

# A Primal-Dual Frank-Wolfe Algorithm for Linear Programming\*

Matthew Hough<sup>†</sup>      Stephen A. Vavasis<sup>‡</sup>

February 28, 2024

## Abstract

We present two first-order primal-dual algorithms for solving saddle point formulations of linear programs, namely FWLP (Frank-Wolfe Linear Programming) and FWLP-P. The former iteratively applies the Frank-Wolfe algorithm to both the primal and dual of the saddle point formulation of a standard-form LP. The latter is a modification of FWLP in which regularizing perturbations are used in computing the iterates. We show that FWLP-P converges to a primal-dual solution with error  $\mathcal{O}(1/\sqrt{k})$  after  $k$  iterations, while no convergence guarantees are provided for FWLP. We also discuss the advantages of using FWLP and FWLP-P for solving very large LPs. In particular, we argue that only part of the matrix  $A$  is needed at each iteration, in contrast to other first-order methods.

## 1 Introduction

In recent years, data science applications have given birth to problems of very large scale. This poses a problem for mature LP solvers that require solving a system of linear equations at each iteration. First-order methods (FoMs) for linear programming aim to solve LPs in such a way that their most expensive operation at each iteration is the product of a matrix and a vector. Their goal is to provide an alternative to the practitioner over LP algorithms such as the simplex method or interior point methods for large-scale problems.

The Frank-Wolfe algorithm [7], also referred to as the conditional gradient algorithm, is a FoM for minimizing a smooth convex objective function over a compact convex set. A major benefit of the Frank-Wolfe algorithm is that each iteration requires only the

---

\*This research was supported in part by a Discovery Grant from the Natural Sciences and Engineering Research Council (NSERC) of Canada.

<sup>†</sup>Department of Combinatorics & Optimization, University of Waterloo, 200 University Ave. W., Waterloo, ON, N2L 3G1, Canada, [mhough@uwaterloo.ca](mailto:mhough@uwaterloo.ca).

<sup>‡</sup>Department of Combinatorics & Optimization, University of Waterloo, 200 University Ave. W., Waterloo, ON, N2L 3G1, Canada, [vavasis@uwaterloo.ca](mailto:vavasis@uwaterloo.ca).

solution of a linear optimization problem over a convex constraint set, a problem which can be solved efficiently over many constraint sets used in practice. It is known that the Frank-Wolfe algorithm converges at a rate of  $\mathcal{O}(1/k)$  [12].

The focus of this paper is on finding optimal solutions to linear programs in standard form, that is, solving the following optimization problem:

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned} \tag{1}$$

where we assume throughout that (1) has an optimal solution. The dual linear program associated with (1) is

$$\begin{aligned} \max \quad & \mathbf{b}^T \mathbf{y} \\ \text{s.t.} \quad & A^T \mathbf{y} \leq \mathbf{c}. \end{aligned} \tag{2}$$

We propose two first-order primal-dual algorithms for simultaneously solving (1) and (2) inspired by the Frank-Wolfe algorithm but using the non-standard step-size of  $1/(k+1)$  analyzed in [8]. We call our algorithms FWLP and FWLP-P. FWLP, first introduced in [11], is derived from iteratively applying the Frank-Wolfe algorithm to the primal and dual problems of a modified saddle-point formulation of (1):

$$\min_{\mathbf{x} \in \Delta} \max_{\mathbf{y} \in \Gamma} \mathcal{L}(\mathbf{x}, \mathbf{y}) := \mathbf{c}^T \mathbf{x} + \mathbf{y}^T (\mathbf{b} - A\mathbf{x}), \tag{3}$$

where we define

$$\Delta = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq \mathbf{0}, \mathbf{e}^T \mathbf{x} \leq \xi\} \text{ and } \Gamma = [-\eta, \eta]^m. \tag{4}$$

Let  $(\mathbf{x}^*, \mathbf{y}^*)$  denote an optimal primal-dual pair of solutions to (1) and (2). The parameters  $\xi, \eta > 0$  are assumed to be chosen large enough to ensure that  $2\xi \geq \|\mathbf{x}^*\|_1$  and  $2\eta \geq \|\mathbf{y}^*\|_\infty$ , thus describing redundant constraints. Despite their redundancy, these constraints are necessary for compactness of the feasible sets corresponding to the primal and dual subproblems solved by FWLP. FWLP-P is a modification of the original FWLP algorithm in which regularizing perturbations are introduced when computing the iterates. FWLP-P is of theoretical importance to FWLP since we are able to prove the convergence of FWLP-P but no convergence proof is known yet for FWLP.

With the ubiquity of very large-scale problems, FoMs for linear programming have seen interest of late [1, 2, 4, 14]. Such methods typically apply existing ideas and algorithms from continuous optimization to linear programming. Our work is related particularly to [9], where Gidel et al. use the Frank-Wolfe algorithm to solve general convex-concave saddle point problems:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}), \tag{5}$$

where  $f$  is a smooth convex-concave function and  $\mathcal{X} \times \mathcal{Y}$  is a convex compact set. In analyzing the convergence of their method, the authors introduce a potential function

bounding the distance to optimality and then show that this potential function decreases to zero at a given rate. We employ this technique in Section 3 to analyze the convergence of FWLP-P. The results of [9] do not directly apply to (3) because their strong-convexity assumption is not satisfied. The FWLP algorithm is also closely related to the Generalized Fictitious Play Algorithm proposed by Hammond in her 1984 PhD thesis [10]. Applying Hammond’s algorithm to (3) yields an algorithm very similar to FWLP, the only difference being that the update of  $\mathbf{y}_k$  is based on  $\mathbf{x}_k$  instead of  $\mathbf{x}_{k+1}$  like in FWLP. For a more in-depth analysis of the similarities between FWLP and Hammond’s Generalized Fictitious Play, see [11, Chapter 6.2].

The algorithms FWLP and FWLP-P are described in Section 2. Iterations of FWLP are simpler and faster than those of FWLP-P, but we do not have a convergence proof of FWLP. Our preliminary computational tests (not reported here) indicate that the two algorithms converge at comparable rates.

Our convergence analysis of FWLP-P, presented in Section 3, first introduces a potential function  $U_k$  which we show that for the iterates of FWLP-P has distance  $\mathcal{O}(1/\sqrt{k})$  from zero after  $k$  iterations. To complete the analysis, we show that as  $U_k \rightarrow 0$ , FWLP-P and FWLP converge to a primal-dual solution of (1). The potential  $U_k$  is similar to a traditional primal-dual optimality gap as discussed in Section 4. A secondary contribution of this paper is the discussion in Section 5 of how FWLP and FWLP-P can be implemented efficiently in a way that only part of the matrix  $A$  is needed at each iteration.

## 2 FWLP and FWLP-P

At each iteration, FWLP performs a Frank-Wolfe update on the primal and the dual of (3) using the step-size  $1/(k+1)$  instead of the standard step-size of  $1/(k+2)$ . The step-size  $1/(k+1)$  was first analyzed by Freund and Grigas in [8]. Notably, FWLP uses the information obtained from the primal update,  $\mathbf{x}_{k+1}$ , in the computation of the dual update:

$$\mathbf{r}_{k+1} := \operatorname{argmin}_{\mathbf{r}} \{ (\mathbf{c} - A^T \mathbf{y}_k)^T \mathbf{r} : \mathbf{r} \in \Delta \}, \quad (6)$$

$$\mathbf{x}_{k+1} := \frac{k}{k+1} \mathbf{x}_k + \frac{1}{k+1} \mathbf{r}_{k+1}, \quad (7)$$

$$\mathbf{s}_{k+1} := \operatorname{argmax}_{\mathbf{s}} \{ (\mathbf{b} - A \mathbf{x}_{k+1})^T \mathbf{s} : \mathbf{s} \in \Gamma \}, \quad (8)$$

$$\mathbf{y}_{k+1} := \frac{k}{k+1} \mathbf{y}_k + \frac{1}{k+1} \mathbf{s}_{k+1}. \quad (9)$$

It is not hard to see that steps (6) and (8) above can be written in closed form (for more detail, see [11, Chapter 6]), giving rise to Algorithm 2.1. In this algorithm and for the remainder of the paper,  $\mathbf{e}_i$  denotes the  $i$ th column of the identity matrix, whose length is determined from the context.

---

**Algorithm 2.1** FWLP: A primal-dual algorithm for (1) based on Frank-Wolfe [11].

---

**Require:** Starting points  $\mathbf{x}_0 \in \mathbb{R}_+^n$ ,  $\mathbf{y}_0 \in \mathbb{R}^m$ , constraint data  $A$  and  $\mathbf{b}$ , and objective  $\mathbf{c}$ .

Parameters:  $\xi, \eta > 0$  such that  $2\|\mathbf{x}^*\| \leq \xi$  and  $2\|\mathbf{y}^*\| \leq \eta$ .

- 1: **for**  $k = 1, 2, \dots$  **do**
- 2:     Determine  $i = \operatorname{argmin}_t [\mathbf{c} - A^\top \mathbf{y}_k]_t$ .
- 3:     **if**  $[\mathbf{c} - A^\top \mathbf{y}_k]_i \geq 0$  **then**
- 4:         Step towards zero in  $\mathbf{x}$ :

$$\mathbf{x}_{k+1} := \frac{k}{k+1} \mathbf{x}_k. \quad (10)$$

- 5:     **else**
- 6:         Step toward  $\xi$  for  $x^{(i)}$ , otherwise step toward zero for  $x^{(j)}$ ,  $j \neq i$ :

$$\mathbf{x}_{k+1} = \frac{k}{k+1} \mathbf{x}_k + \frac{\xi}{k+1} \mathbf{e}_i. \quad (11)$$

- 7:     **end if**
- 8:     Step towards  $\pm\eta$  for each coordinate in  $\mathbf{y}$  according to the sign pattern of  $\mathbf{b} - A\mathbf{x}_{k+1}$ :

$$\mathbf{y}_{k+1} := \frac{k}{k+1} \mathbf{y}_k + \frac{\eta}{k+1} \operatorname{sgn}(\mathbf{b} - A\mathbf{x}_{k+1}). \quad (12)$$

- 9: **end for**
- 

FWLP-P adds the regularizing perturbations  $\|\mathbf{r}\|^2/(2\sqrt{k})$  and  $\|\mathbf{s}\|^2/(2\sqrt{k})$  to steps (10)/(11) and (12) in the above description of FWLP yielding (15) and (17) respectively in Algorithm 2.2. We analyze this algorithm in the next section.

Note that the computation of  $\mathbf{s}_{k+1}$  in FWLP-P is a separable optimization problem allowing a simple median algorithm for each coordinate entry of  $\mathbf{s}_{k+1}$ . The computation of  $\mathbf{r}_{k+1}$  is more elaborate, but it requires only evaluation of the smallest (most negative) elements of  $\mathbf{c} - A^\top \mathbf{y}_k$ . We cover this in more detail in Section 5 but of note here is that FWLP-P retains some of the benefit of FWLP, namely sub-linear time computation per iteration, since only a part of  $A$  is required. Note also that the formulas (15) and (17) for  $\mathbf{r}_{k+1}$  and  $\mathbf{s}_{k+1}$  can be written equivalently as projections:

$$\mathbf{r}_{k+1} := \operatorname{proj}_\Delta(\sqrt{k}(A^\top \mathbf{y}_k - \mathbf{c})), \quad (13)$$

$$\mathbf{s}_{k+1} := \operatorname{proj}_\Gamma(\sqrt{k}(\mathbf{b} - A\mathbf{x}_{k+1})). \quad (14)$$

We discuss how to compute these projections efficiently in Section 5.2.

---

**Algorithm 2.2** FWLP-P: FWLP with perturbations.

---

**Require:** Starting points  $\mathbf{x}_0 \in \mathbb{R}_+^n$ ,  $\mathbf{y}_0 \in \mathbb{R}^m$ , constraint data  $A$  and  $\mathbf{b}$ , and objective  $\mathbf{c}$ .

**Parameters:**  $\xi, \eta > 0$  such that  $2\|\mathbf{x}^*\| \leq \xi$  and  $2\|\mathbf{y}^*\| \leq \eta$ .

1: **for**  $k = 1, 2, \dots$  **do**

$$\mathbf{r}_{k+1} := \operatorname{argmin}_{\mathbf{r}} \left\{ (\mathbf{c} - A^T \mathbf{y}_k)^T \mathbf{r} + \frac{\|\mathbf{r}\|^2}{2\sqrt{k}} : \mathbf{r} \in \Delta \right\}, \quad (15)$$

$$\mathbf{x}_{k+1} := \frac{k}{k+1} \mathbf{x}_k + \frac{1}{k+1} \mathbf{r}_{k+1}, \quad (16)$$

$$\mathbf{s}_{k+1} := \operatorname{argmax}_{\mathbf{s}} \left\{ (\mathbf{b} - A \mathbf{x}_{k+1})^T \mathbf{s} - \frac{\|\mathbf{s}\|^2}{2\sqrt{k}} : \mathbf{s} \in \Gamma \right\}, \quad (17)$$

$$\mathbf{y}_{k+1} := \frac{k}{k+1} \mathbf{y}_k + \frac{1}{k+1} \mathbf{s}_{k+1}. \quad (18)$$

2: **end for**

---

### 3 Convergence analysis of FWLP-P

The forthcoming analysis derives a recursion (30) for potential  $U_k$  below satisfied by the iterates of FWLP-P. The recursion has two perturbation terms  $\delta_{k+1}$  and  $\epsilon_{k+1}$  that are subsequently bounded in (32) and (34). We show that the quantity  $U_k$  goes to 0 like  $1/\sqrt{k}$ .

#### 3.1 Deriving the recursion and potential function

Let us introduce

$$\delta_{k+1} := \mathbf{r}_{k+2}^T (\mathbf{c} - A^T \mathbf{y}_{k+1}) + \frac{1}{2\sqrt{k+1}} \|\mathbf{r}_{k+2}\|^2 - \mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_{k+1}) - \frac{k}{2(k+1)\sqrt{k}} \|\mathbf{r}_{k+1}\|^2. \quad (19)$$

Rearrange this equation:

$$\delta_{k+1} - \mathbf{r}_{k+2}^T (\mathbf{c} - A^T \mathbf{y}_{k+1}) - \frac{1}{2\sqrt{k+1}} \|\mathbf{r}_{k+2}\|^2 = -\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_{k+1}) - \frac{k}{2(k+1)\sqrt{k}} \|\mathbf{r}_{k+1}\|^2. \quad (20)$$

Rewrite the first-term on the RHS of (20), that is,  $-\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_{k+1})$ , using (18):

$$\begin{aligned} -\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_{k+1}) &= -\mathbf{r}_{k+1}^T \left( \mathbf{c} - A^T \left( \frac{k}{k+1} \mathbf{y}_k + \frac{1}{k+1} \mathbf{s}_{k+1} \right) \right) \\ &= -\frac{k}{k+1} \mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{k+1} \mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{s}_{k+1}). \end{aligned} \quad (21)$$

Substitute (21) for the first term on the RHS in (20), moving the second term of (21)  $-\mathbf{r}_{k+1}^T(\mathbf{c} - A^T \mathbf{s}_{k+1})/(k+1)$  to the LHS, thereby rewriting (20) as:

$$\begin{aligned} \delta_{k+1} - \mathbf{r}_{k+2}^T(\mathbf{c} - A^T \mathbf{y}_{k+1}) - \frac{1}{2\sqrt{k+1}} \|\mathbf{r}_{k+2}\|^2 + \frac{1}{k+1} \mathbf{r}_{k+1}^T(\mathbf{c} - A^T \mathbf{s}_{k+1}) \\ = -\frac{k}{k+1} \mathbf{r}_{k+1}^T(\mathbf{c} - A^T \mathbf{y}_k) - \frac{k}{2(k+1)\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 \\ = \frac{k}{k+1} \left( -\mathbf{r}_{k+1}^T(\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 \right). \end{aligned} \quad (22)$$

The point of this algebra is that the 2nd and 3rd terms of the LHS correspond to the RHS advanced from  $k$  to  $k+1$ , thus setting up some of the terms of the recursion.

Next, introduce for  $k \geq 2$

$$\epsilon_{k+1} := \frac{k}{k+1} \left( -\mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_k) + \frac{\sqrt{k+1}}{2k} \|\mathbf{s}_{k+1}\|^2 + \mathbf{s}_k^T(\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 \right). \quad (23)$$

Rewrite this equation as

$$\epsilon_{k+1} + \frac{k}{k+1} \left( \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_k) - \frac{\sqrt{k+1}}{2k} \|\mathbf{s}_{k+1}\|^2 \right) = \frac{k}{k+1} \left( \mathbf{s}_k^T(\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 \right). \quad (24)$$

Rearranging (16) yields  $\mathbf{x}_k = (k+1)\mathbf{x}_{k+1}/k - \mathbf{r}_{k+1}/k$ . This means that the factor  $\mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_k)$  appearing in the LHS of (24) may be rewritten

$$\begin{aligned} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_k) &= \mathbf{s}_{k+1}^T \left( \mathbf{b} - A \left( \frac{k+1}{k} \mathbf{x}_{k+1} - \frac{1}{k} \mathbf{r}_{k+1} \right) \right) \\ &= \frac{k+1}{k} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_{k+1}) - \frac{1}{k} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{r}_{k+1}). \end{aligned} \quad (25)$$

Thus, the LHS of (24) is rewritten:

$$\begin{aligned} \text{LHS of (24)} &= \epsilon_{k+1} + \frac{k}{k+1} \left( \frac{k+1}{k} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_{k+1}) - \frac{1}{k} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{r}_{k+1}) - \frac{\sqrt{k+1}}{2k} \|\mathbf{s}_{k+1}\|^2 \right) \\ &= \epsilon_{k+1} + \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_{k+1}) - \frac{1}{k+1} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{r}_{k+1}) - \frac{1}{2\sqrt{k+1}} \|\mathbf{s}_{k+1}\|^2. \end{aligned} \quad (26)$$

Thus, we have rewritten (24) as

$$\begin{aligned} \epsilon_{k+1} + \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{x}_{k+1}) - \frac{1}{k+1} \mathbf{s}_{k+1}^T(\mathbf{b} - A\mathbf{r}_{k+1}) - \frac{1}{2\sqrt{k+1}} \|\mathbf{s}_{k+1}\|^2 \\ = \frac{k}{k+1} \left( \mathbf{s}_k^T(\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 \right). \end{aligned} \quad (27)$$

Here, we see the correspondence between the 2nd and 4th term on the LHS with the two terms on the RHS (with  $k$  advanced by 1).

Multiply (16) by  $\mathbf{c}$  and rearrange to obtain

$$\mathbf{c}^T \mathbf{x}_{k+1} - \frac{1}{k+1} \mathbf{c}^T \mathbf{r}_{k+1} = \frac{k}{k+1} \mathbf{c}^T \mathbf{x}_k. \quad (28)$$

Similarly, from (18),

$$-\mathbf{b}^T \mathbf{y}_{k+1} + \frac{1}{k+1} \mathbf{b}^T \mathbf{s}_{k+1} = -\frac{k}{k+1} \mathbf{b}^T \mathbf{y}_k. \quad (29)$$

Now add (22), (27), (28), and (29), noting that many quantities on the LHS cancel, while all quantities on the RHS contain the factor  $k/(k+1)$ , to obtain a recursion:

$$\delta_{k+1} + \epsilon_{k+1} + U_{k+1} = \frac{k}{k+1} U_k, \quad (30)$$

where for  $k \geq 2$

$$U_k := -\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 + \mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 + \mathbf{c}^T \mathbf{x}_k - \mathbf{b}^T \mathbf{y}_k. \quad (31)$$

We call  $U_k$  the potential function, and we will use this in our convergence analysis to bound the distance from optimality of (1) and (2).

### 3.2 Bounding the potential function

Our next task is to lower bound  $\delta_{k+1}$  and  $\epsilon_{k+1}$  so that we can use (30) to develop a more useful bound on the potential function.

**Lemma 1.** *Recall  $\epsilon_{k+1}$  defined in (23). We have the bound*

$$\epsilon_{k+1} \geq -\frac{m\eta^2}{6k^2\sqrt{k-1}}.$$

*Proof.* By adding and subtracting multiples of  $\|\mathbf{s}_{k+1}\|^2$  and  $\|\mathbf{s}_k\|^2$  inside (23), we obtain

$$\epsilon_{k+1} = \frac{k}{k+1} (\epsilon'_{k+1} + \epsilon''_{k+1} + \epsilon'''_{k+1}),$$

where

$$\begin{aligned} \epsilon'_{k+1} &:= -\mathbf{s}_{k+1}^T (\mathbf{b} - A\mathbf{x}_k) + \frac{1}{2\sqrt{k-1}} \|\mathbf{s}_{k+1}\|^2 + \mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k-1}} \|\mathbf{s}_k\|^2, \\ \epsilon''_{k+1} &:= \left( \frac{1}{2\sqrt{k-1}} - \frac{1}{2\sqrt{k}} \right) \|\mathbf{s}_k\|^2, \\ \epsilon'''_{k+1} &:= \left( \frac{\sqrt{k+1}}{2k} - \frac{1}{2\sqrt{k-1}} \right) \|\mathbf{s}_{k+1}\|^2. \end{aligned}$$

We see that  $\epsilon'_{k+1} \geq 0$  because  $\mathbf{s}_k$  is the maximizer of  $\mathbf{s}^T(\mathbf{b} - A\mathbf{x}_k) - \|\mathbf{s}\|^2/(2\sqrt{k-1})$  according to the definition (17), while  $\mathbf{s}_{k+1}$  is some other feasible point, and  $\epsilon'_{k+1}$  is the difference between the two objective values. We also see that  $\epsilon''_{k+1} \geq 0$ .

As for  $\epsilon'''_{k+1}$ , use the estimate  $\|\mathbf{s}_k\| \leq \sqrt{m}\eta$  according to (17). Finally, by taking a common denominator and observing that  $(k-1/(3k))^2 \leq k^2-1$ , we obtain a lower bound of  $-1/(6k^2\sqrt{k-1})$  on the parenthesized factor. Thus, adding the three contributions,

$$\epsilon_{k+1} \geq -\frac{m\eta^2}{6k^2\sqrt{k-1}} \cdot \frac{k}{k+1} \geq -\frac{m\eta^2}{6k^2\sqrt{k-1}}. \quad (32)$$

□

**Lemma 2.** Recall  $\delta_{k+1}$  defined in (19). We have the bound

$$\delta_{k+1} \geq -D/k^{3/2},$$

where  $D$  is a positive constant depending on the data defined in (35) below.

*Proof.* Split  $\delta_{k+1}$  into four terms:

$$\delta_{k+1} = \delta'_{k+1} + \delta''_{k+1} + \delta'''_{k+1} + \delta^{\text{iv}}_{k+1}$$

where

$$\begin{aligned} \delta'_{k+1} &:= \mathbf{r}_{k+2}^T(\mathbf{c} - A^T\mathbf{y}_k) + \frac{1}{2\sqrt{k}}\|\mathbf{r}_{k+2}\|^2 - \mathbf{r}_{k+1}^T(\mathbf{c} - A^T\mathbf{y}_k) - \frac{1}{2\sqrt{k}}\|\mathbf{r}_{k+1}\|^2, \\ \delta''_{k+1} &:= (\mathbf{r}_{k+2} - \mathbf{r}_{k+1})^T A^T(\mathbf{y}_k - \mathbf{y}_{k+1}), \\ \delta'''_{k+1} &:= \left( \frac{1}{2\sqrt{k+1}} - \frac{1}{2\sqrt{k}} \right) \|\mathbf{r}_{k+2}\|^2, \\ \delta^{\text{iv}}_{k+1} &:= \left( \frac{1}{2\sqrt{k}} - \frac{k}{2(k+1)\sqrt{k}} \right) \|\mathbf{r}_{k+1}\|^2. \end{aligned}$$

We observe that  $\delta'_{k+1} \geq 0$  because  $\mathbf{r}_{k+1}$  is the minimizer of  $\mathbf{r}^T(\mathbf{c} - A^T\mathbf{y}_k) + \|\mathbf{r}\|^2/(2\sqrt{k})$  according to (15), whereas  $\mathbf{r}_{k+2}$  is some other feasible point.

Next, we turn to  $\delta''_{k+1}$ , a product of three factors. Starting on the first factor,

$$\begin{aligned} \|\mathbf{r}_{k+2} - \mathbf{r}_{k+1}\| &= \|\text{proj}_\Delta(\sqrt{k+1}(A^T\mathbf{y}_{k+1} - \mathbf{c})) - \text{proj}_\Delta(\sqrt{k}(A^T\mathbf{y}_k - \mathbf{c}))\| \\ &\leq \|\sqrt{k+1}(A^T\mathbf{y}_{k+1} - \mathbf{c}) - \sqrt{k}(A^T\mathbf{y}_k - \mathbf{c})\| \\ &= \|(\sqrt{k+1} - \sqrt{k})(A^T\mathbf{y}_{k+1} - \mathbf{c}) + \sqrt{k}(A^T\mathbf{y}_{k+1} - \mathbf{c} - (A^T\mathbf{y}_k - \mathbf{c}))\| \\ &= \|(\sqrt{k+1} - \sqrt{k})(A^T\mathbf{y}_{k+1} - \mathbf{c}) + \sqrt{k}A^T(\mathbf{y}_{k+1} - \mathbf{y}_k)\| \\ &\leq |\sqrt{k+1} - \sqrt{k}| \cdot \|A^T\mathbf{y}_{k+1} - \mathbf{c}\| + \sqrt{k}\|A\| \cdot \|\mathbf{y}_{k+1} - \mathbf{y}_k\|. \end{aligned} \quad (33)$$

Here, the first line follows from (13), the second (that is, (33)) from the fact that the Lipschitz constant of  $\text{proj}_C(\cdot)$  is 1 for any closed nonempty convex set  $C$ , the third line

adds and subtracts the same term, and the last line applies the triangle inequality and submultiplicativity.

Now we use the facts that  $\sqrt{k+1} - \sqrt{k} \leq 1/\sqrt{k}$ , and, from (18),

$$\mathbf{y}_{k+1} - \mathbf{y}_k = \frac{1}{k+1}(\mathbf{s}_{k+1} - \mathbf{y}_k),$$

and finally, the bounds  $\|\mathbf{y}_k\| \leq \sqrt{m}\eta$ ,  $\|\mathbf{s}_k\| \leq \sqrt{m}\eta$  to conclude that  $\|\mathbf{y}_{k+1} - \mathbf{y}_k\| \leq 2\sqrt{m}\eta/(k+1)$  and thus

$$\|\mathbf{r}_{k+2} - \mathbf{r}_{k+1}\| \leq (1/\sqrt{k}) \cdot (\|A\|\sqrt{m}\eta + \|\mathbf{c}\|) + (2/\sqrt{k})\|A\|\sqrt{m}\eta.$$

This takes care of the first factor in  $\delta''_{k+1}$ . The middle factor is bound by  $\|A\|$ , and the third factor  $\|\mathbf{y}_{k+1} - \mathbf{y}_k\|$  is bounded by (see the previous paragraph)  $2\sqrt{m}\eta/(k+1)$ . Thus, overall, we obtain

$$\begin{aligned} \delta''_{k+1} &\geq -\|\mathbf{r}_{k+2} - \mathbf{r}_{k+1}\| \cdot \|A\| \cdot \|\mathbf{y}_k - \mathbf{y}_{k+1}\| \\ &\geq -C/k^{3/2}, \end{aligned}$$

where  $C$  depends on the problem data:

$$C \leq 2\|A\|\sqrt{m}\eta(3\|A\|\sqrt{m}\eta + \|\mathbf{c}\|).$$

For  $\delta'''_{k+1}$ , by finding a common denominator and then multiplying the resulting fraction by  $\sqrt{k} + \sqrt{k+1}$ , we obtain

$$\delta'''_{k+1} \geq -1/(4k^{3/2}) \cdot \|\mathbf{r}_{k+2}\|^2 \geq -1/(4k^{3/2}) \cdot \xi^2.$$

Finally, one sees that  $\delta_{k+1}^{\text{iv}} \geq 0$ . Putting all of these terms together yields:

$$\delta_{k+1} \geq -D/k^{3/2}, \tag{34}$$

where

$$D := 2\|A\|\sqrt{m}\eta(3\|A\|\sqrt{m}\eta + \|\mathbf{c}\|) + \frac{\xi^2}{4}. \tag{35}$$

□

Combining (30), Lemma 1, and Lemma 2, we now have for all  $k \geq 2$

$$U_{k+1} \leq \frac{k}{k+1}U_k + \frac{D}{k^{3/2}} + \frac{m\eta^2}{6k^2\sqrt{k-1}}.$$

In fact, since  $\sqrt{k} \geq 1/\sqrt{k-1}$  for all  $k \geq 2$ , we can enlarge  $D$  to obtain the bound

$$U_{k+1} \leq \frac{k}{k+1}U_k + \frac{\bar{D}}{k^{3/2}}, \tag{36}$$

where

$$\bar{D} := D + \frac{m\eta^2}{6}.$$

With this bound, we prove the following bound on  $U_k$ .

**Theorem 1.** Recall  $U_k$  defined in (31). We have the bound

$$U_{k+1} \leq \frac{F}{\sqrt{k}},$$

for all  $k \geq 2$ , where  $F := \max\{\sqrt{2}U_2, 6\bar{D}\}$ .

*Proof.* The  $k = 2$  case follows immediately from the definition of  $F$ . Suppose inductively that  $U_k \leq F/\sqrt{k}$ . Then we check:

$$\begin{aligned} U_{k+1} &\leq \frac{k}{k+1}U_k + \bar{D}/k^{3/2} \\ &\leq \frac{kF}{(k+1)\sqrt{k}} + \bar{D}/k^{3/2} \\ &= \frac{k^2F + \bar{D}(k+1)}{(k+1)k^{3/2}} \\ &= \frac{1}{\sqrt{k+1}} \cdot \frac{k^2F + \bar{D}(k+1)}{\sqrt{k+1} \cdot k^{3/2}} \\ &\leq \frac{1}{\sqrt{k+1}} \cdot \frac{k^2F + \bar{D}(k+1)}{(\sqrt{k+1}/(3\sqrt{k}))k^{3/2}} \\ &= \frac{1}{\sqrt{k+1}} \cdot \frac{k^2F + \bar{D}(k+1)}{k^2 + k/3} \\ &\leq \frac{1}{\sqrt{k+1}} \cdot \frac{k^2F + kF/3}{k^2 + k/3} \\ &= \frac{1}{\sqrt{k+1}} \cdot F. \end{aligned}$$

Here, we used (36) for the first line, the induction hypothesis for the second line, the inequality  $\sqrt{k+1} \geq \sqrt{k+1}/(3\sqrt{k})$  for  $k \geq 1$  on the 5th line (easy to confirm by squaring) and the assumption  $\bar{D} \leq F/6$  on the 7th line, which implies for all  $k \geq 1$  that  $\bar{D}(k+1) \leq Fk/3$ .

Thus,  $U_k \leq F/\sqrt{k}$ . □

The perturbation terms from (15) and (17) in FWLP-P were essential to this analysis. Without them, as in FWLP, there is no useful bound on how much the primal step  $\mathbf{r}_{k+2}$  can differ from  $\mathbf{r}_{k+1}$  when the index of the most-violated constraint changes from one iteration to the next. With the perturbations, we are able to obtain the inequality (33).

### 3.3 Proving convergence of the potential function implies convergence of the iterates

We have shown that the potential function  $U_k$  decreases at a rate of  $\mathcal{O}(1/\sqrt{k})$  as  $k$  increases. It remains to show that the potential function bounds the distance of the iterates from an optimizer. We first need the following lemmas. Both lemmas have the

same flavor: if an infeasible point for an LP improves on the optimal objective value, then the amount of objective improvement is bounded in terms of the amount of infeasibility.

**Lemma 3.** *Let  $p^*$  denote the optimal value of (1). Assume  $p^*$  is finite, i.e., (1) is feasible and bounded. Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be an arbitrary optimizing primal-dual pair of (1) and (2). Then for an arbitrary  $\hat{\mathbf{x}} \geq \mathbf{0}$ ,*

$$\mathbf{c}^T \hat{\mathbf{x}} \geq p^* - \|\mathbf{y}^*\|_\infty \|\mathbf{b} - A\hat{\mathbf{x}}\|_1. \quad (37)$$

*Proof.* Select an arbitrary optimizing pair  $(\mathbf{x}^*, \mathbf{y}^*)$  for (1). Let  $\mathbf{k} := \mathbf{b} - A\hat{\mathbf{x}}$ . Consider the LP,

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b} - \mathbf{k}, \\ & \mathbf{x} \geq \mathbf{0}. \end{aligned} \quad (38)$$

This LP is clearly feasible since  $\hat{\mathbf{x}}$  satisfies the constraints. It is also bounded, as argued by contradiction: If (38) were feasible and unbounded, then there would exist a certificate of unboundedness, which is a vector  $\mathbf{w}$  such that  $\mathbf{w} \geq \mathbf{0}$ ,  $A\mathbf{w} = \mathbf{0}$ , and  $\mathbf{c}^T \mathbf{w} < 0$ . However, such a  $\mathbf{w}$  would also certify unboundedness of (1), which we have already assumed to be bounded.

Therefore, the dual of (38), which is,

$$\begin{aligned} \max_{\mathbf{y}} \quad & (\mathbf{b} - \mathbf{k})^T \mathbf{y} \\ \text{s.t.} \quad & A^T \mathbf{y} \leq \mathbf{c}, \end{aligned} \quad (39)$$

has an optimal solution, say  $\hat{\mathbf{y}}$ . Since  $\mathbf{y}^*$  is also feasible for (39), we have the following chain of inequalities

$$\begin{aligned} p^* &= \mathbf{b}^T \mathbf{y}^* && \text{(by strong duality of (1))} \\ &= (\mathbf{b} - \mathbf{k})^T \mathbf{y}^* + \mathbf{k}^T \mathbf{y}^* \\ &\leq (\mathbf{b} - \mathbf{k})^T \mathbf{y}^* + \|\mathbf{k}\|_1 \cdot \|\mathbf{y}^*\|_\infty \\ &\leq (\mathbf{b} - \mathbf{k})^T \hat{\mathbf{y}} + \|\mathbf{k}\|_1 \cdot \|\mathbf{y}^*\|_\infty && \text{(since } \hat{\mathbf{y}} \text{ maximizes (39))} \\ &\leq \mathbf{c}^T \hat{\mathbf{x}} + \|\mathbf{k}\|_1 \cdot \|\mathbf{y}^*\|_\infty. && \text{(by weak duality between (38) and (39))} \end{aligned}$$

Recalling  $\mathbf{k} = \mathbf{b} - A\hat{\mathbf{x}}$ , the final line in this chain establishes (37).  $\square$

**Lemma 4.** *Let  $d^*$  denote the optimal value of (2). Assume  $d^*$  is finite, i.e., (2) is feasible and bounded. Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be an arbitrary optimizing primal-dual pair of (1) and (2). Then for an arbitrary  $\hat{\mathbf{y}}$ ,*

$$\mathbf{b}^T \hat{\mathbf{y}} \leq d^* + \|\mathbf{x}^*\|_1 \cdot l, \quad (40)$$

where  $l$  measures the infeasibility of  $\hat{\mathbf{y}}$ , that is,

$$l := \max(0, \max_j \{e_j^T (A^T \hat{\mathbf{y}} - \mathbf{c})\}). \quad (41)$$

*Proof.* This proof is analogous to the previous proof. Let  $(\mathbf{x}^*, \mathbf{y}^*)$  be an arbitrary primal-dual optimizer of (2). Consider the dual-form LP given by

$$\begin{aligned} \max_{\mathbf{y}} \quad & \mathbf{b}^T \mathbf{y} \\ \text{s.t.} \quad & A^T \mathbf{y} \leq \mathbf{c} + l\mathbf{e}, \end{aligned} \tag{42}$$

where  $l$  is from (41) and  $\mathbf{e}$  denotes the vector of all 1's. This LP is clearly feasible since  $\hat{\mathbf{y}}$  satisfies the constraint. Furthermore, it is bounded as argued by contradiction. If (42) were feasible and unbounded, there would exist a certificate of unboundedness, that is, a vector  $\mathbf{z}$  such that  $A^T \mathbf{z} \leq \mathbf{0}$  and  $\mathbf{b}^T \mathbf{z} > 0$ . However, this certificate would also certify the unboundedness of (2), which is assumed to be bounded.

Therefore, the dual of (42), which is

$$\begin{aligned} \min_{\mathbf{x}} \quad & (\mathbf{c} + l\mathbf{e})^T \mathbf{x} \\ \text{s.t.} \quad & A\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned} \tag{43}$$

has an optimal solution which we denote  $\hat{\mathbf{x}}$ . Since  $\mathbf{x}^*$  is also feasible for (43), we have the following chain of inequalities:

$$\begin{aligned} d^* &= \mathbf{c}^T \mathbf{x}^* && \text{(by strong duality of (2))} \\ &= (\mathbf{c} + l\mathbf{e})^T \mathbf{x}^* - l\mathbf{e}^T \mathbf{x}^* \\ &\geq (\mathbf{c} + l\mathbf{e})^T \mathbf{x}^* - l\|\mathbf{x}^*\|_1 \\ &\geq (\mathbf{c} + l\mathbf{e})^T \hat{\mathbf{x}} - l\|\mathbf{x}^*\|_1 && \text{(by the optimality of } \hat{\mathbf{x}} \text{ in (43))} \\ &\geq \mathbf{b}^T \hat{\mathbf{y}} - l\|\mathbf{x}^*\|_1 && \text{(by weak duality between (42) and (43)).} \end{aligned}$$

The final line after rearrangement is (40). □

We are now ready to prove the penultimate theorem that shows convergence of FWLP-P. We show that, assuming a primal-dual optimal LP solution exists and  $\xi, \eta$  have been chosen correctly, the  $U_k$  plus terms that tend to 0 bound the distance from optimality. Note that an analog of Theorem 2 can be derived for FWLP using a similar proof. This analog is not presented here since we are not able to show that  $U_k \rightarrow 0$  for FWLP.

**Theorem 2.** *Suppose  $k \geq 2$ . Then,*

$$\mathbf{x}_k \geq \mathbf{0}, \tag{44}$$

$$\|\mathbf{b} - A\mathbf{x}_k\|_1 \leq \frac{2U_k}{\eta} + \frac{\xi^2}{\eta\sqrt{k}} + \frac{\eta}{\sqrt{k-1}}, \tag{45}$$

$$\max(0, \max_j \{e_j^T (A^T \mathbf{y}_k - \mathbf{c})\}) \leq \frac{2U_k}{\xi} + \frac{\xi}{\sqrt{k}} + \frac{\eta^2}{\xi\sqrt{k-1}}, \tag{46}$$

$$\mathbf{c}^T \mathbf{x}_k - \mathbf{b}^T \mathbf{y}_k \leq U_k, \tag{47}$$

*provided*

$$\xi \geq 2\|\mathbf{x}^*\|_1, \quad \eta \geq 2\|\mathbf{y}^*\|_\infty. \tag{48}$$

*Proof.* It is immediate from the algorithm that (44) holds. Next, notice that  $\mathbf{s}_k$  is a solution to

$$\max_{\mathbf{s} \in \Gamma} \left\{ \mathbf{s}^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k-1}} \|\mathbf{s}\|^2 \right\},$$

so for arbitrary  $\mathbf{s} \in \Gamma$ ,

$$\begin{aligned} \mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 &\geq \mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k-1}} \|\mathbf{s}_k\|^2, \\ &\geq \mathbf{s}^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k-1}} \|\mathbf{s}\|^2. \end{aligned} \quad (49)$$

Setting  $\mathbf{s} = \mathbf{0}$  gives

$$\mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 \geq 0.$$

Similarly,  $\mathbf{r}_{k+1}$  is a solution to

$$\min_{\mathbf{r}} \left\{ \mathbf{r}^T (\mathbf{c} - A^T \mathbf{y}_k) + \frac{1}{2\sqrt{k}} \|\mathbf{r}\|^2 \right\}, \quad (50)$$

so for arbitrary  $\mathbf{r} \in \Delta$ ,

$$\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) + \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 \leq \mathbf{r}^T (\mathbf{c} - A^T \mathbf{y}_k) + \frac{1}{2\sqrt{k}} \|\mathbf{r}\|^2. \quad (51)$$

Setting  $\mathbf{r} = \mathbf{0}$  gives

$$-\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 \geq 0. \quad (52)$$

The above two results imply with (31) that  $U_k \geq \mathbf{c}^T \mathbf{x}_k - \mathbf{b}^T \mathbf{y}_k$ , thus establishing (47).

Setting  $\mathbf{s} := \eta \cdot \text{sgn}(\mathbf{b} - A\mathbf{x}_k) \in \Gamma$  in (49) gives the bound

$$\mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 \geq \eta \|\mathbf{b} - A\mathbf{x}_k\|_1 - \frac{1}{2\sqrt{k-1}} m\eta^2. \quad (53)$$

Similarly, set  $\mathbf{r} := \xi \mathbf{e}_j$  where  $j = \text{argmax}_j \{\mathbf{e}_j^T (A^T \mathbf{y}_k - \mathbf{c})\}$ . Noting  $\mathbf{r} \in \Delta$ , we can use (51) to obtain the bound

$$-\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 \geq \xi \cdot \max_j \{\mathbf{e}_j^T (A^T \mathbf{y}_k - \mathbf{c})\} - \frac{1}{2\sqrt{k}} \xi^2.$$

In fact, from (52) the left-hand side above is nonnegative. It follows that we can strengthen the above bound to

$$-\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 \geq \xi \cdot l - \frac{1}{2\sqrt{k}} \xi^2, \quad (54)$$

where  $l = \max(0, \max_j \{e_j^T (A^T \mathbf{y}_k - \mathbf{c})\})$ . Using Lemma 3, Lemma 4, (31), (53), and (54), we may write

$$\begin{aligned}
U_k &= -\mathbf{r}_{k+1}^T (\mathbf{c} - A^T \mathbf{y}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{r}_{k+1}\|^2 + \mathbf{s}_k^T (\mathbf{b} - A\mathbf{x}_k) - \frac{1}{2\sqrt{k}} \|\mathbf{s}_k\|^2 + \mathbf{c}^T \mathbf{x}_k - \mathbf{b}^T \mathbf{y}_k, \\
&\geq \xi \cdot l - \frac{1}{2\sqrt{k}} \xi^2 + \eta \|\mathbf{b} - A\mathbf{x}_k\|_1 - \frac{1}{2\sqrt{k}-1} m\eta^2 + \mathbf{c}^T \mathbf{x}_k - \mathbf{b}^T \mathbf{y}_k, \\
&\geq \xi \cdot l - \frac{1}{2\sqrt{k}} \xi^2 + \eta \|\mathbf{b} - A\mathbf{x}_k\|_1 - \frac{1}{2\sqrt{k}-1} m\eta^2 + p^* - \|\mathbf{y}^*\|_\infty \|\mathbf{b} - A\mathbf{x}_k\|_1 - d^* - \|\mathbf{x}^*\|_1 \cdot l, \\
&\geq \xi \cdot l - \frac{1}{2\sqrt{k}} \xi^2 + \eta \|\mathbf{b} - A\mathbf{x}_k\|_1 - \frac{1}{2\sqrt{k}-1} m\eta^2 + p^* - d^* - \frac{\eta}{2} \|\mathbf{b} - A\mathbf{x}_k\|_1 - \frac{\xi}{2} \cdot l \\
&= \frac{\xi}{2} \cdot l - \frac{1}{2\sqrt{k}} \xi^2 + \frac{\eta}{2} \|\mathbf{b} - A\mathbf{x}_k\|_1 - \frac{1}{2\sqrt{k}-1} m\eta^2,
\end{aligned}$$

where the fourth line used the assumption (48) and the final line used the fact that  $p^* = d^*$ . We are left with the bound

$$\frac{\xi}{2} \cdot l + \frac{\eta}{2} \|\mathbf{b} - A\mathbf{x}_k\|_1 \leq U_k + \frac{\xi^2}{2\sqrt{k}} + \frac{m\eta^2}{2\sqrt{k}-1},$$

where both terms on the left are nonnegative. It follows that each term must be individually bounded by the right-hand side above. This establishes (45) and (46).  $\square$

**Theorem 3.** *The iterates of FWLP-P converge to an  $\epsilon$ -optimal solution of (1) and its dual (2) after  $\mathcal{O}(1/\epsilon^2)$  iterations.*

*Proof.* This follows immediately by applying the bound from Theorem 1 in Theorem 2.  $\square$

## 4 Relating $U_k$ to the standard primal-dual gap

Consider the general saddle-point problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}).$$

A standard measure of the optimality gap of an iteration  $(\mathbf{x}_k, \mathbf{y}_k)$  used in the analysis of many primal-dual algorithms for saddle-point problems (for example, [6, 9, 15]) is the primal-dual gap:

$$\max \{f(\mathbf{x}_k, \hat{\mathbf{y}}_k) - f(\hat{\mathbf{x}}_k, \mathbf{y}_k) : \hat{\mathbf{x}}_k \in \mathcal{X}, \hat{\mathbf{y}}_k \in \mathcal{Y}\}. \quad (55)$$

Obviously, the solution to (55) is

$$\begin{aligned}
\hat{\mathbf{x}}_k &= \operatorname{argmin}\{f(\mathbf{x}, \mathbf{y}_k) : \mathbf{x} \in \mathcal{X}\}, \\
\hat{\mathbf{y}}_k &= \operatorname{argmax}\{f(\mathbf{x}_k, \mathbf{y}) : \mathbf{y} \in \mathcal{Y}\}.
\end{aligned}$$

It is noted in [2] that for the saddle-point problem associated with LP

$$\min_{\mathbf{x} \geq \mathbf{0}} \max_{\mathbf{y} \in \mathbb{R}^m} \mathbf{c}^T \mathbf{x} + \mathbf{y}^T (\mathbf{b} - A\mathbf{x}),$$

the primal-dual gap can be infinite, since the feasible set  $\mathbb{R}_+^n \times \mathbb{R}^m$  is unbounded. This is not an issue for the modified saddle-point formulation (3) used in our analysis, since the feasible set  $\Delta \times \Gamma$  is bounded by virtue of the redundant constraints.

Recall the definition of  $\mathcal{L}(\mathbf{x}, \mathbf{y})$  from (3). Let

$$\mathcal{M}_k := \max\{\mathcal{L}(\mathbf{x}_k, \mathbf{s}) - \mathcal{L}(\mathbf{r}, \mathbf{y}_k) : (\mathbf{r}, \mathbf{s}) \in \Delta \times \Gamma\},$$

which is the specialization of (55) to our setting. We argue that  $\mathcal{M}_k$  is a perturbation of  $U_k$  via the following bound:

$$\begin{aligned} |\mathcal{M}_k - U_k| &= \left| \max_{\mathbf{s} \in \Gamma} [(\mathbf{b} - A\mathbf{x}_k)^T \mathbf{s} + \mathbf{c}^T \mathbf{x}_k] - \min_{\mathbf{r} \in \Delta} [(\mathbf{c} - A^T \mathbf{y}_k)^T \mathbf{r} + \mathbf{b}^T \mathbf{y}_k] - U_k \right| \\ &\leq \left| \max_{\mathbf{s} \in \Gamma} \left[ (\mathbf{b} - A\mathbf{x}_k)^T \mathbf{s} + \mathbf{c}^T \mathbf{x}_k - \frac{\|\mathbf{s}\|^2}{2\sqrt{k-1}} \right] - \min_{\mathbf{r} \in \Delta} \left[ (\mathbf{c} - A^T \mathbf{y}_k)^T \mathbf{r} + \mathbf{b}^T \mathbf{y}_k + \frac{\|\mathbf{r}\|^2}{2\sqrt{k}} \right] - U_k \right| \\ &\quad + \max_{(\mathbf{r}, \mathbf{s}) \in \Delta \times \Gamma} \left| \frac{\|\mathbf{s}\|^2}{2\sqrt{k-1}} + \frac{\|\mathbf{r}\|^2}{2\sqrt{k}} \right| \\ &\leq \left| \left[ (\mathbf{b} - A\mathbf{x}_k)^T \mathbf{s}_k + \mathbf{c}^T \mathbf{x}_k - \frac{\|\mathbf{s}_k\|^2}{2\sqrt{k-1}} \right] - \left[ (\mathbf{c} - A^T \mathbf{y}_k)^T \mathbf{r}_{k+1} + \mathbf{b}^T \mathbf{y}_k + \frac{\|\mathbf{r}_{k+1}\|^2}{2\sqrt{k}} \right] - U_k \right| \\ &\quad + \frac{m\eta^2}{2\sqrt{k-1}} + \frac{\xi^2}{2\sqrt{k}} \\ &\leq \frac{m\eta^2 + \xi^2}{2\sqrt{k-1}} + \mathcal{O}(k^{-3/2}). \end{aligned}$$

Here, the second line adds and subtracts the same terms and then applies the triangle inequality for  $|\cdot|$ . The third line uses the definitions of  $\mathbf{s}_k$  from (17) and  $\mathbf{r}_k$  from (15). The fourth line uses (31), noting that the terms all cancel out except for the difference between a denominator of  $2\sqrt{k-1}$  versus a denominator of  $2\sqrt{k}$ . This small remainder is written as  $\mathcal{O}(k^{-3/2})$  on the fourth line.

## 5 Efficient implementation of FWLP and FWLP-P

A major advantage of FWLP and FWLP-P is their low computational cost per iteration. Naïve implementations of Algorithms 2.1 and 2.2 have iteration cost bounded by the cost of a full matrix-vector product. We discuss below how this can be significantly improved with an efficient implementation.

### 5.1 FWLP

Consider, e.g., the iteration  $k = 1$ . Suppose we store  $A\mathbf{x}_1$ . In Algorithm 2.1 we can perform an extra step to update  $A\mathbf{x}_2$  by either computing

$$A\mathbf{x}_2 = \frac{k}{2} A\mathbf{x}_1,$$

or

$$A\mathbf{x}_2 = \frac{k}{2}A\mathbf{x}_1 + \frac{\xi}{2}A\mathbf{e}_i,$$

where  $i$  is the index of the most violated dual constraint, computed as  $\operatorname{argmin}\{c_i - \mathbf{e}_i^T A^T \mathbf{y}_1\}$ . Ignoring the cost of computing  $i$ , such an update runs in  $\mathcal{O}(m+n)$  operations per iteration (cost of updating  $\mathbf{x}_k$  and  $A\mathbf{x}_k$ ). Note that  $\mathcal{O}(m+n)$  can be reduced to  $\mathcal{O}(m)$  if we keep track of the product of scaling factors  $k/(k+1)$  of  $\mathbf{x}_{k+1}$  in a separate variable. In contrast, the naïve implementation of the primal update in FWLP, where one instead performs matrix-vector products, runs in  $\mathcal{O}(mn)$  operations per iteration.

The dual update for, e.g.,  $k=1$ , is given by

$$\mathbf{y}_2 = \frac{1}{2}\mathbf{y}_1 + \frac{\eta}{2}\operatorname{sgn}(\mathbf{b} - A\mathbf{x}_2),$$

which also runs in  $\mathcal{O}(m)$  since  $A\mathbf{x}_2$  has already been computed in the primal step. The naïve implementation again takes  $\mathcal{O}(mn)$  operations.

To improve the cost of computing the index  $i$ , first note that in Algorithm 2.1 we only need to compute  $\mathbf{e}_i^T A^T \mathbf{y}_k$  for indices  $i$  such that there is a possibility that  $i = \operatorname{argmin}_j\{c_j - \mathbf{e}_j^T A\mathbf{y}_k\}$ . This could be implemented by a data structure to store the indices  $[n]$  in some order so that only the possibly most infeasible indices need to be considered at each iteration. In more detail, suppose  $j_k$  indexes the most violated dual constraint on iteration  $k$ , i.e.,  $j_k = \operatorname{argmin}_j[\mathbf{c} - A^T \mathbf{y}_k]_j$ . Suppose  $j \in [n]$  is some other index. Based on the value of  $[\mathbf{c} - A^T \mathbf{y}_k]_j - [\mathbf{c} - A^T \mathbf{y}_k]_{j_k}$  and prior knowledge of the stepsize (which tends to 0 with  $k$ ), one knows in advance that constraint  $j$  could not be the most violated constraint prior to iteration  $k'$ , where  $k'$  is a computable index satisfying  $k' > k$ . Then constraint  $j$  does not even have to be considered by the algorithm on all iterations between  $k$  and  $k'$ .

Thus, the algorithm maintains some subset of constraints  $\mathcal{S}_k$  on iteration  $k$  that need to be examined for possibly being the most violated. The computation of  $A^T \mathbf{y}_k$  in the primal update thus runs in  $\mathcal{O}(|\mathcal{S}_k| \cdot m)$  plus the cost of updating the data structure to obtain  $\mathcal{S}_{k+1}$ . Naturally, the estimate  $\mathcal{O}(|\mathcal{S}_k| \cdot m)$  is further reduced if  $A$  is sparse.

The total cost of iteration  $k$  for this implementation of FWLP is thus

$$\mathcal{O}(n + |\mathcal{S}_k| \cdot m),$$

plus the cost of updating the proposed data structure and precomputation. For large problems, one would expect that  $|\mathcal{S}_k| \ll n$ , allowing significant speedup to be achieved by using the proposed efficient implementation of FWLP. And, as mentioned earlier, the “ $n$ ” term may be dropped with a careful implementation for updating  $\mathbf{x}_k$ , and the “ $|\mathcal{S}_k| \cdot m$ ” term is reduced in the presence of sparsity.

## 5.2 FWLP-P

The iterations of FWLP-P differ from FWLP in that the steps (15) and (17) are more computationally involved.

To solve (15), we equivalently consider the projection form (13), which amounts to solving the quadratic programming problem

$$\begin{aligned} \arg \min_{\mathbf{x}} \quad & \frac{1}{2} \|\mathbf{w}_0 - \mathbf{x}\|^2 \\ \text{s.t.} \quad & \mathbf{e}^T \mathbf{x} \leq \xi, \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

where  $\mathbf{w}_0 = \sqrt{k}(A^T \mathbf{y}_k - \mathbf{c})$ . Expanding, rescaling, and dropping constant terms gives the equivalent problem

$$\begin{aligned} \arg \min_{\mathbf{x}} \quad & -\mathbf{w}^T \mathbf{x} + \frac{1}{2} \|\mathbf{x}\|^2 \\ \text{s.t.} \quad & \mathbf{e}^T \mathbf{x} \leq 1, \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned} \tag{56}$$

where  $\mathbf{w} = \mathbf{w}_0/\xi$ . The KKT conditions of (56) are as follows

$$-\mathbf{w} + \mathbf{x} + \mu \mathbf{e} - \mathbf{z} = \mathbf{0}, \tag{57}$$

$$\mathbf{e}^T \mathbf{x} \leq 1, \tag{58}$$

$$\mu(\mathbf{e}^T \mathbf{x} - 1) = 0, \tag{59}$$

$$\mathbf{z}^T \mathbf{x} = 0, \tag{60}$$

$$\mathbf{x}, \mu, \mathbf{z} \geq \mathbf{0}. \tag{61}$$

We first prove the following lemma, which will become useful once we define our algorithm for solving (56).

**Lemma 5.** *Suppose  $(\mathbf{x}, \mu, \mathbf{z})$  form a KKT solution for (56) with  $\mathbf{e}^T \mathbf{x} = 1$ . There must exist an index  $j \in [n]$  such that  $w_j > \mu$ .*

*Furthermore, we have the following cases regardless of whether  $\mathbf{e}^T \mathbf{x} = 1$ .*

(i) *If  $w_j > \mu$  for some  $j \in [n]$ , then  $x_j > 0$  and  $z_j = 0$ .*

(ii) *If  $w_j < \mu$  for some  $j \in [n]$ , then  $x_j = 0$  and  $z_j > 0$ .*

(iii) *If  $w_j = \mu$  for some  $j \in [n]$ , then  $x_j = z_j = 0$ .*

*Proof.* By (57), we have

$$w_i = x_i - z_i + \mu \tag{62}$$

for all  $i \in [n]$ . The condition  $\mathbf{e}^T \mathbf{x} = 1$  along with the nonnegativity of  $\mathbf{x}$  from (61) implies that there must exist an index  $j \in [n]$  such that  $x_j > 0$ . Now, the complementarity condition (60) along with the nonnegativity of  $\mathbf{z}$  from (61) imply that  $z_j = 0$ . It follows from (62) that  $w_j = x_j + \mu > \mu$ .

Now for (i): suppose,  $w_j > \mu$  for some  $j \in [n]$ . By (62),  $x_j - z_j + \mu > \mu$ , which implies  $x_j > z_j$ . But since  $z_j \geq 0$  it follows that  $x_j > 0$ . By applying (60) we see that  $z_j = 0$ . The proof for (ii) is analogous.

For (iii) suppose,  $w_j = \mu$  for some  $j \in [n]$ . By (62),  $x_j - z_j + \mu = \mu$ , which implies  $x_j = z_j$ . By applying (60) we must have  $x_j = z_j = 0$ .  $\square$

We now state Algorithm 5.1 for solving (56). Note that the algorithm computes the KKT multipliers  $\mu, \mathbf{z}$  as well as  $\mathbf{x}$  in order to illustrate its correctness, but in the FWLP-P code, the computation of  $\mathbf{z}$  is not needed.

---

**Algorithm 5.1** An efficient algorithm for finding a KKT solution of (56).

---

**Require:** Linear term coefficient  $\mathbf{w} \in \mathbb{R}^n$ .

- 1: Sort  $\mathbf{w}$  and store result in  $\bar{\mathbf{w}}$ , along with the permutation function  $\sigma : [n] \rightarrow [n]$  which maps the indices of  $\bar{\mathbf{w}}$  back to their original positions in  $\mathbf{w}$ :

$$(\bar{\mathbf{w}}, \sigma) := \text{sort}(\mathbf{w}). \quad (63)$$

- 2: Compute the cumulative sum of  $\bar{\mathbf{w}}$ :

$$S := \text{cumsum}(\bar{\mathbf{w}}). \quad (64)$$

- 3: **for**  $j = 1, 2, \dots, n$  **do**

- 4:     Compute

$$\mu := \frac{S(j) - 1}{j} \quad (65)$$

- 5:     **if**  $\bar{w}_j \geq \mu$  and either  $j = n$ , or  $\bar{w}_{j+1} \leq \mu$  **then**

- 6:         Record index  $j$  in the variable  $J$ .

- 7:         **break**

- 8:     **end if**

- 9: **end for**

- 10: **if**  $\mu \geq 0$  **then**

- 11:     Construct a KKT solution such that  $\mathbf{e}^T \mathbf{x} = 1$ :

$$\begin{aligned} \mathbf{x} &:= \mathbf{0}, \\ \mathbf{x}(\sigma(1 : J)) &:= \bar{\mathbf{w}}(1 : J) - \mu \cdot \mathbf{e}(1 : J), \\ \mathbf{z} &:= \mu \mathbf{e} - \mathbf{w}, \\ \mathbf{z}(\sigma(1 : J)) &:= \mathbf{0}. \end{aligned} \quad (66)$$

- 12:     **return**

- 13: **end if**

- 14: Otherwise, construct a KKT solution such that  $\mathbf{e}^T \mathbf{x} < 1$ :

$$\begin{aligned} \mathbf{x} &:= \max(\mathbf{0}, \mathbf{w}), \\ \mu &:= 0, \\ \mathbf{z} &:= \mathbf{x} - \mathbf{w}. \end{aligned} \quad (67)$$

- 15: **return**
- 

We now prove that Algorithm 5.1 is finite and correct.

**Theorem 4.** *Algorithm 5.1 finds a KKT solution to (56) in finite time.*

*Proof.* It is clear that the algorithm runs in linear time except for the sorting.

Consider the case where  $(\mathbf{x}, \mu, \mathbf{z})$  satisfy the KKT conditions with  $\mathbf{e}^T \mathbf{x} = 1$ . Then by Lemma 5 there must exist an index  $j \in [n]$  such that  $w_j > \mu$ , so  $J$  must be well-defined on line 6 of the algorithm. Moreover, this means that for all  $j \in [J]$ ,  $\bar{w}_j \geq \mu_j$ , and by applying Lemma 5 we get that  $x_j \geq 0$  and  $z_j = 0$  for all  $j \in [J]$ . A similar argument tells us that  $\bar{w}_j < \mu_j$  and thus  $x_j = 0, z_j > 0$  for  $j \in J' := [n] \setminus [J]$ . Using that  $\mathbf{e}^T \mathbf{x} = 1$  and that the entries of  $\mathbf{x}$  must be zero outside of  $[J]$ , as well as  $z_j = 0$  for all  $j \in [J]$ , we may multiply (57) by  $\mathbf{e}_J^T$  and rewrite to obtain:

$$\mathbf{e}_J^T \mathbf{w} = \mathbf{e}_J^T \mathbf{x} - \mathbf{e}_J^T \mathbf{z} + J\mu = 1 + J\mu,$$

where  $\mathbf{e}_J$  is taken to be the vector with ones in entries  $[J]$  and zero everywhere else. Rearranging gives the formula for  $\mu$  used in Algorithm 5.1:

$$\mu = \frac{S(J) - 1}{J}$$

where we note that  $S(J) = \mathbf{e}_J^T \mathbf{w}$ . It follows from assumption of this case  $\mathbf{e}^T \mathbf{x} = 1$  that  $\mu \geq 0$  in line 10 of the algorithm. After equations (66), the algorithm has constructed an  $\mathbf{x}$  such that

$$\mathbf{x}_J = \mathbf{w}_J - \mu \mathbf{e}_J, \quad \mathbf{x}_{J'} = \mathbf{0},$$

and a  $\mathbf{z}$  such that

$$\mathbf{z}_J = \mathbf{0}, \quad \mathbf{z}_{J'} = \mu \mathbf{e}_{J'} - \mathbf{w}_{J'}.$$

It is now easy to check that the  $(\mathbf{x}, \mu, \mathbf{z})$  constructed by the algorithm satisfy the KKT conditions.

Now consider the case where  $(\mathbf{x}, \mu, \mathbf{z})$  satisfy the KKT conditions with  $\mathbf{e}^T \mathbf{x} < 1$ . From (59),  $\mu = 0$  and clearly the  $\mathbf{x}$  and  $\mathbf{z}$  defined in (67) satisfy the remaining KKT conditions.  $\square$

The cost of Algorithm 5.1 is bounded below by the cost of either the matrix-vector product  $A^T \mathbf{y}_k$  in computing  $\mathbf{w}$  or the sort (63), which take  $\mathcal{O}(mn)$  and  $\mathcal{O}(n \log(n))$  operations respectively. Meanwhile, (16) takes  $\mathcal{O}(n)$  operations, so overall the two run in  $\mathcal{O}(n + mn)$  operations if we assume  $m \geq \log(n)$ .

Algorithm 5.1 could be sped up by noticing that in order to construct  $\mathbf{x}$ , we need only knowledge of the most-violated dual constraints, i.e., those for which  $\mathbf{c} - A^T \mathbf{y}_k$  is most negative. Recall from Section 5.2 that the computation of  $A^T \mathbf{y}_k$  would take  $\mathcal{O}(|\mathcal{S}_k| \cdot m)$  operations, plus the cost of updating the data structure storing the indices  $\mathcal{S}_k$ . The sorting computation cost would also be decreased since only the largest entries would need sorting.

We now consider the dual updates (17) and (18). Step (17) can be solved by considering the projection form  $\text{proj}_\Gamma(\sqrt{k}(\mathbf{b} - A\mathbf{x}_{k+1}))$  and noting that this has a well known closed-form solution [5, Lemma 6.26]:

$$\mathbf{s}_{k+1}(i) = \max(-\eta, \min(\eta, \sqrt{k}[\mathbf{b} - A\mathbf{x}_{k+1}]_i)),$$

for each  $i \in [m]$ . The cost of computing  $A\mathbf{x}_{k+1}$  is again  $\mathcal{O}(|\mathcal{S}_k| \cdot m)$  because we can scale  $A\mathbf{x}_{k-1}$  and then add the update  $A\mathbf{r}_k/k$ . The number of nonzero entries in  $\mathbf{r}_k$  is  $|\mathcal{S}_{k+1}|$ , which is the number of entries that partake in the projection.

## 6 Conclusion

We proposed two primal-dual first-order algorithms, namely FWLP and FWLP-P, for solving linear programming problems and discussed how both algorithms can be implemented in such a way that significantly improves their efficiency, especially for large-scale problems. Our convergence analysis of FWLP-P shows that the algorithm converges to a primal-dual solution with error  $\mathcal{O}(1/\sqrt{k})$  after  $k$  iterations. Despite this, no convergence proof is known at this time for the simpler and faster algorithm FWLP, and analysis of this algorithm is a topic for future research.

Another interesting question is how FWLP and FWLP-P can cope with primal or dual infeasibility. It should be noted that our proof that  $U_k \rightarrow 0$  (Theorem 1) did not depend on feasibility nor on the correct choice of  $\xi$  and  $\eta$ , so monitoring  $U_k$  cannot diagnose these conditions. We remark that other first-order algorithms have a theory for infeasibility detection; see, e.g., [3] and [13]. A related question is how the algorithm can detect whether the two parameters  $\xi$  and  $\eta$  of FWLP and FWLP-P have been chosen correctly.

We showed that FWLP-P converges in the sense that the primal and dual infeasibility measures tend to 0, as does the duality gap. However, we did not prove convergence of the iterates. In the case that the LP has multiple optimizers, an open question is whether the algorithms converge to a particular optimizer.

## References

- [1] D. Applegate, M. Díaz, O. Hinder, H. Lu, M. Lubin, B. O’Donoghue, and W. Schudy. Practical large-scale linear programming using primal-dual hybrid gradient. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 20243–20257. Curran Associates, Inc., 2021.
- [2] D. Applegate, O. Hinder, H. Lu, and M. Lubin. Faster first-order primal-dual methods for linear programming using restarts and sharpness. *Mathematical Programming*, 201(1):133–184, Sep 2023.
- [3] David Applegate, Mateo Díaz, Haihao Lu, and Miles Lubin. Infeasibility detection with primal-dual hybrid gradient for large-scale linear programming. *SIAM Journal on Optimization*, 34(1):459–484, 2024.
- [4] K. Basu, A. Ghoting, R. Mazumder, and Y. Pan. ECLIPSE: An extreme-scale linear program solver for web-applications. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 704–714. PMLR, 13–18 Jul 2020.
- [5] Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.

- [6] Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, May 2011.
- [7] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3(1-2):95–110, 1956.
- [8] Robert M. Freund and Paul Grigas. New analysis and results for the Frank–Wolfe method. *Mathematical Programming*, 155(1):199–230, Jan 2016.
- [9] G. Gidel, T. Jebara, and S. Lacoste-Julien. Frank-Wolfe algorithms for saddle point problems. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.
- [10] J.H. Hammond. *Solving Asymmetric Variational Inequality Problems and Systems of Equations with Generalized Nonlinear Programming Algorithms*. PhD thesis, Massachusetts Institute of Technology, 1984.
- [11] Matthew Hough. Solving saddle point formulations of linear programs with Frank-Wolfe. Master’s thesis, University of Waterloo, 2023.
- [12] Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 427–435, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR.
- [13] Tao Jiang, Walaa M Moursi, and Stephen A Vavasis. Range of the displacement operator of PDHG with applications to quadratic and conic programming. *arXiv preprint arXiv:2309.15009*, 2023.
- [14] X. Li, D. Sun, and KC Toh. An asymptotically superlinearly convergent semismooth Newton augmented Lagrangian method for linear programming. *SIAM Journal on Optimization*, 30(3):2410–2440, 2020.
- [15] Arkadi Nemirovski. Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004.