

Formally speaking, if we decide to take a Frank-Wolfe step (See Algorithm 2), we update the vertex set by

$$V(x^{k+1}) = \{s^k\} \text{ if } \alpha_k = 1; \text{ and } V(x^{k+1}) = V(x^k) \cup \{s^k\} \text{ if } \alpha_k \neq 1. \quad (6)$$

We update the weights by

$$\begin{aligned} \lambda_{s^k}(x^{k+1}) &= (1 - \alpha_k)\lambda_{s^k}(x^k) + \alpha_k, \\ \text{and } \lambda_v(x^{k+1}) &= (1 - \alpha_k)\lambda_v(x^k) \quad \text{for } v \in V(x^k) \setminus \{s^k\}. \end{aligned} \quad (7)$$

$$V(x^{k+1}) = V(x^k) \setminus \{v^k\} \text{ if } \alpha_k = \alpha_{\max}; \text{ and } V(x^{k+1}) = V(x^k) \text{ if } \alpha_k \neq \alpha_{\max}. \quad (8)$$

We update the weights by

$$\begin{aligned} \lambda_{v^k}(x^{k+1}) &= (1 + \alpha_k)\lambda_{v^k}(x^k) - \alpha_k, \\ \text{and } \lambda_v(x^{k+1}) &= (1 + \alpha_k)\lambda_v^k \quad \text{for } v \in V(x^k) \setminus \{v^k\}. \end{aligned} \quad (9)$$

This update follows from recent work on away-step Frank-Wolfe methods [31, 42].

3 Computational Guarantees

In this section, we derive computational guarantees of uFW and uAFW. When the objective is smooth, we show that uFW converges to an optimal solution with sublinear rate $O(1/k)$. Additionally, if the objective is strongly convex and the constraint set S is a polyhedron, we show that uAFW converges linearly to an optimal solution.

3.1 Sublinear Rate for uFW

We start by introducing some quantities that appear in the computational guarantees of uFW.

Let X^0 denote the level set of $f(x)$ with maximal value $f(x^0)$, that is,

$$X^0 := \{x \in \mathbb{R}^n \mid f(x) \leq f(x^0)\}. \quad (10)$$

¹Note that other step size rules using the smoothness parameter or backtracking line search [40] also work for uAFW under which linear convergence can be proved (with slightly different arguments). We focus on the exact line-search steps for simplicity.

Algorithm 2 Unbounded Away-Step Frank-Wolfe Method (uAFW)

Initialize. Initialize with $x^0 \in T \oplus S$ such that $\mathcal{P}_T^\perp x^0$ is a vertex of S , vertex set $V(x^0) = \{\mathcal{P}_T^\perp x^0\}$, weight parameters $\lambda_v(x^0) = 1$ for $v \in V(x^0)$, and gradient descent step-size η .

Perform the following updates for iterations $k = 0, 1, 2, \dots$:

- (1) Gradient descent (along subspace T): $y^k = x^k - \eta \mathcal{P}_T \nabla f(x^k)$.
 - (2) Linear oracle: $s^k := \arg \min_{s \in S} \langle \nabla f(y^k), s \rangle$, $v^k := \arg \max_{v \in V(x^k)} \langle \nabla f(y^k), v \rangle$.
 - (3) Find the moving direction:
 - if** $\langle \nabla f(y^k), s^k - \mathcal{P}_T^\perp x^k \rangle < \langle \nabla f(y^k), \mathcal{P}_T^\perp x^k - v^k \rangle$: (Frank-Wolfe step)
 - $d^k := s^k - \mathcal{P}_T^\perp x^k$; $\alpha_{\max} = 1$.
 - else**: (away step)
 - $d^k := \mathcal{P}_T^\perp x^k - v^k$; $\alpha_{\max} = \lambda_{v^k}(x^k)/(1 - \lambda_{v^k}(x^k))$.
 - (4) FW/away step (along subset S): $x^{k+1} = y^k + \alpha_k d^k$, where $\alpha_k \in [0, \alpha_{\max}]$ is chosen by line search.
 - (5) Update vertices $V(x^{k+1})$ and weights $\lambda_v(x^{k+1})$ for $v \in V(x^{k+1})$ according to equations (6) - (9).
-

For the two step-size rules (5) and (4), the iterations x^k and y^k stay in the level set X^0 (see analysis below). Throughout the paper, we assume that f has a bounded level set along the subspace T , i.e.,

$$D_T := \sup_{x, y \in X^0} \|\mathcal{P}_T(x - y)\|_2 < \infty. \quad (11)$$

The traditional Frank-Wolfe method is invariant under affine transformations. [28] introduces a curvature constant of the objective with respect to the bounded constraint set, which is used to improve the convergence rate of the Frank-Wolfe method by using the invariance (of affine transformations). In order to measure the progress of the Frank-Wolfe step in Algorithm 1, we extend the curvature constant in [28] to our setting (3) with an unbounded constraint set:

Definition 3.1. *The curvature constant of f with respect to the constraint set S and level set X^0 is defined as*

$$C_{f, x^0}^S := \sup_{s \in S, x, y \in \mathbb{R}^n, \alpha \in (0, 1]} \left\{ (2/\alpha^2) [f(y) - f(x) - \langle \nabla f(x), y - x \rangle] \mid \mathcal{P}_T^\perp x \in S, x \in X^0, y = x + \alpha(s - \mathcal{P}_T^\perp x) \right\}. \quad (12)$$

The parameter C_{f, x^0}^S quantifies the curvature information of f when moving within the set S , and its value is dependent on the function f , the constraint set $S \oplus T$, and the initial point x^0 . In particular, it is easy to check that C_{f, x^0}^S is smaller than $L' \text{diam}(S)^2$, where L' is the smoothness parameter of f with respect to any norm, and $\text{diam}(S)$ is the diameter of set S with respect to the same chosen norm. A related quantity appears in the standard analysis of the Frank-Wolfe method [28].

To measure the progress of the gradient descent step in Algorithm 1, we introduce below the smoothness constant of f with respect to the subspace T :

Definition 3.2. *The smoothness parameter of f with respect to the subspace T and level set X^0 is defined as the Lipschitz constant of ∇f along subspace T :*

$$L_{f, x^0}^T := \sup_{\substack{x, y \in T \oplus S \\ x \neq y}} \left\{ \frac{\|\mathcal{P}_T \nabla f(x) - \mathcal{P}_T \nabla f(y)\|_2}{\|x - y\|_2} \mid x - y \in T, x, y \in X^0 \right\}. \quad (13)$$

Now we are ready to present the sublinear rate of uFW:

Theorem 3.1. *Consider Algorithm 1 with initial solution x^0 , gradient descent step-size $\eta \leq 1/L_{f, x^0}^T$, and Frank-Wolfe step-size α_k following either the line-search rule (4) or the simple step-size rule (5). Then it*

holds for all $k \geq 1$ that

$$f(x^k) - f^* \leq \frac{2C_{f,x^0}^S + 16D_T^2/\eta}{k+2}$$

where f^* is the optimal objective value for problem (3).

Remark 3.2. We consider two special cases of Theorem 3.1. When $T = \{0\}$, Algorithm 1 recovers the traditional Frank-Wolfe method, and Theorem 3.1 implies $f(x^k) - f^* \leq O(C_{f,x^0}^S/k)$, which recovers the standard convergence rate for the traditional Frank-Wolfe method. When $S = \{0\}$, Algorithm 1 recovers the traditional gradient descent method, and Theorem 3.1 implies $f(x^k) - f^* \leq O(L_{f,x^0}^T D_T^2/k)$ if we take $\eta = 1/L_{f,x^0}^T$. This recovers the standard sublinear rate for gradient descent.

3.1.1 Practical termination rules

It is well known that the FW method (for problem (3) with $T = \{0\}$) naturally produces a primal-dual gap that can be used as a termination criterion (see, e.g. [18]). On the other hand, the gradient descent method (for problem (3) with $S = \{0\}$) can be terminated if the norm of the gradient is below a specified tolerance. In light of these results, we propose to set a termination criteria for Algorithm 1 if both the following quantities

$$G_k := \langle \nabla f(y^k), \mathcal{P}_T^\perp y^k - s^k \rangle, \quad \text{and} \quad H_k := \|\mathcal{P}_T \nabla f(y^k)\|_2 \quad (14)$$

are small. Note that both G_k and H_k are always non-negative, and they can be computed as a by-product of Algorithm 1. The next proposition shows that both $\min_{1 \leq k \leq 2m-1} \{G_k\}$ and $\min_{1 \leq k \leq 2m-1} \{H_k\}$ converge to 0 with a rate $O(1/m)$.

Proposition 3.3. Under the same assumption of Theorem 3.1, it holds for all $m \geq 1$ that

$$\min_{k=1, \dots, 2m-1} G_k \leq \frac{1}{\log(4/3)} \frac{2C_{f,x^0}^S + 8D_T^2/\eta}{m+1}, \quad (15)$$

and

$$\min_{k=1, \dots, 2m-1} H_k \leq \sqrt{\frac{8C_{f,x^0}^S \eta + 32D_T^2}{\log(4/3)}} \left(\frac{\eta^{-1} + L_{f,x^0}^T}{m+1} \right). \quad (16)$$

Finally, note that if some additional problem-specific parameters are known, we can compute upper bounds for the primal optimality gap $f(y^k) - f^*$ based on G_k and H_k .

Proposition 3.4. Under the same assumption of Theorem 3.1, we have

$$f(y^k) - f^* \leq G_k + H_k D_T \quad \forall k \geq 1. \quad (17)$$

If in addition $f(\cdot)$ is μ -strongly convex with respect to $\|\cdot\|_2$ (Definition 3.4), it holds

$$f(y^k) - f^* \leq G_k + H_k^2/(2\mu) \quad \forall k \geq 1. \quad (18)$$

3.2 Proofs of Theorem 3.1 and Proposition 3.3

We first present Lemmas 3.5 and 3.6 before presenting the proofs of Theorem 3.1 and Proposition 3.3:

Lemma 3.5. For any $x, y \in T \oplus S$ with $y - x \in T$ and $x, y \in X^0$, the following holds

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + (L_{f,x^0}^T/2) \|y - x\|_2^2.$$

The proof of Lemma 3.5 is the same as that appearing in the proof of Lemma 1.2.3 in [36] except that we restrict the analysis to the case $y - x \in T$.

For notational convenience we let $z^{k+1} := y^k + (2/(k+2))(s^k - \mathcal{P}_T^\perp x^k)$ for $k \geq 0$.

Lemma 3.6. *Suppose we set $\eta \leq 1/L_{f,x^0}^T$. Then for both step-size rules (4) and (5), we have $x^k \in X^0$ for all $k \geq 0$ and $f(x^{k+1}) \leq f(z^{k+1})$ for all $k \geq 1$.*

Proof: We first note that the conclusions hold true trivially for the line-search step-size rule (4). For the simple step-size rule (5), we first show that $f(x^k) \leq f(x^0)$ for any $k \geq 0$. Suppose there exists a k such that $f(x^k) > f(x^0)$, then without loss of generality, we assume $k \geq 1$ is the smallest index that satisfies $f(x^k) > f(x^0)$, thus it holds that $f(x^{k-1}) \leq f(x^0)$. From the simple step-size rule in (5), if $f(z^k) \leq f(x^0)$, it holds $x^k = z^k$, whereby $f(x^k) = f(z^k) \leq f(x^0)$. Otherwise, $\alpha_{k-1} = 0$, hence it holds $x^k = y^{k-1}$, whereby $f(x^k) = f(y^{k-1}) \leq f(x^{k-1}) \leq f(x^0)$. In either case, we have $f(x^k) \leq f(x^0)$, which contradicts the assumption. Therefore, we have $f(x^k) \leq f(x^0)$ for any $k \geq 0$. Consequently, if $f(z^{k+1}) \leq f(x^0)$, we have $x^{k+1} = z^{k+1}$ whereby $f(x^{k+1}) \leq f(z^{k+1})$. Otherwise we have $f(z^{k+1}) > f(x^0)$, whereby $f(x^{k+1}) \leq f(x^0) < f(z^{k+1})$, therefore it always holds that $f(x^{k+1}) \leq f(z^{k+1})$. This completes the proof of Lemma 3.6. \square

Proof of Theorem 3.1. First, note that

$$\begin{aligned} f(y^k) &\stackrel{(i)}{\leq} f(x^k) + \langle \nabla f(x^k), y^k - x^k \rangle + \|y^k - x^k\|_2^2 / (2\eta) \\ &\stackrel{(ii)}{=} f(x^k) - (\eta/2) \|\mathcal{P}_T \nabla f(x^k)\|_2^2 \leq f(x^k) \leq f(x^0), \end{aligned} \quad (19)$$

where (i) uses Lemma 3.5 and the assumption $\eta \leq 1/L_{f,x^0}^T$, and (ii) is from the update rule $y^k = x^k - \eta \mathcal{P}_T \nabla f(x^k)$. Let x^* be an optimal solution of (3). Recall that we have defined $z^{k+1} := y^k + (2/(k+2))(s^k - \mathcal{P}_T^\perp x^k)$, and shown in Lemma 3.6 that $f(x^{k+1}) \leq f(z^{k+1})$. As a result, we have

$$\begin{aligned} f(x^{k+1}) &\leq f(z^{k+1}) \\ &\leq f(y^k) + \frac{2}{k+2} \langle \nabla f(y^k), s^k - \mathcal{P}_T^\perp x^k \rangle + \frac{2C_{f,x^0}^S}{(k+2)^2} \\ &\leq f(y^k) + \frac{2}{k+2} \langle \nabla f(y^k), \mathcal{P}_T^\perp x^* - \mathcal{P}_T^\perp y^k \rangle + \frac{2C_{f,x^0}^S}{(k+2)^2}, \end{aligned} \quad (20)$$

where the second inequality is from the definition of C_{f,x^0}^S with $s = s^k$, $x = y^k$, $y = z^{k+1}$ and $\alpha = 2/(k+2)$, and the third inequality is by the definition of s^k and $\mathcal{P}_T^\perp x^k = \mathcal{P}_T^\perp y^k$. Since

$$\langle \nabla f(y^k), \mathcal{P}_T^\perp x^* - \mathcal{P}_T^\perp y^k \rangle = \langle \nabla f(y^k), x^* - y^k \rangle + \langle \mathcal{P}_T \nabla f(y^k), y^k - x^* \rangle, \quad (21)$$

we have

$$\begin{aligned} f(x^{k+1}) &\stackrel{(i)}{\leq} \frac{2}{k+2} f(y^k) + \frac{k}{k+2} f(x^k) - \frac{\eta}{2} \frac{k}{k+2} \|\mathcal{P}_T \nabla f(x^k)\|_2^2 + \frac{2C_{f,x^0}^S}{(k+2)^2} \\ &\quad + \frac{2}{k+2} \langle \nabla f(y^k), x^* - y^k \rangle + \frac{2}{k+2} \langle \mathcal{P}_T \nabla f(y^k), y^k - x^* \rangle \\ &\stackrel{(ii)}{\leq} \frac{2}{k+2} f^* + \frac{k}{k+2} f(x^k) + \frac{2C_{f,x^0}^S}{(k+2)^2} \\ &\quad + \frac{2}{k+2} \langle \mathcal{P}_T \nabla f(y^k), y^k - x^* \rangle - \frac{\eta}{2} \frac{k}{k+2} \|\mathcal{P}_T \nabla f(x^k)\|_2^2, \end{aligned} \quad (22)$$

where (i) utilizes (19), (20) and (21), and (ii) is because $f(y^k) + \langle \nabla f(y^k), x^* - y^k \rangle \leq f^*$ due to the convexity of f . Note that

$$\begin{aligned} &\langle \mathcal{P}_T \nabla f(y^k), y^k - x^* \rangle \\ &= \langle \mathcal{P}_T \nabla f(y^k) - \mathcal{P}_T \nabla f(x^k), y^k - x^* \rangle + \langle \mathcal{P}_T \nabla f(x^k), y^k - x^* \rangle \\ &\stackrel{(i)}{\leq} L_{f,x^0}^T \|y^k - x^k\|_2 D_T + \|\mathcal{P}_T \nabla f(x^k)\|_2 D_T \\ &\stackrel{(ii)}{=} (L_{f,x^0}^T \eta + 1) \|\mathcal{P}_T \nabla f(x^k)\|_2 D_T \stackrel{(iii)}{\leq} 2 \|\mathcal{P}_T \nabla f(x^k)\|_2 D_T, \end{aligned} \quad (23)$$

where in the above display, (i) follows from the definitions of D_T and L_{f,x^0}^T in (11) and (13) and the fact $\{y^k, x^*\} \subset X^0$, and (ii) is due to the update rule in step (1) of Algorithm 1, and (iii) utilizes the assumption $\eta \leq 1/L_{f,x^0}^T$. By combining (22) and (23), we arrive at

$$\begin{aligned} f(x^{k+1}) &\leq \frac{2}{k+2}f^* + \frac{k}{k+2}f(x^k) + \frac{2C_{f,x^0}^S}{(k+2)^2} \\ &\quad + \frac{4D_T}{k+2}\|\mathcal{P}_T\nabla f(x^k)\|_2 - \frac{\eta}{2}\frac{k}{k+2}\|\mathcal{P}_T\nabla f(x^k)\|_2^2 \end{aligned} \quad (24)$$

for all $k \geq 0$. For $k \geq 2$, using (24) and the following inequality

$$\frac{8D_T^2}{k(k+2)\eta} + \frac{\eta}{2}\frac{k}{k+2}\|\mathcal{P}_T\nabla f(x^k)\|_2^2 \geq \frac{4D_T}{k+2}\|\mathcal{P}_T\nabla f(x^k)\|_2,$$

(which follows from $a^2 + b^2 \geq 2ab$ for scalars a, b), we have

$$\begin{aligned} f(x^{k+1}) &\leq \frac{2}{k+2}f^* + \frac{k}{k+2}f(x^k) + \frac{2C_{f,x^0}^S}{(k+2)^2} + \frac{8D_T^2}{k(k+2)\eta} \\ &\stackrel{(i)}{\leq} \frac{2}{k+2}f^* + \frac{k}{k+2}f(x^k) + \frac{2C_{f,x^0}^S + 16D_T^2/\eta}{(k+2)^2}. \end{aligned}$$

Note that inequality (i) above makes use of $k \geq 2$.

Let $h_k := f(x^k) - f^*$ and $C := 2C_{f,x^0}^S + 16D_T^2/\eta$, then we have

$$h_{k+1} \leq kh_k/(k+2) + C/(k+2)^2, \quad \forall k \geq 2. \quad (25)$$

To complete the proof, we will show that $h_k \leq C/(k+2)$ by induction. First, it follows from (24) by choosing $k = 0$ that

$$h_1 \leq (1/2)C_{f,x^0}^S + 2D_T\|\mathcal{P}_T\nabla f(x^0)\|_2.$$

Let w^0 be a minimizer of $f(x)$ on the affine subspace $x^0 + T$ (the existence of a minimizer is guaranteed by the assumption that $D_T < \infty$), then it holds that $w^0 \in X^0$, $w^0 - x^0 \in T$ and $\mathcal{P}_T\nabla f(w^0) = 0$. Thus we have

$$\begin{aligned} h_1 &\leq (1/2)C_{f,x^0}^S + 2D_T\|\mathcal{P}_T\nabla f(x^0) - \mathcal{P}_T\nabla f(w^0)\|_2 \\ &\leq (1/2)C_{f,x^0}^S + 2L_{f,x^0}^T D_T^2 \leq C/3. \end{aligned} \quad (26)$$

Again, using (24) with $k = 1$, we have

$$\begin{aligned} h_2 &\leq (1/3)h_1 + (2/9)C_{f,x^0}^S + (4/3)D_T\|\mathcal{P}_T\nabla f(x^1)\|_2 \\ &\leq (1/3)\left((1/2)C_{f,x^0}^S + 2L_{f,x^0}^T D_T^2\right) + (2/9)C_{f,x^0}^S + (4/3)L_{f,x^0}^T D_T^2 \\ &\leq C/4. \end{aligned} \quad (27)$$

Inequalities (26) and (27) imply that $h_k \leq C/(k+2)$ holds for $k = 1, 2$. Now suppose $h_k \leq C/(k+2)$ holds for some $k \geq 2$, then it follows from (25) that $h_{k+1} \leq C/(k+3)$. Therefore, $h_k \leq C/(k+2) \forall k \geq 1$ by induction—this finishes the proof by substituting the value of C . \square

Proof of Proposition 3.3. Let $\tau_k := 2/(k+2)$ and $\tilde{H}_k := \|\mathcal{P}_T\nabla f(x^k)\|_2$ (note that \tilde{H}_k is different from H_k). From (19), (20) and noting that $\mathcal{P}_T x^k = \mathcal{P}_T y^k$, we have

$$f(x^{k+1}) \leq f(x^k) - \tau_k G_k + (\tau_k^2/2)C_{f,x^0}^S - (\eta/2)\tilde{H}_k^2. \quad (28)$$

Summing inequality (28) from $k = m$ to $k = 2m + 1$ and rearranging terms we have

$$\sum_{k=m}^{2m-1} (\tau_k G_k + (\eta/2) \tilde{H}_k^2) \leq f(x^m) - f(x^{2m-1}) + \sum_{k=m}^{2m-1} \frac{\tau_k^2}{2} C_{f,x^0}^S. \quad (29)$$

Note that

$$\sum_{k=m}^{2m-1} \frac{\tau_k^2}{2} = \sum_{k=m}^{2m-1} \frac{2}{(k+2)^2} \leq 2 \sum_{k=m}^{2m-1} \left(\frac{1}{k+1} - \frac{1}{k+2} \right) \leq \frac{2}{m+1}. \quad (30)$$

Combining (29), (30), and noting that $f(x^m) - f(x^{2m-1}) \leq f(x^m) - f^*$, we have

$$\sum_{k=m}^{2m-1} (\tau_k G_k + (\eta/2) \tilde{H}_k^2) \leq f(x^m) - f^* + \frac{2C_{f,x^0}^S}{m+1}. \quad (31)$$

By (31) and the conclusion of Theorem 3.1 we have

$$\sum_{k=m}^{2m-1} (\tau_k G_k + (\eta/2) \tilde{H}_k^2) \leq \frac{2C_{f,x^0}^S + 16D_T^2/\eta}{m+2} + \frac{2C_{f,x^0}^S}{m+1} \leq \frac{4C_{f,x^0}^S + 16D_T^2/\eta}{m+1}. \quad (32)$$

Note that

$$\sum_{k=m}^{2m-1} \tau_k = \sum_{k=m}^{2m-1} \frac{2}{k+2} \geq \int_m^{2m} \frac{2}{t+2} dt = 2 \log \frac{2m+2}{m+2} \geq 2 \log(4/3). \quad (33)$$

So by (32) and (33), there exists some $m \leq j \leq 2m - 1$ satisfying

$$G_j + \eta \tilde{H}_j^2 / (2\tau_j) \leq \left(\sum_{k=m}^{2m-1} \tau_k \right)^{-1} \sum_{k=m}^{2m-1} (\tau_k G_k + (\eta/2) \tilde{H}_k^2) \leq \frac{2C_{f,x^0}^S + 8D_T^2/\eta}{(m+1) \log(4/3)}. \quad (34)$$

The first inequality (15) follows immediately from (34). To prove the second inequality of this proposition, note that

$$\begin{aligned} H_j - \tilde{H}_j &\leq \|\mathcal{P}_T \nabla f(x^j) - \mathcal{P}_T \nabla f(y^j)\|_2 \leq L_{f,x^0}^T \|x^j - y^j\|_2 \\ &\leq L_{f,x^0}^T \eta \|\mathcal{P}_T \nabla f(x^j)\|_2 = L_{f,x^0}^T \eta \tilde{H}_j. \end{aligned} \quad (35)$$

Combining (34) and (35) we have

$$H_j \leq (1 + \eta L_{f,x^0}^T) \tilde{H}_j \leq (1 + \eta L_{f,x^0}^T) \left(\frac{2\tau_j}{\eta} \frac{2C_{f,x^0}^S + 8D_T^2/\eta}{(m+1) \log(4/3)} \right)^{1/2}. \quad (36)$$

Recall that $\tau_j = 2/(j+2) \leq 2/(m+1)$, so by (36) we arrive at (16), which completes the proof. \square

Proof of Proposition 3.4. Let x^* be an optimal solution of (3). By the convexity of $f(\cdot)$, we have

$$\begin{aligned} f(y^k) - f^* &\leq \langle \nabla f(y^k), y^k - x^* \rangle \\ &= \langle \nabla f(y^k), \mathcal{P}_T^\perp y^k - \mathcal{P}_T^\perp x^* \rangle + \langle \mathcal{P}_T \nabla f(y^k), \mathcal{P}_T(y^k - x^*) \rangle \\ &\leq \langle \nabla f(y^k), \mathcal{P}_T^\perp y^k - s^k \rangle + \|\mathcal{P}_T \nabla f(y^k)\|_2 \|\mathcal{P}_T(y^k - x^*)\|_2 \\ &\leq G_k + H_k D_T, \end{aligned} \quad (37)$$

where the second inequality is by the definition of s^k and the third inequality is by the definition of D_T .

If in addition $f(\cdot)$ is μ -strongly convex with respect to $\|\cdot\|_2$, then

$$\begin{aligned}
f(y^k) - f^* &\leq \langle \nabla f(y^k), y^k - x^* \rangle - \frac{\mu}{2} \|y^k - x^*\|_2^2 \\
&= \langle \nabla f(y^k), \mathcal{P}_T^\perp y^k - \mathcal{P}_T^\perp x^* \rangle + \langle \mathcal{P}_T \nabla f(y^k), y^k - x^* \rangle - \frac{\mu}{2} \|y^k - x^*\|_2^2 \\
&\leq \langle \nabla f(y^k), \mathcal{P}_T^\perp y^k - s^k \rangle + \|\mathcal{P}_T \nabla f(y^k)\|_2^2 / (2\mu) \\
&\leq G_k + H_k^2 / (2\mu)
\end{aligned}$$

where the second inequality uses Cauchy-Schwarz inequality and the definition of s^k .

3.3 Linear Rate for uAFW

In this section, we establish the global linear convergence of uAFW (Algorithm 2).

The standard linear convergence result for away-step Frank-Wolfe method relies on a suitable curvature constant of f with respect to an extended constraint set [31] (equivalently, it is the relative smoothness constant to the constraint set presented in [24]). Here we generalize this notion to the setting of (3).

Definition 3.3. We define \bar{C}_{f,x^0}^S , the curvature constant of f with respect to an extended constraint set along the set S as follows:

$$\begin{aligned}
\bar{C}_{f,x^0}^S := & \sup_{\substack{s \in S, x, y \in \mathbb{R}^n \\ \alpha \in [-1, 1] \setminus \{0\}}} \left\{ (2/\alpha^2) (f(y) - f(x) - \langle \nabla f(x), y - x \rangle) \mid \right. \\
& \left. \mathcal{P}_T^\perp x \in S, f(x) \leq f(x^0), y = x + \alpha(s - \mathcal{P}_T^\perp x) \right\}.
\end{aligned}$$

Note that \bar{C}_{f,x^0}^S is different from C_{f,x^0}^S defined in (12). In the definition of \bar{C}_{f,x^0}^S , α is allowed to take values in $[-1, 1] \setminus \{0\}$, while in the definition of C_{f,x^0}^S it takes values in $(0, 1]$. In other words, \bar{C}_{f,x^0}^S quantifies the curvature of f with respect to an extended set $\bar{S} := S + (S - S)$ (i.e., the Minkowski sum of S and $S - S$). In the special case when f is L -smooth over $\bar{S} \oplus T$, if we define $\bar{D} := \text{diam}(\bar{S})$, then we have $\bar{C}_{f,x^0}^S \leq L\bar{D}^2$. Intuitively, \bar{C}_{f,x^0}^S should appear in the computational guarantees for uAFW because the algorithm not only moves towards the vertices but also away from vertices, and \bar{C}_{f,x^0}^S provides an upper bound on the curvature information of f when moving in these directions.

To obtain a linear rate for uAFW we also need to assume strong convexity of the objective function.

Definition 3.4. Let X be a convex set in \mathbb{R}^n . For $\mu > 0$, we say a function f is μ -strongly convex on X with respect to a norm $\|\cdot\|$ if for all $x, y \in X$, we have

$$f(y) - f(x) - \langle \nabla f(x), y - x \rangle \geq (\mu/2) \|y - x\|^2.$$

Finally, following [42], we define the facial distance $\psi(S)$ of the polyhedral set S as follows:

$$\psi(S) := \min_{\substack{F \in \text{faces}(S), \\ \emptyset \neq F \neq S}} \text{dist}(F, \text{conv}(V(S) \setminus F)), \quad (38)$$

where $V(S)$ denotes the set of all vertices of S , $\text{faces}(S)$ the set of all faces of S ; and $V(S) \setminus F$ means $V(S) \cap F^c$. Note $\psi(S)$ appears in the computational guarantees of the classical away-step Frank-Wolfe algorithm in the analysis of [42]. It is strictly positive for a polyhedral S , but may be zero for a general convex set S .

We make the following assumptions to establish the linear convergence rate.

Assumption 3.7. (1) S is a bounded polyhedron.

(2) X^0 is bounded.

(3) f is L -smooth over the extended constraint $\bar{S} \oplus T$ with $\bar{S} := S + (S - S) = \{x + (y - z) \mid x, y, z \in S\}$.

(4) f is μ -strongly convex on $S \oplus T$ with respect to $\|\cdot\|_2$.

As discussed before, condition (3) of Assumption 3.7 guarantees that the parameter \bar{C}_{f,x^0}^S is bounded above. The next theorem presents the linear convergence of Algorithm 2.

Theorem 3.8. Suppose that f is convex and continuously differentiable, and satisfies Assumption 3.7. Let $\{(x^k, y^k)\}_{k \geq 0}$ be the sequence generated by Algorithm 2 by setting $\eta \leq 1/L$; and recall that f^* is the optimal objective value. Let $\sigma = \mu\psi(S)^2/4$, then we have

$$f(y^k) - f^* \leq \left(1 - \min\{1/2, \sigma/\bar{C}_{f,x^0}^S\} \frac{\mu\eta}{4}\right)^{k/2} (f(y^0) - f^*). \quad (39)$$

Remark 3.9. If we set the stepsize $\eta = 1/L$ (recall, f is L -smooth over $\bar{S} \oplus T$), the linear convergence rate in (39) depends on two ratios: $\sigma/\bar{C}_{f,x^0}^S$ and μ/L . The first ratio is in accordance with the linear rate parameter for the classical away-step Frank-Wolfe algorithm, while the second one is the rate for classical gradient descent.

3.3.1 Practical termination rules

Similar to Section 3.1.1, we can use quantities G_k and H_k defined in (14) (with y^k being the iterations produced by uAFW) to derive a practical termination criteria for uAFW. The following proposition which is a variant of Theorem 2 in [31] presents the rate at which G_k and H_k converge to zero.

Proposition 3.10. Under the same assumption of Theorem 3.8, it holds:

(1) $H_k \leq (1 + \eta L_{f,x^0}^T) \sqrt{(2/\eta)(f(x^k) - f^*)}$.

(2) Among the first k iterations of Algorithm 2, there are at least $\lfloor k/2 \rfloor$ iterations satisfying:

$$G_k \leq \min\left\{2(f(y^k) - f^*), \sqrt{2\bar{C}_{f,x^0}^S(f(y^k) - f^*)}\right\}. \quad (40)$$

The proof of Proposition 3.10 is relegated to Appendix B. In light of Theorem 3.8, Proposition 3.10 and noting that $f(x^{k+1}) \leq f(y^k)$, it can be seen that G_k and H_k converge to 0 with a linear rate. Hence we can terminate Algorithm 2 when G_k and H_k are smaller than a tolerance value.

The conclusions of Proposition 3.4 also hold true for Algorithm 2. These results can be used to compute upper bounds on $f(y^k) - f^*$ if D_T or μ is known.

3.4 Proof of Theorem 3.8

We divide the proof of Theorem 3.8 into the following three steps.

Step 1. A reduction to the axis-aligned case: Notice that both Frank-Wolfe and gradient descent are invariant under orthogonal transformations; and so is uAFW. That is, given an orthogonal transformation Q , suppose $\{x^k, y^k\}_{k=0}^\infty$ are the iterations generated by uAFW for the objective function $f(\cdot)$, constraint set $T \oplus S$ and initialization x^0 . Suppose $\{\tilde{x}^k, \tilde{y}^k\}_{k=0}^\infty$ are the iterations generated by uAFW for the objective function $g(\cdot) := f(Q^{-1}\cdot)$, constraint set $Q(T) \oplus Q(S)$ (where $Q(T) := \{Qx \mid x \in T\}$ and $Q(S)$ is defined similarly) and initialization $\tilde{x}^0 = Qx^0$. Then it holds $\tilde{x}^k = Qx^k$ and $\tilde{y}^k = Qy^k$ for all $k \geq 0$. Assumption 3.7 also holds for the function and constraints obtained after the orthogonal transformation. Thus without loss

of generality, we can assume that $\{e_{n-r+1}, e_{n-r+2}, \dots, e_n\}$ is a basis of the r -dimensional subspace T . Then we can convert the original problem (3) to the following form:

$$\min_{u,w} \bar{f}(u, w) \quad \text{s.t.} \quad u \in S \subset \mathbb{R}^{n-r}, w \in \mathbb{R}^r, \quad (41)$$

where $S \subset \mathbb{R}^{n-r}$ is a bounded polyhedral set (with a slight abuse, we are reusing the notation S). We refer to (41) as the *axis-aligned form*. Although the transformed problem (41) is equivalent to the original problem (3), in practice it can be computationally expensive to perform such transformation. We use formulation (41) to help facilitate our proof.

Let $\nabla_u \bar{f}(u, w)$ and $\nabla_w \bar{f}(u, w)$ denote the sub-vectors of $\nabla \bar{f}(u, w)$ corresponding to the derivatives of u and w respectively. Then we can rewrite uAFW (Algorithm 2) as Algorithm 3.

Algorithm 3 uAFW for axis-aligned form problem (41)

Starting from (u^0, w^0) such that u^0 is a vertex of S .

For $k = 0, 1, 2, \dots$:

- (1) $w^{k+1} = w^k - \eta \nabla_w \bar{f}(u^k, w^k)$.
 - (2) $s^k := \arg \min_{s \in S} \langle \nabla_u \bar{f}(u^k, w^{k+1}), s \rangle$, $v^k := \operatorname{argmax}_{v \in V(u^k)} \langle \nabla_u \bar{f}(u^k, w^{k+1}), v \rangle$.
if $\langle \nabla_u \bar{f}(u^k, w^{k+1}), s^k - u^k \rangle < \langle \nabla_u \bar{f}(u^k, w^{k+1}), u^k - v^k \rangle$: (Frank-Wolfe step)
 $d^k := s^k - u^k$; $\alpha_{\max} = 1$.
else: (away step)
 $d^k := u^k - v^k$; $\alpha_{\max} = \lambda_{v^k}(u^k)/(1 - \lambda_{v^k}(u^k))$.
 - (3) $u^{k+1} = u^k + \alpha_k d^k$; where $\alpha_k \in \operatorname{argmin}_{\alpha \in [0, \alpha_{\max}]} f(u^k + \alpha d^k, w^{k+1})$.
 - (4) Update $V(u^{k+1})$ and $\lambda(u^{k+1})$.
-

Note that in Algorithm 3, we use $V(u^{k+1})$ and $\lambda(u^{k+1})$ in place of the corresponding $V(x^{k+1})$ and $\lambda(x^{k+1})$ in Algorithm 2 because their values only depend on u^{k+1} . More specifically, we update $V(u^{k+1})$ and $\lambda(u^{k+1})$ in the following way: For a Frank-Wolfe step, let $V(u^{k+1}) = \{s^k\}$ if $\alpha_k = 1$; otherwise $V(u^{k+1}) = V(u^k) \cup \{s^k\}$. Meanwhile, we set $\lambda_{s^k}(u^{k+1}) := (1 - \alpha_k)\lambda_{s^k}(u^k) + \alpha_k$ and $\lambda_v(u^{k+1}) = (1 - \alpha_k)\lambda_v(u^k)$ for $v \in V(u^k) \setminus \{s^k\}$. For an away step, we update $V(u^{k+1}) = V(u^k) \setminus \{v^k\}$ if $\alpha_k = \alpha_{\max}$; otherwise we set $V(u^{k+1}) = V(u^k)$. We set $\lambda_{v^k}(u^{k+1}) := (1 + \alpha_k)\lambda_{v^k}(u^k) - \alpha_k$ and $\lambda_v(u^{k+1}) := (1 + \alpha_k)\lambda_v^k$ for $v \in V(u^k) \setminus \{v^k\}$.

Below we present the analysis of Algorithm 3 for Problem (41).

Step 2. Reduction in optimality gap via a Frank-Wolfe step: We establish a reduction in optimality gap via the Frank-Wolfe step. This is an adaption of Proposition 9 in [24] to the case of a function $\bar{f}(\cdot, w)$ with a dynamic w . To this end, we define $F : \mathbb{R}^r \rightarrow \mathbb{R}$ as $F(w) = \inf_{u \in S} \bar{f}(u, w)$.

Proposition 3.11. *Consider Algorithm 3 with initial solution (u^0, w^0) . Define $\gamma := \mu\psi(S)^2/(4\bar{C}_{\bar{f}, x^0}^S)$ and $\rho := \min\{1/2, \gamma\}$. Then among iterations $1, \dots, k$, there are at least $k/2$ of them satisfying*

$$\bar{f}(u^{k+1}, w^{k+1}) - f^* \leq (1 - \rho) (\bar{f}(u^k, w^{k+1}) - f^*) + \rho(F(w^{k+1}) - f^*).$$

The proof of Proposition 3.11 is inspired by [24]—it turns out to be a bit technical and we relegate the proof to Appendix A.

Step 3. Piecing together the Frank-Wolfe and gradient descent steps: Here we establish the linear convergence rate of uAFW by combining the contraction of the gradient descent step along the unbounded space T , and the contraction of the Frank-Wolfe step along the bounded set S (presented by Proposition 3.11).

We first present a few technical lemmas that will be used in the proof. Recall that we defined $F : \mathbb{R}^r \rightarrow \mathbb{R}$ as $F(w) = \inf_{u \in S} \bar{f}(u, w)$. Given $w \in \mathbb{R}^r$, we denote $u_w^* \in \operatorname{argmin}_{u \in S} \bar{f}(u, w)$. In addition, let $w^* \in \operatorname{argmin}_{w \in \mathbb{R}^r} F(w)$ and $F^* = f^* = \min_{w \in \mathbb{R}^r} F(w)$.

Lemma 3.12. *Under Assumption 3.7, we have*

(1) F is convex and differentiable over \mathbb{R}^r , with gradient $\nabla F(w) = \nabla_w \bar{f}(u_w^*, w)$.

(2) F is μ -strongly convex on \mathbb{R}^r with respect to the ℓ_2 -norm $\|\cdot\|_2$.

(3) For any $w \in \mathbb{R}^r$, it holds $F(w) - F^* \leq \|\nabla F(w)\|_2^2 / (2\mu)$.

Proof. Part (1) The convexity of F can be found in a standard reference on convex analysis, see e.g. [46]. The differentiability property of F is known as Danskin's theorem (for example, see [7] for a proof).

Part (2) Note that a function $g : \mathbb{R}^m \rightarrow \mathbb{R}$ is μ -strongly convex with respect to $\|\cdot\|_2$ if and only if

$$\langle \nabla g(x) - \nabla g(y), x - y \rangle \geq \mu \|x - y\|_2^2 \quad (42)$$

for all $x, y \in \mathbb{R}^m$ (See e.g. Theorem 2.1.9 of [36]). For any $w^1, w^2 \in \mathbb{R}^r$ we have

$$\begin{aligned} & \langle \nabla F(w^2) - \nabla F(w^1), w^2 - w^1 \rangle \\ \stackrel{(i)}{=} & \langle \nabla_w \bar{f}(u_{w^2}^*, w^2) - \nabla_w \bar{f}(u_{w^1}^*, w^1), w^2 - w^1 \rangle \\ = & \langle \nabla \bar{f}(u_{w^2}^*, w^2) - \nabla \bar{f}(u_{w^1}^*, w^1), (u_{w^2}^*, w^2) - (u_{w^1}^*, w^1) \rangle \\ & - \langle \nabla_u \bar{f}(u_{w^2}^*, w^2) - \nabla_u \bar{f}(u_{w^1}^*, w^1), u_{w^2}^* - u_{w^1}^* \rangle \\ \stackrel{(ii)}{\geq} & \mu \|(u_{w^2}^*, w^2) - (u_{w^1}^*, w^1)\|_2^2 - \langle \nabla_u \bar{f}(u_{w^2}^*, w^2) - \nabla_u \bar{f}(u_{w^1}^*, w^1), u_{w^2}^* - u_{w^1}^* \rangle \\ \geq & \mu \|w^2 - w^1\|_2^2 - \langle \nabla_u \bar{f}(u_{w^2}^*, w^2) - \nabla_u \bar{f}(u_{w^1}^*, w^1), u_{w^2}^* - u_{w^1}^* \rangle . \end{aligned}$$

where inequality (i) follows from the conclusion in part (1), and inequality (ii) follows from the strong-convexity of \bar{f} and the equivalent condition (42). Meanwhile, it follows from $u_{w^1}^* \in \operatorname{argmin}_{u \in S} \{\bar{f}(u, w^1)\}$ and $u_{w^2}^* \in \operatorname{argmin}_{u \in S} \{\bar{f}(u, w^2)\}$ that

$$\langle \nabla_u \bar{f}(u_{w^1}^*, w^1), u_{w^2}^* - u_{w^1}^* \rangle \geq 0, \quad \langle \nabla_u \bar{f}(u_{w^2}^*, w^2), u_{w^1}^* - u_{w^2}^* \rangle \geq 0 .$$

Therefore, we have $\langle \nabla F(w^2) - \nabla F(w^1), w^2 - w^1 \rangle \geq \mu \|w^2 - w^1\|_2^2$, which shows F is μ -strongly convex with respect to $\|\cdot\|_2$.

Part (3) Since F is μ -strongly convex over \mathbb{R}^r , this final claim follows from Theorem 2.1.10 of [36]. \square

Lemma 3.13. [Theorem 2.1.5 of [36]] For any $u^1, u^2 \in S$ and $w^1, w^2 \in \mathbb{R}^r$, we have

$$\begin{aligned} \bar{f}(u^2, w^2) & \geq \bar{f}(u^1, w^1) + \langle \nabla \bar{f}(u^1, w^1), (u^2 - u^1, w^2 - w^1) \rangle \\ & \quad + \frac{1}{2L} \|\nabla \bar{f}(u^2, w^2) - \nabla \bar{f}(u^1, w^1)\|_2^2 . \end{aligned}$$

The next proposition combines the contraction of objective values in the gradient descent and Frank-Wolfe steps:

Proposition 3.14. Suppose in iteration k ($k \geq 1$) the following condition holds:

$$\bar{f}(u^{k+1}, w^{k+1}) - f^* \leq (1 - \rho) (\bar{f}(u^k, w^{k+1}) - f^*) + \rho (F(w^{k+1}) - f^*) . \quad (43)$$

for some $\rho \in (0, 1)$. Then

$$\bar{f}(u^{k+1}, w^{k+2}) - f^* \leq (1 - \frac{\mu\eta}{4}\rho) (\bar{f}(u^k, w^{k+1}) - f^*) . \quad (44)$$

Proof. Let $\tau := \mu\eta/4$. Then $0 < \tau \leq 1/4$. We split the analysis into two cases.

(Case *i*) Here we consider the case when

$$F(w^{k+1}) - f^* \leq (1 - \tau) [\bar{f}(u^k, w^{k+1}) - f^*] . \quad (45)$$

Then we have

$$\begin{aligned} \bar{f}(u^{k+1}, w^{k+2}) - f^* &\leq \bar{f}(u^{k+1}, w^{k+1}) - f^* \\ &\leq (1 - \rho) (\bar{f}(u^k, w^{k+1}) - f^*) + \rho(F(w^{k+1}) - f^*) \\ &\leq (1 - \rho) [\bar{f}(u^k, w^{k+1}) - f^*] + \rho(1 - \tau) [\bar{f}(u^k, w^{k+1}) - f^*] \\ &= (1 - \rho\tau) [\bar{f}(u^k, w^{k+1}) - f^*] , \end{aligned}$$

where the first inequality is from the monotonicity of the iterations, the second inequality is from the contraction condition (43), and the third inequality is from assumption (45).

(Case *ii*) Here we consider the case when

$$F(w^{k+1}) - f^* > (1 - \tau) [\bar{f}(u^k, w^{k+1}) - f^*] \quad (46)$$

which is equivalent to

$$\bar{f}(u^k, w^{k+1}) - F(w^{k+1}) < \tau [\bar{f}(u^k, w^{k+1}) - f^*] . \quad (47)$$

Define $M_k := \bar{f}(u^k, w^{k+1}) - f^*$. Then, it holds that

$$\begin{aligned} \tau M_k &\stackrel{(i)}{>} \bar{f}(u^k, w^{k+1}) - F(w^{k+1}) \\ &\stackrel{(ii)}{\geq} \bar{f}(u^{k+1}, w^{k+1}) - F(w^{k+1}) \\ &\stackrel{(iii)}{\geq} \langle \nabla_u \bar{f}(u_{w^{k+1}}^*, w^{k+1}), u^{k+1} - u_{w^{k+1}}^* \rangle \\ &\quad + (1/(2L)) \|\nabla_w \bar{f}(u^{k+1}, w^{k+1}) - \nabla_w \bar{f}(u_{w^{k+1}}^*, w^{k+1})\|_2^2 \\ &\stackrel{(iv)}{\geq} (1/(2L)) \|\nabla_w \bar{f}(u^{k+1}, w^{k+1}) - \nabla_w \bar{f}(u_{w^{k+1}}^*, w^{k+1})\|_2^2 \\ &= (1/(2L)) \|\nabla_w \bar{f}(u^{k+1}, w^{k+1}) - \nabla F(w^{k+1})\|_2^2 , \end{aligned}$$

where inequality (i) is from (47), (ii) is from monotonicity of the iterations, (iii) is from Lemma 3.13, and (iv) is because of the optimality of $u_{w^{k+1}}^*$. The final equality is from Lemma 3.12 (1). As a result, we have

$$\|\nabla_w \bar{f}(u^{k+1}, w^{k+1}) - \nabla F(w^{k+1})\|_2 \leq \sqrt{2L\tau M_k} ,$$

and hence

$$\begin{aligned} \|\nabla_w \bar{f}(u^{k+1}, w^{k+1})\|_2 &\geq \|\nabla F(w^{k+1})\|_2 - \|\nabla F(w^{k+1}) - \nabla_w \bar{f}(u^{k+1}, w^{k+1})\|_2 \\ &\geq \|\nabla F(w^{k+1})\|_2 - \sqrt{2L\tau M_k} \\ &\stackrel{(i)}{\geq} \sqrt{2\mu(F(w^{k+1}) - f^*)} - \sqrt{2L\tau M_k} \\ &\stackrel{(ii)}{\geq} \sqrt{2\mu(1 - \tau)M_k} - \sqrt{2L\tau M_k} \\ &= \left(\sqrt{2\mu(1 - \tau)} - \sqrt{2L\tau} \right) \sqrt{M_k} , \end{aligned}$$

where (i) uses Lemma 3.12 (3), and (ii) is from inequality (46). As $\tau = \mu\eta/4$, we have

$$\sqrt{2\mu(1 - \tau)} - \sqrt{2L\tau} \geq \sqrt{2\mu(3/4)} - \sqrt{2L\eta(\mu/4)} \geq \left(\sqrt{3/2} - \sqrt{1/2} \right) \sqrt{\mu} > \sqrt{\mu}/2 ,$$

and thus

$$\|\nabla_w \bar{f}(u^{k+1}, w^{k+1})\|_2 \geq (1/2) \sqrt{\mu(\bar{f}(u^k, w^{k+1}) - f^*)}. \quad (48)$$

On the other hand, by the gradient step in the unbounded subspace, we have

$$\bar{f}(u^{k+1}, w^{k+2}) \leq \bar{f}(u^{k+1}, w^{k+1}) - (\eta/2) \|\nabla_w \bar{f}(u^{k+1}, w^{k+1})\|_2^2.$$

Therefore,

$$\bar{f}(u^{k+1}, w^{k+1}) - \bar{f}(u^{k+1}, w^{k+2}) \geq \frac{\eta}{2} \|\nabla_w \bar{f}(u^{k+1}, w^{k+1})\|_2^2 \geq \frac{\eta\mu}{8} (\bar{f}(u^k, w^{k+1}) - f^*),$$

where the last inequality is from (48). Noticing the algorithm is descent, it holds

$$\begin{aligned} \bar{f}(u^k, w^{k+1}) - \bar{f}(u^{k+1}, w^{k+2}) &\geq \bar{f}(u^{k+1}, w^{k+1}) - \bar{f}(u^{k+1}, w^{k+2}) \\ &\geq (\mu\eta/8) (\bar{f}(u^k, w^{k+1}) - f^*). \end{aligned}$$

The above leads to

$$\bar{f}(u^{k+1}, w^{k+2}) - f^* \leq (1 - \frac{\mu\eta}{8}) (\bar{f}(u^k, w^{k+1}) - f^*) \leq \left(1 - \rho \frac{\mu\eta}{4}\right) (\bar{f}(u^k, w^{k+1}) - f^*),$$

where the last inequality is because $\rho = \min\{1/2, \gamma\} \leq 1/2$. \square

Now we are ready to prove Theorem 3.8 by combining Propositions 3.11 and 3.14:

Proof of Theorem 3.8: In view of Proposition 3.11, among the first k iterations, at least $k/2$ of them satisfy condition (43). Proposition 3.14 states that when (43) holds we also have (44), hence (44) holds for at least $k/2$ of the first k iterations. Meanwhile, notice that uAFW is a descent algorithm, therefore we have

$$\bar{f}(u^k, w^{k+1}) - f^* \leq (1 - (\rho\mu\eta/4))^{k/2} (\bar{f}(u^0, w^1) - f^*).$$

We complete the proof of Theorem 3.8 by substituting $\rho = \min\{1/2, \sigma/\bar{C}_{f,x^0}^S\}$ in Proposition 3.11. \square

4 Applications

We discuss how to efficiently apply uFW (Algorithm 1) and uAFW (Algorithm 2) to solve the two motivating applications mentioned in Section 1. Section 4.3 discusses other applications of our framework.

4.1 The ℓ_1 trend filtering problem (of order r)

The ℓ_1 trend filtering problem (1) of order r is an instance of the following optimization problem:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad \|D_n^{(r)} x\|_1 \leq \delta, \quad (49)$$

where f is a smooth convex function; $D_n^{(r)}$ denotes the r -order discrete derivative operator, that is, $D_n^{(1)}$ is a matrix in $\mathbb{R}^{(n-1) \times n}$ with (i, i) -entry being 1 and $(i, i+1)$ -entry being -1 (for $i \in [n-1]$) and all other entries being 0. For $r \geq 2$, it is defined recursively by $D_n^{(r+1)} = D_{n-r}^{(1)} \cdot D_n^{(r)}$.

Popular instances of ℓ_1 trend filtering consider small values of r . ℓ_1 trend filtering with $r = 1$ is also called the fused lasso [49], where the constraint $\|D_n^{(1)} x\|_1 \leq \delta$ tends to induce piecewise constant solutions in x . ℓ_1 trend filtering with $r = 2$ and $r = 3$ favor piecewise linear solutions and piecewise quadratic solutions, respectively.

Trend filtering problems have been extensively studied in the literature. There have been many approaches for solving it, for example, interior point method [30], ADMM type methods [45], and pathwise algorithms [50]. There is also an efficient dynamic programming based method for the case when $r = 1$ [29]. Most of the above approaches work only for a specific form of the objective function f . In the following, we show how to utilize uFW and uAFW to solve this problem with a general convex smooth objective f .

Writing (49) in the form of (3): We can write (49) in the form of (3) with

$$T = \ker(D_n^{(r)}), \quad S = \left\{ x \in \mathbb{R}^n \mid \|D_n^{(r)}x\|_1 \leq \delta, x \in (\ker(D_n^{(r)}))^\perp \right\}.$$

4.1.1 Computing projections onto T and T^\perp

The following lemma gives a characterization of $T = \ker(D_n^{(r)})$.

Lemma 4.1. *Let $U \in \mathbb{R}^{n \times n}$ be an upper triangular matrix with all its upper triangular entries being 1. Then*

$$T = \ker(D_n^{(r)}) = \text{span}\{\mathbf{1}_n, U\mathbf{1}_n, U^2\mathbf{1}_n, \dots, U^{r-1}\mathbf{1}_n\}.$$

Proof. It is straightforward to verify the following identity (e.g., by induction)

$$D_n^{(i)}U^i = [I_{n-i}, \mathbf{0}_{(n-i) \times i}], \quad \forall 1 \leq i \leq r. \quad (50)$$

Hence for all $1 \leq i \leq r-1$, it holds $D_n^{(r)}U^i\mathbf{1}_n = 0$, which implies

$$\text{span}\{\mathbf{1}_n, U\mathbf{1}_n, U^2\mathbf{1}_n, \dots, U^{r-1}\mathbf{1}_n\} \subseteq \ker(D_n^{(r)}).$$

On the other hand, assume that $\sum_{i=0}^{r-1} \alpha_i U^i \mathbf{1}_n = 0$ for some $\alpha_0, \dots, \alpha_{r-1} \in \mathbb{R}$. Multiplying $D_n^{(r-1)}$ in both sides and in view of $D_n^{(i+1)}U^i\mathbf{1}_n = 0$, we have $\alpha_{r-1} = 0$. Similarly, we can prove $\alpha_{r-2} = \dots = \alpha_1 = 0$; and therefore $\{\mathbf{1}_n, U\mathbf{1}_n, U^2\mathbf{1}_n, \dots, U^{r-1}\mathbf{1}_n\}$ are linearly independent. Note also that $\dim(\ker(D_n^{(r)})) = r$, thereby completing the proof of the lemma. \square

Let $A := [\mathbf{1}_n, U\mathbf{1}_n, \dots, U^{r-1}\mathbf{1}_n] \in \mathbb{R}^{n \times r}$. To compute the projection operator \mathcal{P}_T , we perform a QR decomposition of $A = QR$ where Q is a matrix with orthogonal columns and the same size of A , and R is an upper triangular square matrix. Note that the columns of Q form an orthogonal basis of subspace T , and the projections onto T and T^\perp are given by $\mathcal{P}_T x = QQ^\top x$ and $\mathcal{P}_T^\perp x = x - \mathcal{P}_T x$.

4.1.2 Solving the linear programming subproblem

The linear oracle in uFW (Algorithm 1) has the following form

$$\min_x \langle c, x \rangle \quad \text{s.t.} \quad \|D_n^{(r)}x\|_1 \leq \delta, x \in (\ker(D_n^{(r)}))^\perp, \quad (51)$$

where $c = \nabla f(y^k)$ in the k -th iteration. As we discuss below, (51) admits a closed-form solution. To this end, note that (51) can be written as

$$\min_x \langle c, x \rangle \quad \text{s.t.} \quad \|D_n^{(r)}x\|_1 \leq \delta, Q^\top x = 0, \quad (52)$$

where $Q \in \mathbb{R}^{n \times r}$ is computed by the QR decomposition as stated above.

We apply a variable transformation via $x = U^r y$. Now define $\bar{c} := (U^r)^\top c = [\bar{c}_1; \bar{c}_2]$, where $\bar{c}_1 \in \mathbb{R}^{n-r}$ is the first $n-r$ components of \bar{c} and $\bar{c}_2 \in \mathbb{R}^r$ is the last r components of \bar{c} . Similarly, define $y = [z; w]$ where

$z \in \mathbb{R}^{n-r}$ is the first $n-r$ components of y and $w \in \mathbb{R}^r$ is the last r components of y . We let $B_1 \in \mathbb{R}^{r \times (n-r)}$ and $B_2 \in \mathbb{R}^{r \times r}$ be such that $Q^\top U^r = [B_1, B_2]$. Then problem (52) can be converted to

$$\min_{z,w} \langle \bar{c}_1, z \rangle + \langle \bar{c}_2, w \rangle \quad \text{s.t.} \quad \|z\|_1 \leq \delta, \quad B_1 z + B_2 w = 0 .$$

From the equality constraints, we have $w = -B_2^{-1} B_1 z$ and the problem reduces to

$$\min_z \langle \bar{c}_1 - B_1^\top (B_2^{-1})^\top \bar{c}_2, z \rangle \quad \text{s.t.} \quad \|z\|_1 \leq \delta . \quad (53)$$

Indeed, (53) has a closed-form solution. Denote $\tilde{c} := \bar{c}_1 - B_1^\top (B_2^{-1})^\top \bar{c}_2$, and let $j = \operatorname{argmax}_i |\tilde{c}_i|$ be its largest component in absolute value, then the solution of (53) is $z = -\operatorname{sign}(\tilde{c}_j) \delta e_j$, that is, if $\tilde{c}_j \geq 0$, $z = -\delta e_j$; otherwise $z = \delta e_j$. Reversing the above process, we obtain the solution of problem (51).

For uAFW, we need another linear oracle to compute v^k in step (2) of Algorithm 2. This can be computed in a manner similar to that outlined above. Note that the vertex set $V(x^k)$ is a subset of vertices of S . As above, this problem can also be transformed to a similar form as (53) with a fewer number of variables.

The cost of computing the linear oracle above is $O(nr)$. In addition, there is a one-time cost of $O(nr^2)$ for computing the QR decomposition of A and $B_1^\top (B_2^{-1})^\top$.

4.2 Matrix Completion with Generalized Nuclear Norm Constraints

The nuclear norm is often used as a convex surrogate for the rank of a matrix and arises in the matrix completion problem [8, 35, 19]. Here we consider a family of matrix problems with generalized nuclear norm constraints, and the goal is to solve problem (2) where, $P \in \mathbb{R}^{k \times m}$, $Q \in \mathbb{R}^{n \times l}$ are given matrices that encode prior information about the problem under consideration².

In spite of its wide applicability [3, 48, 16], solving (2) is computationally expensive compared to the case when $P = I$ and $Q = I$, which is amenable to the usual Frank-Wolfe method [28, 19]. Below we show that our proposed framework suggests a useful approach for (2)—in particular, we need to compute the largest singular value/vector pair at every iteration, which can be performed efficiently via iterative methods having a low-computational cost per iteration [21, 19].

Writing (2) in the form of (3): Define the linear operator $\varphi_{P,Q} : X \mapsto PXQ$, then the feasible set of problem (2) becomes $\mathcal{F} := \{X \in \mathbb{R}^{m \times n} \mid \|\varphi_{P,Q}(X)\|_* \leq \delta\}$. Thus, we can convert (2) to (3) with

$$T = \ker(\varphi_{P,Q}) \quad \text{and} \quad S = \{X \in \mathbb{R}^{m \times n} \mid X \in \ker(\varphi_{P,Q})^\perp, \|PXQ\|_* \leq \delta\} .$$

4.2.1 Projections onto T and T^\perp

Recall that \mathcal{P}_T and \mathcal{P}_T^\perp denote the projections onto spaces T and T^\perp respectively. Note that the matrix representation of $\varphi_{P,Q}$ is $Q^\top \otimes P$ (Kronecker product). Notice

$$T^\perp = \ker(Q^\top \otimes P)^\perp = \operatorname{range}(Q \otimes P^\top) ,$$

thus the matrix representation of \mathcal{P}_{T^\perp} can be written as

$$(Q \otimes P^\top)(Q \otimes P^\top)^+ = (Q \otimes P^\top)(Q^+ \otimes P^{+\top}) = (QQ^+) \otimes (P^\top P^{+\top}) = (QQ^+) \otimes (P^+ P) ,$$

where P^+ and Q^+ denote the Moore-Penrose inverses of P and Q (respectively). Therefore, we have $\mathcal{P}_T^\perp(X) = P^+ PXQQ^+$ and $\mathcal{P}_T(X) = X - P^+ PXQQ^+$.

²We assume that there exists a finite optimal solution to (2).

4.2.2 Solving the linear subproblem

We show that the linear subproblem in the uFW framework for problem (2) can be solved via computing the leading singular vector of a matrix. Note that in step (2) of the uFW algorithm, we need to solve a problem of the following form

$$\min_X \langle C, X \rangle \quad \text{s.t.} \quad X \in T^\perp, \quad \|PXQ\|_* \leq \delta, \quad (54)$$

where C is a given matrix in $\mathbb{R}^{m \times n}$. Lemma 4.2 shows how to solve (54).

Lemma 4.2. *Let $\bar{C} := (P^+)^\top C (Q^+)^\top$. Suppose the SVD of \bar{C} is $\bar{C} = UDV^\top$, where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices, and D is in $\mathbb{R}^{m \times n}$ with $D_{ii} = \sigma_i$ for $i = 1, \dots, r$ (r is the rank of \bar{C}) and all other entries being 0, satisfying $\sigma_1 \geq \dots \geq \sigma_r > 0$. Let u_1 be the first column of U , v_1 be the first column of V . Then, the rank-one matrix $X^* := -\delta P^+ u_1 v_1^\top Q^+$ is an optimal solution to problem (54).*

Proof. First, we show that X^* satisfies the constraints. Note that

$$T^\perp = \ker(Q^\top \otimes P)^\perp = \text{range}[(Q^\top \otimes P)^+] = \text{range}[(Q^+)^\top \otimes P^+],$$

which means $Y \in T^\perp$ if and only if $\exists Z \in \mathbb{R}^{m \times n}$ such that $Y = P^+ Z Q^+$. Therefore, we have $X^* \in T^\perp$. In addition, we have

$$\|PX^*Q\|_* = \delta \|PP^+ u_1 v_1^\top Q^+ Q\|_* = \delta |v_1^\top Q^+ Q P P^+ u_1|,$$

and hence

$$\|PX^*Q\|_* \leq \delta \|Q^+ Q v_1\|_2 \|P P^+ u_1\|_2 \leq \delta \|v_1\|_2 \|u_1\|_2 = \delta,$$

where the second inequality is because PP^+ and Q^+Q are projection matrices.

Now we verify that X^* is an optimal solution. First, note that

$$\langle C, X^* \rangle = -\delta \langle C, P^+ u_1 v_1^\top Q^+ \rangle = -\delta u_1^\top (P^+)^\top C (Q^+)^\top v_1 = -\delta u_1^\top \bar{C} v_1 = -\delta \sigma_1.$$

Secondly, for any X satisfying $X \in T^\perp$ and $\|PXQ\|_* \leq \delta$, we have $X = \mathcal{P}_T^\perp X = P^+ P X Q Q^+$, and

$$\langle X, C \rangle = \langle P^+ P X Q Q^+, C \rangle = \langle P X Q, (P^+)^\top C (Q^+)^\top \rangle = \langle P X Q, \bar{C} \rangle.$$

Therefore,

$$\langle X, C \rangle \geq -\sigma_1 \|P X Q\|_* \geq -\sigma_1 \delta,$$

where the first inequality is because the spectral norm and the nuclear norm are dual of each other. This finishes the proof. \square

Remark 4.3. *In several applications [16], the matrices P and Q are projection matrices (i.e., $P^+ = P$ and $Q = Q^+$). In such cases (54) can be computed more efficiently. Using the same notations as the statement of Lemma 4.2, let $\bar{C} = PCQ$ and u_1 and v_1 be the leading (left and right) singular vectors of \bar{C} . Then an optimal solution of (54) is given by $X^* = -\delta u_1 v_1^\top$.*

Note that u_1 and v_1 can be computed by iterative procedures [21] for large problems. When P , Q and C are sparse or low-rank, we can further explore these structures via the matrix-vector multiplication steps within the iterative procedure. This makes our algorithm computationally efficient.

4.3 Other Applications

Our presented algorithms (uFW and uAFW) can be generalized to solve the following family of learning problems:

$$\min_x \sum_{i=1}^N h(b^i, a^i, x) \quad \text{s.t.} \quad \varphi(Hx) \leq \delta, \quad (55)$$

where $a^i \in \mathbb{R}^n$ is the feature vector and $b^i \in \mathbb{R}$ is the corresponding response of the i -th sample, $x \in \mathbb{R}^n$ is the model parameter to be optimized, h is a given smooth convex loss function, H is a given matrix and φ is a nonnegative function. When φ is a norm and H does not have full row rank, the constraint set in (55) is unbounded. The constraint set can be written in the form $T \oplus S$ with $T = \ker(H)$ and $S = \{x \in \mathbb{R}^n \mid \varphi(Hx) \leq \delta, x \in (\ker(H))^\perp\}$, hence our algorithms uFW/uAFW can be applied. Concrete examples of this framework include total variation denoising [47], group fused lasso [1] and convex clustering [41].

5 Numerical Experiments

We present numerical experiments with uFW and uAFW for Problems (1) and (2). Our code is available at

<https://github.com/wanghaoyue123/Frank-Wolfe-with-unbounded-constraints>.

5.1 ℓ_1 trend filtering

We first consider the ℓ_1 trend filtering problem.

Data generation. We generate a Gaussian ensemble $A = [a_1, a_2, \dots, a_n]^\top \in \mathbb{R}^{N \times n}$ with iid $N(0, 1)$ entries, and the noise is $\epsilon_i \sim N(0, \sigma^2)$ with variance $\sigma^2 > 0$ for $i \in [N]$. The underlying model coefficient $x^* \in \mathbb{R}^n$ is generated as a piecewise constant vector when $r = 1$ and a piecewise linear vector when $r = 2$. More precisely, for $r = 1$, x^* is piecewise constant with 5 pieces, each of equal length. The value of each piece is generated from the uniform distribution on $[-1/2, 1/2]$ and then normalized such that $\|D_n^{(1)}x^*\|_1 = 1$. For $r = 2$, x^* is piecewise linear with 5 pieces. These 5 pieces have equal lengths, and the slope of each piece is generated from the uniform distribution on $[-1/2, 1/2]$ and then normalized such that $\|D_n^{(2)}x^*\|_1 = 1$. The response is then obtained from the following linear model: $b_i = a_i^T x^* + \epsilon_i$, $i \in [N]$. To estimate x^* , we consider the following problem:

$$\min_x \|b - Ax\|_2^2 \quad \text{s.t.} \quad \|D_n^{(r)}x\|_1 \leq \delta, \quad (56)$$

where the bound δ is taken as $\|D_n^{(r)}x^*\|_1$. We define the *Signal-noise-ratio (SNR)* of the problem as $SNR := \|Ax^*\|^2 / (n\sigma^2)$.

Computational environment. Our algorithms are written in Python 3.7.4, and we compare with MOSEK [2] and ADMM-based solver SCS [37, 38] (SCS for short) by calling the solvers through CVXPY [12]. Computations were performed on MIT Sloan’s engaging cluster with 4 CPUs and 8GB RAM per CPU. The reported results are averaged over three independent trials.

Performance of uFW and uAFW. For the sequence $\{x^k\}$ produced by our algorithms, let f_k be the lowest objective value computed in the first k iterations, and G_k, H_k be the values defined in (14). We let the (relative) G_k -gap be $G_k / \max\{1, |f_k|\}$ and (relative) H_k^2 -gap be $H_k^2 / \max\{1, |f_k|\}$ —they measure algorithm progress and can be used to design termination criteria for our algorithms. Additionally, we measure (relative) *optimality gap* defined as $(f(x^k) - f^*) / \max\{1, |f^*|\}$, where f^* is the optimal objective value computed by MOSEK. The relative measurement adjusts for the scale of the problem, and is commonly used in first-order solvers, such as SCS.

Figure 1 presents the optimality gap, G_k -gap and H_k^2 -gap of uFW and uAFW versus the number of iterations. We consider (56) for $r \in \{1, 2\}$. The data is generated as above with $N = 1000$, $n = 500$ and $SNR = 1$. We compare the following methods:

- *uFW (simple)*: Algorithm 1 with simple step-size rule (5).
- *uFW (linesearch)*: Algorithm 1 with line-search rule (4).

- *uAFW (linesearch)*: Algorithm 2 with line-search rule (4).

The corresponding gradient descent step-size η is chosen as $1/\|A\|^2$ where $\|A\|$ is the operator norm of the data matrix A .

The left panel ($r = 1$) in Figure 1 suggests that the away step variant: uAFW converges linearly to an optimal solution for the ℓ_1 trend filtering problem with $r = 1$, and it can reach very high accuracy. In contrast, uFW algorithms appear to converge at a sub-linear rate. Both these observations are consistent with our theoretical framework. The advantage of uAFW over uFW is less significant in the right panel ($r = 2$) in Figure 1. We believe that this is due to the problem being very ill-conditioned for the case $r = 2$ —uAFW converges linearly, but with an unfavorable parameter for linear convergence that makes it progress slowly. Figure 1 also shows that uFW with line search step-size has a larger optimality gap compared to uFW with simple step-size. This is perhaps because the problem is ill-conditioned, and exact line-search may lead to conservative step-size. However, the G_k -gap and H_k^2 -gap of these methods are comparable. Finally, note that in this example, H_k^2 -gap is much smaller than G_k -gap and optimality gap. This is perhaps because the dimension of T is low and the corresponding gradient steps in T converge faster.

Comparison with Benchmarks. Table 1 compares uFW with two state-of-the-art conic optimization solvers: MOSEK and SCS, for solving the ℓ_1 trend filtering problem (56) with $r = 1$ and $r = 2$. We also compare with the accelerated projection free method for general unbounded constraint set [22] (denoted by APFA).

The data is generated by the procedure discussed above with $SNR = 1$. We use uFW with the simple step-size rule (5), and terminate uFW at iteration k if $G_k/\max\{|f_k|, 1\} < 10^{-4}$ and $H_k^2/\max\{|f_k|, 1\} < 10^{-4}$. We run MOSEK and SCS with their default settings. For reference, Table 1 also reports the (relative) *optimality gap*, denoted as uFW gap. We set the termination tolerance of SCS as 10^{-3} and allow a maximum runtime of 2 hours. For APFA, we use Mosek to solve the linear oracles and the parameters are set following Section 4 of [22]. Let \hat{f} be the objective value computed by uFW upon termination, and \tilde{f}_k be the objective value computed in iteration k of APFA. We terminate APFA at iteration k if $(\tilde{f}_k - \hat{f})/\max\{1, |\hat{f}|\} < 10^{-2}$, with a runtime limit of 2 hours.

For each value of $r \in \{1, 2\}$, Table 1 presents results for three groups of problem sizes when $N \ll n$, $N \approx n$ and $N \gg n$. As we can see, the running time of uFW is significantly better than MOSEK and SCS for all the examples (one or two orders of magnitude better in most of the examples). As the problem sizes grow, the runtimes of uFW appear to be stable. MOSEK and SCS on the other hand, may fail to produce a reasonable solution (within the allocated maximum time-limit of 2hrs and/or memory). We observe that when $G_k/\max\{|f_k|, 1\} < 10^{-4}$ and $H_k^2/\max\{|f_k|, 1\} < 10^{-4}$, the optimality gap of uFW is typically around $10^{-5} \sim 10^{-6}$. Note that uFW takes longer to solve the problem with $r = 2$ than the case $r = 1$. This is because the problem with $r = 2$ has a larger condition number. It can be seen that the runtime of APFA is much longer than the other algorithms even for finding a low-accuracy solution. This is probably due to expensive linear optimization oracles³.

Finally, we note that MOSEK, SCS and uFW/uAFW are three different types of algorithms, each using different termination criteria. MOSEK is an interior-point solver, which naturally produces a high accuracy solution. SCS (based on ADMM) is a primal-dual first-order method, which may produce infeasible solutions, thus potentially negative primal-dual gaps (See Section 3.5 of [38] for a detailed discussion on the termination criteria of SCS). uFW and uAFW are primal algorithms with feasible iterates, and the relative (H_k, G_k) -gaps provide a natural measure of solution quality. The runtime improvements in Table 1 appear to suggest the advantage of uFW over competing methods.

³We need to minimize a linear function over the intersection of $T \oplus S$ and a Euclidean ball—as we are not aware of a simple (closed-form) solution to this, we solve this with Mosek.

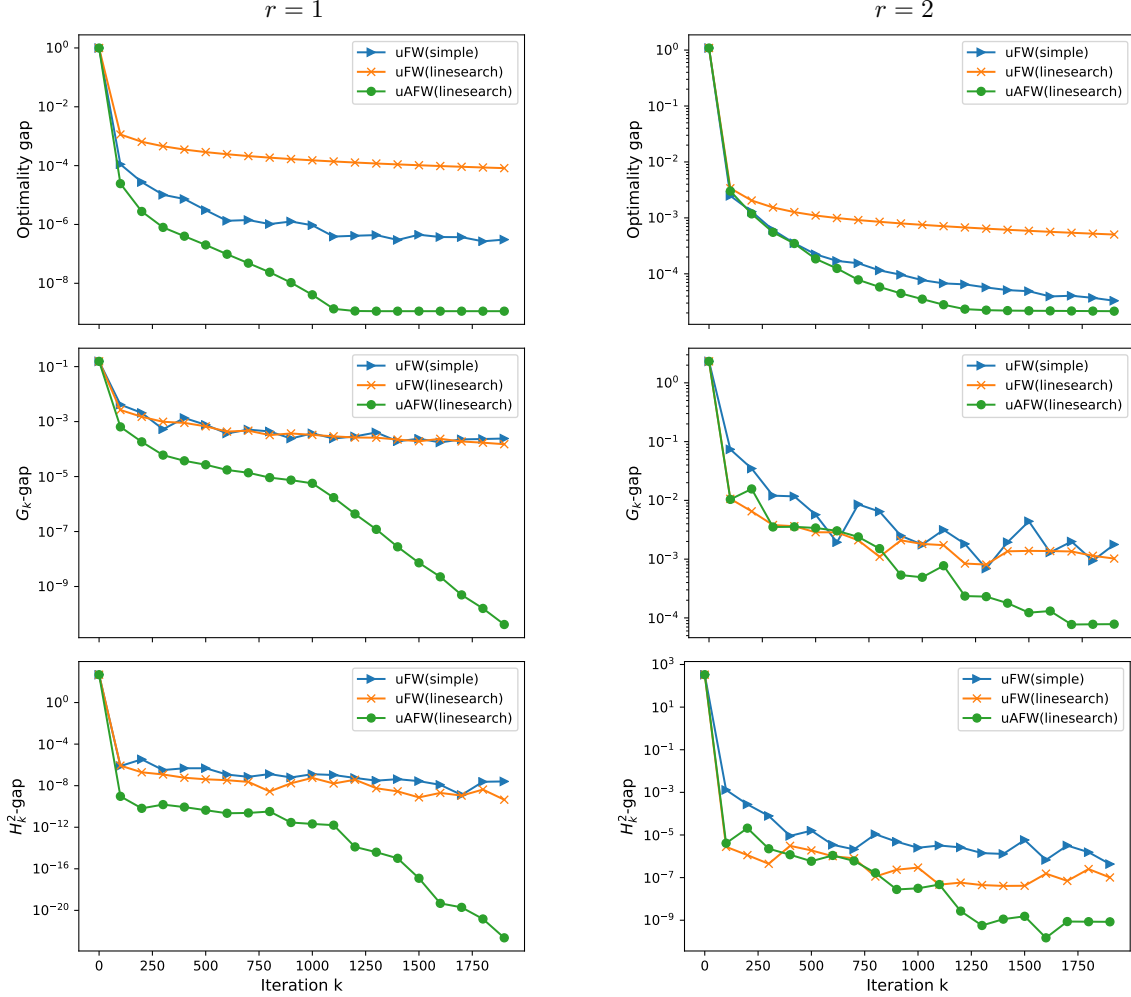


Figure 1: Plots showing the optimality gap, G_k -gap and H_k^2 -gap obtained by uFW and uAFW versus the number of iterations (i.e., denoted by index k in the algorithm descriptions) for solving the ℓ_1 trend filtering problem (56) with $r = 1$ and $r = 2$, respectively.

5.2 Matrix Completion with Side Information

Here we present computational results for the generalized matrix completion problem (2).

Data generation. The matrix of interest B is generated from the model

$$B = P_1 Z^\top + UV^\top + \mathcal{E} \quad (57)$$

where $B \in \mathbb{R}^{m \times n}$, $P_1 \in \mathbb{R}^{m \times r_1}$, $Z \in \mathbb{R}^{r_1 \times n}$, $U \in \mathbb{R}^{m \times r}$, $V \in \mathbb{R}^{n \times r}$ and $\mathcal{E} \in \mathbb{R}^{m \times n}$. Here P_1 corresponds to the known side information about the column space of B , with coefficients Z ; the low-rank component that is not included in this side information is UV^\top , and \mathcal{E} is the noise term. See [16, 15] for more background on this statistical model. The matrices U, V, Z are generated independently with entries (iid) from $\mathcal{N}(0, 1)$. The matrix P_1 has unit-norm orthogonal columns, which is generated by taking the orthogonal basis for the columns of an $m \times r_1$ random matrix with iid standard Gaussian entries. The entries of \mathcal{E} are iid $\mathcal{N}(0, \sigma^2)$.

We study the setting where only a subset of the entries of B are observed. Let Ω be the indices corresponding

Table 1: Table showing runtimes (in secs) for solving (56) by uFW, MOSEK and SCS with $r = 1$ and $r = 2$. For uFW we also display the optimality gap when f^* is available. The symbol “-” for MOSEK means the allocated memory was not sufficient. The symbol “*” means the method did not terminate in 2 hours or the solver reported an error. The symbol “x” means that MOSEK reported an error during execution.

N	n	r=1					r=2				
		uFW time(s)	uFW gap	MOSEK time(s)	SCS time(s)	APFA time(s)	uFW time(s)	uFW gap	MOSEK time(s)	SCS time(s)	APFA time(s)
5K	500	0.27	3.25e-07	3.42	10.70	226.17	1.68	3.02e-06	2.39	53.19	221.08
200K	1K	7.62	7.04e-07	541.82	3879.35	*	28.82	2.98e-06	327.76	4114.53	*
400K	1K	17.48	-	-	*	*	56.50	-	-	*	*
300K	2K	21.32	-	-	*	*	45.69	-	-	*	*
2K	2K	0.37	4.66e-07	7.17	43.14	939.62	0.94	6.20e-06	7.74	97.12	*
10K	10K	2.03	6.03e-07	302.14	2526.14	*	10.24	5.83e-07	448.85	3252.30	*
15K	15K	4.36	-	-	*	*	22.39	-	-	*	*
20K	20K	7.51	-	-	*	*	48.14	-	-	*	*
500	5K	0.23	1.14e-06	4.47	111.60	*	2.97	1.79e-06	4.72	109.12	*
1K	100K	6.29	2.45e-06	262.80	4763.89	*	48.29	-	x	4816.81	*
1K	200K	6.44	-	-	*	*	101.54	-	-	*	*
2K	300K	25.41	-	-	*	*	173.91	-	-	*	*

to the observed entries that are uniformly distributed across the matrix coordinates. For a matrix $A \in \mathbb{R}^{m \times n}$, we let $\mathcal{P}_\Omega(A)$ be a matrix with entries in Ω being the same as A and entries in Ω^c being zeros. We estimate the signal “ $P_1 Z^\top + UV^\top$ ” by solving the following problem (see [16]):

$$\min_X \|\mathcal{P}_\Omega(X - B)\|_F^2 \quad \text{s.t.} \quad \|(I - P_1 P_1^\top)X\|_* \leq \delta, \quad (58)$$

which is a special case of (2) with $P = I_m - P_1 P_1^\top$, $Q = I_n$ and objective function $f(X) = \|\mathcal{P}_\Omega(X - B)\|_F^2$ (here, $\|\cdot\|_F$ denotes the Frobenius norm).

Computational environment. Our code is written in Matlab 2017, and we make use of Matlab built-in function *svds* to compute the leading singular vector with a custom function for matrix-vector multiplication. We compare our proposal with SCS [38] (via CVX). Computations were performed on MIT Sloan’s engaging cluster with 4 CPUs and 20GB RAM for each CPU. Our results are averaged over three independent experiments.

Comparison with SCS. We compare the computation time of uFW and SCS on instances with $r = r_1 = 5$ and different values of m and n , generated from the procedure stated above. We define the following three parameters for the experimental setting:

- *Signal-to-noise Ratio* (SNR): the ratio of signal variance vs noise variance, i.e.,

$$\text{SNR} := \text{var}((P_1 Z^\top + UV^\top)_{i,j}) / \text{var}(\mathcal{E}_{i,j}).$$

- *Non-zero ratio* (nnzr): the percentage of observed entries, i.e., $\text{nnzr} := |\Omega| / (mn)$.
- *Relative diameter* (δ'): the relative value of the parameter δ in (58) with respect to the corresponding value of signal. That is,

$$\delta' := \delta / \|(I_m - P_1 P_1^\top)UV^\top\|_*.$$

We use uFW with the simple step size rule (5). Let f_k be the best objective value obtained in the first k iterations of uFW. We terminate uFW when $G_k / \max\{|f_k|, 1\} \leq 3 \times 10^{-3}$ and $H_k^2 / \max\{|f_k|, 1\} \leq 3 \times 10^{-3}$ and we set the SCS tolerance to be 3×10^{-3} . For instances when SCS is able to output a solution, we set f^* as the objective value of the solution obtained by running SCS with a stricter tolerance 10^{-5} , and define the (relative) *optimality gap* of uFW as $(f_k - f^*) / \max\{1, |f^*|\}$. We also define the optimality gap of SCS as $(\hat{f} - f^*) / \max\{1, |f^*|\}$, where \hat{f} is the objective value for the SCS solution with tolerance 3×10^{-3} . Note

Table 2: Comparison of uFW and SCS for solving (58) with $r = r_1 = 5$. Here, “gap” refers to relative optimality gap as defined in the text. Note that SCS time is the runtime of SCS with tolerance 10^{-3} .

	delta	uFW time	uFW gap	SCS time	SCS gap	SCS feas.
m=700, n=700, nnzr=0.3	0.5	18.18	9.63e-05	694.01	7.00e-05	4.71e-03
	0.8	90.14	5.20e-05	785.85	6.99e-05	4.20e-03
	1.0	194.47	1.01e-04	793.35	1.86e-04	5.34e-03
m=1000, n=500, nnzr=0.3	0.5	13.33	7.81e-05	952.35	2.62e-04	1.68e-02
	0.8	73.34	4.68e-05	1172.53	1.29e-04	6.20e-03
	1.0	180.20	8.48e-05	1261.50	-2.21e-04	2.37e-03
m=300, n=3000, nnzr=0.2	0.5	14.38	1.56e-04	1309.38	-6.96e-04	3.30e-02
	0.8	243.56	2.56e-04	1912.09	8.20e-04	8.97e-03
	1.0	299.84	1.88e-04	1735.94	2.81e-03	1.25e-02

that SCS is an ADMM based solver, so the solution it obtains is not strictly feasible. To this end, we define the value SCS feasibility (*SCS feas.*, in short) to be

$$\text{SCS feas.} = (\|(I_m - P_1 P_1^\top) \tilde{X}\|_* - \delta) / \delta,$$

where \tilde{X} is the the solution obtained upon its termination.

Table 2 presents the comparison of uFW and SCS on runtime and solution accuracy. As presented in Table 2, the computational time of uFW is significantly less than that of SCS. For problems with a smaller δ' (≈ 0.5), uFW is around $30 \sim 80$ times faster than SCS. For problems with larger δ' (≈ 1), even though uFW gets slower, it is still about $3 \sim 5$ times faster than SCS in runtime. Moreover, upon termination, uFW typically finds a solution with (relative) optimality gap less than 3×10^{-4} , which appears to be more accurate than the solution found by SCS with tolerance 3×10^{-3} (even though these two solutions are not strictly comparable) – the latter typically has a *SCS feas.* larger than 3×10^{-3} . Table 2 shows that the “SCS feas.” value is positive (i.e., the solution is not feasible) for all instances. This is the reason why some of the SCS gaps reported in Table 2 are negative.

Table 3 presents more instances with larger values of m and n on which SCS fails to output a solution (for a tolerance of 3×10^{-3} , as above). The termination rule for uFW is the same as mentioned above. Due to its mild per-iteration cost, uFW is able to solve these problems approximately within minutes to hours.

Table 3: Running times of uFW for large instances on which SCS would not run.

	δ	uFW time		δ	uFW time
m=1000, n=1000, nnzr=0.2	0.5	20.64	m=3000, n=300, nnzr=0.2	0.5	16.91
	0.8	100.75		0.8	122.60
	1.0	245.68		1.0	399.60
m=3000, n=3000, nnzr=0.1	0.5	64.68	m=7000, n=7000, nnzr=0.05	0.5	212.86
	0.8	609.70		0.8	3520.89
	1.0	1842.15		1.0	7498.59

6 Acknowledgements

The authors would like to thank the AE and the Reviewers for their comments that led to improvements in the manuscript.

A Proof of Proposition 3.11

The proof of Proposition 3.11 here closely follows the proof of Proposition 9 in [24] with two major differences: (i) we apply the analysis to the function $\tilde{f}(\cdot, w)$ for a dynamically changing w ; and (ii) instead of the relative strong convexity to the constraint, we utilize a lower bound. Since these results do not appear in [24], we present the full proof here.

We first introduce some new notations and establish three auxiliary lemmas (Lemmas A.1–A.3). Let N be the number of vertices of S , and Δ_N be the standard simplex in \mathbb{R}^N , that is, $\Delta_N = \{x \in \mathbb{R}^N : x_i \geq 0 \forall i \in [N], \mathbf{1}_N^\top x = 1\}$. First, we introduce an auxiliary function \tilde{f} . Let $B_1 \in \mathbb{R}^{(n-r) \times N}$ be the matrix whose columns are vertices of S , and let $B := \begin{bmatrix} B_1 & 0 \\ 0 & I_r \end{bmatrix} \in \mathbb{R}^{n \times (N+r)}$. Define $\tilde{f} : \mathbb{R}^{N+r} \rightarrow \mathbb{R} \cup \{\infty\}$ such that $\tilde{f} := \bar{f} \circ B$. Note that $\text{dom}(\tilde{f}) \supseteq \Delta_N \times \mathbb{R}^r$ and recall the definition of $\psi(S)$ appearing in (38).

Before analyzing the Frank-Wolfe steps using the function \tilde{f} , we need a result regarding the strong-convexity of \tilde{f} in certain directions. For $\tilde{u}^0 \in \Delta_N$, we define

$$Z_{B_1, \Delta_N}(\tilde{u}^0) := \{\tilde{u} \in \Delta_N \mid B_1 \tilde{u} = B_1 \tilde{u}^0\} .$$

Lemma A.1. *Suppose f satisfies Assumption 3.7 and \tilde{f} is as defined above. Define $\tilde{\mu} = \mu\psi(S)^2/4$. Then for any $(\tilde{u}^1, w), (\tilde{u}^2, w) \in \Delta_N \times \mathbb{R}^r$ satisfying*

$$\|\tilde{u}^2 - \tilde{u}^1\|_1 = \inf_{\tilde{u} \in Z_{B_1, \Delta_N}(\tilde{u}^2)} \|\tilde{u} - \tilde{u}^1\|_1 , \quad (59)$$

it holds

$$\tilde{f}(\tilde{u}^2, w) \geq \tilde{f}(\tilde{u}^1, w) + \langle \nabla_{\tilde{u}} \tilde{f}(\tilde{u}^1, w), \tilde{u}^2 - \tilde{u}^1 \rangle + (\tilde{\mu}/2) \|\tilde{u}^2 - \tilde{u}^1\|_1^2 .$$

Proof. Let $\tilde{\mu}'$ be the largest possible value such that

$$\tilde{f}(\tilde{u}^2, w) \geq \tilde{f}(\tilde{u}^1, w) + \langle \nabla_{\tilde{u}} \tilde{f}(\tilde{u}^1, w), \tilde{u}^2 - \tilde{u}^1 \rangle + (\tilde{\mu}'/2) \|\tilde{u}^2 - \tilde{u}^1\|_1^2 \quad (60)$$

holds for all \tilde{u}^1 and \tilde{u}^2 satisfying (59). For $\tilde{u}, \tilde{v} \in \Delta_N$, we introduce the notation

$$\|Z_{B_1, \Delta_N}(\tilde{v}) - \tilde{u}\|_1 := \inf_{\tilde{v}^1 \in Z_{B_1, \Delta_N}(\tilde{v})} \|\tilde{v}^1 - \tilde{u}\|_1 .$$

Fix $w \in \mathbb{R}^r$ and consider the functions $\tilde{f}_w(\cdot) = \tilde{f}(\cdot, w) : \Delta_N \rightarrow \mathbb{R}$ and $\bar{f}_w(\cdot) = \bar{f}(\cdot, w) : S \rightarrow \mathbb{R}$, then it is immediate that $\tilde{f}_w = \bar{f}_w \circ B_1$, and (60) is equivalent to

$$\tilde{f}_w(\tilde{u}^2) \geq \tilde{f}_w(\tilde{u}^1) + \langle \nabla \tilde{f}_w(\tilde{u}^1), \tilde{u}^2 - \tilde{u}^1 \rangle + (\tilde{\mu}'/2) \|\tilde{u}^2 - \tilde{u}^1\|_1^2 .$$

It then follows from Theorem 1 of [24] that

$$\tilde{\mu}' \geq \inf_{\substack{\tilde{u}, \tilde{v} \in \Delta_N, \\ \tilde{u} \notin Z_{B_1, \Delta_N}(\tilde{v})}} \frac{\mu \|B_1(\tilde{u} - \tilde{v})\|_2^2}{\|Z_{B_1, \Delta_N}(\tilde{v}) - \tilde{u}\|_1^2} . \quad (61)$$

Furthermore, by Proposition 1 of [25], the right hand side of (61) equals $\mu\psi(S)^2/4 = \tilde{\mu}$. \square

Lemma A.2. *For any $a, b \in \Delta_N$, there exist $p, q \in \Delta_N$ such that*

$$a - b = \|a - b\|_1(p - q)/2 \quad \text{and} \quad \text{supp}(p) \subseteq \text{supp}(a) ,$$

where $\text{supp}(a)$ and $\text{supp}(p)$ denote the indices of nonzero coordinates of a and p respectively.

Proof. For any vector $x \in \mathbb{R}^n$, let x^+ be the vector in \mathbb{R}^n with $x_i^+ = \max\{x_i, 0\}$ and $x^- = x^+ - x$. Assume $a \neq b$ (otherwise the conclusion is trivial). Let

$$p := 2(a-b)^+/\|a-b\|_1 \quad \text{and} \quad q := 2(a-b)^-/\|a-b\|_1.$$

Then we have $a-b = (\|a-b\|_1/2)(p-q)$ and $\text{supp}(p) \subseteq \text{supp}(a)$. Note that

$$1_N^\top p - 1_N^\top q = \frac{2}{\|a-b\|_1}(1_N^\top a - 1_N^\top b) = 0, \quad 1_N^\top p + 1_N^\top q = \frac{2}{\|a-b\|_1}\|a-b\|_1 = 2.$$

As a result, it holds $1_N^\top p = 1_N^\top q = 1$, hence $p, q \in \Delta_N$. \square

Lemma A.3 provides a key inequality that allows us to establish a contraction in optimality gap. Recall that we have defined $F: \mathbb{R}^r \rightarrow \mathbb{R}$ as $F(w) = \inf_{u \in S} \bar{f}(u, w)$.

Lemma A.3. *Let $\{(u^k, w^k, d^k)\}_{k \geq 0}$ be the iterates generated by Algorithm 3. Then*

$$\langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle^2 \geq 2\tilde{\mu}(\bar{f}(u^k, w^{k+1}) - F(w^{k+1})).$$

Proof. Denote $u_{w^{k+1}}^* \in \text{argmin}_{u \in S} \bar{f}(u, w^{k+1})$. Let $\{\tilde{u}^k\}_{k \geq 0}$ be a sequence in Δ_N such that $u^k = B_1 \tilde{u}^k$, and $\tilde{u}_{w^{k+1}}^*$ be a point in \mathbb{R}^N such that

$$\tilde{u}_{w^{k+1}}^* \in \text{argmin}_{\tilde{u} \in \Delta_N} \{\|\tilde{u} - \tilde{u}^k\|_1 \mid B_1 \tilde{u} = u_{w^{k+1}}^*\}.$$

Recall that $V(u^k)$ is the subset of vertices of S corresponding to the support of u^k . Let $\tilde{V}(u^k)$ be the subset of vertices of Δ_N corresponding to $V(u^k)$. Then we know that $\tilde{u}^k = \sum_{\tilde{v} \in \tilde{V}(u^k)} \lambda_{\tilde{v}} \tilde{v}$ where $\lambda_{\tilde{v}} > 0$. (We say that \tilde{u}^k is supported on $\tilde{V}(u^k)$). By Lemma A.2 with $a = \tilde{u}^k$ and $b = \tilde{u}_{w^{k+1}}^*$, there exist $\tilde{p}, \tilde{q} \in \Delta_N$ such that

$$\tilde{u}^k - \tilde{u}_{w^{k+1}}^* = (\beta/2)(\tilde{p} - \tilde{q}), \tag{62}$$

where $\beta := \|\tilde{u}^k - \tilde{u}_{w^{k+1}}^*\|_1$, and the support of \tilde{p} is a subset of $\tilde{V}(u^k)$.

As a result, if we let $p = B_1 \tilde{p}$ and $q = B_1 \tilde{q}$, then we have $p \in \text{conv}(V(u^k))$. In addition, since $1_N^\top \tilde{q} = 1$ and the columns of B_1 correspond to the vertices of S , we have $q \in S$. Multiplying both sides of equality (62) by B_1 , we have

$$u^k - u_{w^{k+1}}^* = (\beta/2)(p - q). \tag{63}$$

With these results at hand, we get

$$\begin{aligned} & \langle \nabla_u \bar{f}(u^k, w^{k+1}), u^k - u_{w^{k+1}}^* \rangle \\ & \stackrel{(i)}{=} (\beta/2) \langle \nabla_u \bar{f}(u^k, w^{k+1}), p - q \rangle \\ & \stackrel{(ii)}{\leq} (\beta/2) \left(\max_{u \in \text{conv}(V(u^k))} \langle \nabla_u \bar{f}(u^k, w^{k+1}), u \rangle - \min_{u \in S} \langle \nabla_u \bar{f}(u^k, w^{k+1}), u \rangle \right) \\ & \stackrel{(iii)}{=} (\beta/2) \langle \nabla_u \bar{f}(u^k, w^{k+1}), v^k - s^k \rangle \end{aligned} \tag{64}$$

where (i) is from (63); (ii) is from the fact $p \in \text{conv}(V(u^k))$ and $q \in S$; (iii) is from the definitions of s^k and v^k in Algorithm 3. From the definition of d^k in Algorithm 3, we have

$$\begin{aligned} \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle & \leq (1/2)[\langle \nabla_u \bar{f}(u^k, w^{k+1}), s^k - u^k \rangle + \langle \nabla_u \bar{f}(u^k, w^{k+1}), u^k - v^k \rangle] \\ & = (1/2) \langle \nabla_u \bar{f}(u^k, w^{k+1}), s^k - v^k \rangle. \end{aligned}$$

Combining with (64), we have

$$\langle \nabla_u \bar{f}(u^k, w^{k+1}), u^k - u_{w^{k+1}}^* \rangle \leq \frac{\beta}{2} \langle \nabla_u \bar{f}(u^k, w^{k+1}), v^k - s^k \rangle \leq -\beta \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle .$$

On the other hand, from the strong convexity of \tilde{f} in the first block (Lemma A.1), we have

$$\begin{aligned} F(w^{k+1}) = \tilde{f}(\tilde{u}_{w^{k+1}}^*, w^{k+1}) &\geq \tilde{f}(\tilde{u}^k, w^{k+1}) + \langle \nabla_{\tilde{u}} \tilde{f}(\tilde{u}^k, w^{k+1}), \tilde{u}_{w^{k+1}}^* - \tilde{u}^k \rangle + \tilde{\mu} \beta^2 / 2 \\ &\stackrel{(i)}{=} \bar{f}(u^k, w^{k+1}) + \langle \nabla_u \bar{f}(u^k, w^{k+1}), u_{w^{k+1}}^* - u^k \rangle + \tilde{\mu} \beta^2 / 2 \\ &\geq \bar{f}(u^k, w^{k+1}) + \beta \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle + \tilde{\mu} \beta^2 / 2 , \end{aligned} \quad (65)$$

where (i) is because $\tilde{f}(\tilde{u}^k, w^{k+1}) = \bar{f}(B_1 \tilde{u}^k, w^{k+1}) = \bar{f}(u^k, w^{k+1})$ and

$$\begin{aligned} \nabla_{\tilde{u}} \tilde{f}(\tilde{u}^k, w^{k+1}) &= \nabla_{\tilde{u}} \bar{f}(B_1 \tilde{u}^k, w^{k+1}) = B_1^\top \nabla_u \bar{f}(u^k, w^{k+1}) \\ \implies \langle \nabla_{\tilde{u}} \tilde{f}(\tilde{u}^k, w^{k+1}), \tilde{u}_{w^{k+1}}^* - \tilde{u}^k \rangle &= \langle \nabla_u \bar{f}(u^k, w^{k+1}), B_1 \tilde{u}_{w^{k+1}}^* - B_1 \tilde{u}^k \rangle \\ &= \langle \nabla_u \bar{f}(u^k, w^{k+1}), u_{w^{k+1}}^* - u^k \rangle . \end{aligned}$$

Therefore, from (65), one has

$$\begin{aligned} -\langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle &\geq (1/\beta)(\bar{f}(u^k, w^{k+1}) - F(w^{k+1})) + \tilde{\mu} \beta / 2 \\ &\geq \sqrt{2\tilde{\mu}(\bar{f}(u^k, w^{k+1}) - F(w^{k+1}))} \end{aligned}$$

where the last step is by Cauchy-Schwarz inequality. \square

Now we are ready to prove Proposition 3.11:

Proof of Proposition 3.11: We prove Proposition 3.11 by discussing two different cases.

(Case a) This case includes two subcases:

(Case a.1) $\alpha_k < \alpha_{\max}$;

(Case a.2) $\alpha_k = \alpha_{\max}$, and iteration k takes a FW step.

When either (Case a.1) or (Case a.2) happens, we have $|V(u^{k+1})| \leq |V(u^k)| + 1$. Furthermore, we claim that

$$\bar{f}(u^{k+1}, w^{k+1}) \leq \min_{\alpha \in [0, 1]} \{ \bar{f}(u^k + \alpha d^k, w^{k+1}) \} . \quad (66)$$

Indeed, when (Case a.2) happens, then $\alpha_{\max} = 1$, so by the updating rule, (66) holds true. When (Case a.1) happens, we have

$$\bar{f}(u^{k+1}, w^{k+1}) = \min_{\alpha \in [0, \alpha_{\max}]} \{ \bar{f}(u^k + \alpha d^k, w^{k+1}) \} = \min_{\alpha \in [0, \infty)} \{ \bar{f}(u^k + \alpha d^k, w^{k+1}) \}$$

where, the first equality is from the updating rule of the algorithm, and the second equality is because $\alpha_k < \alpha_{\max}$ and f is convex. As a result, claim (66) holds true.

From (66) and the definition of \bar{C}_{f, x^0}^S we have

$$\begin{aligned} &\bar{f}(u^{k+1}, w^{k+1}) \\ &\leq \min_{\alpha \in [0, 1]} \left\{ \bar{f}(u^k, w^{k+1}) + \alpha \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle + \bar{C}_{f, x^0}^S \alpha^2 / 2 \right\} . \end{aligned} \quad (67)$$

Let $\bar{\alpha}_k$ be the optimal solution in (67). If $\bar{\alpha}_k < 1$, then (67) implies

$$\begin{aligned}\bar{f}(u^{k+1}, w^{k+1}) &\leq \bar{f}(u^k, w^{k+1}) - \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle^2 / (2\bar{C}_{f,x^0}^S) \\ &\leq \bar{f}(u^k, w^{k+1}) - \tilde{\mu} (\bar{f}(u^k, w^{k+1}) - F(w^{k+1})) / \bar{C}_{f,x^0}^S\end{aligned}\quad (68)$$

where the second inequality is by Lemma A.3. Recall that $\gamma = \mu\psi(S)^2 / (4\bar{C}_{f,x^0}^S) = \tilde{\mu} / \bar{C}_{f,x^0}^S$. As a result, we have

$$\bar{f}(u^{k+1}, w^{k+1}) - f^* \leq (1 - \gamma) (\bar{f}(u^k, w^{k+1}) - f^*) + \gamma (F(w^{k+1}) - f^*). \quad (69)$$

If $\bar{\alpha}_k = 1$, then by taking derivative w.r.t α in (67), we know that

$$-\langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle \geq \bar{C}_{f,x^0}^S.$$

Combining this with (67) again we have

$$\begin{aligned}\bar{f}(u^{k+1}, w^{k+1}) &\leq \bar{f}(u^k, w^{k+1}) + \min_{\alpha \in [0,1]} \{(\alpha - \alpha^2/2) \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle\} \\ &= \bar{f}(u^k, w^{k+1}) + (1/2) \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle,\end{aligned}\quad (70)$$

where the last equality is because $\langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle \leq 0$. On the other hand,

$$\begin{aligned}\langle \nabla_u \bar{f}(u^k, w^{k+1}), u_{w^{k+1}}^* - u^k \rangle &\stackrel{(i)}{\geq} \langle \nabla_u \bar{f}(u^k, w^{k+1}), s^k - u^k \rangle \\ &\stackrel{(ii)}{\geq} \langle \nabla_u \bar{f}(u^k, w^{k+1}), d^k \rangle\end{aligned}\quad (71)$$

where (i) is from the definition of s^k , and (ii) is from the definition of d^k in Algorithm 3. Combining (71) and (70), we have

$$\begin{aligned}\bar{f}(u^{k+1}, w^{k+1}) &\leq \bar{f}(u^k, w^{k+1}) + (1/2) \langle \nabla_u \bar{f}(u^k, w^{k+1}), u_{w^{k+1}}^* - u^k \rangle \\ &\leq \bar{f}(u^k, w^{k+1}) + (1/2) (F(w^{k+1}) - \bar{f}(u^k, w^{k+1})),\end{aligned}$$

where the last inequality is from the convexity of f . As a result, we arrive at

$$\bar{f}(u^{k+1}, w^{k+1}) - f^* \leq (1/2) [\bar{f}(u^k, w^{k+1}) - f^*] + (1/2) [F(w^{k+1}) - f^*]. \quad (72)$$

Recall that $\rho = \min\{\gamma, 1/2\}$ and $F(w^{k+1}) \leq \bar{f}(u^{k+1}, w^{k+1})$, so combining (69) and (72) we have

$$\bar{f}(u^{k+1}, w^{k+1}) - f^* \leq (1 - \rho) (\bar{f}(u^k, w^{k+1}) - f^*) + \rho (F(w^{k+1}) - f^*). \quad (73)$$

(Case b): $\alpha_k = \alpha_{\max}$ and iteration k takes an away step. In this case, we have $|V(u^{k+1})| \leq |V(u^k)| - 1$. Since $|V(u^0)| = 1$ and $|V(u^k)| \geq 1$ for $k \geq 1$, it follows that for each iteration when (Case b) occurs, there must have been at least one earlier iteration index at which (Case a) occurred. Hence in the first k iterations, (Case b) occurs at most $k/2$ times.

Combining the discussions in (Case a) and (Case b), we reach the conclusion.

B Proof of Proposition 3.10

First, by (19) (this inequality also holds for iterations generated by Algorithm 2), we know

$$\|\mathcal{P}_T \nabla f(x^k)\|_2 \leq \sqrt{(2/\eta)(f(x^k) - f(y^k))} \leq \sqrt{(2/\eta)(f(x^k) - f^*)}. \quad (74)$$

Combining the above with (35) (and recalling, $\tilde{H}_k = \|\mathcal{P}_T \nabla f(x^k)\|_2$), we have

$$H_k \leq (1 + \eta L_{f,x^0}^T) \|\mathcal{P}_T \nabla f(x^k)\|_2 \leq (1 + \eta L_{f,x^0}^T) \sqrt{(2/\eta)(f(x^k) - f^*)}. \quad (75)$$

This completes the proof of Part (1).

To prove Part (2), similar to the proof of Theorem 3.8, we reduce the arguments to the axis-aligned case (Algorithm 3). Note that in this case, $x^k = (u^k, w^k)$ and $y^k = (u^k, w^{k+1})$. In the proof of Proposition 3.11, we show that in the first k iterations, there are at least $k/2$ iterations such that (Case a) happens. If (Case a) happens with $\bar{\alpha}_k < 1$, then we have

$$G_k \leq -\langle \nabla f(y^k), d^k \rangle \leq \sqrt{2\bar{C}_{f,x^0}^S(f(y^k) - f(x^{k+1}))} \leq \sqrt{2\bar{C}_{f,x^0}^S(f(y^k) - f^*)} \quad (76)$$

where the first inequality is by the definition of d^k , and the second is by the first inequality in (68). If (Case a) happens with $\bar{\alpha}_k = 1$, then we have

$$G_k \leq -\langle \nabla f(y^k), d^k \rangle \leq 2(f(y^k) - f(w^{k+1})) \leq 2(f(y^k) - f^*) \quad (77)$$

where the second inequality is by (70). Combining (76) and (77), the proof of Part (2) is complete.

References

- [1] C. M. Alaíz, A. Barbero, and J. Dorransoro. Group fused lasso. In *International Conference on Artificial Neural Networks*, pages 66–73. Springer, 2013.
- [2] E. D. Andersen and K. D. Andersen. The mosek interior point optimizer for linear programming: an implementation of the homogeneous algorithm. In *High performance optimization*, pages 197–232. Springer, 2000.
- [3] R. Angst, C. Zach, and M. Pollefeys. The generalized trace-norm and its application to structure-from-motion problems. In *2011 International Conference on Computer Vision*, pages 2502–2509. IEEE, 2011.
- [4] A. Beck and S. Shtern. Linearly convergent away-step conditional gradient for non-strongly convex functions. *Mathematical Programming*, 164(1-2):1–27, 2017.
- [5] Immanuel M Bomze, Francesco Rinaldi, and Samuel Rota Buló. First-order methods for the impatient: Support identification in finite time with convergent frank–wolfe variants. *SIAM Journal on Optimization*, 29(3):2211–2226, 2019.
- [6] Immanuel M Bomze, Francesco Rinaldi, and Damiano Zeffiro. Active set complexity of the away-step frank–wolfe algorithm. *SIAM Journal on Optimization*, 30(3):2470–2500, 2020.
- [7] J. F. Bonnans and A. Shapiro. Optimization problems with perturbations: A guided tour. *SIAM review*, 40(2):228–264, 1998.
- [8] E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717, 2009.
- [9] K. Chiang, C. Hsieh, and I. Dhillon. Matrix completion with noisy side information. In *Advances in Neural Information Processing Systems*, pages 3447–3455, 2015.
- [10] K. L. Clarkson. Coresets, sparse greedy approximation, and the frank-wolfe algorithm. *ACM Transactions on Algorithms (TALG)*, 6(4):1–30, 2010.
- [11] V. Demyanov and A. Rubinov. *Approximate methods in optimization problems*.
- [12] S. Diamond and S. Boyd. Cvxpy: A python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(1):2909–2913, 2016.
- [13] J. C. Dunn. Rates of convergence for conditional gradient algorithms near singular and nonsingular extremals. *SIAM Journal on Control and Optimization*, 17(2):187–211, 1979.
- [14] J. C. Dunn and S. Harshbarger. Conditional gradient algorithms with open loop step size rules. *Journal of Mathematical Analysis and Applications*, 62(2):432–444, 1978.

- [15] A. Eftekhari, D. Yang, and M. B. Wakin. Weighted matrix completion and recovery with prior subspace information. *IEEE Transactions on Information Theory*, 64(6):4044–4071, 2018.
- [16] W. Fithian and R. Mazumder. Flexible low-rank statistical modeling with missing data and side information. *Statistical Science*, 33 No.2:238–260.
- [17] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 1956.
- [18] R. M. Freund and P. Grigas. New analysis and results for the frank-wolfe method. *Mathematical Programming*, 155:199–230, 2016.
- [19] R. M. Freund, P. Grigas, and R. Mazumder. An extended frank-wolfe method with in-face directions, and its application to low-rank matrix completion. *SIAM Journal on Optimization*, 27:319–346, 2013.
- [20] Dan Garber and Elad Hazan. A linearly convergent variant of the conditional gradient algorithm under strong convexity, with applications to online and stochastic optimization. *SIAM Journal on Optimization*, 26(3):1493–1528, 2016.
- [21] G. Golub and C. Van Loan. *Matrix Computations (3rd Ed.)*. Johns Hopkins University Press, Baltimore, MD, USA, 1996.
- [22] Max LN Gonçalves, Jefferson G Melo, and Renato DC Monteiro. Projection-free accelerated method for convex optimization. *Optimization Methods and Software*, pages 1–27, 2020.
- [23] J. Guélat and P. Marcotte. Some comments on wolfe’s ‘away step’. *Mathematical Programming*, 35(1):110–119, 1986.
- [24] D. H. Gutman and J. F. Pena. The condition number of a function relative to a set. *Mathematical Programming*, pages 1–40, 2020.
- [25] David H Gutman and Javier F Pena. The condition of a function relative to a polytope. *arXiv preprint arXiv:1802.00271*, 2018.
- [26] Z. Harchaoui, A. Juditsky, and A. Nemirovski. Conditional gradient algorithms for norm-regularized smooth convex optimization. *Mathematical Programming*, 152(1-2):75–112, 2015.
- [27] E. Hazan. Sparse approximate solutions to semidefinite programs. In *Latin American symposium on theoretical informatics*, pages 306–316. Springer, 2008.
- [28] M. Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th international conference on machine learning*, pages 427–435, 2013.
- [29] N. Johnson. A dynamic programming algorithm for the fused lasso and l_0 -segmentation. *Journal of Computational and Graphical Statistics*, 22(2):246–260, 2013.
- [30] S. Kim, K. Koh, S. Boyd, and D. Gorinevsky. l_1 trend filtering. *Siam Review*, 51. No. 2:339–360.
- [31] S. Lacoste-Julien and M. Jaggi. On the global linear convergence of frank-wolfe optimization variants. In *Advances in neural information processing systems*, pages 496–504, 2015.
- [32] Guanghui Lan. The complexity of large-scale convex programming under a linear optimization oracle. *arXiv preprint arXiv:1309.5550*, 2013.
- [33] Guanghui Lan. *First-order and Stochastic Optimization Methods for Machine Learning*. Springer, 2020.
- [34] Guanghui Lan and Yi Zhou. Conditional gradient sliding for convex optimization. *SIAM Journal on Optimization*, 26(2):1379–1409, 2016.
- [35] R. Mazumder, T. Hastie, and R. Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *Journal of Machine Learning Research*, 11:2287–2322.
- [36] Y. Nesterov. *Introductory Lectures on Convex Optimization*. Kluwer Academic Publishers, 2003.
- [37] B. O’Donoghue, E. Chu, N. Parikh, and S. Boyd. SCS: Splitting conic solver, version 2.1.2. <https://github.com/cvxgrp/scs>, November 2019.
- [38] B. O’Donoghue, E. Chu, N. Parikh, and S. Boyd. Conic optimization via operator splitting and homogeneous self-dual embedding. *Journal of Optimization Theory and Applications*, 169(3):1042–1068, 2016.
- [39] N. Parikh, S. Boyd, et al. Proximal algorithms. *Foundations and Trends® in Optimization*, 1(3):127–239, 2014.

- [40] Fabian Pedregosa, Geoffrey Negiar, Armin Askari, and Martin Jaggi. Linearly convergent frank-wolfe with backtracking line-search. In *International Conference on Artificial Intelligence and Statistics*, pages 1–10. PMLR, 2020.
- [41] K. Pelckmans, J. De Brabanter, J. Suykens, and B. De Moor. Convex clustering shrinkage. In *PASCAL Workshop on Statistics and Optimization of Clustering Workshop*, 2005.
- [42] J. Pena and D. Rodriguez. Polytope conditioning and linear convergence of the frank-wolfe algorithm. *Mathematics of Operations Research*, 44(1):1–18, 2019.
- [43] J. Pena, D. Rodríguez, and N. Soheili. On the von neumann and frank-wolfe algorithms with away steps. *SIAM Journal on Optimization*, 26(1):499–512, 2016.
- [44] B. T. Polyak. Introduction to optimization. *Optimization software, Inc., New York*, 1, 1987.
- [45] A. Ramdas and R. J. Tibshirani. Fast and flexible admm algorithms for trend filtering. *Journal of Computational and Graphical Statistics*, 25(3):839–858, 2016.
- [46] R. T. Rockafellar. *Convex analysis*. Princeton university press, 1970.
- [47] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992.
- [48] N. Srebro and R. R. Salakhutdinov. Collaborative filtering in a non-uniform world: Learning with the weighted trace norm. In *Advances in Neural Information Processing Systems*, pages 2056–2064, 2010.
- [49] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B*, 67:91–108.
- [50] R. Tibshirani and J. Taylor. The solution path of the generalized lasso. *Annals of Statistics*, 39(3):1335–1371.
- [51] R. J. Tibshirani. Adaptive piecewise polynomial estimation via trend filtering. *Annals of Statistics*, 42(1):285–323.
- [52] B. Wahlberg, S. Boyd, M. Annergren, and Y. Wang. An admm algorithm for a class of total variation regularized estimation problems. *IFAC Proceedings Volumes*, 45(16):83–88, 2012.
- [53] P. Wolfe. Convergence theory in nonlinear programming. *Integer and nonlinear programming*, pages 1–36, 1970.