

SCALABLE ROBUST MATRIX RECOVERY: FRANK-WOLFE MEETS PROXIMAL METHODS

CUN MU*, YUQIAN ZHANG†, JOHN WRIGHT†, AND DONALD GOLDFARB*

Abstract. Recovering matrices from compressive and grossly corrupted observations is a fundamental problem in robust statistics, with rich applications in computer vision and machine learning. In theory, under certain conditions, this problem can be solved in polynomial time via a natural convex relaxation, known as Compressive Principal Component Pursuit (CPCP). However, many existing provably convergent algorithms for CPCP suffer from superlinear per-iteration cost, which severely limits their applicability to large-scale problems. In this paper, we propose provably convergent, scalable and efficient methods to solve CPCP with (essentially) linear per-iteration cost. Our method combines classical ideas from Frank-Wolfe and proximal methods. In each iteration, we mainly exploit Frank-Wolfe to update the low-rank component with rank-one SVD and exploit the proximal step for the sparse term. Convergence results and implementation details are discussed. We demonstrate the practicability and scalability of our approach with numerical experiments on visual data.

Key words. robust matrix recovery, compressive principal component pursuit, Frank-Wolfe, conditional gradient, proximal methods, scalability

AMS subject classifications. 90C06, 90C25, 90C52

1. Introduction. Suppose that a matrix $\mathbf{M}_0 \in \mathbb{R}^{m \times n}$ is of the form $\mathbf{M}_0 = \mathbf{L}_0 + \mathbf{S}_0 + \mathbf{N}_0$, where \mathbf{L}_0 is a low-rank matrix, \mathbf{S}_0 is a sparse error matrix, and \mathbf{N}_0 is a dense noise matrix. Linear measurements

$$(1.1) \quad \mathbf{b} = \mathcal{A}[\mathbf{M}_0] = (\langle \mathbf{A}_1, \mathbf{M}_0 \rangle, \langle \mathbf{A}_2, \mathbf{M}_0 \rangle, \dots, \langle \mathbf{A}_p, \mathbf{M}_0 \rangle)^\top \in \mathbb{R}^p$$

are collected, where $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$ is the sensing operator, \mathbf{A}_k is the sensing matrix for the k -th measurement and $\langle \mathbf{A}_k, \mathbf{M}_0 \rangle \doteq \text{Tr}(\mathbf{M}_0^\top \mathbf{A}_k)$. Can we, in a tractable way, recover \mathbf{L}_0 and \mathbf{S}_0 from \mathbf{b} , given \mathcal{A} ?

One natural approach is to solve the optimization combining the fidelity term and the structural terms:

$$(1.2) \quad \min_{\mathbf{L}, \mathbf{S}} \frac{1}{2} \|\mathbf{b} - \mathcal{A}[\mathbf{L} + \mathbf{S}]\|_2^2 + \lambda_L \text{rank}(\mathbf{L}) + \lambda_S \|\mathbf{S}\|_0.$$

Here, λ_L and λ_S are regularization parameters, and $\|\mathbf{S}\|_0$ denotes the number of nonzero entries in \mathbf{S} .

Unfortunately, problem (1.2) is nonconvex, and hence is not directly tractable. However, by replacing the ℓ_0 norm $\|\mathbf{S}\|_0$ with the ℓ_1 norm $\|\mathbf{S}\|_1 \doteq \sum_{i=1}^m \sum_{j=1}^n |S_{ij}|$, and replacing the rank $\text{rank}(\mathbf{L})$ with the nuclear norm $\|\mathbf{L}\|_*$ (defined as the sum of the singular values of \mathbf{L}), we obtain a natural, tractable, convex relaxation of (1.2),

$$(1.3) \quad \min_{\mathbf{L}, \mathbf{S}} \frac{1}{2} \|\mathbf{b} - \mathcal{A}[\mathbf{L} + \mathbf{S}]\|_2^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1.$$

This convex surrogate is sometimes referred to as *compressive principal component pursuit (CPCP)* [1]. Equivalently, since

$$\{ \mathbf{M} \in \mathbb{R}^{m \times n} \mid \mathbf{b} = \mathcal{A}[\mathbf{M}] \} = \{ \mathbf{M} \in \mathbb{R}^{m \times n} \mid \mathcal{P}_Q[\mathbf{M}] = \mathcal{P}_Q[\mathbf{M}_0] \},$$

*Department of Industrial Engineering and Operations Research, Columbia University (cm3052@columbia.edu, goldfarb@columbia.edu). DG was funded by NSF Grants DMS-1016571 and CCF-1527809.

†Department of Electrical Engineering, Columbia University, (yq2409@cs.columbia.edu, johnwright@ee.columbia.edu). JW was funded by ONR-N00014-13-0492.

where $\mathcal{Q} \subseteq \mathbb{R}^{m \times n}$ is a linear subspace spanned by the set of sensing matrices $\{\mathbf{A}_i\}_{i=1}^p$, and $\mathcal{P}_{\mathcal{Q}}$ denotes the projection operator onto that subspace, we can rewrite problem (1.3) in the (possibly) more compact form, *

$$(1.4) \quad \min_{\mathbf{L}, \mathbf{S}} f(\mathbf{L}, \mathbf{S}) \doteq \frac{1}{2} \|\mathcal{P}_{\mathcal{Q}}[\mathbf{L} + \mathbf{S} - \mathbf{M}_0]\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1.$$

Recently, CPCP and its close variants have been studied for different sensing operators \mathcal{A} (or equivalently different subspaces \mathcal{Q}). In specific, [2, 3, 4, 5, 6] consider the case where a subset $\Omega \subseteq \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$ of the entries of \mathbf{M}_0 is observed. Then CPCP can be reduced to

$$(1.5) \quad \min_{\mathbf{L}, \mathbf{S}} \frac{1}{2} \|\mathcal{P}_{\Omega}[\mathbf{L} + \mathbf{S} - \mathbf{M}_0]\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1,$$

where $\mathcal{P}_{\Omega}[\cdot]$ denotes the orthogonal projection onto the linear space of matrices supported on Ω , i.e., $\mathcal{P}_{\Omega}[\mathbf{M}_0](i, j) = (\mathbf{M}_0)_{ij}$ if $(i, j) \in \Omega$ and $\mathcal{P}_{\Omega}[\mathbf{M}_0](i, j) = 0$ otherwise. [1] studies the case where each \mathbf{A}_k is an i.i.d. $\mathcal{N}(0, 1)$ matrix, which is equivalent (in distribution) to saying that we choose a linear subspace \mathcal{Q} uniformly at random from the set of all p -dimensional subspaces of $\mathbb{R}^{m \times n}$ and observe $\mathcal{P}_{\mathcal{Q}}[\mathbf{M}_0]$. Accordingly, all the above provide theoretical guarantees for CPCP, under fairly mild conditions, to produce accurate estimates of \mathbf{L}_0 and $\mathcal{P}_{\Omega}[\mathbf{S}_0]$ (or \mathbf{S}_0), even when the number of measurements p is substantially less than mn .

Inspired by these theoretical results, researchers from different fields have leveraged CPCP to solve many practical problems, including video background modeling [3], batch image alignment [7], face verification [8], photometric stereo [9], dynamic MRI [10], topic modeling [11], latent variable graphical model learning [12] and outlier detection and robust Principal Component Analysis [3], to name just a few.

Living in the era of *big data*, most of these applications involve large datasets and high dimensional data spaces. Therefore, to fully realize the benefit of the theory, we need *provably convergent* and *scalable* algorithms for CPCP. This has motivated much research into the development of first-order methods for problem (1.4) and its variants; e.g see [13, 14, 15, 16, 17, 18]. These methods, in essence, all exploit a closed-form expression for the proximal operator of the nuclear norm, which involves the singular value decomposition (SVD). Hence, the dominant cost in each iteration is computing an SVD of the same size as the input data. This is substantially more scalable than off-the-shelf interior point solvers such as SDPT3 [19]. Nevertheless, the superlinear cost of each iteration has limited the practical applicability of these first-order methods to problems involving several thousands of data points and several thousands of dimensions. The need to compute a sequence of full or partial SVDs is a serious bottleneck for truly large-scale applications.

As a remedy, in this paper, we design more scalable algorithms to solve CPCP that compute only a rank-one SVD in each iteration. Our approach leverages two classical and widely studied ideas – Frank-Wolfe iterations to handle the nuclear norm, and proximal steps to handle the ℓ_1 norm. This turns out to be an ideal combination of techniques to solve large-scale CPCP problems. In particular, it yields algorithms that

*To transform problem (1.3) into problem (1.4), simple procedures like Gram-chmidt might be invoked. Despite being equivalent, one formulation might be preferred over the other in practice, depending on the specifications of the sensing operator $\mathcal{A}[\cdot]$. In this paper, we will mainly focus on solving problem (1.4) and its variants. Our methods, however, are not restrictive to (1.4) and can be easily extended to problem (1.3).

are substantially *more scalable* than prox-based first-order methods such as ISTA and FISTA [20], and converge *much faster* in practice than a straightforward application of Frank-Wolfe.

The remainder of this paper is organized as follows. Section 2 reviews the general properties of the Frank-Wolfe algorithm, and describes several basic building blocks that we will use in our algorithms. Section 3 and Section 4 respectively describe how to modify the Frank-Wolfe algorithm to solve CPCP's *norm constrained* version

$$(1.6) \quad \min_{\mathbf{L}, \mathbf{S}} l(\mathbf{L}, \mathbf{S}) \doteq \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}_0]\|_F^2 \quad \text{s.t.} \quad \|\mathbf{L}\|_* \leq \tau_L, \|\mathbf{S}\|_1 \leq \tau_S,$$

and the penalized version, i.e. problem (1.4), by incorporating proximal regularization to more effectively handle the ℓ_1 norm. Convergence results and our implementation details are also discussed. Section 5 presents numerical experiments on large datasets that demonstrate the scalability of our proposed algorithms. In Section 6, we summarize our contributions and discuss potential future works.

2. Preliminaries.

2.1. Frank-Wolfe method. The Frank-Wolfe (FW) method [21], also known as the conditional gradient method [22], applies to the general problem of minimizing a differentiable convex function h over a compact, convex domain $\mathcal{D} \subseteq \mathbb{R}^n$:

$$(2.1) \quad \text{minimize} \quad h(\mathbf{x}) \quad \text{subject to} \quad \mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^n.$$

Here, ∇h is assumed to be L -Lipschitz:

$$(2.2) \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D}, \quad \|\nabla h(\mathbf{x}) - \nabla h(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|.$$

Throughout, we let $D = \max_{\mathbf{x}, \mathbf{y} \in \mathcal{D}} \|\mathbf{x} - \mathbf{y}\|$ denote the diameter of the feasible set \mathcal{D} .

In its simplest form, the Frank-Wolfe algorithm proceeds as follows. At each iteration k , we linearize the objective function h about the current point \mathbf{x}^k :

$$(2.3) \quad h(\mathbf{v}) \approx h(\mathbf{x}^k) + \langle \nabla h(\mathbf{x}^k), \mathbf{v} - \mathbf{x}^k \rangle.$$

We minimize the linearization over the feasible set \mathcal{D} to obtain

$$(2.4) \quad \mathbf{v}^k \in \arg \min_{\mathbf{v} \in \mathcal{D}} \langle \nabla h(\mathbf{x}^k), \mathbf{v} \rangle,$$

and then take a step in the feasible descent direction $\mathbf{v}^k - \mathbf{x}^k$:

$$(2.5) \quad \mathbf{x}^{k+1} = \mathbf{x}^k + \frac{2}{k+2}(\mathbf{v}^k - \mathbf{x}^k).$$

This yields a very simple procedure, which we summarize as Algorithm 1. The particular step size, $\frac{2}{k+2}$, comes from the convergence analysis of the algorithm, which we discuss in more details below.

First proposed in [21], FW-type methods have been frequently revisited in different fields. Recently, they have experienced a resurgence in statistics, machine learning and signal processing, due to their ability to yield highly scalable algorithms for optimization with structure-encouraging norms such as the ℓ_1 norm and nuclear norm. In particular, if \mathbf{x} is a matrix and $\mathcal{D} = \{\mathbf{x} \mid \|\mathbf{x}\|_* \leq \beta\}$ is a nuclear norm ball, the subproblem

$$(2.6) \quad \min_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}) \rangle$$

Algorithm 1 Frank-Wolfe method for problem (2.1)

```

1: Initialization:  $\mathbf{x}^0 \in \mathcal{D}$ ;
2: for  $k = 0, 1, 2, \dots$  do
3:    $\mathbf{v}^k \in \operatorname{argmin}_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}^k) \rangle$ ;
4:    $\gamma = \frac{2}{k+2}$ ;
5:    $\mathbf{x}^{k+1} = \mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)$ ;
6: end for

```

can be solved using only the singular vector pair corresponding to the single leading singular value of the matrix $\nabla h(\mathbf{x})$. Thus, at each iteration, we only have to compute a rank-one partial SVD. This is substantially cheaper than the full/partial SVD exploited in proximal methods [23, 24]. We recommend [25] as a comprehensive survey of the latest developments in FW-type methods.

Algorithm 2 Frank-Wolfe method for problem (2.1) with general updating scheme

```

1: Initialization:  $\mathbf{x}^0 \in \mathcal{D}$ ;
2: for  $k = 0, 1, 2, \dots$  do
3:    $\mathbf{v}^k \in \operatorname{argmin}_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}^k) \rangle$ ;
4:    $\gamma = \frac{2}{k+2}$ ;
5:   Update  $\mathbf{x}^{k+1}$  to some point in  $\mathcal{D}$  such that  $h(\mathbf{x}^{k+1}) \leq h(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k))$ ;
6: end for

```

In the past five decades, numerous variants of Algorithm 1 have been proposed and implemented. Many modify Algorithm 1 by replacing the simple updating rule (2.5) with more sophisticated schemes, e.g.,

$$(2.7) \quad \mathbf{x}^{k+1} \in \operatorname{argmin}_{\mathbf{x}} h(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in \operatorname{conv}\{\mathbf{x}^k, \mathbf{v}^k\}$$

or

$$(2.8) \quad \mathbf{x}^{k+1} \in \operatorname{argmin}_{\mathbf{x}} h(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in \operatorname{conv}\{\mathbf{x}^k, \mathbf{v}^k, \mathbf{v}^{k-1}, \dots, \mathbf{v}^{k-j}\}.$$

The convergence of these schemes can be analyzed simultaneously, using the fact that they produce iterates \mathbf{x}^{k+1} whose objective is no greater than that produced by the original Frank-Wolfe update scheme:

$$h(\mathbf{x}^{k+1}) \leq h(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)).$$

Algorithm 2 states a general version of Frank-Wolfe, whose update is only required to satisfy this relationship. It includes as special cases the updating rules (2.5), (2.7) and (2.8). This flexibility will be crucial for effectively handling the sparse structure in the CPCP problems (1.4) and (1.6).

The convergence of Algorithm 2 can be proved using well-established techniques [24, 25, 26, 27, 28, 29, 30, 31]. Using these ideas, one can show that it converges at a rate of $O(1/k)$ in function value:

THEOREM 2.1. *Let \mathbf{x}^* be an optimal solution to (2.1). For $\{\mathbf{x}^k\}$ generated by Algorithm 2, we have for $k = 0, 1, 2, \dots$,*

$$(2.9) \quad h(\mathbf{x}^k) - h(\mathbf{x}^*) \leq \frac{2LD^2}{k+2}.$$

Proof. For $k = 0, 1, 2, \dots$, we have

$$\begin{aligned}
 h(\mathbf{x}^{k+1}) &\leq h(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)) \\
 &\leq h(\mathbf{x}^k) + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \frac{L\gamma^2}{2} \|\mathbf{v}^k - \mathbf{x}^k\|^2 \\
 (2.10) \quad &\leq h(\mathbf{x}^k) + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \frac{\gamma^2 LD^2}{2} \\
 &\leq h(\mathbf{x}^k) + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{x}^* - \mathbf{x}^k \rangle + \frac{\gamma^2 LD^2}{2}
 \end{aligned}$$

$$(2.11) \quad \leq h(\mathbf{x}^k) + \gamma(h(\mathbf{x}^*) - h(\mathbf{x}^k)) + \frac{\gamma^2 LD^2}{2},$$

where the second inequality holds since $\nabla h(\cdot)$ is L -Lipschitz continuous; the third line follows because D is the diameter for the feasible set \mathcal{D} ; the fourth inequality follows from $\mathbf{v}^k \in \operatorname{argmin}_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}^k) \rangle$ and $\mathbf{x}^* \in \mathcal{D}$; the last one holds since $h(\cdot)$ is convex.

Rearranging terms in (2.11), one obtains that for $k = 0, 1, 2, \dots$,

$$(2.12) \quad h(\mathbf{x}^{k+1}) - h(\mathbf{x}^*) \leq (1 - \gamma)(h(\mathbf{x}^k) - h(\mathbf{x}^*)) + \frac{\gamma^2 LD^2}{2}.$$

Therefore, by mathematical induction, it can be verified that

$$h(\mathbf{x}^k) - h(\mathbf{x}^*) \leq \frac{2LD^2}{k+2}, \quad \text{for } k = 1, 2, 3, \dots$$

□

REMARK 1. Note that the constant in the rate of convergence depends on the Lipschitz constant L of h and the diameter D .

While Theorem 2.1 guarantees that Algorithm 2 converges at a rate of $O(1/k)$, in practice it is useful to have a more precise bound on the suboptimality at iterate k . The surrogate duality gap

$$(2.13) \quad d(\mathbf{x}^k) = \langle \mathbf{x}^k - \mathbf{v}^k, \nabla h(\mathbf{x}^k) \rangle,$$

provides a useful upper bound on the suboptimality $h(\mathbf{x}^k) - h(\mathbf{x}^*)$:

$$\begin{aligned}
 h(\mathbf{x}^k) - h(\mathbf{x}^*) &\leq -\langle \mathbf{x}^* - \mathbf{x}^k, \nabla h(\mathbf{x}^k) \rangle \\
 (2.14) \quad &\leq -\min_{\mathbf{v}} \langle \mathbf{v} - \mathbf{x}^k, \nabla h(\mathbf{x}^k) \rangle = \langle \mathbf{x}^k - \mathbf{v}^k, \nabla h(\mathbf{x}^k) \rangle = d(\mathbf{x}^k).
 \end{aligned}$$

This was first proposed in [21] and later [25] showed that $d(\mathbf{x}^k) = O(1/k)$. Next, we provide a refinement of this result, using ideas from [25, 30]:

THEOREM 2.2. Let $\{\mathbf{x}^k\}$ be the sequence generated by Algorithm 2. Then for any $K \geq 1$, there exists $1 \leq k \leq K$ such that

$$(2.15) \quad d(\mathbf{x}^{\hat{k}}) \leq \frac{6LD^2}{K+2}.$$

Proof. For notational convenience, we denote $h^k \doteq h(\mathbf{x}^k)$, $\Delta^k \doteq h(\mathbf{x}^k) - h(\mathbf{x}^*)$, $d^k \doteq d(\mathbf{x}^k)$, $C \doteq 2LD^2$, $B \doteq K+2$, $\hat{k} \doteq \lceil \frac{1}{2}B \rceil - 1$, $\mu \doteq \lceil \frac{1}{2}B \rceil / B$.

Suppose on the contrary that

$$(2.16) \quad d^k > \frac{3C}{B}, \quad \text{for all } k \in \left\{ \lceil \frac{1}{2}B \rceil - 1, \lceil \frac{1}{2}B \rceil, \dots, K \right\}.$$

From (2.10), we know that for any $k \geq 1$

$$(2.17) \quad \Delta^{k+1} \leq \Delta^k + \gamma \langle \nabla h(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \frac{\gamma^2 L D^2}{2} = \Delta^k - \frac{2d^k}{k+2} + \frac{C}{(k+2)^2}.$$

Therefore, by using (2.17) repeatedly, one has

$$\begin{aligned} \Delta^{K+1} &\leq \Delta^{\hat{k}} - \sum_{k=\hat{k}}^K \frac{2d^k}{k+2} + \sum_{k=\hat{k}}^K \frac{C}{(k+2)^2} \\ &< \Delta^{\hat{k}} - \frac{6C}{B} \sum_{k=\hat{k}}^K \frac{1}{k+2} + C \sum_{k=\hat{k}}^K \frac{1}{(k+2)^2} \\ &= \Delta^{\hat{k}} - \frac{6C}{B} \sum_{k=\hat{k}+2}^B \frac{1}{k} + C \sum_{k=\hat{k}+2}^B \frac{1}{k^2} \\ &\leq \frac{C}{\mu B} - \frac{6C}{B} \cdot \frac{B - \hat{k} - 1}{B} + C \cdot \frac{B - \hat{k} - 1}{B(\hat{k}+1)} \\ &= \frac{C}{\mu B} - \frac{6C}{B} (1 - \mu) + \frac{C}{B} \frac{1 - \mu}{\mu} \\ (2.18) \quad &= \frac{C}{\mu B} (2 - 6\mu(1 - \mu) - \mu) \end{aligned}$$

where the second line is due to our assumption (2.16); the fourth line holds since $\Delta^{\hat{k}} \leq \frac{C}{\hat{k}+2}$ by Theorem 1, and $\sum_{k=a}^b \frac{1}{k^2} \leq \frac{b-a+1}{b(a-1)}$ for any $b \geq a > 1$.

Now define $\phi(x) = 2 - 6x(1-x) - x$. Clearly $\phi(\cdot)$ is convex. Since $\phi(\frac{1}{2}) = \phi(\frac{2}{3}) = 0$, we have $\phi(x) \leq 0$ for any $x \in [\frac{1}{2}, \frac{2}{3}]$. As $\mu = \lceil \frac{1}{2}B \rceil / B \in [\frac{1}{2}, \frac{2}{3}]$, from (2.18), we have

$$\Delta^{K+1} = h(\mathbf{x}^{K+1}) - h(\mathbf{x}^*) < \frac{C}{\mu B} \phi(\mu) \leq 0,$$

which is a contradiction. \square

REMARK 2. *The convergence rate for the duality gap matches the one for $h(\mathbf{x}^k) - h(\mathbf{x}^*)$ (see (2.9)), which suggests that the upper bound $d(\mathbf{x}^k)$ can serve as a practical stopping criterion.*

For our problem, the main computational burden in Algorithms 1 and 2 will be solving the linear subproblem $\min_{\mathbf{v} \in \mathcal{D}} \langle \mathbf{v}, \nabla h(\mathbf{x}^k) \rangle$,[†] i.e. minimizing linear functions over the unit balls for $\|\cdot\|_*$ and $\|\cdot\|_1$. Fortunately, both of these operations have simple closed-form solutions, which we will describe in the next section.

2.2. Optimization oracles. We now describe several optimization oracles involving the ℓ_1 norm and the nuclear norm, which serve as the main building blocks for our methods. These oracles have computational costs that are (essentially) linear in the size of the input.

[†]In some situations, we can significantly reduce this cost by solving this problem inexactly [27, 25]. Our algorithms and results can also tolerate inexact step calculations; we omit the discussion here for simplicity.

Minimizing a linear function over the nuclear norm ball. Since the dual norm of the nuclear norm is the operator norm, i.e., $\|\mathbf{Y}\| = \max_{\|\mathbf{X}\|_* \leq 1} \langle \mathbf{Y}, \mathbf{X} \rangle$, the optimization problem

$$(2.19) \quad \text{minimize}_{\mathbf{X}} \langle \mathbf{Y}, \mathbf{X} \rangle \quad \text{subject to } \|\mathbf{X}\|_* \leq 1$$

has optimal value $-\|\mathbf{Y}\|$. One minimizer is the rank-one matrix $\mathbf{X}^* = -\mathbf{u}\mathbf{v}^\top$, where \mathbf{u} and \mathbf{v} are the left- and right- singular vectors corresponding to the leading singular value of \mathbf{Y} , and can be efficiently computed (e.g. using power method).

Minimizing a linear function over the ℓ_1 ball. Since the dual norm of the ℓ_1 norm is the ℓ_∞ norm, i.e., $\|\mathbf{Y}\|_\infty := \max_{(i,j)} |Y_{ij}| = \max_{\|\mathbf{X}\|_1 \leq 1} \langle \mathbf{Y}, \mathbf{X} \rangle$, the optimization problem

$$(2.20) \quad \text{minimize}_{\mathbf{X}} \langle \mathbf{Y}, \mathbf{X} \rangle \quad \text{subject to } \|\mathbf{X}\|_1 \leq 1$$

has optimal value $-\|\mathbf{Y}\|_\infty$. One minimizer is the one-sparse matrix

$$\mathbf{X}^* = -\text{sgn}(Y_{i^*j^*})\mathbf{e}_{i^*}\mathbf{e}_{j^*}^\top,$$

where $(i^*, j^*) \in \arg \max_{(i,j)} |Y_{ij}|$; i.e. \mathbf{X}^* has exactly one nonzero element.

Projection onto the ℓ_1 -ball. To effectively handle the sparse term in the norm constrained problem (1.6), we will need to modify the Frank-Wolfe algorithm by incorporating additional projection steps. For any $\mathbf{Y} \in \mathbb{R}^{m \times n}$ and $\beta > 0$, the projection onto the ℓ_1 -ball:

$$(2.21) \quad \mathcal{P}_{\|\cdot\|_1 \leq \beta}[\mathbf{Y}] = \arg \min_{\|\mathbf{X}\|_1 \leq \beta} \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2,$$

can be easily solved with $O(mn(\log m + \log n))$ cost [32]. Moreover, a divide and conquer algorithm, achieving linear cost in expectation to solve (2.21), has also been proposed in [32].

Proximal mapping of ℓ_1 norm. To effectively handle the sparse term arising in problem (1.4), we will need to modify the Frank-Wolfe algorithm by incorporating additional proximal steps. For any $\mathbf{Y} \in \mathbb{R}^{m \times n}$ and $\lambda > 0$, the proximal mapping of ℓ_1 norm has the following closed-form expression

$$(2.22) \quad \mathcal{T}_\lambda[\mathbf{Y}] = \arg \min_{\mathbf{X} \in \mathbb{R}^{m \times n}} \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 + \lambda \|\mathbf{X}\|_1,$$

where $\mathcal{T}_\lambda : \mathbb{R} \rightarrow \mathbb{R}$ denotes the soft-thresholding operator $\mathcal{T}_\lambda(x) = \text{sgn}(x) \max\{|x| - \lambda, 0\}$, and extension to matrices is obtained by applying the scalar operator $\mathcal{T}_\lambda(\cdot)$ to each element.

3. FW-P Method for Norm Constrained Problem. In this section, we develop scalable algorithms for the norm-constrained compressive principal component pursuit problem,

$$(3.1) \quad \min_{\mathbf{L}, \mathbf{S}} l(\mathbf{L}, \mathbf{S}) = \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 \quad \text{s.t.} \quad \|\mathbf{L}\|_* \leq \tau_L, \|\mathbf{S}\|_1 \leq \tau_S.$$

We first describe a straightforward application of the Frank-Wolfe method to this problem. We will see that although it has relatively cheap iterations, it converges very slowly on typical numerical examples, because it only makes a one-sparse update to the sparse term \mathbf{S} at a time. We will show how to remedy this problem by augmenting the FW iteration with an additional proximal step (essentially a projected gradient step) in each iteration, yielding a new algorithm which updates \mathbf{S} much more efficiently. Because it combines Frank-Wolfe and projection steps, we will call this new algorithm Frank-Wolfe-Projection (FW-P).

Properties of the objective and constraints. To apply Frank-Wolfe to (3.1), we first note that the objective $l(\mathbf{L}, \mathbf{S})$ in (3.1) is differentiable, with

$$(3.2) \quad \nabla_{\mathbf{L}} l(\mathbf{L}, \mathbf{S}) = \mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]$$

$$(3.3) \quad \nabla_{\mathbf{S}} l(\mathbf{L}, \mathbf{S}) = \mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}].$$

Moreover, the following lemma shows that the gradient map $\nabla l(\mathbf{L}, \mathbf{S}) = (\nabla_{\mathbf{L}} l, \nabla_{\mathbf{S}} l)$ is 2-Lipschitz:

LEMMA 3.1. *For all (\mathbf{L}, \mathbf{S}) and $(\mathbf{L}', \mathbf{S}')$, we have $\|\nabla l(\mathbf{L}, \mathbf{S}) - \nabla l(\mathbf{L}', \mathbf{S}')\|_F \leq 2\|(\mathbf{L}, \mathbf{S}) - (\mathbf{L}', \mathbf{S}')\|_F$.*

Proof. From (3.2) and (3.3), we have

$$\begin{aligned} \|\nabla l(\mathbf{L}, \mathbf{S}) - \nabla l(\mathbf{L}', \mathbf{S}')\|_F^2 &= 2\|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}] - \mathcal{P}_Q[\mathbf{L}' + \mathbf{S}' - \mathbf{M}]\|_F^2 \\ &= 2\|\mathcal{P}_Q[\mathbf{L} + \mathbf{S}] - \mathcal{P}_Q[\mathbf{L}' + \mathbf{S}']\|_F^2 \\ &\leq 2\|\mathbf{L} + \mathbf{S} - \mathbf{L}' - \mathbf{S}'\|_F^2 \\ &\leq 4\|\mathbf{L} - \mathbf{L}'\|_F^2 + 4\|\mathbf{S} - \mathbf{S}'\|_F^2 \\ &= 4\|(\mathbf{L}, \mathbf{S}) - (\mathbf{L}', \mathbf{S}')\|_F^2, \end{aligned}$$

which implies the result. \square

The feasible set in (3.1) is compact. The following lemma bounds its diameter D :

LEMMA 3.2. *The feasible set $\mathcal{D} = \{(\mathbf{L}, \mathbf{S}) \mid \|\mathbf{L}\|_* \leq \tau_L, \|\mathbf{S}\|_1 \leq \tau_S\}$ has diameter $D \leq 2\sqrt{\tau_L^2 + \tau_S^2}$.*

Proof. For any $\mathbf{Z} = (\mathbf{L}, \mathbf{S})$ and $\mathbf{Z}' = (\mathbf{L}', \mathbf{S}') \in \mathcal{D}$,

$$\begin{aligned} \|\mathbf{Z} - \mathbf{Z}'\|_F^2 &= \|\mathbf{L} - \mathbf{L}'\|_F^2 + \|\mathbf{S} - \mathbf{S}'\|_F^2 \leq (\|\mathbf{L}\|_F + \|\mathbf{L}'\|_F)^2 + (\|\mathbf{S}\|_F + \|\mathbf{S}'\|_F)^2 \\ (3.4) \quad &\leq (\|\mathbf{L}\|_* + \|\mathbf{L}'\|_*)^2 + (\|\mathbf{S}\|_1 + \|\mathbf{S}'\|_1)^2 \leq 4\tau_L^2 + 4\tau_S^2. \end{aligned}$$

\square

3.1. Frank-Wolfe for problem (3.1). Since (3.1) asks us to minimize a convex, differentiable function with Lipschitz gradient over a compact convex domain, the Frank-Wolfe method in Algorithm 1 applies. It generates a sequence of iterates $\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k)$. Using the expression for the gradient in (3.2)-(3.3), at each iteration, the step direction $\mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k)$ is generated by solving the linearized subproblem

$$(3.5) \quad \begin{pmatrix} \mathbf{V}_L^k \\ \mathbf{V}_S^k \end{pmatrix} \in \arg \min \left\langle \begin{pmatrix} \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \\ \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \end{pmatrix}, \begin{pmatrix} \mathbf{V}_L \\ \mathbf{V}_S \end{pmatrix} \right\rangle$$

s.t. $\|\mathbf{V}_L\|_* \leq \tau_L, \|\mathbf{V}_S\|_1 \leq \tau_S,$

which decouples into two independent subproblems:

$$\mathbf{V}_L^k \in \arg \min_{\|\mathbf{V}_L\|_* \leq \tau_L} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_L \rangle,$$

$$\mathbf{V}_S^k \in \arg \min_{\|\mathbf{V}_S\|_1 \leq \tau_S} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_S \rangle.$$

These subproblems can be easily solved by exploiting the linear optimization oracles introduced in Section 2.2. In particular,

$$(3.6) \quad \mathbf{V}_L^k = -\tau_L \mathbf{u}^k (\mathbf{v}^k)^\top,$$

$$(3.7) \quad \mathbf{V}_S^k = -\tau_S \cdot \delta_{i^* j^*}^k \cdot \mathbf{e}_{i^*}^k (\mathbf{e}_{j^*}^k)^\top,$$

where \mathbf{u}^k and \mathbf{v}^k are leading left- and right- singular vectors of $\mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]$ and (i^*, j^*) is the of the largest element of $\mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]$ in magnitude and $\delta_{ij}^k := \text{sgn}[(\mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}])_{ij}]$. Algorithm 3 gives the Frank-Wolfe method specialized to problem (3.1).

Algorithm 3 Frank-Wolfe method for problem (3.1)

```

1: Initialization:  $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$ ;
2: for  $k = 0, 1, 2, \dots$  do
3:    $\mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle$ ;  $\mathbf{V}_L^k = \tau_L \mathbf{D}_L^k$ ;
4:    $\mathbf{D}_S^k \in \arg \min_{\|\mathbf{D}_S\|_1 \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S \rangle$ ;  $\mathbf{V}_S^k = \tau_S \mathbf{D}_S^k$ ;
5:    $\gamma = \frac{2}{k+2}$ ;
6:    $\mathbf{L}^{k+1} = \mathbf{L}^k + \gamma(\mathbf{V}_L^k - \mathbf{L}^k)$ ;
7:    $\mathbf{S}^{k+1} = \mathbf{S}^k + \gamma(\mathbf{V}_S^k - \mathbf{S}^k)$ ;
8: end for
```

The major advantage of Algorithm 3 lies in the simplicity of the update rules (3.6)-(3.7). Both have closed form, and both can be computed in time (essentially) linear in the size of the input. Because \mathbf{V}_L^k is rank-one, the algorithm can be viewed as performing a sequence of rank one updates.

The major disadvantage of Algorithm 3 is that \mathbf{S} has only a one-sparse update at each iteration, since $\mathbf{V}_S^k = -\tau_S \mathbf{e}_{i^*}^k (\mathbf{e}_{j^*}^k)^\top$ has only one nonzero entry. This is a significant disadvantage in practice, as the optimal \mathbf{S}^* may have a relatively large number of nonzero entries. Indeed, in theory, the CPCP relaxation works even when a constant fraction of the entries in \mathbf{S}_0 are nonzero. In applications such as foreground-background separation, the number of nonzero entries in the target sparse term can be quite large. The dashed curves in Figure 1 show the effect of this on the practical convergence of the algorithm, on a simulated example of size $1,000 \times 1,000$, in which about 1% of the entries in the target sparse matrix \mathbf{S}_0 are nonzero. As shown, the progress is quite slow.

3.2. FW-P algorithm: combining Frank-Wolfe and projected gradient.

To overcome the drawback of the naive Frank-Wolfe algorithm described above, we propose incorporating an additional gradient projection step after each Frank-Wolfe update. This additional step updates the sparse term \mathbf{S} only, with the goal of accelerating convergence in these variables. At iteration k , let $(\mathbf{L}^{k+1/2}, \mathbf{S}^{k+1/2})$ be the result produced by Frank-Wolfe. To produce the next iterate, we retain the low rank term $\mathbf{L}^{k+1/2}$, but set

$$(3.8) \quad \mathbf{S}^{k+1} = \mathcal{P}_{\|\cdot\|_1 \leq \tau_S} \left[\mathbf{S}^{k+\frac{1}{2}} - \nabla_{\mathbf{S}} l(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}) \right]$$

$$(3.9) \quad = \mathcal{P}_{\|\cdot\|_1 \leq \tau_S} \left[\mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}] \right];$$

i.e. we simply take an additional projected gradient step in the sparse term \mathbf{S} . The resulting algorithm is presented as Algorithm 4 below. We call this method the FW-P algorithm, as it combines Frank-Wolfe steps and projections. In Figure 1, we compare Algorithms 3 and 4 on synthetic data. In this example, the FW-P method is clearly more efficient in recovering \mathbf{L}_0 and \mathbf{S}_0 .

The convergence of Algorithm 4 can be analyzed by recognizing it as a specific instance of the generalized Frank-Wolfe iteration in Algorithm 2. This projection step

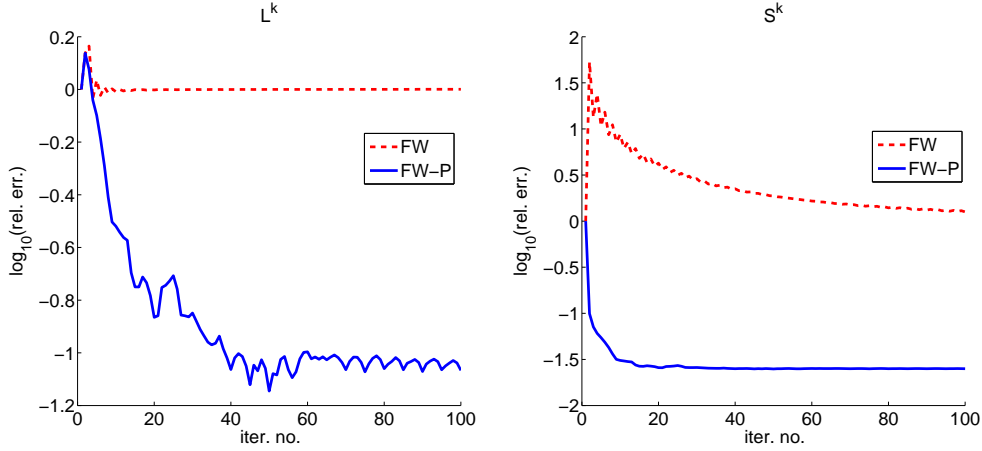


FIG. 1. *Comparisons between Algorithms 3 and 4 for problem (3.1) on synthetic data.* The data are generated in Matlab as $m = 1000$; $n = 1000$; $r = 5$; $L_0 = \text{randn}(m, r) * \text{randn}(r, n)$; $\Omega = \text{ones}(m, n)$; $S_0 = 100 * \text{randn}(m, n) * (\text{rand}(m, n) < 0.01)$; $M = L_0 + S_0 + \text{randn}(m, n)$; $\tau_L = \text{norm_nuc}(L_0)$; $\tau_S = \text{norm}(\text{vec}(S_0), 1)$; The left figure plots $\log_{10}(\|L^k - L_0\|_F / \|L_0\|_F)$ versus the iteration number k . The right figure plots $\log_{10}(\|S^k - S_0\|_F / \|S_0\|_F)$ versus k . The FW-P method is clearly more efficient than the straightforward FW method in recovering L_0 and S_0 .

Algorithm 4 FW-P method for problem (3.1)

- 1: **Initialization:** $L^0 = S^0 = 0$;
 - 2: **for** $k = 0, 1, 2, \dots$ **do**
 - 3: $D_L^k \in \arg \min_{\|D_L\|_* \leq 1} \langle \mathcal{P}_Q[L^k + S^k - M], D_L \rangle$; $V_L^k = \tau_L D_L^k$;
 - 4: $D_S^k \in \arg \min_{\|D_S\|_1 \leq 1} \langle \mathcal{P}_Q[L^k + S^k - M], D_S \rangle$; $V_S^k = \tau_S D_S^k$;
 - 5: $\gamma = \frac{2}{k+2}$;
 - 6: $L^{k+\frac{1}{2}} = L^k + \gamma(V_L^k - L^k)$;
 - 7: $S^{k+\frac{1}{2}} = S^k + \gamma(V_S^k - S^k)$;
 - 8: $S^{k+1} = \mathcal{P}_{\|\cdot\|_1 \leq \tau_S} [S^{k+\frac{1}{2}} - \mathcal{P}_Q[L^{k+\frac{1}{2}} + S^{k+\frac{1}{2}} - M]]$;
 - 9: $L^{k+1} = L^{k+\frac{1}{2}}$;
 - 10: **end for**
-

(3.9) can be regarded as a proximal step to set S^{k+1} as

$$\arg \min_{\|S\|_1 \leq \tau_S} \hat{l}^{k+\frac{1}{2}}(S) := l(L^{k+\frac{1}{2}}, S^{k+\frac{1}{2}}) + \langle \nabla_S l(L^{k+\frac{1}{2}}, S^{k+\frac{1}{2}}), S - S^{k+\frac{1}{2}} \rangle + \frac{1}{2} \|S - S^{k+\frac{1}{2}}\|_F^2.$$

It can then be easily verified that

$$(3.10) \quad \hat{l}^{k+\frac{1}{2}}(S^{k+\frac{1}{2}}) = l(L^{k+\frac{1}{2}}, S^{k+\frac{1}{2}}), \quad \text{and} \quad \hat{l}^{k+\frac{1}{2}}(S) \geq l(L^{k+\frac{1}{2}}, S) \quad \text{for any } S,$$

since $\nabla_S l(L, S)$ is 1-Lipschitz. This implies that the FW-P algorithm chooses a next iterate whose objective is no worse than that produced by the Frank-Wolfe step:

$$l(L^{k+1}, S^{k+1}) = l(L^{k+\frac{1}{2}}, S^{k+1}) \leq \hat{l}^{k+\frac{1}{2}}(S^{k+1}) \leq \hat{l}^{k+\frac{1}{2}}(S^{k+\frac{1}{2}}) = l(L^{k+\frac{1}{2}}, S^{k+\frac{1}{2}}).$$

This is precisely the property that is required to invoke Algorithm 2 and Theorems 2.1 and 2.2. Using Lemmas 4.1 and 4.2 to estimate the Lipschitz constant of ∇l and the diameter of \mathcal{D} , we obtain the following result, which shows that FW-P retains the $O(1/k)$ convergence rate of the original FW method:

THEOREM 3.3. *Let l^* be the optimal value to problem (3.1), $\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k)$ and $\mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k)$ be the sequence produced by Algorithm 4. Then we have*

$$(3.11) \quad l(\mathbf{L}^k, \mathbf{S}^k) - l^* \leq \frac{16(\tau_L^2 + \tau_S^2)}{k + 2}.$$

Moreover, for any $K \geq 1$, there exists $1 \leq \tilde{k} \leq K$ such that the surrogate duality gap (defined in (2.13)) satisfies

$$(3.12) \quad d(\mathbf{x}^{\tilde{k}}) = \langle \mathbf{x}^{\tilde{k}} - \mathbf{v}^{\tilde{k}}, \nabla l(\mathbf{x}^{\tilde{k}}) \rangle \leq \frac{48(\tau_L^2 + \tau_S^2)}{K + 2}.$$

Proof. Substituting $L = 2$ (Lemma 3.1) and $D \leq 2\sqrt{\tau_L^2 + \tau_S^2}$ (Lemma 3.2) into Theorems 2.1 and 2.2, we can easily obtain the above result. \square

4. FW-T Method for Penalized Problem. In this section, we develop a scalable algorithm for the penalized version of the CPCP problem,

$$(4.1) \quad \min_{\mathbf{L}, \mathbf{S}} \quad f(\mathbf{L}, \mathbf{S}) \doteq \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1.$$

In Section 4.1, we reformulate problem (4.1) into the form of (2.1) so that the Frank-Wolfe method can be applied. In Section 4.2, we apply the Frank-Wolfe method directly to the reformulated problem, achieving linear per-iteration cost and $O(1/k)$ convergence in function value. However, because it updates the sparse term one element at a time, it converges very slowly on typical numerical examples. In Section 4, we introduce our FW-T method, which resolves this issue. Our FW-T method essentially exploits the Frank-Wolfe step to handle the nuclear norm and a proximal gradient step to handle the ℓ_1 -norm, while keeping iteration cost low and retaining convergence guarantees.

4.1. Reformulation as smooth, constrained optimization. Note that problem (4.1) has a non-differentiable objective function and an unbounded feasible set. To apply the Frank-Wolfe method, we exploit a two-step reformulation to transform (4.1) into the form of (2.1). First, we borrow ideas from [24] and work with the epigraph reformulation of (4.1),

$$(4.2) \quad \begin{aligned} \min \quad & g(\mathbf{L}, \mathbf{S}, t_L, t_S) \doteq \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L t_L + \lambda_S t_S \\ \text{s.t.} \quad & \|\mathbf{L}\|_* \leq t_L, \quad \|\mathbf{S}\|_1 \leq t_S, \end{aligned}$$

obtained by introducing auxiliary variables t_L and t_S . Now the objective function $g(\mathbf{L}, \mathbf{S}, t_L, t_S)$ is differentiable, with

$$(4.3) \quad \nabla_{\mathbf{L}} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \nabla_{\mathbf{S}} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}],$$

$$(4.4) \quad \nabla_{t_L} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \lambda_L, \quad \nabla_{t_S} g(\mathbf{L}, \mathbf{S}, t_L, t_S) = \lambda_S.$$

A calculation, which we summarize in the following lemma, shows that the gradient $\nabla g(\mathbf{L}, \mathbf{S}, t_L, t_S) = (\nabla_{\mathbf{L}} g, \nabla_{\mathbf{S}} g, \nabla_{t_L} g, \nabla_{t_S} g)$ is 2-Lipschitz:

LEMMA 4.1. *For all $(\mathbf{L}, \mathbf{S}, t_L, t_S)$ and $(\mathbf{L}', \mathbf{S}', t'_L, t'_S)$ feasible to (4.2),*
(4.5)
$$\|\nabla g(\mathbf{L}, \mathbf{S}, t_L, t_S) - \nabla g(\mathbf{L}', \mathbf{S}', t'_L, t'_S)\|_F \leq 2 \|(\mathbf{L}, \mathbf{S}, t_L, t_S) - (\mathbf{L}', \mathbf{S}', t'_L, t'_S)\|_F.$$

Proof. Based on (4.3) and (4.4), it follows directly that

$$\begin{aligned} \|\nabla g(\mathbf{L}, \mathbf{S}, t_L, t_S) - \nabla g(\mathbf{L}', \mathbf{S}', t'_L, t'_S)\|_F^2 &\leq 4 \|\mathbf{L} - \mathbf{L}'\|_F^2 + 4 \|\mathbf{S} - \mathbf{S}'\|_F^2 \\ &\leq 4 \|(\mathbf{L}, \mathbf{S}, t_L, t_S) - (\mathbf{L}', \mathbf{S}', t'_L, t'_S)\|_F^2, \end{aligned}$$

which implies the result. \square

However, the Frank-Wolfe method still cannot deal with (4.2), since its feasible region is unbounded. If we could somehow obtain upper bounds on the optimal values of t_L and t_S : $U_L \geq t_L^*$ and $U_S \geq t_S^*$, then we could solve the equivalent problem

$$(4.6) \quad \begin{aligned} \min \quad & \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L t_L + \lambda_S t_S \\ \text{s.t.} \quad & \|\mathbf{L}\|_* \leq t_L \leq U_L, \quad \|\mathbf{S}\|_1 \leq t_S \leq U_S, \end{aligned}$$

which now has a compact and convex feasible set. One simple way to obtain such U_L, U_S is as follows. One trivial feasible solution to problem (4.2) is $\mathbf{L} = \mathbf{0}, \mathbf{S} = \mathbf{0}, t_L = 0, t_S = 0$. This solution has objective value $\frac{1}{2} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2$. Hence, the optimal objective value is no larger than this. This implies that for any optimal t_L^*, t_S^* ,

$$(4.7) \quad t_L^* \leq \frac{1}{2\lambda_L} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2, \quad t_S^* \leq \frac{1}{2\lambda_S} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2.$$

Hence, we can always choose

$$(4.8) \quad U_L = \frac{1}{2\lambda_L} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2, \quad U_S = \frac{1}{2\lambda_S} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2$$

to produce a valid, bounded feasible region. The following lemma bounds its diameter D :

LEMMA 4.2. *The feasible set $\mathcal{D} = \{(\mathbf{L}, \mathbf{S}, t_L, t_S) \mid \|\mathbf{L}\|_* \leq t_L \leq U_L, \|\mathbf{S}\|_1 \leq t_S \leq U_S\}$ has diameter $D \leq \sqrt{5} \cdot \sqrt{U_L^2 + U_S^2}$.*

Proof. Since for any $\mathbf{Z} = (\mathbf{L}, \mathbf{S}, t_L, t_S), \mathbf{Z}' = (\mathbf{L}', \mathbf{S}', t'_L, t'_S) \in \mathcal{D}$, we have

$$\begin{aligned} \|\mathbf{Z} - \mathbf{Z}'\|_F^2 &= \|\mathbf{L} - \mathbf{L}'\|_F^2 + \|\mathbf{S} - \mathbf{S}'\|_F^2 + (t_L - t'_L)^2 + (t_S - t'_S)^2 \\ &\leq (\|\mathbf{L}\|_F + \|\mathbf{L}'\|_F)^2 + (\|\mathbf{S}\|_F + \|\mathbf{S}'\|_F)^2 + (t_L - t'_L)^2 + (t_S - t'_S)^2 \\ &\leq (\|\mathbf{L}\|_* + \|\mathbf{L}'\|_*)^2 + (\|\mathbf{S}\|_1 + \|\mathbf{S}'\|_1)^2 + (t_L - t'_L)^2 + (t_S - t'_S)^2 \\ &\leq (U_L + U_L)^2 + (U_S + U_S)^2 + U_L^2 + U_S^2 \\ &= 5(U_L^2 + U_S^2), \end{aligned}$$

which implies the result. \square

With these modifications, we can apply Frank-Wolfe directly to obtain a solution $(\widehat{\mathbf{L}}, \widehat{\mathbf{S}}, \widehat{t}_L, \widehat{t}_S)$ to (4.6), and hence to produce a solution $(\widehat{\mathbf{L}}, \widehat{\mathbf{S}})$ to the original problem (4.1). In subsection 4.2, we describe how to do this. Unfortunately, this straightforward solution has two main disadvantages. First, as in the norm constrained case, it produces only one-sparse updates to \mathbf{S} , which results in slow convergence. Second, the exact primal convergence rate in Theorem 2.1 depends on the diameter of the feasible set, which in turn depends on the accuracy of our (crude) upper bounds U_L and U_S . In subsection 4.3, we show how to remedy both issues, yielding a Frank-Wolfe-Thresholding method that performs significantly better in practice.

4.2. Frank-Wolfe for problem (4.6). Applying the Frank-Wolfe method in Algorithm 1 generates a sequence of iterates $\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$. Using the expressions for the gradient in (4.3) and (4.4), at each iteration, $\mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k, V_{t_L}^k, V_{t_S}^k)$ is generated by solving the linearized subproblem

$$(4.9) \quad \mathbf{v}^k \in \arg \min_{\mathbf{v} \in \mathcal{D}} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_L + \mathbf{V}_S \rangle + \lambda_L V_{t_L} + \lambda_S V_{t_S},$$

which can be decoupled into two independent subproblems,

$$(4.10) \quad (\mathbf{V}_L^k, V_{t_L}^k) \in \arg \min_{\|\mathbf{V}_L\|_* \leq V_{t_L} \leq U_L} g_L(\mathbf{V}_L, V_{t_L}) \doteq \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_L \rangle + \lambda_L V_{t_L}$$

$$(4.11) \quad (\mathbf{V}_S^k, V_{t_S}^k) \in \arg \min_{\|\mathbf{V}_S\|_1 \leq V_{t_S} \leq U_S} g_S(\mathbf{V}_S, V_{t_S}) \doteq \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{V}_S \rangle + \lambda_S V_{t_S}.$$

Let us consider problem (4.10) first. Set

$$(4.12) \quad \mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \hat{g}_L(\mathbf{D}_L) \doteq \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle + \lambda_L.$$

Because $g_L(\mathbf{V}_L, V_{t_L})$ is a homogeneous function, i.e., $g_L(\alpha \mathbf{V}_L, \alpha V_{t_L}) = \alpha g_L(\mathbf{V}_L, V_{t_L})$, for any $\alpha \in \mathbb{R}$, its optimal value $g(\mathbf{V}_L^k, V_{t_L}^k) = V_{t_L}^k \hat{g}_L(\mathbf{D}_L^k)$. Hence $V_{t_L}^k = U_L$ if $\hat{g}_L(\mathbf{D}_L^k) < 0$, and $V_{t_L}^k = 0$ if $\hat{g}_L(\mathbf{D}_L^k) > 0$. From this observation, it can be easily verified (see also [24, Lemma 1] for a more general result) that

$$(4.13) \quad (\mathbf{V}_L^k, V_{t_L}^k) \in \begin{cases} \{(\mathbf{0}, 0)\} & \text{if } \hat{g}_L(\mathbf{D}_L^k) > 0 \\ \text{conv}\{(\mathbf{0}, 0), U_L(\mathbf{D}_L^k, 1)\} & \text{if } \hat{g}_L(\mathbf{D}_L^k) = 0 \\ \{U_L(\mathbf{D}_L^k, 1)\} & \text{if } \hat{g}_L(\mathbf{D}_L^k) < 0. \end{cases}$$

In a similar manner, we can update $(\mathbf{V}_S^k, V_{t_S}^k)$. This leads fairly directly to the implementation of the Frank-Wolfe method for problem (4.6), described in Algorithm 5. As a direct corollary of Theorem 2.1, using parameters calculated in Lemmas 4.1 and 4.2, we have

COROLLARY 4.3. *Let $\mathbf{x}^* = (\mathbf{L}^*, \mathbf{S}^*, t_L^*, t_S^*)$ be an optimal solution to (4.6). For $\{\mathbf{x}^k\}$ generated by Algorithm 5, we have for $k = 0, 1, 2, \dots$,*

$$(4.14) \quad g(\mathbf{x}^k) - g(\mathbf{x}^*) \leq \frac{20(U_L^2 + U_S^2)}{k+2}.$$

Proof. Applying Theorem 2.1 with parameters calculated in Lemmas 4.1 and 4.2, we directly have

$$(4.15) \quad g(\mathbf{x}^k) - g(\mathbf{x}^*) \leq \frac{2 \cdot 2 \cdot \left(\sqrt{5(U_L^2 + U_S^2)} \right)^2}{k+2} = \frac{20(U_L^2 + U_S^2)}{k+2}.$$

A more careful calculation below slightly improves the constant in (4.15).

$$(4.16) \quad \begin{aligned} g(\mathbf{x}^{k+1}) &= g(\mathbf{x}^k + \gamma(\mathbf{v}^k - \mathbf{x}^k)) \\ &\leq g(\mathbf{x}^k) + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + \gamma^2 \|\mathbf{V}_L^k - \mathbf{L}^k\|_F^2 + \gamma^2 \|\mathbf{V}_S^k - \mathbf{S}^k\|_F^2 \\ &\leq g(\mathbf{x}^k) + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + 4\gamma^2(U_L^2 + U_S^2), \end{aligned}$$

Algorithm 5 Frank-Wolfe method for problem (4.6)

```

1: Initialization:  $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$ ;  $t_L^0 = t_S^0 = 0$ ;
2: for  $k = 0, 1, 2, \dots$  do
3:    $\mathbf{D}_L^k \in \arg \min_{\|\mathbf{D}_L\|_* \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L \rangle$ ;
4:    $\mathbf{D}_S^k \in \arg \min_{\|\mathbf{D}_S\|_1 \leq 1} \langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S \rangle$ ;
5:   if  $\lambda_L \geq -\langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_L^k \rangle$  then
6:      $\mathbf{V}_L^k = \mathbf{0}$ ;  $V_{t_L}^k = 0$ 
7:   else
8:      $\mathbf{V}_L^k = U_L \mathbf{D}_L^k$ ,  $V_{t_L}^k = U_L$ ;
9:   end if
10:  if  $\lambda_S \geq -\langle \mathcal{P}_Q[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}], \mathbf{D}_S^k \rangle$  then
11:     $\mathbf{V}_S^k = \mathbf{0}$ ;  $V_{t_S}^k = 0$ ;
12:  else
13:     $\mathbf{V}_S^k = U_S \mathbf{D}_S^k$ ,  $V_{t_S}^k = U_S$ ;
14:  end if
15:   $\gamma = \frac{2}{k+2}$ ;
16:   $\mathbf{L}^{k+1} = (1 - \gamma)\mathbf{L}^k + \gamma\mathbf{V}_L^k$ ,  $t_L^{k+1} = (1 - \gamma)t_L^k + \gamma V_{t_L}^k$ ;
17:   $\mathbf{S}^{k+1} = (1 - \gamma)\mathbf{S}^k + \gamma\mathbf{V}_S^k$ ,  $t_S^{k+1} = (1 - \gamma)t_S^k + \gamma V_{t_S}^k$ ;
18: end for

```

where the second line holds by noting that g is only linear in t_L and t_S ; the last line holds as

$$\begin{aligned} \|\mathbf{V}_L^k - \mathbf{L}^k\|_F^2 &\leq (\|\mathbf{V}_L^k\|_F + \|\mathbf{L}^k\|_F)^2 \leq (U_L + U_L)^2 = 4U_L^2, \quad \text{and} \\ \|\mathbf{V}_S^k - \mathbf{S}^k\|_F^2 &\leq (\|\mathbf{V}_S^k\|_F + \|\mathbf{S}^k\|_F)^2 \leq (U_S + U_S)^2 = 4U_S^2. \end{aligned}$$

Following the arguments in the proof of Theorem 1 with (2.10) replaced by (4.16), we can easily obtain that

$$g(\mathbf{x}^k) - g(\mathbf{x}^*) \leq \frac{16(U_L^2 + U_S^2)}{k + 2}.$$

□

In addition to the above convergence result, another major advantage of Algorithm 5 is the simplicity of the update rules (lines 3-4 in Algorithm 5). Both have closed-form solutions that can be computed in time (essentially) linearly dependent on the size of the input.

However, two clear limitations substantially hinder Algorithm 5's efficiency. First, as in the norm constrained case, \mathbf{V}_S^k has only one nonzero entry, so \mathbf{S} has a one-sparse update in each iteration. Second, the exact rate of convergence relies on our (crude) guesses of U_L and U_S (Corollary 4.3). In the next subsection, we present remedies to resolve both issues.

4.3. FW-T algorithm: combining Frank-Wolfe and proximal methods.

To alleviate the difficulties faced by Algorithm 5, we propose a new algorithm called Frank-Wolfe-Thresholding (FW-T) (Algorithm 6), that combines a modified FW step with a proximal gradient step. Below we highlight the key features of FW-T.

Proximal gradient step for \mathbf{S} . To update \mathbf{S} in a more efficient way, we incorporate an additional proximal gradient step for \mathbf{S} . At iteration k , let $(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}})$

be the result produced by Frank-Wolfe step. To produce the next iterate, we retain the low-rank term $\mathbf{L}^{k+\frac{1}{2}}$, but execute a proximal gradient step for the function $f(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S})$ at the point $\mathbf{S}^{k+\frac{1}{2}}$, i.e.

$$(4.17) \quad \begin{aligned} \mathbf{S}^{k+1} &\in \arg \min_{\mathbf{S}} \left\langle \nabla_{\mathbf{S}} f(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}), \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\rangle + \frac{1}{2} \left\| \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\|_F^2 + \lambda_S \|\mathbf{S}\|_1 \\ &= \arg \min_{\mathbf{S}} \left\langle \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}], \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\rangle + \frac{1}{2} \left\| \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\|_F^2 + \lambda_S \|\mathbf{S}\|_1 \end{aligned}$$

which can be easily computed using the soft-thresholding operator:

$$(4.18) \quad \mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S} \left[\mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}] \right].$$

Exact line search. For the Frank-Wolfe step, instead of choosing the fixed step length $\frac{2}{k+2}$, we implement an exact line search by solving a two-dimensional quadratic problem (4.20), as in [24]. This modification turns out to be crucial to achieve a primal convergence result that only weakly depends on the tightness of our guesses U_L and U_S .

Adaptive updates of U_L and U_S . We initialize U_L and U_S using the crude bound (4.8). Then, at the end of the k -iteration, we respectively update

$$(4.19) \quad U_L^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_L, \quad U_S^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_S.$$

This scheme maintains the property that $U_L^{k+1} \geq t_L^*$ and $U_S^{k+1} \geq t_S^*$. Moreover, we prove (Lemma 4.4) that g is non-increasing through our algorithm, and so this scheme produces a sequence of tighter upper bounds for U_L^* and U_S^* . Although this dynamic scheme does not improve the theoretical convergence result, some acceleration is empirically exhibited.

Convergence analysis. Since both the FW step and the proximal gradient step do not increase the objective value, we can easily recognize FW-T method as a descent algorithm:

LEMMA 4.4. *Let $\{(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)\}$ be the sequence of iterates produced by the FW-T algorithm. For each $k = 0, 1, 2, \dots$,*

$$(4.21) \quad g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1}) \leq g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}}) \leq g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k).$$

Proof. Since $(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$ is always feasible to the quadratic program (4.20),

$$(4.22) \quad g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}}) \leq g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k).$$

Based on (4.17), the threshold step (line 6 in Algorithm 3) can be written as

$$\begin{aligned} \mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} \quad & \hat{g}^{k+\frac{1}{2}}(\mathbf{S}) \doteq \frac{1}{2} \left\| \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}] \right\|_F^2 + \lambda_L t_L^{k+\frac{1}{2}} + \lambda_S \|\mathbf{S}\|_1 \\ & + \left\langle \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}], \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\rangle + \frac{1}{2} \left\| \mathbf{S} - \mathbf{S}^{k+\frac{1}{2}} \right\|_F^2. \end{aligned}$$

The following properties of $\hat{g}^{k+\frac{1}{2}}(\cdot)$ can be easily verified

$$\hat{g}^{k+\frac{1}{2}}(\mathbf{S}^{k+\frac{1}{2}}) = g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, \|\mathbf{S}^{k+\frac{1}{2}}\|_1) \leq g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}});$$

Algorithm 6 FW-T method for problem (4.1)

1: **Input:** data matrix $\mathbf{M} \in \mathbb{R}^{m \times n}$; weights $\lambda_L, \lambda_S > 0$; max iteration number T ;
2: **Initialization:** $\mathbf{L}^0 = \mathbf{S}^0 = \mathbf{0}$; $t_L^0 = t_S^0 = 0$; $U_L^0 = g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0)/\lambda_L$; $U_S^0 = g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0)/\lambda_S$;
3: **for** $k = 0, 1, 2, \dots, T$ **do**
4: *same as lines 3-14 in Algorithm 5*;
5: $\left(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}} \right)$ is computed as an optimizer to

$$(4.20) \quad \min \quad \frac{1}{2} \|\mathcal{P}_Q[\mathbf{L} + \mathbf{S} - \mathbf{M}]\|_F^2 + \lambda_L t_L + \lambda_S t_S$$

s.t. $\begin{pmatrix} \mathbf{L} \\ t_L \end{pmatrix} \in \text{conv} \left\{ \begin{pmatrix} \mathbf{L}^k \\ t_L^k \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{t_L}^k \\ V_{t_L}^k \end{pmatrix} \right\}$
 $\begin{pmatrix} \mathbf{S} \\ t_S \end{pmatrix} \in \text{conv} \left\{ \begin{pmatrix} \mathbf{S}^k \\ t_S^k \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{t_S}^k \\ V_{t_S}^k \end{pmatrix} \right\};$

6: $\mathbf{S}^{k+1} = \mathcal{T}[\mathbf{S}^{k+\frac{1}{2}} - \mathcal{P}_Q[\mathbf{L}^{k+\frac{1}{2}} + \mathbf{S}^{k+\frac{1}{2}} - \mathbf{M}], \lambda_S]$;
7: $\mathbf{L}^{k+1} = \mathbf{L}^{k+\frac{1}{2}}, t_L^{k+1} = t_L^{k+\frac{1}{2}}; t_S^{k+1} = \|\mathbf{S}^{k+1}\|_1$;
8: $U_L^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_L$;
9: $U_S^{k+1} = g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1})/\lambda_S$;
10: **end for**

$$\hat{g}^{k+\frac{1}{2}}(\mathbf{S}) \geq g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}, t_L^{k+\frac{1}{2}}, \|\mathbf{S}\|_1), \quad \text{for any } \mathbf{S}.$$

Therefore, we have

$$(4.23) \quad \begin{aligned} g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1}) &= g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+1}, t_L^{k+\frac{1}{2}}, t_S^{k+1}) \leq \hat{g}^{k+\frac{1}{2}}(\mathbf{S}^{k+1}) \\ &\leq \hat{g}^{k+\frac{1}{2}}(\mathbf{S}^{k+\frac{1}{2}}) \leq g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}}) \end{aligned}$$

Combining (4.22) and (4.23), we obtain

$$g(\mathbf{L}^{k+1}, \mathbf{S}^{k+1}, t_L^{k+1}, t_S^{k+1}) \leq g(\mathbf{L}^{k+\frac{1}{2}}, \mathbf{S}^{k+\frac{1}{2}}, t_L^{k+\frac{1}{2}}, t_S^{k+\frac{1}{2}}) \leq g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k).$$

□

Moreover, we can establish primal convergence (almost) independent of U_L^0 and U_S^0 :

THEOREM 4.5. *Let r_L^* and r_S^* be the smallest radii such that*

$$(4.24) \quad \left\{ (\mathbf{L}, \mathbf{S}) \mid f(\mathbf{L}, \mathbf{S}) \leq g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0) = \frac{1}{2} \|\mathcal{P}_Q[\mathbf{M}]\|_F^2 \right\} \subseteq \overline{B(r_L^*)} \times \overline{B(r_S^*)},$$

where $\overline{B(r)} \doteq \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \|\mathbf{X}\|_F \leq r\}$ for any $r \geq 0$.[‡] Then for the sequence $\{(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)\}$ generated by Algorithm 6, we have

$$(4.25) \quad \begin{aligned} g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k) - g(\mathbf{L}^*, \mathbf{S}^*, t_L^*, t_S^*) \\ \leq \frac{\min\{4(t_L^* + r_L^*)^2 + 4(t_S^* + r_S^*)^2, 16(U_L^0)^2 + 16(U_S^0)^2\}}{k+2}. \end{aligned}$$

[‡]Since the objective function in problem (4.1) is coercive, i.e. $\lim_{k \rightarrow +\infty} f(\mathbf{L}^k, \mathbf{S}^k) = +\infty$ for any sequence $(\mathbf{L}^k, \mathbf{S}^k)$ such that $\lim_{k \rightarrow +\infty} \|(\mathbf{L}^k, \mathbf{S}^k)\|_F = +\infty$, clearly $r_L^* \geq 0$ and $r_S^* \geq 0$ exist.

Proof. For notational convenience, we denote

$$\mathbf{x}^k = (\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k), \mathbf{x}^* = (\mathbf{L}^*, \mathbf{S}^*, t_L^*, t_S^*) \text{ and } \mathbf{v}^k = (\mathbf{V}_L^k, \mathbf{V}_S^k, \mathbf{V}_{t_L}^k, \mathbf{V}_{t_S}^k).$$

For any point $\mathbf{x} = (\mathbf{L}, \mathbf{S}, t_L, t_S) \in \mathbb{R}^{m \times n} \times \mathbb{R}^{m \times n} \times \mathbb{R} \times \mathbb{R}$, we adopt the notation that $\mathbf{L}[\mathbf{x}] = \mathbf{L}$, $\mathbf{S}[\mathbf{x}] = \mathbf{S}$, $t_L[\mathbf{x}] = t_L$ and $t_S[\mathbf{x}] = t_S$.

Since $g(\mathbf{x}^k) - g(\mathbf{x}^*) \leq \frac{16(U_L^0)^2 + 16(U_S^0)^2}{k+2}$ can be easily established following the proof of Corollary 4.3, below we will focus on the other part that $g(\mathbf{x}^k) - g(\mathbf{x}^*) \leq \frac{4(t_L^* + r_L^*)^2 + 4(t_S^* + r_S^*)^2}{k+2}$.

Let us first make two simple observations.

Since $f(\mathbf{L}^*, \mathbf{S}^*) \leq g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$, we have

$$(4.26) \quad U_L^k = g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)/\lambda_L \geq t_L^* \quad \text{and} \quad U_S^k = g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)/\lambda_S \geq t_S^*.$$

Therefore, our U_L^k and U_S^k always bound t_L^* and t_S^* from above.

From Lemma 4.4, $g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$ is non-increasing,

$$f(\mathbf{L}^k, \mathbf{S}^k) \leq g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k) \leq g(\mathbf{L}^0, \mathbf{S}^0, t_L^0, t_S^0),$$

which implies that $(\mathbf{L}^k, \mathbf{S}^k) \subseteq \overline{B(r_L^*)} \times \overline{B(r_S^*)}$, i.e. $\|\mathbf{L}^k\|_F \leq r_L^*$ and $\|\mathbf{S}^k\|_F \leq r_S^*$.

Let us now consider the k -th iteration. Similar to the proof in [24], we introduce the auxiliary point $\mathbf{v}_+^k = (\frac{t_L^*}{U_L^k} \mathbf{V}_L^k, \frac{t_S^*}{U_S^k} \mathbf{V}_S^k, \frac{t_L^*}{U_L^k} \mathbf{V}_{t_L}^k, \frac{t_S^*}{U_S^k} \mathbf{V}_{t_S}^k)$. Then based on our argument for (4.13), it can be easily verified that

$$(4.27) \quad (\mathbf{L}[\mathbf{v}_+^k], t_L[\mathbf{v}_+^k]) \in \arg \min_{\|\mathbf{V}_L\|_* \leq V_{t_L} \leq t_L^*} g_L(\mathbf{V}_L, V_{t_L})$$

$$(4.28) \quad (\mathbf{S}[\mathbf{v}_+^k], t_S[\mathbf{v}_+^k]) \in \arg \min_{\|\mathbf{V}_S\|_1 \leq V_{t_S} \leq t_S^*} g_S(\mathbf{V}_S, V_{t_S}).$$

Recall $\gamma = \frac{2}{k+2}$. We have

$$\begin{aligned} & g(\mathbf{x}^{k+\frac{1}{2}}) \\ & \leq g(\mathbf{x}^k + \gamma(\mathbf{v}_+^k - \mathbf{x}^k)) \\ & \leq g(\mathbf{x}^k) + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}_+^k - \mathbf{x}^k \rangle + \gamma^2 \left(\|\mathbf{L}[\mathbf{v}_+^k] - \mathbf{L}[\mathbf{x}^k]\|_F^2 + \|\mathbf{S}[\mathbf{v}_+^k] - \mathbf{S}[\mathbf{x}^k]\|_F^2 \right) \\ & \leq g(\mathbf{x}^k) + \gamma (g_L(\mathbf{L}[\mathbf{v}_+^k] - \mathbf{x}^k], t_L[\mathbf{v}_+^k] - \mathbf{x}^k]) + g_S(\mathbf{S}[\mathbf{v}_+^k] - \mathbf{x}^k], t_S[\mathbf{v}_+^k] - \mathbf{x}^k]) \\ & \quad + \gamma^2 ((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2) \\ & \leq g(\mathbf{x}^k) + \gamma (g_L(\mathbf{L}[\mathbf{x}^* - \mathbf{x}^k], t_L[\mathbf{x}^* - \mathbf{x}^k]) + g_S(\mathbf{S}[\mathbf{x}^* - \mathbf{x}^k], t_S[\mathbf{x}^* - \mathbf{x}^k])) \\ & \quad + \gamma^2 ((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2) \\ & = g(\mathbf{x}^k) + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{x}^* - \mathbf{x}^k \rangle + \gamma^2 ((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2) \\ & \leq g(\mathbf{x}^k) + \gamma (g(\mathbf{x}^*) - g(\mathbf{x}^k)) + \gamma^2 ((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2), \end{aligned}$$

where the first inequality holds since $\mathbf{x}^k + \gamma(\mathbf{v}_+^k - \mathbf{x}^k)$ is feasible to the quadratic program (4.20) while $\mathbf{x}^{k+\frac{1}{2}}$ minimizes it; the third inequality is due to the facts that

$$\begin{aligned} \|\mathbf{L}[\mathbf{v}_+^k] - \mathbf{L}[\mathbf{x}^k]\|_F & \leq \|\mathbf{L}[\mathbf{v}_+^k]\|_F + \|\mathbf{L}[\mathbf{x}^k]\|_F \leq \|\mathbf{L}[\mathbf{v}_+^k]\|_* + \|\mathbf{L}[\mathbf{x}^k]\|_F \leq t_L^* + r_L^* \\ \|\mathbf{S}[\mathbf{v}_+^k] - \mathbf{S}[\mathbf{x}^k]\|_F & \leq \|\mathbf{S}[\mathbf{v}_+^k]\|_F + \|\mathbf{S}[\mathbf{x}^k]\|_F \leq \|\mathbf{S}[\mathbf{v}_+^k]\|_1 + \|\mathbf{S}[\mathbf{x}^k]\|_F \leq t_S^* + r_S^*; \end{aligned}$$

the fourth inequality holds as $(\mathbf{L}[\mathbf{x}^*], t_L[\mathbf{x}^*])$ and $(\mathbf{S}[\mathbf{x}^*], t_S[\mathbf{x}^*])$ are respectively feasible to (4.27) and (4.28) while $(\mathbf{L}[\mathbf{v}_+^k], t_L[\mathbf{v}_+^k])$ and $(\mathbf{S}[\mathbf{v}_+^k], t_S[\mathbf{v}_+^k])$ respectively minimize (4.27) and (4.28);

Therefore, we obtain

$$g(\mathbf{x}^{k+\frac{1}{2}}) - g(\mathbf{x}^*) \leq (1 - \gamma) (g(\mathbf{x}^k) - g(\mathbf{x}^*)) + \gamma^2 ((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2).$$

Moreover, by Lemma 4.4, we have

$$g(\mathbf{x}^{k+1}) \leq g(\mathbf{x}^{k+\frac{1}{2}}).$$

Thus, we obtain the recurrence

$$g(\mathbf{x}^{k+1}) - g(\mathbf{x}^*) \leq (1 - \gamma) (g(\mathbf{x}^k) - g(\mathbf{x}^*)) + \gamma^2 ((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2).$$

Applying mathematical induction, one can easily obtain that

$$g(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k) - g(\mathbf{L}^*, \mathbf{S}^*, t_L^*, t_S^*) \leq \frac{4((t_L^* + r_L^*)^2 + (t_S^* + r_S^*)^2)}{k + 2}.$$

□

Since U_L^0 and U_S^0 are quite crude upper bounds for t_L^* and t_S^* , $16(U_L^0)^2 + 16(U_S^0)^2$ could be much larger than $4(t_L^* + r_L^*)^2 + 4(t_S^* + r_S^*)^2$. Therefore, this primal convergence results depend on U_L^0 and U_S^0 in a very weak manner.

However, the convergence result of the surrogate duality gap $d(\mathbf{x}^k)$ still hinges upon the upper bounds:

THEOREM 4.6. *Let \mathbf{x}^k denote $(\mathbf{L}^k, \mathbf{S}^k, t_L^k, t_S^k)$ generated by Algorithm 6. Then for any $K \geq 1$, there exists $1 \leq \tilde{k} \leq K$ such that*

$$(4.29) \quad g(\mathbf{x}^{\tilde{k}}) - g(\mathbf{x}^*) \leq d(\mathbf{x}^{\tilde{k}}) \leq \frac{48((U_L^0)^2 + (U_S^0)^2)}{K + 2}.$$

Proof. Define $\Delta^k = g(\mathbf{x}^k) - g(\mathbf{x}^*)$. Following (4.16), we have

$$(4.30) \quad \Delta^{k+1} \leq \Delta^k + \gamma \langle \nabla g(\mathbf{x}^k), \mathbf{v}^k - \mathbf{x}^k \rangle + 4\gamma^2 ((U_L^0)^2 + (U_S^0)^2).$$

Then following the arguments in the proof of Theorem 2 with (2.17) replaced by (4.30), we can easily obtain the result. □

Stopping criterion. Compared to the convergence of $g(\mathbf{x}^k)$ (Theorem 4.5), the convergence result for $d(\mathbf{x}^k)$ can be much slower (Theorem 4.6). Therefore, here the surrogate duality gap $d(\cdot)$ is not that suitable to serve as a stopping criterion. Consequently, in our implementation, we terminate Algorithm 6 if

$$(4.31) \quad |g(\mathbf{x}^{k+1}) - g(\mathbf{x}^k)| / g(\mathbf{x}^k) \leq \varepsilon,$$

for five consecutive iterations.

5. Numerical Experiments. In this section, we report numerical results obtained by applying our FW-T method (Algorithm 6) to problem (1.5) with real data arising from applications considered in [3]: *foreground/background separation in surveillance videos*, and *shadow and specular removal from face images*.

Given observations $\{\mathbf{M}_0(i, j) \mid (i, j) \in \Omega\}$, where $\Omega \subseteq \{1, \dots, m\} \times \{1, \dots, n\}$ is the index set of the observable entries in $\mathbf{M}_0 \in \mathbb{R}^{m \times n}$, we assigned weights

$$\lambda_L = \delta \rho \|\mathcal{P}_\Omega[\mathbf{M}_0]\|_F \quad \text{and} \quad \lambda_S = \delta \sqrt{\rho} \|\mathcal{P}_\Omega[\mathbf{M}_0]\|_F / \sqrt{\max(m, n)}$$

to problem (1.5),[§] where $\rho = |\Omega|/mn$ and δ is chosen as 0.001 for the surveillance problem and 0.01 for the face problem.

We compared our FW-T method with the popular first-order methods *iterative soft-thresholding algorithm* (ISTA) and *fast iterative soft-thresholding algorithm* (FISTA) [20], both of whose implementations used partial *singular value decomposition* (SVD). In subsection 5.1, we provided detailed descriptions and implementations of ISTA and FISTA.

We set $\varepsilon = 10^{-3}$ in FW-T's stopping criterion (4.31),[¶] and terminated ISTA and FISTA whenever they reached the objective value returned by the FW-T method.^{||} All the experiments were conducted on a computer with Intel Xeon E5-2630 Processor (12 cores at 2.4 GHz), and 64GB RAM running MATLAB R2012b (64 bits).

5.1. ISTA & FISTA for problem (1.5). *Iterative soft-thresholding algorithm* (ISTA), is an efficient way to tackle unconstrained nonsmooth optimization problem especially at large scale. ISTA follows the general idea by iteratively minimizing an upper bound of the original objective. In particular, when applied to problem (1.5) of our interest, ISTA updates (\mathbf{L}, \mathbf{S}) for the k -th iteration by solving

$$(5.1) \quad (\mathbf{L}^{k+1}, \mathbf{S}^{k+1}) = \arg \min_{\mathbf{L}, \mathbf{S}} \left\langle \begin{pmatrix} \nabla_{\mathbf{L}} l(\mathbf{L}^k, \mathbf{S}^k) \\ \nabla_{\mathbf{S}} l(\mathbf{L}^k, \mathbf{S}^k) \end{pmatrix}, \begin{pmatrix} \mathbf{L} - \mathbf{L}^k \\ \mathbf{S} - \mathbf{S}^k \end{pmatrix} \right\rangle + \frac{L_f}{2} \left\| \begin{pmatrix} \mathbf{L} \\ \mathbf{S} \end{pmatrix} - \begin{pmatrix} \mathbf{L}^k \\ \mathbf{S}^k \end{pmatrix} \right\|_F^2 + \lambda_L \|\mathbf{L}\|_* + \lambda_S \|\mathbf{S}\|_1.$$

Here $L_f = 2$ denotes the Lipschitz constant of $\nabla l(\mathbf{L}, \mathbf{S})$ with respect to (\mathbf{L}, \mathbf{S}) , and $\nabla_{\mathbf{L}} l(\mathbf{L}^k, \mathbf{S}^k) = \nabla_{\mathbf{S}} l(\mathbf{L}^k, \mathbf{S}^k) = \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]$. Since \mathbf{L} and \mathbf{S} are decoupled in (5.1), equivalently we have

$$(5.2) \quad \mathbf{L}^{k+1} = \arg \min_{\mathbf{L}} \left\| \mathbf{L} - \left(\mathbf{L}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right) \right\|_F^2 + \lambda_L \|\mathbf{L}\|_*,$$

$$(5.3) \quad \mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} \left\| \mathbf{S} - \left(\mathbf{S}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right) \right\|_F^2 + \lambda_S \|\mathbf{S}\|_1.$$

The solution to problem (5.3) can be given explicitly in terms of the proximal mapping of $\|\cdot\|_1$ as introduced in Section 2.2, i.e.,

$$\mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S/2} \left[\mathbf{S}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right].$$

For a matrix \mathbf{X} and any $\tau \geq 0$, let $\mathcal{D}_\tau(\mathbf{X})$ denote the singular value thresholding operator $\mathcal{D}_\tau(\mathbf{X}) = \mathbf{U} \mathcal{T}_\tau(\mathbf{\Sigma}) \mathbf{V}^\top$, where $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ is the singular value decomposition

[§]The ratio $\lambda_L/\lambda_S = \sqrt{\rho \max(m, n)}$ follows the suggestion in [3]. For applications in computer vision at least, our choices in λ_L and λ_S seem to be quite robust, although it is possible to improve the performance by making slight adjustments to our current settings of λ_L and λ_S .

[¶]As discussed in [33, 34], with noisy data, solving optimization problems to high accuracy does not necessarily improve the recovery quality. Consequently, we set ε to a modest value.

^{||}All codes are available at: <https://sites.google.com/site/mucun1988/publi>

of \mathbf{X} . It is not difficult to show [35, 36] that the solution to problem (5.2) can be given explicitly by

$$\mathbf{L}^{k+1} = \mathcal{D}_{\lambda_L/2} \left[\mathbf{L}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}] \right].$$

Algorithm 7 summarizes our ISTA implementation for problem (1.5).

Algorithm 7 ISTA for problem (1.5)

```

1: Initialization:  $\mathbf{L}^0 = \mathbf{0}$ ,  $\mathbf{S}^0 = \mathbf{0}$ ;
2: for  $k = 0, 1, 2, \dots$  do
3:    $\mathbf{L}^{k+1} = \mathcal{D}_{\lambda_L/2} [\mathbf{L}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]]$ ;
4:    $\mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S/2} [\mathbf{S}^k - \frac{1}{2} \mathcal{P}_\Omega[\mathbf{L}^k + \mathbf{S}^k - \mathbf{M}]]$ ;
5: end for
```

Regarding ISTA's speed of convergence, it can be proved that $f(\mathbf{L}^k, \mathbf{S}^k) - f^* = O(1/k)$, where f^* denotes the optimal value of problem (1.5).

Fast iterative soft-thresholding algorithm (FISTA) introduced in [20], is an accelerated version of ISTA, which incorporate a momentum step borrowed from Nesterov's optimal gradient scheme [37]. For FISTA, a better convergence result, $f(\mathbf{L}^k, \mathbf{S}^k) - f^* = O(1/k^2)$, can be achieved with a cost per iteration that is comparable to ISTA. Algorithm 8 summarizes our FISTA implementation for problem (1.5).

Algorithm 8 FISTA for problem (1.5)

```

1: Initialization:  $\hat{\mathbf{L}}^0 = \mathbf{L}^0 = \mathbf{0}$ ,  $\hat{\mathbf{S}}^0 = \mathbf{S}^0 = \mathbf{0}$ ,  $t_0 = 1$ ;
2: for  $k = 0, 1, 2, \dots$  do
3:    $\mathbf{L}^{k+1} = \mathcal{D}_{\lambda_L/2} [\hat{\mathbf{L}}^k - \frac{1}{2} \mathcal{P}_\Omega[\hat{\mathbf{L}}^k + \hat{\mathbf{S}}^k - \mathbf{M}]]$ ;
4:    $\mathbf{S}^{k+1} = \mathcal{T}_{\lambda_S/2} [\hat{\mathbf{S}}^k - \frac{1}{2} \mathcal{P}_\Omega[\hat{\mathbf{L}}^k + \hat{\mathbf{S}}^k - \mathbf{M}]]$ ;
5:    $t^{k+1} = \frac{1 + \sqrt{1 + 4(t^k)^2}}{2}$ ;
6:    $\hat{\mathbf{L}}^{k+1} = \mathbf{L}^{k+1} + \frac{t^k - 1}{t^{k+1}} (\mathbf{L}^{k+1} - \mathbf{L}^k)$ ;
7:    $\hat{\mathbf{S}}^{k+1} = \mathbf{S}^{k+1} + \frac{t^k - 1}{t^{k+1}} (\mathbf{S}^{k+1} - \mathbf{S}^k)$ ;
8: end for
```

Partial SVD. In each iteration of either ISTA or FISTA, we only need those singular values that are larger than $\lambda_S/2$ and their corresponding singular vectors. Therefore, a partial SVD can be utilized to reduce the computational burden of a full SVD. Since most partial SVD software packages (e.g. PROPACK [38]) require specifying in advance the number of top singular values and singular vectors to compute, we heuristically determine this number (denoted as sv^k at iteration k). Specifically, let $d = \min\{m, n\}$, and svp^k denote the number of computed singular values that were larger than $\lambda_S/2$ in the k -th iteration. Similar to [17], in our implementation, we start with $sv^0 = d/10$, and adjust sv^k dynamically as follows:

$$sv^{k+1} = \begin{cases} \min\{svp^k + 1, d\} & \text{if } svp^k < sv^k \\ \min\{svp^k + \text{round}(0.05d), d\} & \text{otherwise.} \end{cases}$$

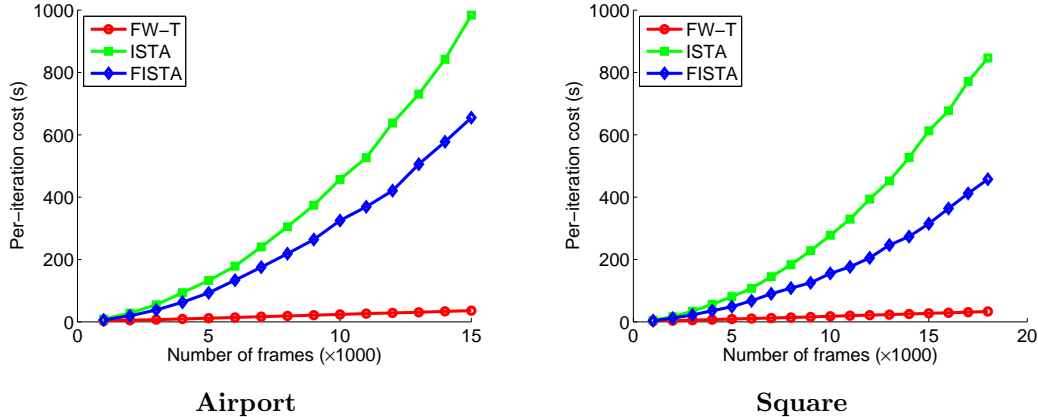


FIG. 2. *Per-iteration cost vs. the number of frames in Airport and Square videos with full observation.* The per-iteration cost of our FW-T method grows linearly with the size of data, in contrast with the superlinear per-iteration cost of ISTA and FISTA. That makes the FW-T method more advantageous or may even be the only feasible choice for large problems.

5.2. Foreground-background separation in surveillance video. In surveillance videos, due to the strong correlation between frames, it is natural to model the background as low rank; while foreground objects, such as cars or pedestrians, that normally occupy only a fraction of the video, can be treated as sparse. So, if we stack each frame as a column in the data matrix \mathbf{M}_0 , it is reasonable to assume that $\mathbf{M}_0 \approx \mathbf{L}_0 + \mathbf{S}_0$, where \mathbf{L}_0 captures the background and \mathbf{S}_0 represents the foreground movements. Here, we solved problem (1.5) for videos introduced in [39] and [40]. The observed entries were sampled uniformly with ratio ρ chosen respectively as 1, 0.8 and 0.6.

Table 1 summarizes the numerical performances of FW-T, ISTA and FISTA in terms of the iteration number and running time (in seconds). As can be observed, our FW-T method is more efficient than ISTA and FISTA, and the advantage becomes more prominent as the size of the data grows and the observations are more compressed (with smaller sampling ratio ρ). Even though the FW-T method took more iterations than FISTA and in many cases than ISTA, it took less time in many cases but one due to its low per-iteration cost. To illustrate this more clearly, in Figure 2, we plot the per-iteration cost of these three methods on the Airport and Square videos as a function of the number of frames. The computational cost of FW-T scales linearly with the size of the data, whereas the cost of the other methods increases superlinearly. Another observation is that as the number of measurements decreases, the iteration numbers of both ISTA and FISTA methods grow substantially, while those of the FW-T method remain quite stable. This explains the more favorable behavior of the FW-T method when ρ is small. In Figure 3, frames of the original videos, the backgrounds and the foregrounds produced by the FW-T method are presented, and the separation achieved is quite satisfactory.

5.3. Shadow and specular removal from face images. Images taken under varying illumination can also be modeled as the superposition of low-rank and sparse components. Here, the data matrix \mathbf{M}_0 is again formed by stacking each image as a column. The low-rank term \mathbf{L}_0 captures the smooth variations [41], while the sparse term \mathbf{S}_0 represents cast shadows and specularities [42, 8]. CPCP can

TABLE 1

Comparisons of FW-T, ISTA and FISTA on surveillance video data. The advantage of our FW-T method becomes prominent when the data are at large scale and compressed (i.e. the small ρ scenarios).

Data	ρ	FW-T		ISTA		FISTA	
		iter.	time	iter.	time	iter.	time
Lobby (20480 \times 1000)	1.0	96	1.94e+02	144	3.64e+02	41	1.60e+02
	0.8	104	2.33e+02	216	1.03e+03	52	3.55e+02
	0.6	133	3.12e+02	380	1.67e+03	74	5.10e+02
Campus (20480 \times 1439)	1.0	45	1.56e+02	78	1.49e+03	23	4.63e+02
	0.8	44	1.57e+02	122	2.34e+03	30	6.45e+02
	0.6	41	1.39e+02	218	4.27e+03	43	1.08e+03
Escalator (20800 \times 3417)	1.0	81	7.40e+02	58	4.19e+03	25	2.18e+03
	0.8	80	7.35e+02	90	8.18e+03	32	3.46e+03
	0.6	82	7.68e+02	162	1.83e+04	43	5.73e+03
Mall (81920 \times 1286)	1.0	38	4.70e+02	110	5.03e+03	35	1.73e+03
	0.8	35	4.58e+02	171	7.32e+03	44	2.34e+03
	0.6	44	5.09e+02	308	1.31e+04	62	3.42e+03
Restaurant (19200 \times 3055)	1.0	70	5.44e+02	52	3.01e+03	20	1.63e+03
	0.8	74	5.51e+02	81	4.84e+03	26	1.82e+03
	0.6	76	5.73e+02	144	9.93e+03	38	3.31e+03
Hall (25344 \times 3584)	1.0	60	6.33e+02	52	2.98e+03	21	1.39e+03
	0.8	62	6.52e+02	81	6.45e+03	28	2.90e+03
	0.6	70	7.43e+02	144	1.42e+04	39	4.94e+03
Airport (25344 \times 15730)	1.0	130	6.42e+03	29	2.37e+04	14	1.37e+04
	0.8	136	6.65e+03	45	6.92e+04	18	4.27e+04
	0.6	154	7.72e+03	77	1.78e+05	24	7.32e+04
Square (19200 \times 28181)	1.0	179	1.24e+04	29	3.15e+04	13	1.51e+04
	0.8	181	1.26e+04	44	1.04e+05	17	6.03e+04
	0.6	191	1.31e+04	78	2.63e+05	22	9.88e+05

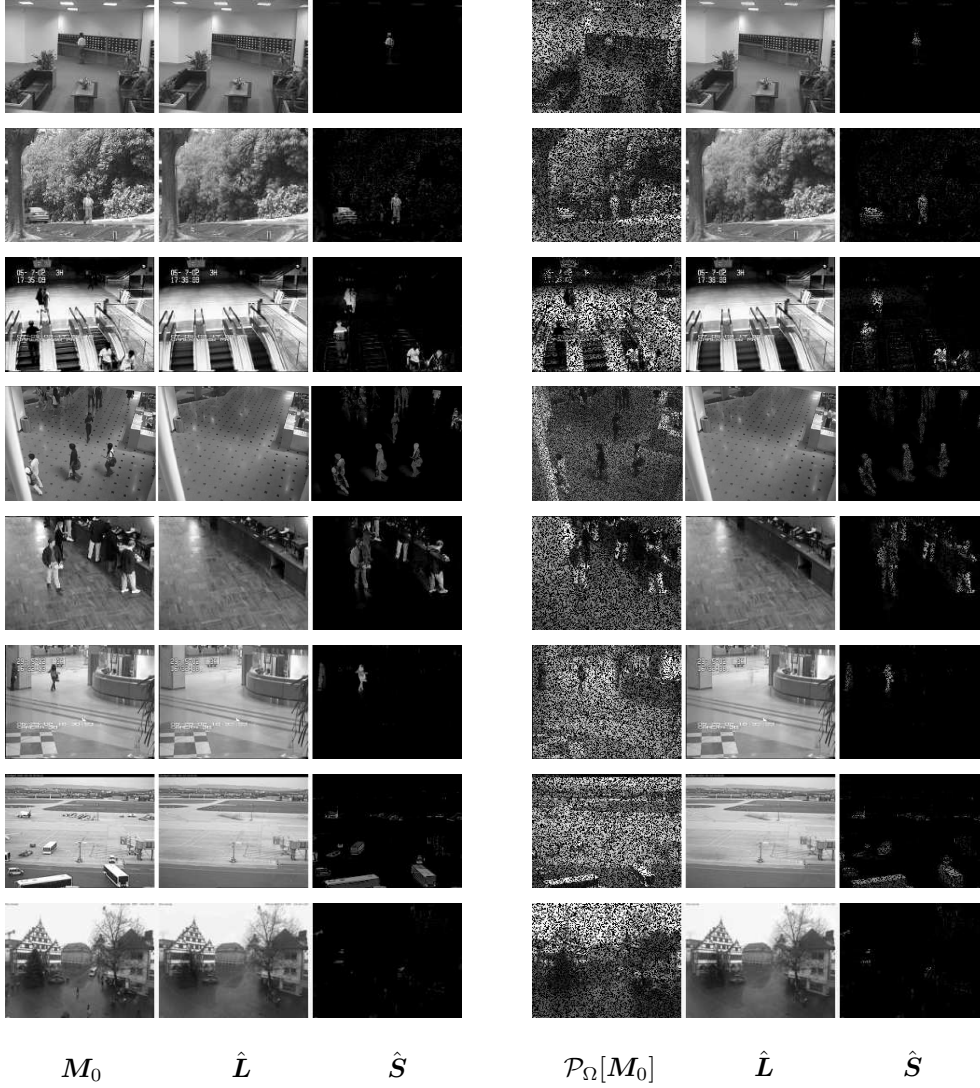


FIG. 3. **Surveillance videos.** The videos from top to bottom are respectively Lobby, Campus, Escalator, Mall, Restaurant, Hall, Airport and Square. The left panel presents videos with full observation ($\rho = 1$) and the right one presents videos with partial observation ($\rho = 0.6$). Visually, the low-rank component successfully recovers the background and the sparse one captures the moving objects (e.g. vehicles, pedestrians) in the foreground.

be used to remove the shadows and specularities [3, 8]. Here, we solved problem (1.4) for YaleB face images [43]. Table 2 summarizes the numerical performances of FW-T, ISTA and FISTA. Similar to the observation made regarding the above surveillance video experiment, the number of iterations required by ISTA and FISTA grows much faster than it does for the FW-T method when ρ decreases. However, unlike in those tests, where the number of frames in each dataset was at least several thousand, the number of frames here is just 65. This prevents the FW-T method from significantly benefiting from its linear per-iteration cost and consequently, while FW-T still outperforms ISTA for values of $\rho \leq 0.7$, the FISTA method is always the

TABLE 2

Comparisons of FW-T, ISTA and FISTA on YaleB face data. The number of frames, 65, is relatively small for this application. This disables the FW-T method to significantly benefit from its linear per-iteration cost and consequently the FISTA method consistently has a better performance.

Data	ρ	FW-T		ISTA		FISTA	
		iter.	time	iter.	time	iter.	time
YaleB01 (32256 \times 65)	1.0	65	34.0	49	21.4	17	8.69
	0.9	68	35.6	59	23.9	19	8.62
	0.8	79	42.2	76	35.3	22	10.9
	0.7	76	39.9	97	44.0	25	11.1
	0.6	71	37.5	127	50.2	29	12.9
	0.5	80	40.5	182	77.9	35	15.2
YaleB02 (32256 \times 65)	1.0	64	34.6	51	19.2	18	7.31
	0.9	64	26.8	61	22.6	20	7.32
	0.8	71	33.9	78	27.7	22	8.61
	0.7	71	31.3	99	36.6	26	11.0
	0.6	73	36.6	132	53.7	30	12.4
	0.5	63	28.0	177	64.6	35	13.4
YaleB03 (32256 \times 65)	1.0	62	26.0	49	16.6	18	6.00
	0.9	71	27.5	62	20.3	20	6.43
	0.8	69	30.0	78	26.0	22	8.32
	0.7	78	31.5	101	32.9	26	9.00
	0.6	73	28.7	132	40.4	30	10.6
	0.5	70	28.0	181	60.3	36	12.8
YaleB04 (32256 \times 65)	1.0	63	28.5	47	16.6	17	6.35
	0.9	67	28.7	58	23.1	19	7.98
	0.8	68	31.7	72	26.3	23	9.39
	0.7	69	30.7	92	35.9	26	9.84
	0.6	71	29.4	124	40.0	29	10.1
	0.5	74	29.4	174	67.3	36	14.3

fastest. In Figure 4, the original images, the low-rank and the sparse parts produced by the FW-T method are presented. Visually, the recovered low-rank component is smoother and better conditioned for face recognition than the original image, while the sparse component corresponds to shadows and specularities.

6. Discussion. In this paper, we have proposed scalable algorithms called Frank-Wolfe-Projection (FW-P) and Frank-Wolfe-Thresholding (FW-T) for norm constrained and penalized versions of CPCP. Essentially, these methods combine classical ideas in Frank-Wolfe and Proximal methods to achieve linear per-iteration cost, $O(1/k)$ convergence in function value and practical efficiency in updating the sparse component. Extensive numerical experiments were conducted on computer vision related

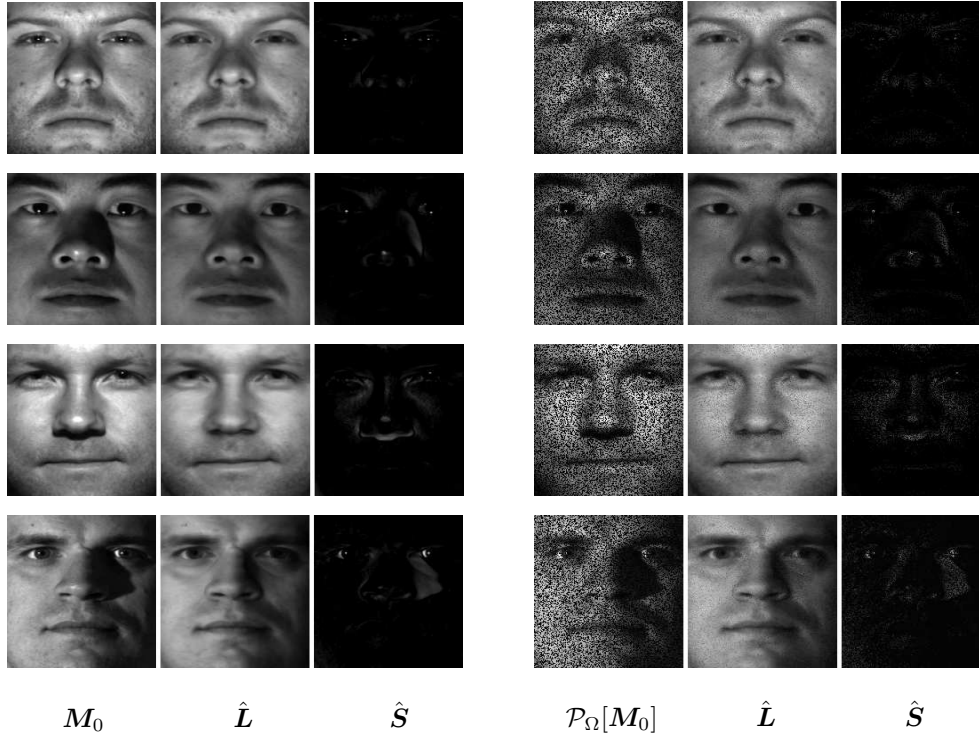


FIG. 4. **Face images.** The pictures from top to bottom are respectively YaleB01, YaleB02, YaleB03 and YaleB04 face images. The left panel presents the case with full observation ($\rho = 1$), while the right panel presents the case with partial observation ($\rho = 0.6$). Visually, the recovered low-rank component is smoother and better conditioned for face recognition than the original image, while the sparse component corresponds to shadows and specularities.

applications of CPCP, which demonstrated the great potential of our methods for dealing with problems of very large scale. Moreover, the general idea of leveraging different methods to deal with different functions may be valuable for other demixing problems.

We are also aware that though our algorithms are extremely efficient in the beginning iterations and quickly arrive at an approximate solution of practical significance, they become less competitive in solutions of very high accuracy, due to the nature of Frank-Wolfe. This suggests further hybridization under our framework (e.g. using nonconvex approaches to handle the nuclear norm) might be utilized in certain applications (see [44] for research in that direction).

Acknowledgements. We are grateful to the associate editor Chen Greif and three anonymous reviewers for their helpful suggestions and comments that substantially improve the paper.

REFERENCES

- [1] J. Wright, A. Ganesh, K. Min, and Y. Ma, “Compressive principal component pursuit,” *Information and Inference*, vol. 2, no. 1, pp. 32–68, 2013.
- [2] V. Chandrasekaran, S. Sanghavi, P. Parrilo, and A. Willsky, “Rank-sparsity incoherence for matrix decomposition,” *SIAM Journal on Optimization*, vol. 21, no. 2, pp. 572–596, 2011.

- [3] E. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?,” *Journal of the ACM (JACM)*, vol. 58, no. 3, pp. 11:1–11:37, 2011.
- [4] Z. Zhou, X. Li, J. Wright, E. Candès, and Y. Ma, “Stable principal component pursuit,” in *ISIT*, 2010.
- [5] D. Hsu, S. Kakade, and T. Zhang, “Robust matrix decomposition with sparse corruptions,” *IEEE Transactions on Information Theory*, vol. 57, no. 11, pp. 7221–7234, 2011.
- [6] A. Agarwal, S. Negahban, and M. Wainwright, “Noisy matrix decomposition via convex relaxation: Optimal rates in high dimensions,” *The Annals of Statistics*, vol. 40, no. 2, pp. 1171–1197, 2012.
- [7] Y. Peng, A. Ganesh, J. Wright, W. Xu, and Y. Ma, “Rasl: Robust alignment by sparse and low-rank decomposition for linearly correlated images,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 11, pp. 2233–2246, 2012.
- [8] Y. Zhang, C. Mu, H. Kuo, and J. Wright, “Towards guaranteed illumination models for non-convex objects,” in *ICCV*, 2013.
- [9] L. Wu, A. Ganesh, B. Shi, Y. Matsushita, Y. Wang, and Y. Ma, “Robust photometric stereo via low-rank matrix completion and recovery,” in *ACCV*, 2011.
- [10] R. Otazo, E. Candès, and D. K. Sodickson, “Low-rank plus sparse matrix decomposition for accelerated dynamic MRI with separation of background and dynamic components,” *Magnetic Resonance in Medicine*, 2014.
- [11] K. Min, Z. Zhang, J. Wright, and Y. Ma, “Decomposing background topics from keywords by principal component pursuit,” in *CIKM*, 2010.
- [12] V. Chandrasekaran, P. Parrilo, and A. Willsky, “Latent variable graphical model selection via convex optimization,” *Annals of Statistics*, vol. 40, no. 4, pp. 1935–1967, 2012.
- [13] Z. Lin, A. Ganesh, J. Wright, L. Wu, M. Chen, and Y. Ma, “Fast convex optimization algorithms for exact recovery of a corrupted low-rank matrix,” in *CAMSAP*, 2009.
- [14] Z. Lin, M. Chen, and Y. Ma, “The augmented lagrange multiplier method for exact recovery of corrupted low-rank matrices,” *arXiv preprint arXiv:1009.5055*, 2010.
- [15] X. Yuan and J. Yang, “Sparse and low-rank matrix decomposition via alternating direction methods,” *preprint*, 2009.
- [16] N. S. Aybat, D. Goldfarb, and G. Iyengar, “Fast first-order methods for stable principal component pursuit,” *arXiv preprint arXiv:1105.2126*, 2011.
- [17] M. Tao and X. Yuan, “Recovering low-rank and sparse components of matrices from incomplete and noisy observations,” *SIAM Journal on Optimization*, vol. 21, no. 1, pp. 57–81, 2011.
- [18] N. S. Aybat, D. Goldfarb, and S. Ma, “Efficient algorithms for robust and stable principal component pursuit problems,” *Computational Optimization and Applications*, pp. 1–29, 2012.
- [19] R. Tütüncü, K. Toh, and M. Todd, “Solving semidefinite-quadratic-linear programs using sdpt3,” *Mathematical Programming*, vol. 95, no. 2, pp. 189–217, 2003.
- [20] A. Beck and M. Teboulle, “A fast iterative shrinkage-thresholding algorithm for linear inverse problems,” *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [21] M. Frank and P. Wolfe, “An algorithm for quadratic programming,” *Naval Research Logistics Quarterly*, vol. 3, no. 1-2, pp. 95–110, 1956.
- [22] E. Levitin and B. Polyak, “Constrained minimization methods,” *USSR Computational Mathematics and Mathematical Physics*, vol. 6, no. 5, pp. 1–50, 1966.
- [23] M. Jaggi and M. Sulovsk, “A simple algorithm for nuclear norm regularized problems,” in *ICML*, 2010.
- [24] Z. Harchaoui, A. Juditsky, and A. Nemirovski, “Conditional gradient algorithms for norm-regularized smooth convex optimization,” *Mathematical Programming*, pp. 1–38, 2014.
- [25] M. Jaggi, “Revisiting Frank-Wolfe: Projection-free sparse convex optimization,” in *ICML*, 2013.
- [26] V. F. Demianov and A. M. Rubinov, *Approximate methods in optimization problems*. Modern analytic and computational methods in science and mathematics, American Elsevier Pub. Co., 1970.
- [27] J. C. Dunn and S. Harshbarger, “Conditional gradient algorithms with open loop step size rules,” *Journal of Mathematical Analysis and Applications*, vol. 62, no. 2, pp. 432 – 444, 1978.
- [28] M. Patriksson, “Partial linearization methods in nonlinear programming,” *Journal of Optimization Theory and Applications*, vol. 78, no. 2, pp. 227–246, 1993.
- [29] T. Zhang, “Sequential greedy approximation for certain convex optimization problems,” *IEEE Transactions on Information Theory*, vol. 49, no. 3, pp. 682–691, 2003.
- [30] K. Clarkson, “Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm,” *ACM Trans. Algorithms*, vol. 6, no. 4, pp. 63:1–63:30, 2010.

- [31] R. M. Freund and P. Grigas, “New analysis and results for the frank–wolfe method,” *Mathematical Programming*, vol. 155, no. 1-2, pp. 199–230, 2016.
- [32] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra, “Efficient projections onto the l_1 -ball for learning in high dimensions,” in *ICML*, 2008.
- [33] J. Yang and Y. Zhang, “Alternating direction algorithms for ℓ_1 -problems in compressive sensing,” *SIAM Journal on Scientific Computing*, vol. 33, no. 1, pp. 250–278, 2011.
- [34] J. Yang and X. Yuan, “Linearized augmented lagrangian and alternating direction methods for nuclear norm minimization,” *Mathematics of Computation*, vol. 82, no. 281, pp. 301–329, 2013.
- [35] J. Cai, E. Candès, and Z. Shen, “A singular value thresholding algorithm for matrix completion,” *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010.
- [36] S. Ma, D. Goldfarb, and L. Chen, “Fixed point and Bregman iterative methods for matrix rank minimization,” *Mathematical Programming*, vol. 128, no. 1-2, pp. 321–353, 2011.
- [37] Y. Nesterov, “A method of solving a convex programming problem with convergence rate $O(1/k^2)$,” in *Soviet Mathematics Doklady*, vol. 27, pp. 372–376, 1983.
- [38] R. M. Larsen, “Propack-software for large and sparse svd calculations,” *Available online. URL <http://sun.stanford.edu/rmunk/PROPACK>*, pp. 2008–2009, 2004.
- [39] L. Li, W. Huang, I. Y. Gu, and Q. Tian, “Statistical modeling of complex backgrounds for foreground object detection,” *IEEE Transactions on Image Processing*, vol. 13, no. 11, pp. 1459–1472, 2004.
- [40] N. Jacobs, N. Roman, and R. Pless, “Consistent temporal variations in many outdoor scenes,” in *CVPR*, 2007.
- [41] R. Basri and D. Jacobs, “Lambertian reflectance and linear subspaces,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 2, pp. 218–233, 2003.
- [42] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma, “Robust face recognition via sparse representation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 2, pp. 210–227, 2009.
- [43] A. Georgiades, P. Belhumeur, and D. Kriegman, “From few to many: Illumination cone models for face recognition under variable lighting and pose,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 6, pp. 643–660, 2001.
- [44] S. Laue, “A hybrid algorithm for convex semidefinite optimization,” in *ICML*, 2012.