

Rescaled proximal methods for linearly constrained convex problems

Paulo J. S. Silva*

Carlos Humes Jr.†

Instituto de Matemática e Estatística, Universidade de São Paulo

In honor of Clóvis Gonzaga 60th birthday

September 15, 2004

Abstract

We present an inexact interior point proximal method to solve linearly constrained convex problems. In fact, we derive a primal-dual algorithm to solve the KKT conditions of the optimization problem using a modified version of the rescaled proximal method. We also present a pure primal method.

The proposed proximal method has as distinctive feature the possibility of allowing inexact inner steps even for Linear Programming. This is achieved by using an error criterion that bounds the subgradient of the regularized function, instead of using ϵ -subgradients of the original objective function. Quadratic convergence for LP is also proved using a more stringent error criterion.

1 Introduction

The idea of proximal methods can be traced back to Martinet [9, 10] and Rockafellar [12]. In the context of optimization, the classical proximal method replaces a minimization problem by a sequence of better behaved problems with a quadratic regularization term added to the objective function.

Many generalizations of this classical proximal algorithm were proposed in the last years. One of the main objectives was to replace the squared Euclidean norm by coercive regularizations that were able to implicitly deal with simple constraints, giving raise to interior point proximal methods [4, 15, 5, 8, 7, 16, 3, 2, 14, 13]. This effort followed the success of interior point methods for optimization, particularly Linear Programming (LP). For a classical survey on interior point methods for LP see [6].

Initially, interior point proximal methods were formulated for box constrained problems. More recently these methods were generalized for linearly constrained problems [1, 7, 16, 19]. In particular, exact versions of the proximal algorithm were applied to solve Linear Programming (LP) problems resulting on quadratic convergent methods.

*Partially supported by CNPQ, grant 304691/2002-0, and PRONEX - Optimization.

†Partially supported by CNPQ, grant 307105/2003-2, and PRONEX - Optimization.

Our intent in this paper is to extend the rescaled proximal method from [14, 13] to linearly constrained optimization problems. This will be done in such a way that allows for inexact proximal subproblems even for LP, maintaining quadratic convergence. This objective is attained in four steps. First, we observe that the Karush-Kuhn-Tucker (KKT) conditions of a linearly constrained problem have a special structure that can be exploited by a specialized proximal method. Second, we present the detailed proximal method and prove its convergence. Next, we apply the method to linearly constrained problems. Finally, we show the convergence rate result.

2 Preliminaries

Consider the following convex problem constrained to a polyhedron in the standard form:

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & x \in \mathcal{F} \stackrel{\text{def}}{=} \{x \in \mathbb{R}^n \mid Ax = b, x \geq 0\}, \end{aligned} \tag{1}$$

where f is a proper lower semi-continuous convex function.

The associated KKT conditions are:

$$\begin{aligned} 0 &\in \partial f(x) - A'y + u \\ Ax &= b \\ x &\geq 0 \\ u &\leq 0 \\ \langle u, x \rangle &= 0. \end{aligned} \tag{2}$$

Or, using the variational inequalities notation,

$$0 \in \begin{bmatrix} \partial f(x) - A'y \\ Ax - b \end{bmatrix} + \begin{bmatrix} N_{\mathbb{R}_+^n}(x) \\ 0 \end{bmatrix}. \tag{3}$$

The variational inequality formulation opens up the opportunity to use a proximal method based only on box constraints. In this paper, we are specially interested on applying the Rescaling Proximal Method for Variational Inequalities (RPMVI) described in [14]. However, a naive, direct application of this idea may raise questions as:

1. It is not clear why it is necessary to add a regularization term for the dual variables, y , since they are unconstrained.
2. It may be difficult to prove that [14, Assumption 2.3.2] holds and therefore the convergence theorem from [14] may not apply. This assumption is known to hold only if the proximal method is applied to an optimization problem directly [14], instead of applying it to KKT conditions, or if the chosen regularizations can guarantee Fejér or quasi-Fejér convergence of the iterates [13].

Although there may be sensible answers to these points, we present below a detour that is able to avoid these possible issues.

3 The Modified Rescaled Proximal Method

Let us consider a variational problem inspired by (3),

$$0 \in T(x, y) + \begin{bmatrix} N_B(x) \\ 0 \end{bmatrix}, \quad (4)$$

where T is a (possibly set-valued) maximal monotone operator on \mathbb{R}^{n+m} and B is a box, that is $B \stackrel{\text{def}}{=} ([a_1, b_1] \times \dots \times [a_n, b_n]) \cap \mathbb{R}^n$, where $-\infty \leq a_i < b_i \leq +\infty$ for all $i = 1, \dots, n$.

As discussed above, we will present a variation of the RPMVI [14] that regularizes only the constrained variable x . Since we are interested in coercive proximal methods, we must assume from now on that:

Assumption 3.1

$$\text{dom } T \cap \text{int } B \times \mathbb{R}^m \neq \emptyset.$$

This condition also guarantees that the sum in (4) is a maximal monotone operator.

Before presenting the modified proximal algorithm, we introduce the class of coercive distances it employs:

Definition 3.2 For $i = 1, \dots, n$, a function $d_i : \mathbb{R} \times (a_i, b_i) \rightarrow (-\infty, \infty]$ is called a rescaled distance if it presents the following properties:

3.2.1. For all $z_i \in (a_i, b_i)$, $d_i(\cdot, z_i)$ is closed and strictly convex, with its minimum at z_i . Moreover, $\text{int dom } d_i(\cdot, z_i) = (a_i, b_i)$.

3.2.2. d_i is differentiable with respect to the first variable on $(a_i, b_i) \times (a_i, b_i)$, and this partial derivative is continuous at all points of the form $(x_i, z_i) \in (a_i, b_i) \times (a_i, b_i)$. Moreover, we will use the notation

$$d'_i(x_i, z_i) \stackrel{\text{def}}{=} \frac{\partial d_i}{\partial x_i}(x_i, z_i).$$

3.2.3. For all $z_i \in (a_i, b_i)$, $d_i(\cdot, z_i)$ is essentially smooth [11, Chapter 26].

3.2.4. There exist $L, \epsilon > 0$ such that if either $-\infty < a_i < z_i \leq x_i < a_i + \epsilon$ or $b_i - \epsilon < x_i \leq z_i < b_i < +\infty$, then $|d'_i(x_i, z_i)| \leq L|x_i - z_i|$.

This definition has appeared in [13] and is a slight generalization of [14, Assumption 2.1], as it does not explicitly require twice-differentiable distances. The relation between both assumptions becomes clear if we identify the above distances with the ones in [14, Assumption 2.1] divided by the convenient second derivative.

We can now present the modified proximal algorithm:

Modified Rescaling Proximal Method (MRPM)

1. **Initialization:** Let $k = 0$. Choose a scalar $\alpha > 0$, $\{\beta_k\}$ a positive real sequence converging to zero, and an initial iterate $(x^0, y^0) \in \text{int } B \times \mathbb{R}^m$. Finally, let D be a separable distance defined by

$$D(x, x^k) \stackrel{\text{def}}{=} \sum_{i=1}^n d_i(x_i, x_i^k), \quad (5)$$

where each d_i is a rescaled distance.

2. Iteration:

(a) Choose α^k such that each $\alpha_i^k \in [\underline{\alpha}, \infty)$, $i = 1, \dots, n$.

(b) Find x^{k+1} and e^{k+1} such that

$$\begin{bmatrix} e^{k+1} \\ 0 \end{bmatrix} \in T(x^{k+1}, y^{k+1}) + \begin{bmatrix} \text{diag}(\alpha^k)^{-1} \nabla_1 D(x^{k+1}, x^k) \\ 0 \end{bmatrix}, \quad (6)$$

where

$$|e_i^k| \leq \frac{1}{\alpha_i^{k-1}} |d'_i(x_i^k, x_i^{k-1})| + \beta_k. \quad (7)$$

(c) Let $k = k + 1$, and repeat the iteration. □

The algorithm above takes into account that a component of the original variational inequality may be an easily solvable equation. In such cases, it dismisses the need for a regularization in the respective variables. This structure appears in (3), where the lower part of the system is linear. We will further explore this fact in the next section.

However, let us first turn our attention to analyze the algorithm convergence. The first step is to adapt [14, Assumption 2.3.2] to our present context.

Assumption 3.3 *Define*

$$\gamma^k = \begin{bmatrix} \gamma_x^k \\ \gamma_y^k \end{bmatrix} \stackrel{\text{def}}{=} \begin{bmatrix} e^k - \text{diag}(\alpha^{k-1})^{-1} \nabla_1 D(x^k, x^{k-1}) \\ 0 \end{bmatrix}. \quad (8)$$

Let \bar{x} be an accumulation point of the x component in a sequence computed by the MRPM, i.e., there is an infinite set $\mathcal{K} \subseteq \mathbb{N}$ such that $x^k \rightarrow_{\mathcal{K}} \bar{x}$. Denote by $\gamma_{x_i}^k$ the i -th component of γ_x^k . Then, either $\gamma_{x_i}^k \rightarrow_{\mathcal{K}} 0$ or there is an infinite set $\mathcal{K}' \subseteq \mathcal{K}$ such that $x_i^{k-1} \rightarrow_{\mathcal{K}'} \bar{x}_i$.

Using the above definitions, we can now follow the convergence proof of the RPMVI from [14] to show the convergence of MRPM. We must now modify and prove [14, Lemmas 2.4 and 2.5, and Theorem 2.6]. This task is made particularly easy since the proof of the key result [14, Lemma 2.4] already considers each component of γ^k isolatedly.

Lemma 3.4 [14, Lemma 2.4] *Suppose that Assumption 3.3 hold. Let $\bar{x} \in \mathbb{R}^n$ be a limit point of $\{x^k\}$, i.e., $x^k \rightarrow_{\mathcal{K}} \bar{x}$ for some infinite set $\mathcal{K} \subseteq \mathbb{N}$. Then, for $i = 1, \dots, n$,*

$$\begin{aligned} \lim_{k \rightarrow_{\mathcal{K}} \infty} \gamma_{x_i}^k &= 0 && \text{if } \bar{x}_i \in (a_i, b_i) \\ \liminf_{k \rightarrow_{\mathcal{K}} \infty} \gamma_{x_i}^k &\geq 0 && \text{if } \bar{x}_i = a_i \\ \limsup_{k \rightarrow_{\mathcal{K}} \infty} \gamma_{x_i}^k &\leq 0 && \text{if } \bar{x}_i = b_i \\ \lim_{k \rightarrow_{\mathcal{K}} \infty} \gamma_y^k &= 0. \end{aligned} \quad (9)$$

Proof. The limit for $\{\gamma_y^k\}$ holds trivially, since this sequence is constant equal to 0. The component $\{\gamma_x^k\}$ is analyzed coordinate by coordinate following the the proof of [14, Lemma 2.4].

Indeed, let us consider the case $x_i \in (a_i, b_i)$. If we assume, for the sake of a contradiction that $\gamma_{x_i}^k \not\rightarrow 0$, Assumption 3.3 asserts that there is an infinite index set $\mathcal{K}' \subset \mathcal{K}$ and $\zeta > 0$ such that $|\gamma_{x_i}^k| > \zeta$ and $x_i^{k-1} \rightarrow_{\mathcal{K}'} \bar{x}_i$. Then,

$$\begin{aligned}
|\gamma_{x_i}^k| &= |e_i^k - 1/\alpha_i^k d'_i(x_i^k, x_i^{k-1})| \\
&\leq |e_i^k| + 1/\alpha_i^k |d'_i(x_i^k, x_i^{k-1})| \\
&\leq 2/\alpha_i^k |d'_i(x_i^k, x_i^{k-1})| + \beta_k && \text{[Error criterion]} \\
&\leq 2/\alpha_i^k |d'_i(x_i^k, x_i^{k-1})| + \beta_k && [\alpha_i^k > \alpha] \\
&\rightarrow_{\mathcal{K}'} d'_i(\bar{x}_i, \bar{x}_i) && \text{[Assumption 3.2.2]} \\
&= 0, && \text{[Assumption 3.2.1]}
\end{aligned}$$

a contradiction with $|\gamma_{x_i}^k| > \zeta > 0$.

Now, we consider $\bar{x}_i = a_i > \infty$ and assume that $\liminf_{k \rightarrow \infty} \gamma_{x_i}^k < 0$. Then, using Assumption 3.3, there must be an infinite index set $\mathcal{K}' \subset \mathcal{K}$ and $\zeta > 0$ such that $\gamma_{x_i}^k < -\zeta$ and $x_i^{k-1} \rightarrow_{\mathcal{K}'} \bar{x}_i$. Hence

$$\begin{aligned}
\zeta &\leq |\gamma_{x_i}^k| \\
&= |e_i^k - 1/\alpha_i^k d'_i(x_i^k, x_i^{k-1})| \\
&\leq |e_i^k| + 1/\alpha_i^k |d'_i(x_i^k, x_i^{k-1})| \\
&\leq 2/\alpha_i^k |d'_i(x_i^k, x_i^{k-1})| + \beta_k && \text{[Error criterion]} \\
&\leq 2/\alpha_i^k |d'_i(x_i^k, x_i^{k-1})| + \beta_k. && [\alpha_i^k > \alpha]
\end{aligned}$$

Now, let ϵ be as in Assumption 3.2.4, if there is an infinite index set $\mathcal{K}'' \subset \mathcal{K}'$ such that $a_i < x_i^{k-1} \leq x_i^k < a_i + \epsilon$, this assumption would imply that

$$\zeta \leq 2L/\alpha |x_i^k - x_i^{k-1}| + \beta_k \rightarrow 0, \quad \text{[Error criterion and Assumption 3.2.4]}$$

a contradiction with $\zeta > 0$. Hence, for sufficiently large $k \in \mathcal{K}'$, $x_i^k < x_i^{k-1}$.

Using Assumption 3.2.1 and $x_i^k < x_i^{k-1}$, we conclude that $d'_i(x_i^k, x_i^{k-1}) < 0$. Therefore,

$$\begin{aligned}
-\zeta &> \gamma_{x_i}^k \\
&= e_i^k - 1/\alpha_i^k d'_i(x_i^k, x_i^{k-1}) \\
&\geq 1/\alpha_i^k d'_i(x_i^k, x_i^{k-1}) - \beta_k - 1/\alpha_i^k d'_i(x_i^k, x_i^{k-1}) && \text{[Error criterion]} \\
&= \beta_k \rightarrow 0.
\end{aligned}$$

Once again, a contradiction with $\zeta > 0$.

The case $\bar{x}_i = b_i$, $b_i < \infty$ is analogous. □

Just as in the above lemma, the reasoning used to prove the boundedness of γ^k in [14, Lemma 2.5] still holds, if we remember that the y coordinate of γ^k is identically 0. We now have sufficient tools to prove the (subsequence) convergence of the MRPM.

Theorem 3.5 *Let $\{(x^k, y^k)\}$ be a sequence computed by the MRPM with Assumptions 3.1 and 3.3 holding, then all of its limit points are solutions to the variational inequality problem (4).*

Proof. The proof is analogous to the one in [14, Theorem 2.6]. It relies upon the boundedness of $\{\gamma^k\}$, discussed above, and upon the signal structure given by Lemma 3.4. □

4 Application to linearly constrained convex programming

Once the convergence of the MRPM is proved, we can analyze the original problem (1). We will derive a primal-dual interior proximal method to solve its KKT conditions based on the MRPM. Afterwards, we will show how to derive a pure primal method for a problem with special constraint structure.

4.1 A primal-dual method

In order to apply the results above to solve the linearly constrained problem (1), we define L as its Lagrangian considering only the equality constraints. Identifying $T(x, y)$ with $(\partial_x L(x, y), -\partial_y L(x, y))$, the KKT conditions for (1) are clearly a particular case of the variational inequality (4). Additionally, it is a well known result that such operator is maximal monotone [11, Corollary 37.5.2].

Applying the MRPM to find KKT pairs, each iteration needs to find (x^{k+1}, y^{k+1}) such that

$$\begin{aligned} e^{k+1} &\in \partial f(x^{k+1}) - A'y^{k+1} + \text{diag}(\alpha^k)^{-1} \nabla_1 D(x^{k+1}, x^k) \\ Ax^{k+1} &= b, \end{aligned} \quad (10)$$

where,

$$|e_i^{k+1}| \leq \frac{1}{\alpha_i^k} |d'_i(x_i^{k+1}, x_i^k)|. \quad (11)$$

Note that the error criterion above is equivalent to (7) will β_k identically zero. It is also possible to choose more permissive $\{\beta_k\}$ [14, Section 3], however it would clutter the proofs without presenting any particular advantage in the present context.

It is easy to recognize the above iteration as a variation of a classical primal-dual interior point method for convex programming under linear constraints[18]. However, the above algorithm presents the following important distinctions:

1. The barrier is replaced by the coercive generalized distance D .
2. The update of the barrier parameter is replaced by the update of the distance center, typical of proximal methods and, possibly, α^k .
3. The error criterion that must be achieved to change centers is now controlled by (11) instead of using a measure of the distance to the central path.

Let us verify the conditions required by Theorem 3.5, the convergence theorem for the MRPM. First, considering (3), it can be easily seen that Assumption 3.1 is implied by a Slater type constraint qualification, that is obligatory in interior point methods:

Assumption 4.1

$$\exists x \in \text{ri dom } f, \text{ such that } Ax = b, x > 0.$$

Since the constraint set is a polyhedron, this assumption also asserts that the KKT conditions are necessary and sufficient for optimality [11, Corollary 28.3.1].

On the other hand, Assumption 3.3 automatically holds for linearly constrained optimization problems. We prove this fact with the following three lemmas.

Lemma 4.2 For all feasible x and for all iterations of the MRPM,

$$f(x) \geq f(x^k) + \langle \gamma_x^k, x - x^k \rangle,$$

i.e., γ_x^k can be seen as a subgradient of the objective function among the feasible points.

Proof. Remembering (8), the definition of γ_x^k , we see that $\gamma_x^k \in \partial(f - \langle A'y^k, \cdot \rangle)(x^k)$. Therefore

$$\begin{aligned} f(x) &= f(x) + \langle y^k, b - Ax \rangle && [Ax = b] \\ &= f(x) - \langle A'y^k, x \rangle + \langle b, y^k \rangle \\ &\geq f(x^k) - \langle A'y^k, x^k \rangle + \langle \gamma_x^k, x - x^k \rangle + \langle b, y^k \rangle \\ &= f(x^k) + \langle y^k, b - Ax^k \rangle + \langle \gamma_x^k, x - x^k \rangle \\ &= f(x^k) + \langle \gamma_x^k, x - x^k \rangle. && [Ax^k = b] \end{aligned}$$

□

This allow us to prove the following extensions to [14, Lemma 3.4 and 3.6].

Lemma 4.3 In all iterations of the MRPM, $\gamma_{x_i}^k(x_i^{k-1} - x_i^k) \geq 0$.

Proof. We have

$$\begin{aligned} \gamma_{x_i}^k(x_i^{k-1} - x_i^k) &= \left(e_i^k - \frac{1}{\alpha_i^{k-1}} d'_i(x_i^k, x_i^{k-1}) \right) (x_i^{k-1} - x_i^k) \\ &\geq \underbrace{-\frac{1}{\alpha_i^{k-1}} d'_i(x_i^k, x_i^{k-1})(x_i^{k-1} - x_i^k) - |e_i^k| |x_i^{k-1} - x_i^k|}_{\geq 0} \\ &= \frac{1}{\alpha_i^{k-1}} (|d'_i(x_i^k, x_i^{k-1})| - \alpha_i^{k-1} |e_i^k|) |x_i^{k-1} - x_i^k| \\ &\geq \frac{1}{\alpha_i^{k-1}} |0| |x_i^{k-1} - x_i^k| \quad [\text{Due to (11)}] \\ &= 0. \end{aligned}$$

□

Lemma 4.4 $\{f(x^k)\}$ is non-increasing and hence convergent if f is bounded below on \mathcal{F} . In this case, we also have

$$|\gamma_{x_i}^k| |x_i^{k-1} - x_i^k| \rightarrow 0 \quad \forall i = 1, \dots, n.$$

Proof. From Lemma 4.2,

$$f(x^{k-1}) \geq f(x^k) + \langle \gamma_x^k, x^{k-1} - x^k \rangle = f(x^k) + \sum_{i=1}^n \gamma_i^k(x_i^{k-1} - x_i^k),$$

whence the result follows from the non-negativity of $\gamma_{x_i}^k(x_i^{k-1} - x_i^k)$.

□

Note that, if f is bounded below on \mathcal{F} , Assumption 3.3 is direct a consequence of the last lemma. Thus, Theorem 3.5 implies the optimality of all accumulation points of the sequence $\{(x^k, y^k)\}$ defined by (10)-(11). Actually, we can strength this result as:

Theorem 4.5 *Suppose that the problem (1) conforms Assumption 4.1 and that f is bounded below on the feasible set \mathcal{F} . Let $\{(x^k, y^k)\}$ be a sequence computed by (10)-(11). Then, if $\{x^k\}$ has a limit point, $\{f(x^k)\}$ converges to the optimal value and all its limit points will be minimizers of (1). A condition that ensures the boundedness of $\{x^k\}$ is the boundedness of the solution set, or any other level set of f restricted to \mathcal{F} .*

Proof. Let \bar{x} be a limit point of $\{x^k\}$, i.e. $x^k \rightarrow_{\mathcal{K}} \bar{x}$, for some infinite set $\mathcal{K} \subseteq \mathbb{N}$. Since $\{\gamma^k\}$ is bounded, as discussed before Theorem 3.5, we can assume without loss of generality, that there is $\bar{\gamma}_x \in \mathbb{R}^n$ such that $\gamma_x^k \rightarrow_{\mathcal{K}} \bar{\gamma}_x$.

Using Lemma 4.2 we have that

$$\forall x \in \mathcal{F}, f(x) \geq f(x^k) + \langle \gamma_x^k, x - x^k \rangle.$$

Taking limits, and remembering that f is l.s.c. we have:

$$\begin{aligned} \forall x \in \mathcal{F}, f(x) &\geq \lim_{k \rightarrow \infty} f(x^k) + \langle \bar{\gamma}_x, x - \bar{x} \rangle \\ &\geq f(\bar{x}), \end{aligned}$$

where the last inequality comes from Lemma 3.4. Therefore \bar{x} is a minimizer of f in \mathcal{F} . Letting $x = \bar{x}$ above, it also follows that, $f(x^k) \rightarrow f(\bar{x})$.

Finally, as Lemma 4.4 states that $\{f(x^k)\}$ is non-increasing, $\{x^k\}$ is included in a level set of f plus the indicator function of $\{x \geq 0 \mid Ax = b\}$. Since the boundedness of this set, or any other level set of the sum, is equivalent to the boundedness of the optimal solution set, the result follows. \square

We end this section addressing an important issue. The theorems that guarantee that the proximal iterations are well defined are usually based on the fact that all variables are being regularized. This is not the case in the MRPM. Therefore we must prove that the MRPM iterations can be done, at least in the linearly constrained optimization case:

Proposition 4.6 *Suppose that (1) conforms to Assumption 4.1. Let (\bar{x}, \bar{y}) be a KKT pair of*

$$\begin{aligned} \min \quad & f(x) + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i, x_i^k) \\ \text{s.t.} \quad & Ax = b. \end{aligned} \tag{12}$$

Then, (\bar{x}, \bar{y}) solves (10) with zero error. Moreover, such KKT pairs always exist if f is bounded below on \mathcal{F} .

Proof. If we write down the KKT conditions for (12) we immediately recognize (10).

Moreover, since the constraints are all affine and Assumption 4.1, asserts that there is a point in $\text{ri dom}(f(\cdot) + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(\cdot, x_i^k))$ that is feasible, KKT is necessary and sufficient for optimality [11, Corollary 28.3.1].

Now, let l be a lower bound of f on F . Since each d_i is only finite at the positive orthant for a give $\zeta \in \mathbb{R}$ the level set:

$$\begin{aligned} \left\{ x \in \mathbb{R}^n \mid Ax = b, f(x) + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i, x_i^k) \leq \zeta \right\} &= \\ &= \left\{ x \geq 0 \mid Ax = b, f(x) + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i, x_i^k) \leq \zeta \right\} \\ &\subset \left\{ x \geq 0 \mid Ax = b, \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i, x_i^k) \leq \zeta - l \right\}. \end{aligned}$$

The last level set is a level set of $\sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i, x_i^k)$ which is bounded as this function attains unique minimum at x^k by Assumption 3.2.1. Therefore (12) admits solutions and hence KKT pairs. \square

4.2 A pure primal method

Using the primal-dual method it is also possible derive a pure primal method for problem with a special constraint structure. Let us consider a problem in the form.

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & Ax \geq b \\ & x \geq 0 \end{aligned} \tag{13}$$

Adding slack variables, s , the above problem is equivalent to

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & Ax - s = b \\ & x \geq 0 \\ & s \geq 0 \end{aligned}$$

Which is constrained by a standard polyhedron and therefore is in the form described by (1). We may then apply the primal-dual method described in the last section. The respective iteration asks to solve approximately

$$0 \in \begin{bmatrix} \partial f(x^{k+1}) & -A'y^{k+1} & + \text{diag}(\alpha_x^k)^{-1} \nabla D(x^{k+1}, x^k) \\ 0 & +y^{k+1} & + \text{diag}(\alpha_s^k)^{-1} \nabla \tilde{D}(s^{k+1}, s^k) \\ & b - Ax^{k+1} + s^{k+1} & \end{bmatrix}.$$

where D regularizes the original variables and \tilde{D} the slack variables. It is interesting to note that the last two groups of equations can be forced to hold by defining $s^{k+1} \stackrel{\text{def}}{=} Ax^{k+1} - b$ and $y^{k+1} \stackrel{\text{def}}{=} \text{diag}(\alpha_s^k)^{-1} \nabla \tilde{D}(s^{k+1}, s^k)$. Substituting this definitions back in the first group, we need to solve approximately

$$0 \in \partial f(x^{k+1}) + \text{diag}(\alpha_s^k)^{-1} A' \nabla \tilde{D}(Ax^{k+1} - b, Ax^k - b) + \text{diag}(\alpha_x^k)^{-1} \nabla D(x^{k+1}, x^k).$$

This is equivalent o minimize approximately

$$f_k(x) \stackrel{\text{def}}{=} f(x) + \sum_{j=1}^m \frac{1}{(\alpha_s^k)_j} \tilde{d}_j((Ax - b)_j, (Ax^k - b)_j) + \sum_{i=1}^n \frac{1}{(\alpha_x^k)_i} d_i(x_i, x_i^k).$$

The acceptance criteria (11) turns out to be

$$|(\partial f_k(x^{k+1}))_i| \leq \frac{1}{(\alpha_x^k)_i} |d'_i(x_i^{k+1}, x_i^k)|. \quad (14)$$

The method above can be viewed as a variation of the algorithm proposed in [2, Section 4]. This is not the case as its error criterion is based on the subgradient of f and not in ϵ -subgradients. Error criteria based on ϵ -subgradients are particularly useful in situations where it is difficult to compute f or its subgradients. Good examples are multiplier methods or situations where f is non-smooth. On the other hand, if the subgradients of the objective function are readily available, a criterion based on the subdifferential itself may be preferable. This is clearly the case in Linear Programming, as the only ϵ -subgradient of f is actually its gradient. Under such conditions, a criterion as in [2, Section 4] would require exact solutions of the proximal subproblems, while (11) or (14) still allow for approximate steps.

Finally, it should be clear that any Linear Programming problem can be reduced to the form (13) using standard textbook techniques.

5 Rate of convergence for Linear Programming

The convergence rate of proximal methods for LP has been studied by many authors. Rockafellar proved in [12, Proposition 8] that a classical proximal algorithm, based on the squared Euclidean norm, converges in finite many steps. Tseng and Bertsekas showed in [17] that the exponential multiplier method with rescaling is super-linearly convergent. This last result was somewhat extended for general φ -divergences by Iusem and Teboulle [7, Section 4]. In [1, 16], Auslender, Haddou, and Teboulle have proved that methods based on a special φ -divergence may be quadratically convergent if the stepsizes converge fast enough to zero. In [2, Section 6], Auslender *et. al.* proved quadratic convergence if the regularizations used based on second order homogeneous kernels with an extra quadratic term. Finally, Silva and Eckstein [13] generalized the quadratic convergence for double regularizations, encompassing kernels based on Bregman distances. All these results required the exact solution of the proximal subproblems.

In this section, we demonstrate that the primal-dual method described by (10)-(11) is also quadratic convergent for Linear Programming even with inexact proximal steps. In order to do this, we will use two extra assumptions. The first one limits the class of generalized distances, imposing an upper bound on their value. The second strengthens the acceptance criterion (11), but still allows for inexact subproblems.

We assume from now on that problem (1) have the form

$$\begin{aligned} \min \quad & \langle c, x \rangle \\ \text{s.t.} \quad & Ax = b \\ & x \geq 0. \end{aligned} \quad (15)$$

Assumption 5.1 For $i = 1, \dots, n$, let d_i be a rescaled distance. There must be $C > 0$ such that for all $x_i, z_i > 0$,

$$d_i(x_i, y_i) \leq C(x_i - y_i)^2.$$

Two examples of distances that obey the above condition, with C equal 1 and 2 respectively, are the rescaled Bregman distances [14, Section 2.2.1] based on the kernels $x \log(x)$ and $x - \sqrt{x}$.

The second assumption aims to enforce that at each iteration the objective function decreases “enough”. As stated in Proposition 4.6, if (10) is solved exactly, x^{k+1} would minimize

$$f(x) + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i, x_i^k)$$

among all the feasible solutions. Hence, between iterations the objective function would decrease $\sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i^{k+1}, x_i^k)$. Assumption 5.2 requires that at least a fraction of this bound should be attained in any inexact iteration.

Assumption 5.2 Let $\beta \in [0, 1)$. For a given k , define $\epsilon_k(x) \stackrel{\text{def}}{=} \beta \sum_i 1/\alpha_i^k d_i(x_i, x_i^k)$. The pair (x^{k+1}, y^{k+1}) computed by the primal dual method (10)-(11) should also conform to

$$0 \in c - A'y^{k+1} + \text{diag}(\alpha_k)^{-1} \partial_\epsilon D(x^{k+1}, x^k), \quad \epsilon \leq \epsilon_k(x^{k+1}). \quad (16)$$

Note that in the above assumption the ϵ -subgradient operator is applied to the generalized distance, rather than to the objective function. Thus inexact proximal steps are allowed.

We proceed to prove quadratic convergence rate for Linear Programming. We begin stating an auxiliary lemma.

Lemma 5.3 Let x^* be any solution of (15). Let $\{(x^k, y^k)\}$ be a sequence computed by the primal-dual method (10)-(11) with Assumptions 4.1, 5.1, and 5.2 holding. Then,

$$\langle c, x^{k+1} \rangle + (1 - \beta) \sum_i \frac{1}{\alpha_i^k} d_i(x_i^{k+1}, x_i^k) \leq \langle c, x^* \rangle + \sum_i \frac{1}{\alpha_i^k} d_i(x_i^*, x_i^k)$$

Proof. Equation (16) states that x^{k+1} is at least an $\epsilon_k(x^{k+1})$ -minimum of $h(x) \stackrel{\text{def}}{=} \langle c, x \rangle - \langle A'y^{k+1}, x \rangle + \sum_i \frac{1}{\alpha_i^k} d_i(x, x_i^k)$, therefore:

$$\begin{aligned} \langle c, x^{k+1} \rangle - \langle A'y^{k+1}, x^{k+1} \rangle + \sum_i \frac{1}{\alpha_i^k} d_i(x^{k+1}, x_i^k) &\leq \\ \langle c, x^* \rangle - \langle A'y^{k+1}, x^* \rangle + \sum_i \frac{1}{\alpha_i^k} d_i(x_i^*, x_i^k) + \epsilon_k(x^{k+1}). \end{aligned}$$

As x^* and x^{k+1} are both feasible, $\langle A'y^{k+1}, x^* \rangle = \langle A'y^{k+1}, x^{k+1} \rangle$, and so

$$\begin{aligned} \langle c, x^{k+1} \rangle + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i^{k+1}, x_i^k) &\leq \\ &\leq \langle c, x^* \rangle + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i^*, x_i^k) + \epsilon_k(x^{k+1}) \\ &= \langle c, x^* \rangle + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i^*, x_i^k) + \beta \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i^{k+1}, x_i^k). \end{aligned}$$

□

We can now prove the following convergence rate result:

Theorem 5.4 Let $X^* \neq \emptyset$ be the solution set of (15) and f^* the optimal value. Let $\{(x^k, y^k)\}$ a sequence computed by the primal-dual method (10)-(11) with Assumptions 4.1, 5.1, and 5.2 holding. If $\{x^k\}$ has a limit point, then the distance of x^k to the solution set converges at least quadratically to zero and $\langle c, x^k \rangle$ converges at least quadratically to f^* .

Proof. Theorem 4.5 already proves convergence. We need to focus only on the convergence rate estimates. Let $d(x, X^*)$ denote the distance between a point $x \in \mathbb{R}^n$ and the solution set X^* .

Lemma 6.1 from [2] shows that there is a $\mu > 0$, depending only on the problem data, such that

$$d(x^k, X^*) \leq \mu(\langle c, x^k \rangle - f^*), \quad \forall k = 1, 2, \dots \quad (17)$$

Let \bar{x}^k be the point in X^* that realizes $d(x^k, X^*)$. Lemma 5.3 asserts that

$$\begin{aligned} \langle c, x^{k+1} \rangle &\leq \langle c, x^{k+1} \rangle + (1 - \beta) \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(x_i^{k+1}, x_i^k) \\ &\leq f^* + \sum_{i=1}^n \frac{1}{\alpha_i^k} d_i(\bar{x}_i^k, x_i^k) \\ &\leq f^* + \sum_{i=1}^n \frac{1}{\underline{\alpha}} d_i(\bar{x}_i^k, x_i^k) \end{aligned}$$

Using Assumption 5.1 and (17), it follows that:

$$\begin{aligned} \langle c, x^{k+1} \rangle - f^* &\leq \frac{C}{\underline{\alpha}} \sum_i (x_i^k - \bar{x}_i^k)^2 \\ &\leq \frac{\mu^2 C}{\underline{\alpha}} (\langle c, x^k \rangle - f^*)^2. \end{aligned} \quad (18)$$

Therefore $\langle c, x^k \rangle$ converges to f^* at least quadratically. Moreover, as for any $x \in X^*$ we have

$$\langle c, x^k \rangle - f^* = \langle c, x^k - x \rangle \leq \|c\|_2 \|x^k - x\|_2$$

We may then combine (17) and (18) to conclude that:

$$d(x^{k+1}, X^*) \leq \frac{\mu^3 C}{\underline{\alpha}} \|c\|_2^2 d(x^k, X^*)^2.$$

□

This new result on the quadratic convergence of an inexact proximal method induces future computational experiments. Particularly, the search for efficient practical procedures to ensure Assumption 5.2 should be investigated.

References

- [1] A. Auslender and M. Haddou. An interior proximal method for convex linearly constrained problems. *Mathematical Programming*, 71:77–100, 1995.

- [2] A. Auslender, M. Teboulle, and S. Ben-Tiba. Interior proximal and multiplier methods based on second order homogeneous kernels. *Mathematics of Operations Research*, 24:645–668, 1999.
- [3] A. Auslender, M. Teboulle, and S. Ben-Tiba. A logarithmic-quadratic proximal method for variational inequalities. *Computational Optimization and Applications*, 12:31–40, 1999.
- [4] Y. Censor and J Zenios. The proximal minimization algorithms with D-functions. *Journal of Optimization: Theory and Applications*, 73:451–464, 1992.
- [5] J. Eckstein. Nonlinear proximal point algorithms using Bregman functions, with applications to convex programming. *Mathematics of Operations Research*, 18:202–226, 1993.
- [6] C. C. Gonzaga. Path-following methods for linear programming. *SIAM Review*, 34(2):167–224, 1992.
- [7] A. N. Iusem and M. Teboulle. Convergence rate analysis of nonquadratic proximal methods for convex and linear programming. *Mathematics of Operations Research*, 20(3):657–677, 1995.
- [8] A. N. Iusem, M. Teboulle, and B. Svaiter. Entropy-like proximal methods in convex programming. *Mathematics of Operations Research*, 19(4):790–814, November 1994.
- [9] B. Martinet. Regularisation d’inéquations variationnelles par approximations successives. *Rev. Française Inf. Rech. Oper.*, pages 154–159, 1970.
- [10] B. Martinet. Determination approach d’un point fixe d’une application pseudo-contractante. *C.R. Acad. Sci. Paris*, 274A:163–165, 1972.
- [11] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [12] R. T Rockafellar. Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization*, 14:887–898, August 1976.
- [13] P. J. S. Silva and J. Eckstein. Double regularization proximal methods, with complementarity applications. Technical Report RR 29-03, Rutgers Center for Operations Research, 2003. submitted.
- [14] P. J. S. Silva, J. Eckstein, and C. Humes Jr. Rescaling and stepsize selection in proximal methods using generalized distances. *SIAM Journal on Optimization*, 12(1):238–261, 2001.
- [15] M. Teboulle. Entropic proximal methods with applications to nonlinear programming. *Mathematics of Operations Research*, 17:670–690, 1992.
- [16] M. Teboulle. Convergence of proximal-like algorithms. *SIAM Journal on Optimization*, 7(4):1069–1083, 1997.
- [17] P. Tseng and D. Bertsekas. On the convergence of the exponential multiplier method for convex programming. *Mathematical Programming*, 60:1–19, 1993.
- [18] Stephen J. Wright. *Primal-Dual Interior-Point Methods*. SIAM, 1997.
- [19] N. Yamashita, C. Kanzow, T. Morimoto, and M. Fukushima. An infeasible interior proximal method for convex programming problems with linear constraints. *Journal of Nonlinear and Convex Analysis*, 2(2):139–156, 2001.