

An Inexact Interior-Point Lagrangian Decomposition Algorithm with Inexact Oracles

Deyi Liu · Quoc Tran-Dinh*

Received: date / Accepted: date

Abstract We develop a new inexact interior-point Lagrangian decomposition method to solve a wide range class of constrained composite convex optimization problems. Our method relies on four techniques: Lagrangian dual decomposition, self-concordant barrier smoothing, path-following, and proximal-Newton technique. It also allows one to approximately compute the solution of the primal subproblems (called *the slave problems*), which leads to inexact oracles (i.e., inexact gradients and Hessians) of the smoothed dual problem (called *the master problem*). The smoothed dual problem is nonsmooth, we propose to use an inexact proximal-Newton method to solve it. By appropriately controlling the inexact computation at both levels: the slave and master problems, we still estimate a polynomial-time iteration-complexity of our algorithm as in standard short-step interior-point methods. We also provide a strategy to recover primal solutions and establish complexity to achieve an approximate primal solution. We illustrate our method through two numerical examples on well-known models with both synthetic and real data and compare it with some existing state-of-the-art methods.

Keywords Interior-point Lagrangian decomposition · barrier smoothing · inexact oracle · proximal-Newton method · constrained convex optimization

Mathematics Subject Classification (2000) 90C25 · 90-08

1 Introduction

The Lagrangian dual decomposition framework is a classical technique to handle constrained convex optimization problems with separable structures such as conic, multi-stage stochastic, network, and distributed optimization problems [3, 4, 9, 10, 25]. This approach has been incorporated with interior-point

Deyi Liu · Quoc Tran-Dinh

Department of Statistics and Operations Research, The University of North Carolina at Chapel Hill

333 Hanes Hall, UNC-Chapel Hill, NC27599

Emails: {deyi@live.unc.edu, quoctd@email.unc.edu}

*Corresponding author.

methods to obtain a dual decomposition interior-point framework in early 1990s [16]. Since then, many researchers have regularly applied this approach to different problems. For example, [13] exploited this idea to develop a dual decomposition algorithm for semidefinite programming, and [39] considered this method for general convex and multi-stage stochastic programming. The authors in [18] further investigated the method from [39] to solve a more general class of problems and obtained more intensive and rigorous complexity guarantees. The work [34] studied this framework under the effect of inexact oracle computed by inexactly solving the primal subproblems up to a given accuracy. Other related theoretical results include [5, 12, 14, 15, 24, 27, 37]. In particular, [24] solved loosely coupled problems using message passing, and [5] applied it to multi-agent optimization problems. However, none of these works has investigated general constrained composite convex optimization settings involving linear operators and allows both inexactness in the slave problems and master problem altogether. In addition, existing methods do not handle directly nonsmooth objectives but often introduce auxiliary variables to reformulate the underlying problem into a smooth problem which may significantly increase problem size and loose their theoretical guarantee.

Motivation and goals: Although the Lagrangian decomposition method is classical, it is very useful to handle large-scale constrained convex problems with separable structure by means of parallel and distributed computational architectures. In this paper, we conduct an intensive study on the interior-point Lagrangian decomposition (IPLD) framework considered in many existing works, especially [18, 34, 39], from the following aspects.

- (a) Firstly, we consider a more general problem class than [16, 18, 34, 39] by handling directly a nonsmooth composite convex function with a linear operator (see (2)) instead of couple linear equality constraints as in existing methods by means of proximal Newton-type methods (see Subsection 5.1).
- (b) Secondly, our method works with inexact oracles of the dual problem arising from inexactly solving the primal subproblems (the slave problems). We explicitly describe the range of accuracies to flexibly control the tolerance of the subproblems (see Subsection 4.2).
- (c) Thirdly, we also exploit inexact proximal-Newton method to handle general nonsmooth terms of the dual problems.
- (d) Fourthly, we provide a thorough analysis for both the primal and dual problems and derive concrete iteration-complexity bounds for our method.
- (e) Finally, we incorporate our approach with a recent concept called “generalized self-concordance” developed in [29] to handle new applications.

We are interested in the class of constrained composite convex problems where g is smooth and satisfies some additional properties so that existing methods often do not have a theoretical convergence guarantee. For instance, the objective function does not have Lipschitz gradient or is not “tractably proximal”. We also consider a generic convex set where the projection onto it may not be tractable to compute such as general polyhedra. Under such assumptions, our problem setting covers a wide range class of applications ranging from optimal control, operations research, and networks to machine learning, statistics,

and signal processing [2, 7]. It also covers standard conic programming such as linear programming, second-order cone programming, and semidefinite programming.

Our contribution: We exploit the approach from [16, 18, 34, 39] to develop a new algorithm for solving a class of constrained convex optimization problems. The main idea is to smooth the dual problem using a self-concordant barrier function [22] associated with the constraint set, and apply a path-following scheme to solve the smoothed dual problem. While [16, 18, 39] exactly follow this main stream, [34] proposed another path-following scheme and analyzed its convergence under inexact computation. It also provides a strategy to recover an approximate primal solution from its approximate dual solution. Compared to [34], this work studies a much more general problem class than [34]. In addition, it is different from existing works, including [34], in several aspects as previously mentioned. To this end, we can summarize our contribution as follows:

- (a) We exploit the approach in [16, 18, 34, 39] and combine it with recent new mathematical tools in [23, 29] to develop a new algorithm. The new mathematical tools allow us to cover much broader class of models than [16, 18, 34, 39], and to analyze polynomial-time iteration-complexity. In addition, we handle a more general class of problems than [16, 18, 34, 39] by allowing general composite convex objectives involving linear operators (see (2)).
- (b) We propose a new inexact interior-point Lagrangian decomposition algorithm to solve this class of problems. Our algorithm can deal with inexact oracles of the dual problems arising from approximating the primal subproblem solutions. It also uses an inexact proximal-Newton scheme to approximate the search direction in the dual problem. We characterize explicitly the choice of all related parameters and accuracies based on our analysis.
- (c) We establish a polynomial-time iteration-complexity estimate of our method to find an approximate optimal solution. Our algorithm can be viewed as a short-step interior-point methods for general convex problems involving Nesterov and Nemirovskii's self-concordance structures. Our complexity bound is the same as in standard interior-point methods (up to a constant factor), while it is able to directly handle nonsmooth objective by means of proximal operator.

In addition to the above main contribution, let us highlight some technical contribution of our methods. Firstly, unlike other methods involving inexact oracles in the literature [11], our inexact oracle is rendered from inexact solution of the subproblem. The accuracy level can be adaptively chosen instead fixing as in existing methods to flexibly trade-off the computation cost by choosing rough accuracy at the early iterations and decrease it in the last iterations. Secondly, solving the primal subproblem (slave problem) is reduced to solve a nonlinear equation instead of a general convex problem as in some existing decomposition methods. As a result, we can characterize an implementable criterion to control the inexactness of the primal subproblems by using Newton-type schemes. Thirdly, instead of using unspecified parameters

such as the radius of quadratic convergence region and contraction factor, we compute these parameters explicitly using the theory of self-concordant barriers as often seen in interior-point methods [19, 22]. Finally, combining inexact oracle and inexact methods make our algorithm practical since this computation is unavoidable in iterative methods, especially, in decomposition approaches when handling complex models.

Paper organization: The rest of this paper is organized as follows. Section 2 states the problem of interest, basic assumptions, and its dual form. Section 3 recalls some preliminary results on (generalized) self-concordance and self-concordant barriers [22]. Section 4 focuses on barrier smoothing techniques and inexact oracles. Section 5 presents our main algorithm and its complexity analysis as well as convergence guarantees. Section 6 provides two numerical examples to verify the theoretical results. For the sake of presentation, we move all the technical proofs to the appendix.

2 Problem statement, basic assumptions, and dual formulation

Notation and terminologies: We work with finite dimensional vector space \mathbb{R}^p or \mathbb{R}^n endowed with standard inner product $x^\top y$ or $\langle x, y \rangle$ and Euclidean norm $\|x\|_2 := \sqrt{x^\top x}$. We denote by \mathbb{S}_+^p (resp., \mathbb{S}_{++}^p) the set of symmetric positive semidefinite matrices (resp., symmetric positive definite matrices). Given $H \in \mathbb{S}_{++}^p$, we define a weighted norm $\|u\|_H := (u^\top H u)^{1/2}$ and its dual norm $\|v\|_H^* := (v^\top H^{-1} v)^{1/2}$ for any vectors u and v in \mathbb{R}^p . For $X, Y \in \mathbb{S}_+^p$, $X \preceq Y$ means that $Y - X \in \mathbb{S}_+^p$ and $X \succeq Y$ stands for $X - Y \in \mathbb{S}_+^p$.

Given a three-time differentiable and strictly convex function f , we define the following local norms for any u and v in \mathbb{R}^p :

$$\|u\|_x := (u^\top \nabla^2 f(x) u)^{1/2}, \quad \text{and} \quad \|v\|_x^* := (v^\top \nabla^2 f(x)^{-1} v)^{1/2}. \quad (1)$$

They also satisfy the Cauchy-Schwarz inequality, i.e. $u^\top v \leq \|u\|_x \|v\|_x^*$. We say that f is μ_f -strongly convex if $f(\cdot) - (\mu_f/2)\|\cdot\|^2$ remains convex. We also often use the following two convex functions: $\omega(\tau) := \tau - \ln(1 + \tau)$ for $\tau \geq 0$, and $\omega_*(\tau) := -\tau - \ln(1 - \tau)$ for $\tau \in [0, 1)$. These functions are smooth and strictly convex. We also use $\mathcal{O}(\cdot)$ to denote big-O complexity notion.

2.1 The primal problem and basic assumptions

Consider the following constrained composite convex optimization problem:

$$P^* := \min_{x \in \mathcal{K}} \left\{ P(x) := g(x) + \phi(Ax) \right\}, \quad (2)$$

where $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is a smooth and convex function, $\phi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper, closed, and convex function, $A \in \mathbb{R}^{n \times p}$, and \mathcal{K} is a nonempty, closed, and convex set in \mathbb{R}^p . As a special case of (2), if we choose $\phi := \delta_{\mathcal{C}}$, the indicator of a nonempty, closed, and convex set \mathcal{C} in \mathbb{R}^n , then (2) reduces to the following general constrained convex problem:

$$g^* := \min_{x \in \mathcal{K}} \left\{ g(x) \quad \text{s.t.} \quad Ax \in \mathcal{C} \right\}. \quad (3)$$

Without loss of generality, we can also assume that g and \mathcal{K} possess a separable structure as follows:

$$g(x) := \sum_{i=1}^N g_i(x_i) \quad \text{and} \quad \mathcal{K} := \mathcal{K}_1 \times \cdots \times \mathcal{K}_N, \quad (4)$$

for $N \geq 1$, where $x_i \in \mathbb{R}^{p_i}$, $\mathcal{K}_i \subseteq \mathbb{R}^{p_i}$, and $\sum_{i=1}^N p_i = p$ for $i = 1, \dots, N$.

Note that the separable structure (4) frequently appears in graph and network optimization. It is also a natural structure in conic programming such as linear programming and monotropic programming [26]. Another example is convex empirical minimization models in statistical learning, which can also be reformulated into (2) by duplicating variables.

Basic assumptions: Our approach relies on the following assumptions:

Assumption 2.1 *The optimal solution set \mathcal{X}^* of (2) is nonempty, and hence the optimal value P^* is finite. The following Slater condition holds:*

$$0 \in \text{ri}(\text{dom}(\phi) - A(\text{dom}(g) \cap \mathcal{K})), \quad (5)$$

where $\text{ri}(\mathcal{Z})$ is the relative interior of \mathcal{Z} , and $\text{dom}(\cdot)$ is the domain of (\cdot) .

Assumption 2.2 *The function g is standard self-concordant as in Definition 3.1. \mathcal{K} is endowed with a ν_f -self-concordant barrier f as in Definition 3.2 and A is full-row rank.*

Note that Assumption 2.1 is standard and required in any primal-dual optimization method to guarantee strong duality. Assumption 2.2 is also not restrictive. First, the self-concordance of g can be relaxed to a broader class called generalized self-concordant function as shown in Proposition 3.1 with additional structures. Next, the full-row rankness of A can always be obtained by eliminating redundant rows. Finally, the self-concordant barrier of \mathcal{K} is always guaranteed under mild condition as discussed in [22].

Throughout this paper, we assume that both Assumptions 2.1 and 2.2 hold without recalling them in the sequel.

2.2 Dual problem and optimality condition

The dual problem associated with (2) can be written as

$$D^* := \min_{y \in \mathbb{R}^n} \left\{ D(y) := \underbrace{\max_{x \in \mathcal{K}} \left\{ \langle Ax, y \rangle - g(x) \right\}}_{d(y)} + \phi^*(-y) \right\}, \quad (6)$$

where $\phi^*(\cdot) := \sup_u \{ \langle \cdot, u \rangle - \phi(u) \}$ is the Fenchel conjugate of ϕ . Under the separable structure (4), we can decompose the dual function d into N functions d_i on smaller spaces \mathbb{R}^{p_i} . That is

$$d(y) := \sum_{i=1}^N d_i(y) \quad \text{with} \quad d_i(y) := \max_{x_i \in \mathcal{K}_i} \left\{ \langle A_i x_i, y \rangle - g_i(x_i) \right\}.$$

This computation can be carried out in parallel. Moreover, under Assumption 2.1, the dual optimal solution set \mathcal{Y}^* of (6) is nonempty, and the strong

duality holds, i.e. $P^* + D^* = 0$. The optimality condition of the primal problem (2) can be written as

$$\begin{cases} 0 \in \nabla g(x^*) - A^\top y^* + \mathcal{N}_{\mathcal{K}}(x^*), & \text{(primal optimality)} \\ 0 \in y^* + \partial\phi(Ax^*) & \text{(dual optimality)} \\ x^* \in \mathcal{K}, & \text{(primal feasibility)}. \end{cases} \quad (7)$$

Under Assumption 2.1, (7) is the necessary and sufficient condition for $x^* \in \mathcal{X}^*$ to be a primal optimal solution of (2), and $y^* \in \mathcal{Y}^*$ to be a dual optimal solution of (6). Note that $0 \in y^* + \partial\phi(Ax^*)$ can be written as

$$0 \in Ax^* - \partial\phi^*(-y^*) \equiv \nabla d(y^*) - \partial\phi^*(-y^*). \quad (8)$$

This is exactly the optimality condition of the dual problem (6). Our goal is to approximate a primal-dual solution of (2) and (6) in the sense of Definition 4.2.

3 Generalized self-concordance and self-concordant barriers

Let us review the theory of generalized self-concordant functions [29] and self-concordant barriers [19, 22], which will be used in the sequel.

Generalized self-concordance and standard self-concordance: Assume that $f : \text{dom}(f) \subseteq \mathbb{R}^p \rightarrow \mathbb{R}$ is a three-time continuously differentiable convex function, i.e. $f \in \mathcal{C}^3(\text{dom}(f))$, we use $\nabla^3 f(x)[u]$ to denote the third order derivative of f at $x \in \text{dom}(f)$ along a direction $u \in \mathbb{R}^p$. We recall the following definition [29].

Definition 3.1 ([29]) A \mathcal{C}^3 -convex function f is said to be (M_f, θ) -generalized self-concordant with the parameter $M_f \geq 0$, and order $\theta > 0$ if

$$|\langle \nabla^3 f(x)[v]u, u \rangle| \leq M_f \|u\|_x^2 \|v\|_x^{\theta-2} \|u\|_x^{3-\theta}, \quad x \in \text{dom}(f), \quad u, v \in \mathbb{R}^p, \quad (9)$$

where we use the convention $\frac{0}{0} = 0$ for the case $\theta < 2$ and $\theta > 3$. If $\theta = 3$, then f reduces to the self-concordant function defined by Nesterov and Nemirovskii in [22]. If $\theta = 3$ and $M_f = 2$, then f is said to be standard self-concordant.

Basic properties: Basic and fundamental properties as well as examples of generalized self-concordant functions can be found in [29]. We recall the following Legendre conjugate of a generalized self-concordant function. Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be an (M_f, θ) -generalized self-concordant function, we define

$$f^*(y) := \sup_{x \in \text{dom}(f)} \{-y^\top x - f(x)\}, \quad (10)$$

the Legendre conjugate of f (i.e. $f^*(-y)$ is the Fenchel conjugate of f). For generalized self-concordant functions and their conjugates, we have the following result.

Proposition 3.1 (a) *If f is (M_f, θ) -generalized self-concordant with $\theta \in (0, 3)$ and μ_f -strongly convex w.r.t. the Euclidean norm $\|\cdot\|_2$, then f is \hat{M}_f -self-concordant with $\hat{M}_f := \mu_f^{\frac{\theta-3}{2}} M_f$.*

- (b) If f is an (M_f, θ) -generalized self-concordant function with $\theta \in [3, 6)$, then its Legendre conjugate $f^*(\cdot)$ is also (M_f, θ_*) -generalized self-concordant with $\theta_* := 6 - \theta$.
- (c) Assume that f is M_f -self-concordant on $\text{dom}(f)$ and g is nonlinear and (M_g, θ) -generalized self-concordant on $\text{dom}(g)$ with $\theta \in (0, 3]$. If $\text{dom}(f) \cap \text{dom}(g)$ is nonempty, closed, and bounded, then $h := f + g$ is M_h -self-concordant with $M_h := \max \left\{ M_f, \hat{M}_g \right\}$, where $\hat{M}_g := \mu_g^{\frac{\theta-3}{2}} M_g$ and

$$\mu_g := \min \left\{ \lambda_{\min}(\nabla^2 g(x)) \mid x \in \text{dom}(f) \cap \text{dom}(g) \right\} \in (0, +\infty), \quad (11)$$

if $\theta < 3$, and $\hat{M}_g := M_g$ if $\theta = 3$.

Proof The proof of statements (a) and (b) can be found in [29, Propositions 4 and 6]. If $\text{dom}(f) \cap \text{dom}(g)$ is nonempty, closed, and bounded, then g is also μ_g -strongly convex on $\text{dom}(f) \cap \text{dom}(g)$ with μ_g defined by (11). Applying statement (a) to the strongly convex function g , we obtain statement (c). \square

Discussion: Proposition 3.1 shows that the class of self-concordant functions can be extended to cover at least three classes of smooth convex functions. The first one is the class of smooth and strongly convex functions that is also generalized self-concordant as studied in [29]. In the case it is not strongly convex, one can add a small quadratic regularizer to obtain this property. The second class is the conjugate of generalized self-concordant functions with Lipschitz continuous gradient. The third class of functions is generalized self-concordant functions on bounded domain. We believe that these three classes of functions cover a sufficiently large class of applications, see [29] for more detailed examples and additional properties.

Standard self-concordant barriers: Next, we recall the class of standard self-concordant barriers, and its properties.

Definition 3.2 Given a nonempty, closed, and convex set \mathcal{K} in \mathbb{R}^p , we say that f is a ν_f -self-concordant barrier of \mathcal{K} if f is standard self-concordant on $\text{dom}(f) \equiv \text{int}(\mathcal{K})$, $f(x) \rightarrow +\infty$ as x approaches the boundary $\partial\mathcal{K}$ of \mathcal{K} , and

$$\sup_{u \in \mathbb{R}^p} \left\{ \nabla f(x)^\top u - \|u\|_x^2 \right\} \leq \nu_f, \quad \forall x \in \text{dom}(f). \quad (12)$$

The self-concordant barrier f is said to be a logarithmically homogeneous self-concordant barrier if $f(\tau x) = f(x) - \nu_f \ln(\tau)$ for any $\tau > 0$ and $x \in \text{dom}(f)$.

Given a self-concordant barrier of \mathcal{K} , we define $x_f^* := \underset{x \in \mathcal{K}}{\text{argmin}} f(x)$ the analytical center of \mathcal{K} if x_f^* exists. Clearly, if \mathcal{K} is bounded, then x_f^* exists. In addition to these properties, we also have $\|x - x_f^*\|_{x_f^*} \leq \rho_f$ for any $x \in \text{dom}(f)$, where $\rho_f := \nu_f + 2\sqrt{\nu_f}$ for general self-concordant barrier f and $\rho_f := \nu_f$ if f is logarithmically homogeneous.

4 Barrier smoothing technique and inexact oracles

In this section, we describe a barrier smoothing technique for (2) which has been used in [16, 18, 20, 34, 39]. Without loss of generality, we can assume that $M_g = 2$, since any self-concordant function g with the parameter $M_g > 0$, $(M_g^2/4)g$ is standard self-concordant.

4.1 Smoothed dual problem

Under Assumption 2.2, we consider the following self-concordant barrier smoothed dual problem of (2) (shortly, smoothed dual problem):

$$\bar{D}_t^* := \min_{y \in \mathbb{R}^n} \left\{ \bar{D}_t(y) := \underbrace{\max_{x \in \text{int}(\mathcal{K})} \left\{ y^\top Ax - g(x) - tf(x) \right\}}_{\bar{d}_t(y)} + \underbrace{\phi^*(-y)}_{\bar{h}(y)} \right\}. \quad (13)$$

Note that $g(\cdot) + tf(\cdot)$ is self-concordant with the parameter $M_t := \max \left\{ 2, \frac{2}{\sqrt{t}} \right\}$ on $\text{dom}(f) \cap \text{dom}(g)$. To make it standard self-concordant, we rescale (13) as follows:

$$D_t^* := \min_{y \in \mathbb{R}^n} \left\{ D_t(y) := \underbrace{\frac{M_t^2}{4} \bar{d}_t(y)}_{d_t(y)} + \underbrace{\frac{M_t^2}{4} \bar{h}(y)}_{h_t(y)} \right\}. \quad (14)$$

From [18] or [39], d_t is standard self-concordant. Clearly, if $t \in (0, 1]$, then $M_t = \frac{2}{\sqrt{t}}$. In this case, we have $\bar{d}_t(y) = td_t(y)$ and $\bar{h}(y) = th_t(y)$.

To evaluate the (normalized) smoothed dual function d_t and its derivative, we consider the following standard self-concordant function:

$$\psi_t(x; y) := \frac{M_t^2}{4} [g(x) + tf(x) - y^\top Ax]. \quad (15)$$

Primal local norms: Note that $\nabla^2 \psi_t(x; y) = \frac{M_t^2}{4} [\nabla^2 g(x) + t \nabla^2 f(x)] = \nabla^2 \psi_t(x)$ is symmetric positive definite on $\text{dom}(g) \cap \text{dom}(f)$ and independent of y . Therefore, we define the following local norms on the primal space:

$$|u|_{x,t} := (u^\top \nabla^2 \psi_t(x) u)^{1/2}, \quad \text{and} \quad |v|_{x,t}^* := (v^\top \nabla^2 \psi_t(x)^{-1} v)^{1/2}, \quad (16)$$

for any $u, v \in \mathbb{R}^p$. If $t \in (0, 1]$, then $|u|_{x,t} = (u^\top [\nabla^2 f(x) + \frac{1}{t} \nabla^2 g(x)] u)^{1/2}$.

Exact oracles of the dual function d_t : We can summarize the properties of d_t defined in (14) into the following proposition which we omit the proof.

Proposition 4.1 *Under Assumption 2.2, $\psi_t(\cdot; y)$ defined by (15) and $d_t(\cdot)$ defined by (14) are standard self-concordant. Moreover, if the following primal subproblem has optimal solution*

$$x_t^*(y) := \arg \min_{x \in \text{int}(\mathcal{K})} \left\{ \psi_t(x; y) := \frac{M_t^2}{4} [g(x) + tf(x) - y^\top Ax] \right\}, \quad (17)$$

then its solution is unique. The optimality condition of this subproblem is

$$\nabla \psi_t(x_t^*(y); y) \equiv \frac{M_t^2}{4} [\nabla g(x_t^*(y)) + t \nabla f(x_t^*(y)) - A^\top y] = 0, \quad (18)$$

which is necessary and sufficient for $x_t^*(y)$ to be optimal to (17). The function value and derivatives of d_t in (14) can be evaluated as (see [22])

$$(\mathbf{Exact\ oracles}): \begin{cases} d_t(y) &= -\psi_t(x_t^*(y); y), \\ \nabla d_t(y) &= \frac{M_t^2}{4} A x_t^*(y), \\ \nabla^2 d_t(y) &= \frac{M_t^4}{16} A \nabla^2 \psi_t(x_t^*(y))^{-1} A^\top. \end{cases} \quad (19)$$

Dual local norms: Since $\nabla^2 d_t(y) \succ 0$, we can define the following local norms in the dual space:

$$\|u\|_{y,t} := (u^\top \nabla^2 d_t(y) u)^{1/2} \quad \text{and} \quad \|v\|_{y,t}^* := (v^\top \nabla^2 d_t(y)^{-1} v)^{1/2}. \quad (20)$$

4.2 Inexact oracles of the smoothed dual function

When g and \mathcal{K} are not trivial, solving the smoothed slave subproblem (17) exactly is impractical. We can only approximately solve (17) or (18) up to a given accuracy as defined in the following.

Definition 4.1 Let $x_t^*(y)$ be the exact solution of (17) at $y \in \mathbb{R}^n$. We call $\tilde{x}_t^*(y)$ a δ -(approximate) solution of (17) if $\delta_t(y) := |\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y),t} \leq \delta$, where $|\cdot|_{x,t}$ is defined by (16).

Given an inexact solution $\tilde{x}_t^*(y)$ of (17) as defined in Definition 4.1, we define an inexact oracle of d_t as follows:

$$(\mathbf{Inexact\ oracles}): \begin{cases} \tilde{d}_t(y) &= -\psi_t(\tilde{x}_t^*(y); y), \\ \tilde{\nabla} d_t(y) &= \frac{M_t^2}{4} A \tilde{x}_t^*(y), \\ \tilde{\nabla}^2 d_t(y) &= \frac{M_t^4}{16} A \nabla^2 \psi_t(\tilde{x}_t^*(y))^{-1} A^\top. \end{cases} \quad (21)$$

Since $\nabla^2 \psi_t(\cdot)$ is positive definite and A is full-row rank, $\tilde{\nabla}^2 d_t(y)$ is positive definite. Now we define the following local norms using inexact oracles:

$$\|u\|_{y,t} := (u^\top \tilde{\nabla}^2 d_t(y) u)^{1/2} \quad \text{and} \quad \|v\|_{y,t}^* := (v^\top \tilde{\nabla}^2 d_t(y)^{-1} v)^{1/2}. \quad (22)$$

We first prove some properties of inexact solution $\tilde{x}_t^*(y)$ and inexact oracles of d_t defined by (21) in the following proposition, whose proof can be found in Appendix A.1.

Proposition 4.2 For any $\delta \in [0, 1)$, we have:

$$\text{if } |\nabla \psi_t(\tilde{x}_t^*(y); y)|_{\tilde{x}_t^*(y),t}^* \leq \frac{\delta}{1+\delta} \text{ then } \delta_t(y) := |\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y),t} \leq \delta. \quad (23)$$

In addition, for d_t and its derivatives defined by (19), and its inexact oracle defined by (21), the following properties hold

$$\begin{cases} 0 \leq \omega\left(\frac{\delta_t(y)}{1+\delta_t(y)}\right) \leq d_t(y) - \tilde{d}_t(y) \leq \omega_*\left(\frac{\delta_t(y)}{1-\delta_t(y)}\right), \\ (1 - \delta_t(y))^2 \tilde{\nabla}^2 d_t(y) \preceq \nabla^2 d_t(y) \preceq (1 - \delta_t(y))^{-2} \tilde{\nabla}^2 d_t(y), \\ \|\tilde{\nabla} d_t(y) - \nabla d_t(y)\|_{y,t}^* \leq \delta_t(y), \end{cases} \quad (24)$$

where $\omega(\tau) := \tau - \ln(1+\tau)$ for $\tau \geq 0$ and $\omega_*(\tau) := -\tau - \ln(1-\tau)$ for $\tau \in [0, 1)$.

Discussion: The first estimate (23) shows that to obtain an approximate solution $\tilde{x}_t^*(y)$ such that $|\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y),t} \leq \delta$, we need to solve the slave problem (17) such that

$$|\nabla g(\tilde{x}_t^*(y)) + t\nabla f(\tilde{x}_t^*(y)) - A^\top y|_{\tilde{x}_t^*(y),t}^* \leq \frac{4\delta}{M_t^2(1+\delta)}. \quad (25)$$

This condition is implementable, e.g., when we apply a Newton-type method to solve the nonlinear system (18). The estimates in (24) show us how the inexact oracles in (21) approximate the exact ones in (19).

Approximate primal-dual solutions: Given an accuracy $\varepsilon > 0$, our goal is to compute an ε -approximate primal-dual solution $(\tilde{x}^*, \tilde{y}^*)$ to (x^*, y^*) of (7) in the following sense:

Definition 4.2 A pair $(\tilde{x}^*, \tilde{y}^*)$ is called an ε -approximate primal-dual solution to an exact primal-dual one (x^*, y^*) of (7) if

$$\begin{cases} |A^\top \tilde{y}^* - \nabla g(\tilde{x}^*)|_{\tilde{x}^*,t}^* \leq \varepsilon & (\varepsilon\text{-primal optimality}), \\ r \in \tilde{y}^* + \partial\phi(A\tilde{x}^* + e) & (\varepsilon\text{-dual optimality}), \\ \tilde{x}^* \in \text{int}(\mathcal{K}) & (\text{primal feasibility}), \\ \|e\|_{\tilde{y}^*,t}^* \leq \varepsilon \quad \text{and} \quad |A^\top r|_{\tilde{x}^*,t}^* \leq \varepsilon. \end{cases} \quad (26)$$

Here, the errors are measured through local norms in primal and dual spaces defined in (16) and (22). These norms are computable since they are defined through \tilde{x}^* and \tilde{y}^* . In addition, since A is full-row rank, $A^\top r = 0$ if and only if $r = 0$. Because $\tilde{x}^* \in \text{int}(\mathcal{K})$, we have $\mathcal{N}_{\mathcal{K}}(\tilde{x}^*) = \{\mathbf{0}\}$. Therefore the first line of (26) can approximate the first line of (7). Similarly, the second line of (26) approximates the second line of (7), i.e. $0 \in y^* + \partial\phi(Ax^*)$. Therefore, Definition 4.2 is consistent with the optimality condition (7).

5 Inexact IPLD Method with Inexact Oracles

We develop an inexact interior-point Lagrangian decomposition method to solve (2) by using the inexact oracles (21).

5.1 Inexact proximal-Newton method for (14)

The optimality condition of (14): Recall the smoothed dual problem (14), its optimality condition is

$$0 \in \nabla d_t(y) + \partial h_t(y) = \frac{M_t^2}{4} Ax_t^*(y) + \partial h_t(y). \quad (27)$$

Any y_t^* satisfies (27) is an optimal solution of (14). The sequence $\{(x_t^*(y_t^*), y_t^*)\}_{t \geq 0}$ forms a central path, which converges to (x^*, y^*) a primal-dual solution of (2).

Exact Proximal Newton scheme: Suppose that we are currently at y^k , since d_t is twice differentiable, we will apply proximal-Newton method to compute \bar{y}^{k+1} , which leads to

$$0 \in \tilde{\nabla}^2 d_{t_{k+1}}(y^k)(\bar{y}^{k+1} - y^k) + \tilde{\nabla} d_{t_{k+1}}(y^k) + \partial h_{t_{k+1}}(\bar{y}^{k+1}). \quad (28)$$

If we define

$$Q_{t_{k+1}}(y) := \langle \tilde{\nabla} d_{t_{k+1}}(y^k), y - y^k \rangle + \frac{1}{2} \langle \tilde{\nabla}^2 d_{t_{k+1}}(y^k)(y - y^k), y - y^k \rangle + h_{t_{k+1}}(y), \quad (29)$$

then we can write $\bar{y}^{k+1} := \operatorname{argmin}_y Q_{t_{k+1}}(y)$. Introducing the notation $\operatorname{prox}_{h_t}^{\tilde{\nabla}^2 d_t}(\cdot)$, we can write (28) in the following form (see [35] for a concrete definition)

$$\bar{y}^{k+1} := \operatorname{prox}_{h_{t_{k+1}}}^{\tilde{\nabla}^2 d_{t_{k+1}}(y^k)} \left(y^k - \tilde{\nabla}^2 d_{t_{k+1}}(y^k)^{-1} \tilde{\nabla} d_{t_{k+1}}(y^k) \right). \quad (30)$$

Inexact Proximal Newton scheme: Similarly, we can also approximately solve (28) up to a given accuracy as.

$$y^{k+1} \approx \operatorname{prox}_{h_{t_{k+1}}}^{\tilde{\nabla}^2 d_{t_{k+1}}(y^k)} \left(y^k - \tilde{\nabla}^2 d_{t_{k+1}}(y^k)^{-1} \tilde{\nabla} d_{t_{k+1}}(y^k) \right). \quad (31)$$

Here, the approximation “ \approx ” is defined in the following sense:

Definition 5.1 For a given $\epsilon \geq 0$ and $Q_{t_{k+1}}$ defined by (29), a vector y^{k+1} given in (31) is said to be an ϵ -approximate solution to \bar{y}^{k+1} of (28) if

$$Q_{t_{k+1}}(y^{k+1}) - Q_{t_{k+1}}(\bar{y}^{k+1}) \leq \frac{\epsilon^2}{2}. \quad (32)$$

Note that (32) implies $\|y^{k+1} - \bar{y}^{k+1}\|_{y^k, t_{k+1}} \leq \epsilon$. There exists several convex optimization methods to compute y^{k+1} in (31). For example, we can apply accelerated proximal gradient methods such as FISTA [1, 21] to compute this point. We can also apply semi-smooth Newton-CG augmented Lagrangian methods in [17, 40] to solve this problem. We will discuss the computation of y^{k+1} in detail in Section 6.

Generalized gradient mapping: Now let us define the following inexact generalized gradient mapping

$$\tilde{G}_t(y) := \tilde{\nabla}^2 d_t(y) \left(y - \operatorname{prox}_{h_t}^{\tilde{\nabla}^2 d_t(y)}(y - \tilde{\nabla}^2 d_t(y)^{-1} \tilde{\nabla} d_t(y)) \right). \quad (33)$$

Using $\tilde{\nabla}^2 d_t(\cdot)$ defined by (21), we further define the following quantity:

$$\lambda_t(y) := \|\tilde{G}_t(y)\|_{y,t}^* = \langle \tilde{\nabla}^2 d_t(y)^{-1} \tilde{G}_t(y), \tilde{G}_t(y) \rangle^{1/2}. \quad (34)$$

We call $\lambda_t(y)$ the inexact proximal-Newton decrement. In Subsection 5.4 we can show that this quantity can be used to characterize the optimality condition (7).

5.2 The algorithm

From the above analysis, we can combine all the steps together and describe an algorithm to solve (2) as in Algorithm 1. In this algorithm, we explicitly show how to choose the accuracy of inexact oracles and inexact proximal-Newton direction, and how to update the penalty parameter t .

Note that we have not specified how to find a starting point (x^0, y^0) to guarantee (35) and how to set k_{\max} in Algorithm 1. In Subsection 5.5, we will show that such an (x^0, y^0) can be found in finite steps. In Subsection 5.4, we show how to set k_{\max} to get an ε -approximate primal-dual solution of (2).

Algorithm 1 (*Inexact Interior-Point Lagrangian Decomposition Algorithm*)

- 1: **Phase 1: Find an initial point.** Given any value $t_0 \in (0, 1]$ and $\beta \in (0, \frac{1}{10}]$, find starting points $y^0 \in \mathbb{R}^n$ and $x^0 \in \mathbb{R}^p$ such that

$$\|\tilde{G}_{t_0}(y^0)\|_{y^0, t_0}^* \leq \beta \quad \text{and} \quad |\nabla \psi_{t_0}(x^0; y^0)|_{x^0, t_0}^* \leq \frac{\tilde{\delta}_0}{1 + \tilde{\delta}_0}, \quad (35)$$

by using Algorithm 2 below, for any predefined accuracy $\tilde{\delta}_0 \in (0, \frac{\beta}{100}]$.

- 2: **Phase 2: Main iteration.** For $k = 0$ to k_{\max} , perform
- 3: Update t_k as $t_{k+1} := \sigma t_k$, where $\sigma \in (0, 1)$ is defined by (45) below.
- 4: Solve approximately (18) at $y = y^k$ up to an accuracy $\delta_k \in (0, \frac{\beta}{100}]$ to get $x^{k+1} := \tilde{x}_{t_{k+1}}^*(y^k)$, i.e.:

$$|\nabla \psi_{t_{k+1}}(x^{k+1}; y^k)|_{x^{k+1}, t_{k+1}}^* \leq \frac{\delta_k}{1 + \delta_k}.$$

- 5: (**Inexact oracles**): Evaluate inexact gradient and Hessian of d_t as

$$\begin{cases} \tilde{\nabla} d_{t_{k+1}}(y^k) &:= \frac{M_{t_{k+1}}^2}{4} A x^{k+1}, \\ \tilde{\nabla}^2 d_{t_{k+1}}(y^k) &:= \frac{M_{t_{k+1}}^4}{16} A \nabla^2 \psi_{t_{k+1}}(x^{k+1})^{-1} A^\top. \end{cases} \quad (36)$$

- 6: (**Inexact proximal-Newton step**): Compute y^{k+1} up to an accuracy $\epsilon_k \in (0, \frac{\beta}{100}]$, i.e.:

$$y^{k+1} \approx \text{prox}_{\mu_{t_{k+1}}}^{\tilde{\nabla}^2 d_{t_{k+1}}(y^k)} \left(y^k - \tilde{\nabla}^2 d_{t_{k+1}}(y^k)^{-1} \tilde{\nabla} d_{t_{k+1}}(y^k) \right).$$

7: **End.**

5.3 Convergence analysis

Our analysis consists of several steps and is organized as follows:

- Lemma 5.1 provides an estimate between $\lambda_{t_{k+1}}(y^{k+1})$ and $\lambda_{t_{k+1}}(y^k)$ in (34).
- Lemma 5.2 bounds $\lambda_{t_{k+1}}(y^k)$ in terms of $\tilde{\Delta}_{t_k}$, $\tilde{\Delta}_{t_{k+1}}$ and $\lambda_{t_k}(y^k)$, where $\tilde{\Delta}_{t_k}$ and $\tilde{\Delta}_{t_{k+1}}$ measure the distances between $\tilde{x}_{t_k}^*(y^k)$ and $\tilde{x}_{t_{k+1}}^*(y^k)$.
- Lemma 5.3 shows how to upper bound $\tilde{\Delta}_{t_k}$ and $\tilde{\Delta}_{t_{k+1}}$.
- The main result of this section is Theorem 5.1 which provides an update rule of t to maintain the point y^k in the neighborhood of the central path. The proof of this theorem is obtained by combining all the above lemmas.

Firstly, we state the main estimate of the inexact Newton-type step at Step 6 of Algorithm 1 in Lemma 5.1, whose proof is given in Appendix A.2.1.

Lemma 5.1 *Let $\{y^k\}$ be generated by Algorithm 1, and*

$$|\tilde{x}_{t_{k+1}}^*(y^{k+1}) - x_{t_{k+1}}^*(y^{k+1})|_{\tilde{x}_{t_{k+1}}^*(y^{k+1}), t_{k+1}}^* \leq \tilde{\delta}_{k+1}.$$

Then

$$\begin{aligned} \lambda_{t_{k+1}}(y^{k+1}) &\leq \tilde{\delta}_{k+1} + \frac{1}{(1 - \tilde{\delta}_{k+1})(1 - \delta_k - \lambda_{t_{k+1}}(y^k) - \epsilon_k)} \left[3\epsilon_k + \delta_k \right. \\ &\quad \left. + \sqrt{4\delta_k - 2\delta_k^2}(\lambda_{t_{k+1}}(y^k) + \epsilon_k) + \frac{(\lambda_{t_{k+1}}(y^k) + \epsilon_k)^2}{(1 - \lambda_{t_{k+1}}(y^k) - \delta_k - \epsilon_k)} \right]. \end{aligned} \quad (37)$$

In particular, if $\tilde{\delta}_{k+1} = 0$, $\delta_k = 0$, and $\epsilon_k = 0$, then (37) reduces to

$$\lambda_{t_{k+1}}(y^{k+1}) \leq \frac{\lambda_{t_{k+1}}(y^k)^2}{(1 - \lambda_{t_{k+1}}(y^k))^2}. \quad (38)$$

Note that if we solve both the slave problem at Step 4 and the master problem at Step 6 exactly, then we could obtain the estimate (38), which is the same as in standard interior-point path-following methods [19]. Next, we show a relation between $\lambda_{t_{k+1}}(y^k)$ and $\lambda_{t_k}(y^k)$, whose proof is in Appendix A.2.2.

Lemma 5.2 *Let t_k be updated as $t_{k+1} := \sigma t_k$ for given $\sigma \in (0, 1)$. Define*

$$\begin{cases} \tilde{\Delta}_{t_k} &:= |\tilde{x}_{t_{k+1}}^*(y^k) - \tilde{x}_{t_k}^*(y^k)|_{\tilde{x}_{t_k}^*(y^k), t_k}, \\ \tilde{\Delta}_{t_{k+1}} &:= |\tilde{x}_{t_{k+1}}^*(y^k) - \tilde{x}_{t_k}^*(y^k)|_{\tilde{x}_{t_{k+1}}^*(y^k), t_{k+1}}. \end{cases} \quad (39)$$

Then, the following estimate holds

$$\lambda_{t_{k+1}}(y^k) \leq \tilde{\Delta}_{t_{k+1}} + \left[\frac{1 + \sqrt{1 - 2\sigma(1 - \tilde{\Delta}_{t_k})^2 + \sigma}}{\sigma(1 - \tilde{\Delta}_{t_k})} \right] \lambda_{t_k}(y^k). \quad (40)$$

The following lemma shows how to bound $\tilde{\Delta}_{t_k}$ and $\tilde{\Delta}_{t_{k+1}}$, the distances between $\tilde{x}_{t_k}^*(y^k)$ and $\tilde{x}_{t_{k+1}}^*(y^k)$, whose proof is given in Appendix A.2.3.

Lemma 5.3 *Let $\tilde{\Delta}_{t_k}$ and $\tilde{\Delta}_{t_{k+1}}$ be defined by (39), and $t_{k+1} := \sigma t_k$ for some $\sigma \in (0, 1)$. We define the following quantities:*

$$\begin{cases} \hat{\delta}_{t_k} &:= |\nabla \psi_{t_k}(\tilde{x}_{t_k}^*(y^k); y^k)|_{\tilde{x}_{t_k}^*(y^k), t_k}^* \\ \hat{\delta}_{t_{k+1}} &:= |\nabla \psi_{t_{k+1}}(\tilde{x}_{t_{k+1}}^*(y^k); y^k)|_{\tilde{x}_{t_{k+1}}^*(y^k), t_{k+1}}^*. \end{cases} \quad (41)$$

Then, we have

$$\begin{cases} \frac{\tilde{\Delta}_{t_k}^2}{1 + \tilde{\Delta}_{t_k}} &\leq \tilde{\Delta}_{t_k} \hat{\delta}_{t_k} + \left(\sigma \hat{\delta}_{t_{k+1}} + (1 - \sigma) \sqrt{\nu_f} \right) \tilde{\Delta}_{t_{k+1}} \\ \frac{\tilde{\Delta}_{t_{k+1}}^2}{1 + \tilde{\Delta}_{t_{k+1}}} &\leq \tilde{\Delta}_{t_{k+1}} \hat{\delta}_{t_{k+1}} + \left(\frac{\hat{\delta}_{t_k}}{\sigma} + \left(\frac{1 - \sigma}{\sigma} \right) \sqrt{\nu_f} \right) \tilde{\Delta}_{t_k}, \end{cases} \quad (42)$$

where ν_f is the barrier parameter of f . In particular, for fixed $\delta \in [0, 1]$, if we choose $\hat{\delta}_{t_k} \leq \delta$ and $\hat{\delta}_{t_{k+1}} \leq \delta$, then

$$\begin{cases} \frac{\tilde{\Delta}_{t_k}^2}{1 + \tilde{\Delta}_{t_k}} \leq \delta \cdot \tilde{\Delta}_k + c_\nu(\sigma) \cdot \tilde{\Delta}_{t_{k+1}} \\ \frac{\tilde{\Delta}_{t_{k+1}}^2}{1 + \tilde{\Delta}_{t_{k+1}}} \leq \delta \cdot \tilde{\Delta}_{t_{k+1}} + c_\nu(\sigma) \cdot \tilde{\Delta}_{t_k}, \end{cases} \quad (43)$$

where $c_\nu(\sigma) := \frac{\delta}{\sigma} + \left(\frac{1-\sigma}{\sigma}\right)\sqrt{\nu_f}$ is a decreasing function of σ on $(0, 1]$. As a consequence, we also have

$$\tilde{\Delta}_{t_k} \leq \frac{\delta + c_\nu(\sigma)}{1 - \delta - c_\nu(\sigma)} \quad \text{and} \quad \tilde{\Delta}_{t_{k+1}} \leq \frac{\delta + c_\nu(\sigma)}{1 - \delta - c_\nu(\sigma)}. \quad (44)$$

Utilizing the results of Lemma 5.1, Lemma 5.2 and Lemma 5.3, we can prove the following main result on the iteration-complexity of Algorithm 1.

Theorem 5.1 *Let us choose $\beta \in (0, \frac{1}{10}]$. Suppose that we choose $\tilde{\delta}_0, \delta_k, \tilde{\delta}_{k+1}, \epsilon_k \in [0, \frac{\beta}{100}]$ and update t_k in Algorithm 1 as $t_{k+1} := \sigma t_k$ with*

$$\sigma := 1 - \frac{0.29\beta}{0.3\beta + \sqrt{\nu_f}} \in (0, 1). \quad (45)$$

In addition, if $y^0 \in \mathbb{R}^n$ and $x^0 \in \mathbb{R}^p$ satisfy (35), then for all $k \geq 0$, we have

$$\lambda_{t_k}(y^k) \leq \beta.$$

Consequently, the number of iterations to obtain $t_k \leq \hat{\varepsilon}$ for a given $\hat{\varepsilon} > 0$ and $\lambda_{t_k}(y^k) \leq \beta$ does not exceed:

$$k_{\max} := \left\lceil \frac{\ln\left(\frac{t_0}{\hat{\varepsilon}}\right)}{-\ln(\sigma)} \right\rceil = \mathcal{O}\left(\sqrt{\nu_f} \ln\left(\frac{t_0}{\hat{\varepsilon}}\right)\right), \quad (46)$$

where ν_f is the barrier parameter of f and $t_0 \in (0, 1]$.

Proof Let us first assume that $\lambda_{t_{k+1}}(y^k) \leq 2.1\beta$. Using $\delta_k, \tilde{\delta}_{k+1}, \epsilon_k \in [0, 10^{-2}\beta]$, after a few elementary calculations, we can overestimate (37) in Lemma 5.1 as

$$\begin{aligned} \lambda_{t_{k+1}}(y^{k+1}) &\leq \frac{\beta}{100} + \frac{0.04\beta + 0.1\sqrt{4\beta - 0.02\beta^2}(2.11\beta)}{(1 - 10^{-2}\beta)(1 - 2.12\beta)} \\ &\quad + \frac{(2.11\beta)^2}{(1 - 10^{-2}\beta)(1 - 2.12\beta)^2} \\ &\leq \beta, \end{aligned} \quad (47)$$

when $\beta \in (0, \frac{1}{10}]$. Now, we prove that $\lambda_{t_{k+1}}(y^k) \leq 2.1\beta$ is always satisfied. Indeed, since

$$\begin{aligned} |\nabla\psi_{t_k}(\tilde{x}_{t_k}^*(y^k); y^k)|_{\tilde{x}_{t_k}^*(y^k), t_k}^* &\leq \frac{\tilde{\delta}_k}{1 + \delta_k} \leq \frac{\beta}{100} \\ |\nabla\psi_{t_{k+1}}(\tilde{x}_{t_{k+1}}^*(y^k); y^k)|_{\tilde{x}_{t_{k+1}}^*(y^k), t_{k+1}}^* &\leq \frac{\delta_k}{1 + \delta_k} \leq \frac{\beta}{100}, \end{aligned}$$

we can choose δ in Lemma 5.3 to be $\frac{\beta}{100}$. In addition, from $\sigma := 1 - \frac{0.29\beta}{0.3\beta + \sqrt{\nu_f}}$, we can show that

$$c_\nu(\sigma) := \frac{\delta}{\sigma} + \frac{1-\sigma}{\sigma} \sqrt{\nu_f} = \frac{10^{-2}\beta}{\sigma} + \frac{1-\sigma}{\sigma} \sqrt{\nu_f} = 0.3\beta.$$

Next, using (44), we get

$$\tilde{\Delta}_{t_k} \leq \frac{0.01\beta + 0.3\beta}{1 - 10^{-2}\beta - 0.3\beta} \leq 0.4493\beta \quad \text{and} \quad \tilde{\Delta}_{t_{k+1}} \leq \frac{10^{-2}\beta + 0.3\beta}{1 - 10^{-2}\beta - 0.3\beta} \leq 0.4493\beta.$$

Finally, combining these estimates and (40) we can show that

$$\lambda_{t_{k+1}}(y^k) \leq 0.4493\beta + \frac{1 + \sqrt{1 - 2\sigma(1 - 0.4493\beta)^2 + \sigma}}{\sigma(1 - 0.4493\beta)} \beta \stackrel{(45)}{\leq} 2.1\beta,$$

when $\beta \in (0, \frac{1}{10}]$. Since $t_k := \sigma^k t_0 = \left(1 - \frac{0.29\beta}{0.3\beta + \sqrt{\nu_f}}\right)^k t_0$, to guarantee $t_k \leq \hat{\varepsilon}$, we impose $\sigma^k t_0 = \left(1 - \frac{0.29\beta}{0.3\beta + \sqrt{\nu_f}}\right)^k t_0 \leq \hat{\varepsilon}$. Note that $-\ln\left(1 - \frac{0.29\beta}{0.3\beta + \sqrt{\nu_f}}\right) \sim \frac{0.29\beta}{0.3\beta + \sqrt{\nu_f}} \sim \frac{1}{\sqrt{\nu_f}}$, we have

$$k \geq \left\lceil \frac{\ln\left(\frac{t_0}{\hat{\varepsilon}}\right)}{-\ln(\sigma)} \right\rceil = \mathcal{O}\left(\sqrt{\nu_f} \ln\left(\frac{t_0}{\hat{\varepsilon}}\right)\right),$$

as stated in (46). Here, \sim means that two quantities can be approximated by the same order. \square

The worst-case iteration complexity: Theorem 5.1 shows that for any $\hat{\varepsilon} > 0$, the number of iterations k to obtain y^k such that $\lambda_{t_k}(y^k) \leq \beta$ and $t_k \leq \hat{\varepsilon}$ does not exceed

$$\mathcal{O}\left(\sqrt{\nu_f} \ln\left(\frac{t_0}{\hat{\varepsilon}}\right)\right),$$

which is the same as in standard interior-point methods [19, 22] up to a constant factor. It depends on $\sqrt{\nu_f}$, where ν_f is the barrier parameter of f . Note that the parameter β in Algorithm 1 represents the radius of the central path neighborhood as in standard path-following methods. While the range of β in standard exact path-following methods [19] is $(0, \frac{3-\sqrt{5}}{2}]$, it is $[0, \frac{1}{10}]$ in our method. Clearly, the latter is much smaller than the former one. However, this range was roughly estimated in our analysis and it is affected by the inexactness in our algorithm.

As we will show in Subsection 5.4, the conditions $\lambda_{t_k}(y^k) \leq \beta$ and $t_k \leq \hat{\varepsilon}$ imply an approximate solution of (2) and (6).

5.4 Optimality certification

Our goal is to compute an approximate solution of the primal problem (2). The following theorem shows how we can find this approximate solution for both the primal and dual problem.

Theorem 5.2 *Let $\{(x^{k+1}, y^k)\}$ be the sequence generated by Algorithm 1. Then, for $t_{k+1} \in (0, 1]$ we have the following guarantees:*

$$\begin{cases} x^{k+1} \in \text{int}(\mathcal{K}), \\ |A^\top y^k - \nabla g(x^{k+1})|_{x^{k+1}, t_{k+1}}^* \leq \left(\sqrt{\nu_f} + \frac{\delta_k}{1+\delta_k}\right) t_{k+1}, \\ r^k \in y^k + \partial\phi(Ax^{k+1} + e^k), \\ \|e^k\|_{y^k, t_{k+1}}^* \leq t_{k+1} \lambda_{t_{k+1}}(y^k), \text{ and } |A^\top r^k|_{x^{k+1}, t_{k+1}}^* \leq t_{k+1} \lambda_{t_{k+1}}(y^k). \end{cases} \quad (48)$$

Consequently, the number of iterations to obtain an ε -primal-dual solution (x^{k+1}, y^k) in the sense of Definition 4.2 does not exceed:

$$k_{\max} := \mathcal{O}\left(\sqrt{\nu_f} \ln\left(\frac{\sqrt{\nu_f} t_0}{\varepsilon}\right)\right), \quad (49)$$

where $t_0 \in (0, 1]$ and ν_f is the barrier parameter of f .

Proof From Step 4 of Algorithm 1, we can see that $x^{k+1} := \tilde{x}_{t_{k+1}}^*(y^k) \in \text{int}(\mathcal{K})$. Moreover, Step 4 also leads to

$$\begin{aligned} |\nabla g(x^{k+1}) - A^\top y^k|_{x^{k+1}, t_{k+1}}^* &\leq |\nabla g(x^{k+1}) - A^\top y^k + t_{k+1} \nabla f(x^{k+1})|_{x^{k+1}, t_{k+1}}^* \\ &\quad + t_{k+1} |\nabla f(x^{k+1})|_{x^{k+1}, t_{k+1}}^* \\ &\leq \frac{\delta_k t_{k+1}}{1+\delta_k} + t_{k+1} |\nabla f(x^{k+1})|_{x^{k+1}, t_{k+1}}^*. \end{aligned} \quad (50)$$

Next, for $t \in (0, 1]$, it is obvious that $\nabla^2 \psi_t(x; y) = \frac{M_t^2}{4} [\nabla^2 g(x) + t \nabla^2 f(x)] = \nabla^2 f(x) + \frac{1}{t} \nabla^2 g(x)$. Consequently, one has $\nabla^2 \psi_t(x; y) \succeq \nabla^2 f(x)$. Using this fact, we can easily show that

$$|\nabla f(x^{k+1})|_{x^{k+1}, t_{k+1}}^* \leq \|\nabla f(x^{k+1})\|_{x^{k+1}}^* \leq \sqrt{\nu_f}.$$

Combining this inequality and (50), we obtain the second estimate of (48).

Now, from (28), we have

$$-\tilde{\nabla}^2 d_{t_{k+1}}(y^k)(\bar{y}^{k+1} - y^k) - \tilde{\nabla} d_{t_{k+1}}(y^k) \in \partial h_{t_{k+1}}(\bar{y}^{k+1}).$$

Using (21) and the definition of h_t , the last estimate becomes

$$-t_{k+1} \tilde{\nabla}^2 d_{t_{k+1}}(y^k)(\bar{y}^{k+1} - y^k) \in Ax^{k+1} - \partial\phi^*(-\bar{y}^{k+1}).$$

If we define $r^k := y^k - \bar{y}^{k+1}$ and $e^k := t_{k+1} \tilde{\nabla}^2 d_{t_{k+1}}(y^k)(\bar{y}^{k+1} - y^k)$, then the last expression leads to

$$-e^k \in Ax^{k+1} - \partial\phi^*(-y^k + r^k) \Leftrightarrow r^k \in y^k + \partial\phi(Ax^{k+1} + e^k).$$

It is obvious to show that

$$\|e^k\|_{y^k, t_{k+1}}^* = t_{k+1} \|y^k - \bar{y}^{k+1}\|_{y^k, t_{k+1}} = t_{k+1} \lambda_{t_{k+1}}(y^k),$$

which is the first statement in the last line of (48).

Now, from (36) and $t_{k+1} \in (0, 1]$, we have $\tilde{\nabla}^2 d_{t_{k+1}}(y^k) = \frac{1}{t_{k+1}^2} A(\nabla^2 f(x^{k+1}) + \frac{1}{t_{k+1}} \nabla^2 g(x^{k+1}))^{-1} A^\top$. This implies that

$$\begin{aligned} \lambda_{t_{k+1}}(y^k)^2 &= (\bar{y}^{k+1} - y^k)^\top \tilde{\nabla}^2 d_{t_{k+1}}(y^k) (\bar{y}^{k+1} - y^k) \\ &= \frac{1}{t_{k+1}^2} (r^k)^\top A(\nabla^2 f(x^{k+1}) + \frac{1}{t_{k+1}} \nabla^2 g(x^{k+1}))^{-1} A^\top r^k \\ &= \frac{1}{t_{k+1}^2} (|A^\top r^k|_{x^{k+1}, t_{k+1}}^*)^2. \end{aligned}$$

Therefore, we have $|A^\top r^k|_{x^{k+1}, t_{k+1}}^* = t_{k+1} \lambda_{t_{k+1}}(y^k)$, which proves the second statement in the last line of (48).

From (48), to obtain an ε -primal-dual solution (x^{k+1}, y^k) in the sense of Definition 4.2, we need to set $(\sqrt{\nu_f} + \frac{\delta_k}{1+\delta_k}) t_{k+1} \leq \varepsilon$ and $t_{k+1} \lambda_{t_{k+1}}(y^k) \leq \varepsilon$. Since $\lambda_{t_{k+1}}(y^k) \leq 2.1\beta$ (see the proof in Theorem 5.1) and $\delta_k \leq \frac{\beta}{100}$, we can set $t_{k+1} \leq \hat{\varepsilon}$ such that

$$\varepsilon \geq \hat{\varepsilon}(\sqrt{\nu_f} + 1) \geq \hat{\varepsilon} \max \left\{ \sqrt{\nu_f} + \frac{0.01\beta}{1 + 0.01\beta}, 2.1\beta \right\},$$

i.e., $\hat{\varepsilon} \leq \frac{\varepsilon}{(1+\sqrt{\nu_f})}$. Combining this expression and (46), we can show that the number of iterations to obtain an ε -primal-dual solution does not exceed $\mathcal{O}\left(\sqrt{\nu_f} \ln\left(\frac{\sqrt{\nu_f} t_0}{\varepsilon}\right)\right)$, which is exactly (49). \square

Discussion: Theorem 5.2 estimates the maximum iterations k_{\max} to obtain an ε -primal-dual solution (x^{k+1}, y^k) of (2) and (6). It shows that such a number of iterations remains the same as in standard path-following methods [19] up to a constant factor. Although the norms in (48) are local norms, but this is the standard metric used in general interior-point methods [19, 22].

5.5 Finding an initial point in Algorithm 1

We need to find (x^0, y^0) such that the condition (35) holds. As in standard interior-point methods, we need to perform a damped proximal-Newton method. Such a method can be found in, e.g. [31, 32], but since we use inexact oracles, we need to customize this method in our context. More specifically, we describe this routine in Algorithm 2.

We terminate Algorithm 2 if we find $x^0 := \hat{x}^{j_{\max}}$ and $y^0 := \hat{y}^{j_{\max}}$ such that (35) holds. Since the constraint of x^0 in (35) is always satisfied from Step 3 of Algorithm 2, we only need to guarantee that $\lambda_{t_0}(y^0) \leq \beta$.

The following theorem estimates the number of iterations to obtain (x^0, y^0) satisfying (35).

Algorithm 2 (*Find an initial point x^0, y^0*)

- 1: **Initialization.** Choose an initial point $\hat{y}^0 \in \mathbb{R}^n$ and fix a value $t_0 \in (0, 1]$.
- 2: **Main iteration.** For $j = 0$ to j_{\max} , perform
- 3: Solve approximately (18) at $y = \hat{y}^j$ up to an accuracy $\delta_j \in (0, \frac{\beta}{100}]$ to get $\hat{x}^j := \tilde{x}_{t_0}^*(\hat{y}^j)$, i.e.:

$$|\nabla \psi_{t_0}(\hat{x}^j; \hat{y}^j)|_{\hat{x}^j, t_0}^* \leq \frac{\delta_j}{1 + \delta_j}.$$

- 4: (**Inexact oracles**): Evaluate inexact gradient and Hessian of d_{t_0} as

$$\begin{cases} \tilde{\nabla} d_{t_0}(\hat{y}^j) &:= \frac{M_{t_0}^2}{4} A \hat{x}^j, \\ \tilde{\nabla}^2 d_{t_0}(\hat{y}^j) &:= \frac{M_{t_0}^4}{16} A \nabla^2 \psi_{t_0}(\hat{x}^j)^{-1} A^\top. \end{cases} \quad (51)$$

- 5: (**Inexact damped-step proximal-Newton step**): Compute \hat{s}^j up to an accuracy $\epsilon_j \in (0, \frac{\beta}{100}]$ and update \hat{y}^j , i.e.:

$$\begin{cases} \hat{s}^j & \approx s^j := \text{prox}_{h_{t_0}}^{\tilde{\nabla}^2 d_{t_0}(\hat{y}^j)} \left(\hat{y}^j - \tilde{\nabla}^2 d_{t_0}(\hat{y}^j)^{-1} \tilde{\nabla} d_{t_0}(\hat{y}^j) \right) \\ \hat{y}^{j+1} &:= (1 - \alpha_j) \hat{y}^j + \alpha_j \hat{s}^j, \end{cases}$$

where $\alpha_j := \frac{(\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)^2}{((1 - \delta_j)(\hat{\lambda}_j - \epsilon_j - \delta_j) + 1)\hat{\lambda}_j} \in (0, 1)$ and $\hat{\lambda}_j := \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0}$.

- 6: **End.**
-

Theorem 5.3 *Let us define $\hat{\lambda}_j := \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0}$ and $\lambda_j := \|s^j - \hat{y}^j\|_{\hat{y}^j, t_0}$. Let $\{(\hat{x}^j, \hat{y}^j)\}$ be the sequence generated by Algorithm 2, where we choose $\delta_j, \epsilon_j \in (0, \frac{\beta}{100}]$ and the step-size*

$$\alpha_j := \frac{(\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)^2}{[1 + (1 - \delta_j)(\hat{\lambda}_j - \epsilon_j - \delta_j)]\hat{\lambda}_j} \in (0, 1). \quad (52)$$

Then, after at most finite number of iterations j_{\max} as

$$j_{\max} := \left\lfloor \frac{D_{t_0}(\hat{y}^0) - D_{t_0}(y_{t_0}^*)}{\omega(0.97\beta(1 - 10^{-2}\beta))} \right\rfloor + 1, \quad (53)$$

we obtain $y^0 := \hat{y}^{j_{\max}}$ and $x^0 := \hat{x}^{j_{\max}}$ such that $\lambda_{t_0}(y^0) \leq \beta$ and (35) holds, where $y_{t_0}^$ is the optimal solution of (14) at $t := t_0$.*

Proof Note that at each iteration j of Algorithm 2, we always have $\lambda_j > \beta$. By the triangle inequality and the choice of ϵ_j , we can easily show that

$$\hat{\lambda}_j \geq \|s^j - \hat{y}^j\|_{\hat{y}^j, t_0} - \|s^j - \hat{s}^j\|_{\hat{y}^j, t_0} \geq \lambda_j - \epsilon_j > (1 - 10^{-2})\beta.$$

In addition, from Lemma A.2 in Appendix A.3, we have

$$D_{t_0}(\hat{y}^{j+1}) \leq D_{t_0}(\hat{y}^j) - \omega\left((\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)\right),$$

where $\omega(\tau) := \tau - \ln(1 + \tau) \geq 0$. Using $\epsilon_j \leq 10^{-2}\beta$, $\delta_j \leq 10^{-2}\beta$, and $\hat{\lambda}_j \geq (1 - 10^{-2})\beta$ in the above inequality, we get

$$D_{t_0}(\hat{y}^{j+1}) \leq D_{t_0}(\hat{y}^j) - \omega(0.97\beta(1 - 10^{-2}\beta)).$$

Summing up this inequality from $j = 0$ to $j = j_{\max}$, we obtain

$$j_{\max}\omega(0.97\beta(1 - 10^{-2}\beta)) \leq D_{t_0}(\hat{y}^0) - D_{t_0}(\hat{y}^{j_{\max}}) \leq D_{t_0}(\hat{y}^0) - D_{t_0}(y_{t_0}^*),$$

which implies $j_{\max} \leq \frac{D_{t_0}(\hat{y}^0) - D_{t_0}(y_{t_0}^*)}{\omega(0.97\beta(1 - 10^{-2}\beta))}$. Consequently, we obtain (53). \square

Discussion: Theorem 5.3 shows that the number of iterations to obtain a starting point (x^0, y^0) is finite even with inexact oracles and inexact proximal-Newton methods. However, the convergence rate of Algorithm 2 is sublinear in j . If t_0 is large (i.e., close to 1), Algorithm 2 often requires a small number of iterations. Another possibility is to apply a path-following procedure as in [33] to obtain a new variant with linear convergence rate. Note that the per-iteration complexity of Algorithm 2 is essentially the same as in Algorithm 1 since the computation of $\hat{\lambda}_j$ is neglectable. In particular, if we choose $\epsilon_j = \delta_j = 0$, the steps size α_j will become the standard damped Newton step-size $\frac{1}{1+\lambda_j}$ in the theory of self-concordant function [22].

6 Numerical Experiments

We provide two numerical examples to illustrate our algorithm and compare it with some existing methods. We choose SDPT3 [30] as a common used conic solver, and Chambolle-Pock's (CP) primal-dual method [8] as one of the most powerful first-order methods that can handle our problem. The first example is the well-known network utility maximization (NUM) problem, and the second one is the spectrum management problem for multi-user DSL networks studied in [36]. Our method and the CP method are implemented in Matlab 2018b, running on a Linux server with 3.4GHz Intel Xeon E5 and 16Gb memory.

6.1 Implementation remarks

We discuss how we implement two main steps of Algorithm 1 as follows. First, we need to solve the slave problem at Step 4 up to a given accuracy δ_k such that $\delta_k \leq 10^{-2}\beta$. Solving this problem is equivalent to solving the nonlinear equation $\nabla\psi_{t_{k+1}}(x; y^k) = 0$ in x . Since $\psi_{t_{k+1}}(\cdot; y^k)$ is standard self-concordant, we can apply a damped-step Newton method to solve it. Combining this method and a warm-start strategy, we can solve this equation efficiently. Second, if $\phi = \delta_{\{b\}}$ in (2) for a given $b \in \mathbb{R}^n$, then the master problem at Step 6 reduces to a positive definite linear system $\tilde{\nabla}^2 d_{t_{k+1}}(y^k)(y - y^k) = -\tilde{\nabla} d_{t_{k+1}}(y^k) + 0.25M_{t_{k+1}}^2 b$, which can be efficiently solved by, e.g., preconditioned conjugate gradient methods. However, since ϕ usually does not have such a simple form, we need to apply iterative methods such as accelerated proximal gradient method [1, 19] to solve this problem which has a linear convergence rate. Note that we can also apply a semi-smooth Newton-type methods as in [38] to solve this problem efficiently. In our numerical test, we use FISTA which seems working well.

6.2 Network Utility Maximization

Consider a network consisting of a finite set \mathcal{S} of N nodes and a finite set \mathcal{E} of undirected capacitated edges. Let x_{ij} denote the rate of sending data from node i to node j . We assume that such a flow f_{ij} from node i to node j is fixed and unique (we usually choose f_{ij} to be the shortest path from i to j).

Assume that each node i is associated with a utility function $u_i(x_i) := \log(d_i^\top x_i + \mu_i)$, where $x_i := (x_{i1}, \dots, x_{iN})^\top$, $d_i := (d_{i1}, \dots, d_{iN})^\top$ and μ_i is a scalar. Since we ignore self-links from node i to itself, we set $d_{ii} = 0$ and $f_{ii} = \emptyset$. We further assume that the rate x_{ij} is constrained to lie in a given interval $[0, M]$, where the scalar M denotes the maximum capacity of flows.

Under this setting, we formulate the problem of interest into the following constrained convex optimization problem called NUM:

$$\begin{cases} \max_x \left\{ \sum_{i \in \mathcal{S}} \ln(d_i^\top x_i + \mu_i) - \frac{\rho}{2} \|x_i - r_i\|^2 \right\} \\ \text{s.t. } L_e \leq \sum_{e \in f_{ij}} x_{ij} \leq U_e, \quad \forall e \in \mathcal{E}, \\ 0 \leq x_{ij} \leq M, \quad \forall i, j \in \mathcal{S}. \end{cases} \quad (54)$$

Here, L_e and U_e are the lower bound and upper bound capacity of each edge, respectively, r_{ij} is the initial designed rate from node i to node j and we do not want to have the rate x_{ij} to be far away from our target r_{ij} , and ρ is the corresponding penalty parameter to control the distance from x_{ij} to r_{ij} . By defining $g(x) := -\sum_{i \in \mathcal{S}} \ln(d_i^\top x_i + \mu_i) + \frac{\rho}{2} \|x_i - r_i\|^2$, $Ax = \sum_{e \in f_{ij}} x_{ij}$, $\phi(\cdot) := \delta_{[L_e, U_e]}(\cdot)$, and $\mathcal{K} := [0, M]$, we can reformulate (54) into (2). Clearly, this problem satisfies Assumptions 2.1 and 2.2.

We implement Algorithm 1 using Algorithm 2 to find an initial point using $t_0 := 0.25$. We also implement the Chambolle-Pock method in [8] and use SDPT3 to solve (54) as our competitors. Note that SDPT3 can directly handle log-terms in g compared to other interior-point solvers such as SeDuMi, SDPA, or Mosek. To avoid solving subproblems in the Chambolle-Pock method, we reformulate (54) by introducing auxiliary variables $z_i := d_i^\top x_i + \mu_i$ for $i \in \mathcal{S}$. Since the Chambolle-Pock method has two step-sizes τ and σ , we tune τ for each run and let $\sigma := 0.99/(\tau \|K\|^2)$, where K is the linear operator obtained from reformulating (54) into a composite form. The best values of τ we found are between 10^{-6} and 10^{-7} depending on problem.

All algorithms are terminated when both infeasibility and relative duality gap reach 10^{-7} accuracy or the maximum number of iterations $k_{\max} := 20,000$ is exceeded. In the first case, we certify that the problem is “solved”, while in the second case, we mark it by “*”. If problem is too big to solve by our computer, we also mark it by “*”.

We use the “tech-router-rf” dataset from <http://networkrepository.com/tech-routers-rf.php> from [28], where we have approximately 2000 nodes and 6000 edges. In this network, each node is either a router or a computer IP. Each computer IP has to go through one or multiple routers to send data to another computer IP. For larger networks, we use the “tech-ggp” dataset from

<http://networkrepository.com/tech-pgp.php> from [6], which is a social network with approximately 11000 nodes and 24000 edges. Given a network structure, we generate the input data as follows. The initial designed rate r_i are generated from a uniform distribution $\mathcal{U}(0, 1)$ between 0 and 1. The upper and lower bounds of capacity are generated as $L_e := (1 - \mathcal{U}(0, 0.5))\bar{b}$ and $U_e := (1 + \mathcal{U}(0, 0.5))\bar{b}$, where $\bar{b} := \sum_{i \in \mathcal{S}} A_i r_i$. The maximum limit of rate M is 1 and the penalty parameter ρ is chosen to be 0.01. Both d_i and μ_i are generated randomly using $\mathcal{U}(0, 1)$. To have different problem instances, we use different sub-networks of the original one.

We run three algorithms on 10 problems instances of different sizes. The results are reported in Table 1, where n is the number of linear inequality constraints, p is the number of variables in (54), IPLD is Algorithm 1, and CP is the Chambolle-Pock method in [8].

Table 1 Numerical results of three solvers on 10 problem instances of (54).

Problem size		CPU time [s]			Feasibility violation			Objective value f^*		
n	p	IPLD	CP	SDPT3	IPLD	CP	SDPT3	IPLD	CP	SDPT3
96	17,686	0.70	3.80	4.24	4.747e-09	9.986e-08	0.000e+00	-60.1375	-60.1375	-60.1375
188	29,502	1.31	4.44	9.59	5.126e-09	9.974e-08	0.000e+00	-116.8216	-116.8216	-116.8216
239	38,050	1.85	6.64	11.75	7.227e-10	9.983e-08	0.000e+00	-171.4066	-171.4066	-171.4066
306	53,048	2.43	9.45	227.32	2.266e-10	9.995e-08	0.000e+00	-228.6001	-228.6001	-228.6001
242	72,016	2.78	10.25	809.57	9.055e-09	9.982e-08	0.000e+00	-288.6970	-288.6970	-272.1732
324	125,848	5.61	26.98	*	1.107e-08	9.987e-08	*	-569.2405	-569.2405	*
658	243,936	19.37	78.58	*	5.478e-09	9.987e-08	*	-1133.8747	-1133.8747	*
833	432,218	47.86	203.95	*	8.007e-08	9.999e-08	*	-2124.3265	-2124.3265	*
1,383	1,194,500	206.88	571.17	*	9.473e-08	9.983e-08	*	-3236.1724	-3236.1724	*
1,619	2,389,000	556.34	1297.86	*	4.422e-09	9.994e-08	*	-6474.3812	-6474.3812	*

From Table 1, we observe the following facts:

- IPLD can solve large-scale problems with huge variables and moderate number of couple linear inequality constraints relatively fast and accurate. IPLD outperforms SDPT3 and CP in a majority of problems in terms of CPU time and achieves the same accuracy in the objective value and constraint violation.
- It is not surprising that CP can also achieve high accuracy but requires very large number of iterations. The CP algorithm requires from 6500 to 15200 iterations to achieve our specified accuracy depending on problem.
- SDPT3 is quickly prohibited to handle larger instances due to the increase of variables and constraints when transforming it into a conic and log form. Therefore, the problem cannot be fit into our computer memory.

In summary, we believe that our method, IPLD, can potentially solve large-scale convex problems of the form (2) as long as they satisfy Assumptions 2.1 and 2.2. It can often achieve high accuracy within reasonably computational effort and can be easily parallelized. While primal-dual first-order methods require to tune the step-size to obtain good performance, our method is relatively robust to inexact oracles and inexact Newton-type methods as well as the choice of parameter $t_0 \in (0, 1]$.

6.3 Spectrum management of multi-user DSL networks

We consider the spectrum management problem of multi-user DSL networks studied in [36], which can be cast into the following constrained problem:

$$\begin{cases} \min_{x \in \mathbb{R}^m} \left\{ g(x) := -\sum_{i=1}^M [a_i^\top x_i - c_i^\top \ln(H_i x_i + g_i)] \right\} \\ \text{s.t.} \quad \sum_{i=1}^M x_i \leq b, \\ \quad \quad 0 \leq x_i \leq L, \quad i = 1, \dots, M. \end{cases} \quad (55)$$

where $x_i \in \mathbb{R}^m$, $a_i \in \mathbb{R}^m$, $c_i \in \mathbb{R}_+^m$, $b \in \mathbb{R}^m$, $L \in \mathbb{R}_{++}^m$, $g_i \in \mathbb{R}^m$, and $H_i \in \mathbb{R}^{m \times m}$. Here, m is the number of users, and M is the number of channels. For the detail explanation of this model, we refer the reader to [36]. Clearly, (55) can be cast into (2), where g is self-concordant, $Ax = \sum_{i=1}^M x_i$, $\phi := \delta_{(-\infty, b]}$ the indicator of $(-\infty, b]$, and $\mathcal{K} := [0, L]^M$.

Our goal in this example is to verify the performance of Algorithm 1 using different accuracy levels both for inexact oracles and inexact proximal-Newton method. For this purpose, we use two real datasets to test our algorithm. More precisely, we first fix the tolerance δ_k of the inexact oracles at 10^{-5} and change the tolerance ϵ_k of the inexact proximal-Newton method from 10^{-2} to 10^{-11} . Then, we fix the tolerance ϵ_k at 10^{-5} in the inexact proximal-Newton scheme and vary δ_k in the inexact oracles between 10^{-2} and 10^{-8} . In all these cases, we terminate our algorithm whenever the feasibility violation is below 10^{-5} and the relative gap is below 10^{-6} .

In the first test, we use a 7-user asymmetric ADSL downstream dataset, where $m = 7$ and $M = 224$. Figure 1 shows how the number of iterations and the normalized CPU time depend on the tolerances, where the normalized CPU time is computed by $(T - T_{\min})/(T_{\max} - T_{\min})$ with the time T .

We can see from the top row of Figure 1 that with $\delta_k = 10^{-5}$ fixed and $\epsilon_k \leq 10^{-4}$, the number of iterations is almost stable and the computational time does not decrease significantly. This suggests that the accuracy $\epsilon_k = 10^{-4}$ is sufficiently for computing proximal-Newton direction in the dual problem. If $\epsilon_k > 10^{-4}$, then the number of iterations and CPU time increase significantly. Similarly, if we fix $\epsilon_k = 10^{-5}$ and increase δ_k from 10^{-8} to 10^{-2} , then we can observe from the bottom row of Figure 1 that $\delta_k \leq 10^{-4}$ is sufficient to accommodate the inexact oracles.

To confirm our above statement, we again test our algorithm with the second dataset, 12-user VDSL upstream dataset, where $n = 12$ and $M = 1147$. Figure 2 provides the number of iterations and normalized CPU time by rescaling it between $[0, 1]$ as in Figure 1. We again observe very similar behavior in both situations, but since the problem is relatively larger than that of the first dataset, the computational time increases significantly when we decrease the accuracy δ_k of the inexact oracles.

Acknowledgements This work was partly supported by the National Science Foundation (NSF), awarded number: DMS-1619884.

A Appendix: The proof of technical results in the main text

We provide all the missing proofs in the main text.

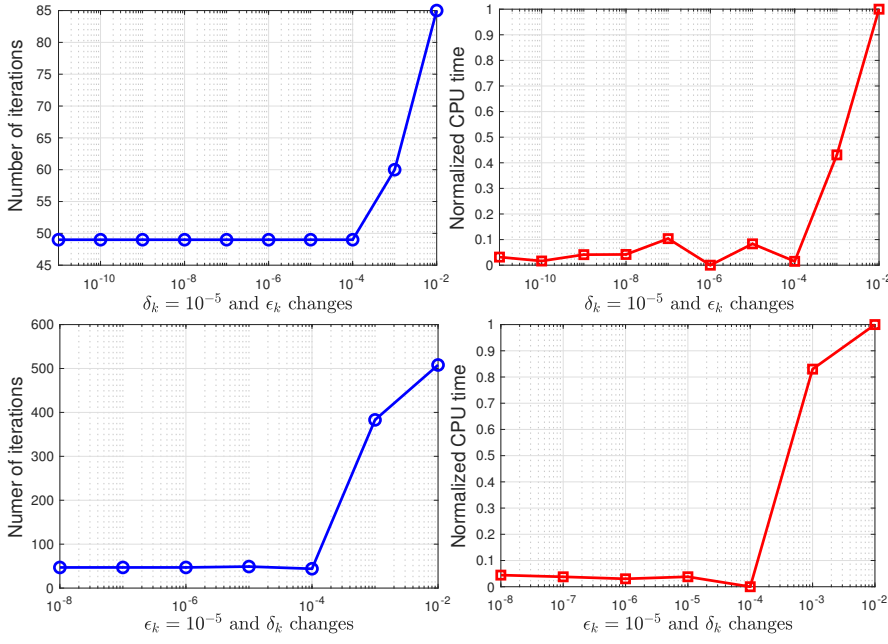


Fig. 1 The number of iterations and normalized CPU time of Algorithm 1 on the 7-users dataset. The first row shows the number of iterations and normalized CPU time when δ_k is fixed at 10^{-5} and ϵ_k changes from 10^{-12} to 10^{-2} , while the second row is for the case $\epsilon_k = 10^{-5}$ and δ_k changes from 10^{-8} to 10^{-2} .

A.1 The proof of Proposition 4.2: Properties of inexact oracles.

Since $x_t^*(y)$ is the exact solution of (18), we have

$$\nabla \psi_t(x_t^*(y); y) \equiv -\frac{M_t^2}{4} [\nabla g(x_t^*(y)) + t \nabla f(x_t^*(y)) - A^\top y] = 0.$$

Therefore, using the standard self-concordance of ψ_t , we can show that

$$\begin{aligned} \langle \nabla \psi_t(\tilde{x}_t^*(y); y), \tilde{x}_t^*(y) - x_t^*(y) \rangle &= \langle \nabla \psi_t(\tilde{x}_t^*(y); y) - \nabla \psi_t(x_t^*(y); y), \tilde{x}_t^*(y) - x_t^*(y) \rangle \\ &\geq \frac{|\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y), t}^2}{1 + |\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y), t}}. \end{aligned}$$

By the Cauchy-Schwarz inequality, we have

$$\langle \nabla \psi_t(\tilde{x}_t^*(y); y), \tilde{x}_t^*(y) - x_t^*(y) \rangle \leq |\nabla \psi_t(\tilde{x}_t^*(y); y)|_{\tilde{x}_t^*(y), t}^* |\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y), t}.$$

Combining the last two inequalities, we eventually get

$$\frac{|\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y), t}}{1 + |\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y), t}} \leq |\nabla \psi_t(\tilde{x}_t^*(y); y)|_{\tilde{x}_t^*(y), t}^* \leq \frac{\delta}{1 + \delta}.$$

This implies that $|\tilde{x}_t^*(y) - x_t^*(y)|_{\tilde{x}_t^*(y), t} \leq \delta$.

Next, using (19) and (21), we have $d_t(y) - \tilde{d}_t(y) = \psi_t(\tilde{x}_t^*(y); y) - \psi_t(x_t^*(y); y)$. Therefore, applying [19, Theorems 4.1.7 and 4.1.8] respectively, we obtain the first estimate of (24).

Note that since $\nabla^2 d_t(y) = \frac{M_t^4}{16} A \nabla^2 \psi_t(x_t^*(y); y)^{-1} A^\top$ and $\tilde{\nabla}^2 d_t(y) = \frac{M_t^4}{16} A \nabla^2 \psi_t(\tilde{x}_t^*(y); y)^{-1} A^\top$, using [19, Theorem 4.1.6], we obtain the second estimate of (24).

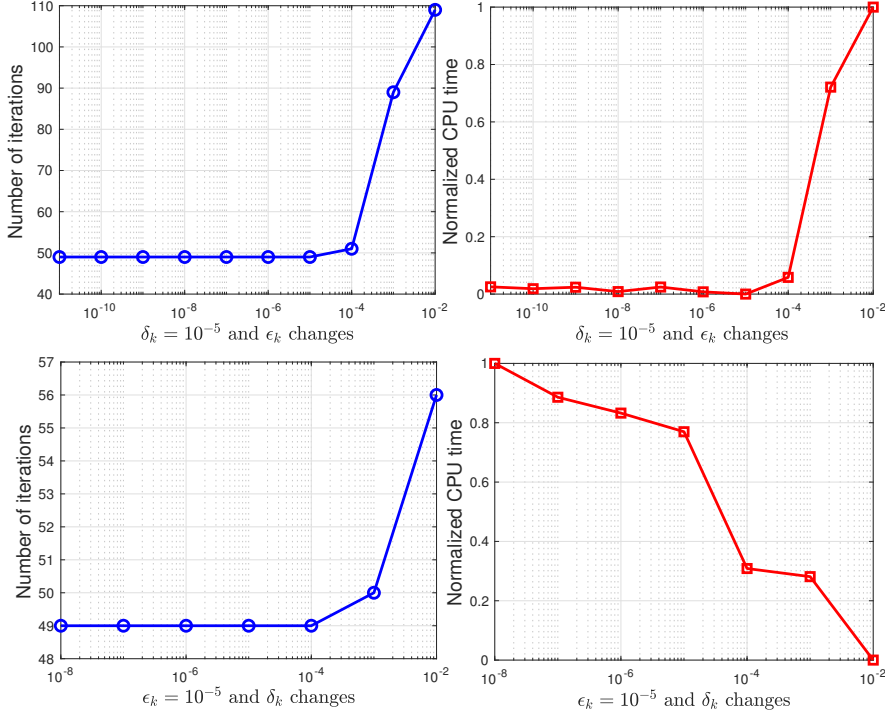


Fig. 2 The number of iterations and normalized CPU time of Algorithm 1 on the 12-users dataset. The first row shows the number of iterations and normalized CPU time when δ_k is fixed at 10^{-5} and ϵ_k changes from 10^{-12} to 10^{-2} , while the second row is for the case $\epsilon_k = 10^{-5}$ and δ_k changes from 10^{-8} to 10^{-2} .

Finally, since $\nabla d_t(y) - \tilde{\nabla} d_t(y) = \frac{M_t^2}{4} A(x_t^*(y) - \tilde{x}_t^*(y))$, we have

$$\begin{aligned}
 & [\|\nabla d_t(y) - \tilde{\nabla} d_t(y)\|_{y,t}^*]^2 \\
 &= \frac{M_t^4}{16} (x_t^*(y) - \tilde{x}_t^*(y)) A^\top \left(\frac{M_t^4}{16} A \nabla^2 \psi_t(\tilde{x}_t^*(y); y)^{-1} A^\top \right)^{-1} A (x_t^*(y) - \tilde{x}_t^*(y)) \\
 &\leq \frac{M_t^4}{16} (x_t^*(y) - \tilde{x}_t^*(y))^\top \frac{16}{M_t^4} \nabla^2 \psi_t(\tilde{x}_t^*(y); y) (x_t^*(y) - \tilde{x}_t^*(y)) \\
 &= |x_t^*(y) - \tilde{x}_t^*(y)|_{\tilde{x}_t^*(y), t}^2.
 \end{aligned}$$

In the last inequality, we use $A^\top (AQ^{-1}A^\top)^{-1}A \preceq Q$ for any symmetric positive definite matrix Q and any full-row rank matrix A . Hence, we obtain the third estimate of (24). \square

A.2 The technical proofs of Subsection 5.3: Convergence analysis

We provide the proof of technical results in Subsection 5.3.

A.2.1 The proof of Lemma 5.1: Key estimate of the inexact PN scheme (31).

For simplicity of presentation, we redefine $t := t_k$, $t_+ := t_{k+1}$, $y := y^k$, $y_+ := y^{k+1}$, and $\bar{y}_+ := \bar{y}^{k+1}$, where \bar{y}^{k+1} is defined by (28) or (30). Using these new notations, we also denote $\delta := \delta_{t_+}(y)$, $\delta_+ := \delta_{t_+}(y_+)$, $\epsilon := \|y_+ - \bar{y}_+\|_{y, t_+}$, and $\hat{\lambda} := \|y_+ - y\|_{y, t_+}$ to make our analysis more clean.

If we define $r_{t_+}(y) := \tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y) + \tilde{\nabla} d_{t_+}(y)$, then from (28), we have

$$-r_{t_+}(y) \in \partial h_{t_+}(\bar{y}_+),$$

which is equivalent to

$$\bar{y}_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} r_{t_+}(y) \in \bar{y}_+ + \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} \partial h_{t_+}(\bar{y}_+).$$

Utilizing the scaled proximal operator defined by (30), we can write the last statement as

$$\bar{y}_+ = \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y_+)}(\bar{y}_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} r_{t_+}(y)). \quad (56)$$

Using the definition of $\tilde{G}_t(y)$ in (33) and of $\lambda_t(y)$ in (34), we can derive

$$\begin{aligned} \lambda_{t_+}(y_+) &= \|\tilde{G}_{t_+}(y_+)\|_{y_+, t_+}^* \\ &= \left\| \tilde{\nabla}^2 d_{t_+}(y_+) \left[y_+ - \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y_+)}(y_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} \tilde{\nabla} d_{t_+}(y_+)) \right] \right\|_{y_+, t_+}^* \\ &= \left\| y_+ - \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y_+)}(y_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} \tilde{\nabla} d_{t_+}(y_+)) \right\|_{y_+, t_+} \\ &\leq \|y_+ - \bar{y}_+\|_{y_+, t_+} + \left\| \bar{y}_+ - \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y_+)}(y_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} \tilde{\nabla} d_{t_+}(y_+)) \right\|_{y_+, t_+} \\ &\stackrel{(56)}{=} \|y_+ - \bar{y}_+\|_{y_+, t_+} + \left\| \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y_+)}(\bar{y}_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} r_{t_+}(y)) \right. \\ &\quad \left. - \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y_+)}(y_+ - \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} \tilde{\nabla} d_{t_+}(y_+)) \right\|_{y_+, t_+}. \end{aligned}$$

By the non-expansiveness of $\text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}}(\cdot)$, see [32], we can further estimate this term as

$$\begin{aligned} \lambda_{t_+}(y_+) &\leq \|y_+ - \bar{y}_+\|_{y_+, t_+} + \left\| \bar{y}_+ - y_+ + \tilde{\nabla}^2 d_{t_+}(y_+)^{-1} (\tilde{\nabla} d_{t_+}(y_+) - r_{t_+}(y)) \right\|_{y_+, t_+} \\ &\leq 2\|y_+ - \bar{y}_+\|_{y_+, t_+} + \|\tilde{\nabla} d_{t_+}(y_+) - \tilde{\nabla} d_{t_+}(y) - \tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y)\|_{y_+, t_+}^*. \end{aligned} \quad (57)$$

Next, we decompose the following term $R_{t_+}(y)$ as

$$\begin{aligned} R_{t_+}(y) &:= \tilde{\nabla} d_{t_+}(y_+) - \tilde{\nabla} d_{t_+}(y) - \tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y) \\ &= [\tilde{\nabla} d_{t_+}(y_+) - \nabla d_{t_+}(y_+)] - [\tilde{\nabla} d_{t_+}(y) - \nabla d_{t_+}(y)] \\ &\quad - [\tilde{\nabla}^2 d_{t_+}(y) - \nabla^2 d_{t_+}(y)](y_+ - y) - \tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y_+) \\ &\quad + [\nabla d_{t_+}(y_+) - \nabla d_{t_+}(y) - \nabla^2 d_{t_+}(y)(y_+ - y)]. \end{aligned} \quad (58)$$

Before we estimate the five terms of $R_{t_+}(y)$, we recall the following inequalities, which will be repeatedly used in our proof.

$$\frac{1}{1 - \|y_+ - y\|_{y, t_+}} \stackrel{(24)}{\leq} \frac{1}{1 - \frac{1}{1 - \delta_{t_+}(y)} \|y_+ - y\|_{y, t_+}} = \frac{1 - \delta}{1 - \delta - \hat{\lambda}}. \quad (59)$$

$$\|\cdot\|_{y_+, t_+}^* \stackrel{(24)}{\leq} \frac{\|\cdot\|_{y_+, t_+}^*}{1 - \delta_{t_+}(y_+)} \leq \frac{\|\cdot\|_{y, t_+}^*}{(1 - \delta_{t_+}(y_+))(1 - \|y_+ - y\|_{y, t_+})} \stackrel{(59)}{\leq} \frac{(1 - \delta)\|\cdot\|_{y, t_+}^*}{(1 - \delta_+)(1 - \delta - \hat{\lambda})}. \quad (60)$$

Here, the second last inequality of (60) is from [19, Theorem 4.1.6]. Note that (60) also holds for $\|\cdot\|_{y_+, t_+}$ and $\|\cdot\|_{y, t_+}$.

Using (60), we have

$$\|\cdot\|_{y_+, t_+}^* \stackrel{(60)}{\leq} \frac{(1 - \delta)\|\cdot\|_{y, t_+}^*}{(1 - \delta_+)(1 - \delta - \hat{\lambda})} \stackrel{(24)}{\leq} \frac{\|\cdot\|_{y, t_+}^*}{(1 - \delta_+)(1 - \delta - \hat{\lambda})}. \quad (61)$$

Note that (61) also holds for $\|\cdot\|_{y_+, t_+}$ and $\|\cdot\|_{y, t_+}$.

Now, we estimate the first term in $R_{t_+}(y)$ of (58) as

$$\|\tilde{\nabla} d_{t_+}(y_+) - \nabla d_{t_+}(y_+)\|_{y_+, t_+}^* \stackrel{(24)}{\leq} \delta_{t_+}(y_+) = \delta_{t_+}. \quad (62)$$

For the second term of (58), we have

$$\begin{aligned} \|\tilde{\nabla} d_{t_+}(y) - \nabla d_{t_+}(y)\|_{y_+, t_+}^* &\stackrel{(61)}{\leq} \frac{1}{(1-\delta_+)(1-\delta-\hat{\lambda})} \|\tilde{\nabla} d_{t_+}(y) - \nabla d_{t_+}(y)\|_{y, t_+}^* \\ &\stackrel{(24)}{\leq} \frac{\delta_{t_+}(y)}{(1-\delta_+)(1-\delta-\hat{\lambda})} \\ &= \frac{\delta}{(1-\delta_+)(1-\delta-\hat{\lambda})}. \end{aligned} \quad (63)$$

To estimate the third term of (58), let $S(y) := [\tilde{\nabla}^2 d_{t_+}(y) - \nabla^2 d_{t_+}(y)](y_+ - y)$. We have

$$\|S(y)\|_{y_+, t_+}^* \stackrel{(60)}{\leq} \frac{(1-\delta)\|S(y)\|_{y, t_+}^*}{(1-\delta_+)(1-\delta-\hat{\lambda})} \quad (64)$$

However, $\|S(y)\|_{y, t_+}^*$ can be estimated as

$$\begin{aligned} \left[\|S(y)\|_{y, t_+}^*\right]^2 &= (y_+ - y)^\top [\tilde{\nabla}^2 d_{t_+}(y) - \nabla^2 d_{t_+}(y)] \nabla^2 d_{t_+}(y)^{-1} [\tilde{\nabla}^2 d_{t_+}(y) - \nabla^2 d_{t_+}(y)] (y_+ - y) \\ &= (y_+ - y)^\top \tilde{\nabla}^2 d_{t_+}(y) \nabla^2 d_{t_+}(y)^{-1} \tilde{\nabla}^2 d_{t_+}(y) (y_+ - y) \\ &\quad - 2(y_+ - y)^\top \tilde{\nabla}^2 d_{t_+}(y) (y_+ - y) + (y_+ - y)^\top \nabla^2 d_{t_+}(y) (y_+ - y) \\ &\stackrel{(24)}{\leq} \left[\frac{2}{(1-\delta_{t_+}(y))^2} - 2 \right] (y_+ - y)^\top \tilde{\nabla}^2 d_{t_+}(y) (y_+ - y) \\ &= \frac{4\delta-2\delta^2}{(1-\delta)^2} \|y_+ - y\|_{y, t_+}^2 = \frac{4\delta-2\delta^2}{(1-\delta)^2} \hat{\lambda}^2 \end{aligned}$$

Using this estimate into (64), we finally get

$$\begin{aligned} \|\tilde{\nabla}^2 d_{t_+}(y) - \nabla^2 d_{t_+}(y)\|_{y_+, t_+}^* &\leq \frac{1-\delta}{(1-\delta_+)(1-\hat{\lambda}-\delta)} \frac{\sqrt{4\delta-2\delta^2}}{(1-\delta)} \hat{\lambda} \\ &= \frac{\sqrt{4\delta-2\delta^2} \hat{\lambda}}{(1-\delta_+)(1-\hat{\lambda}-\delta)}. \end{aligned} \quad (65)$$

For the fourth term $\tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y_+)$ of (58), we have

$$\begin{aligned} \|\tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y_+)\|_{y_+, t_+}^* &\stackrel{(61)}{\leq} \frac{1}{(1-\delta_+)(1-\hat{\lambda}-\delta)} \|\tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y_+)\|_{y, t_+}^* \\ &= \frac{1}{(1-\delta_+)(1-\hat{\lambda}-\delta)} \|\bar{y}_+ - y_+\|_{y, t_+} \\ &= \frac{\epsilon}{(1-\delta_+)(1-\hat{\lambda}-\delta)}. \end{aligned} \quad (66)$$

Finally, we estimate last terms $\mathcal{T}_5 := \|\nabla d_{t_+}(y_+) - \nabla d_{t_+}(y) - \nabla^2 d_{t_+}(y)(y_+ - y)\|_{y_+, t_+}^*$.

Note that

$$\begin{aligned} \mathcal{T}_5 &\stackrel{(60)}{\leq} \frac{1-\delta}{(1-\delta_+)(1-\delta-\hat{\lambda})} \|\nabla d_{t_+}(y_+) - \nabla d_{t_+}(y) - \nabla^2 d_{t_+}(y)(y_+ - y)\|_{y, t_+}^* \\ &\leq \frac{1-\delta}{(1-\delta_+)(1-\delta-\hat{\lambda})} \left(\frac{\|y_+ - y\|_{y, t_+}^2}{1 - \|y_+ - y\|_{y, t_+}} \right) \\ &\stackrel{(59)}{\leq} \frac{(1-\delta)^2}{(1-\delta_+)(1-\delta-\hat{\lambda})^2} \|y_+ - y\|_{y, t_+}^2 \\ &\stackrel{(24)}{\leq} \frac{\hat{\lambda}^2}{(1-\delta_+)(1-\delta-\hat{\lambda})^2}, \end{aligned} \quad (67)$$

where the second inequality follows from [35, Theorem 1].

Plugging (62), (63), (65), (66), and (67) into (58), we can estimate

$$\begin{aligned} \|R_{t_+}(y)\|_{y_+, t_+}^* &\leq \delta_+ + \frac{\delta}{(1-\delta_+)(1-\hat{\lambda}-\delta)} + \frac{\sqrt{4\delta-2\delta^2}\hat{\lambda}}{(1-\delta_+)(1-\hat{\lambda}-\delta)} \\ &\quad + \frac{\epsilon}{(1-\delta_+)(1-\hat{\lambda}-\delta)} + \frac{\hat{\lambda}^2}{(1-\delta_+)(1-\delta-\hat{\lambda})^2}. \end{aligned} \quad (68)$$

Note that

$$\|y_+ - \bar{y}_+\|_{y_+, t_+} \stackrel{(61)}{\leq} \frac{\|y_+ - \bar{y}_+\|_{y, t_+}}{(1-\delta_+)(1-\delta-\hat{\lambda})} = \frac{\epsilon}{(1-\delta_+)(1-\delta-\hat{\lambda})}.$$

Substituting this estimate and (68) into (57), we finally obtain

$$\begin{aligned} \lambda_{t_+}(y) &\leq \frac{3\epsilon}{(1-\delta_+)(1-\delta-\hat{\lambda})} + \delta_+ + \frac{\delta}{(1-\delta_+)(1-\hat{\lambda}-\delta)} \\ &\quad + \frac{\sqrt{4\delta-2\delta^2}\hat{\lambda}}{(1-\delta_+)(1-\hat{\lambda}-\delta)} + \frac{\hat{\lambda}^2}{(1-\delta_+)(1-\delta-\hat{\lambda})^2}. \end{aligned} \quad (69)$$

However, from the definition of $\lambda_{t_+}(y)$, we have

$$\begin{aligned} \lambda_{t_+}(y) &:= \left\| \tilde{\nabla}^2 d_{t_+}(y) \left(y - \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y)} (y - \tilde{\nabla}^2 d_{t_+}(y)^{-1} \tilde{\nabla} d_{t_+}(y)) \right) \right\|_{y, t_+}^* \\ &= \| \bar{y}_+ - y \|_{y, t_+} \\ &\geq \| y_+ - y \|_{y, t_+} - \| y_+ - \bar{y}_+ \|_{y, t_+} \\ &= \hat{\lambda} - \epsilon. \end{aligned}$$

This implies $\hat{\lambda} := \|y_+ - y\|_{y, t_+} \leq \lambda_{t_+}(y) + \epsilon$. Substituting this estimate into (69), we obtain (37). In particular, if $\delta = \delta_+ = \epsilon = 0$, then we can simplify (37) to obtain (38). \square

A.2.2 The proof of Lemma 5.2: The relationship between $\lambda_{t_+}(y)$ and $\tilde{\Delta}$.

We again redefine $t := t_k$, $t_+ := t_{k+1}$, $y := y^k$, $y_+ := y^{k+1}$, and $\bar{y}_+ := \bar{y}^{k+1}$ as in Lemma 5.1. In addition, we also define $\bar{u} := \text{prox}_{h_t}^{\tilde{\nabla}^2 d_t(y)} (y - \tilde{\nabla}^2 d_t(y)^{-1} \tilde{\nabla} d_t(y))$.

First, we show that $\lambda_{t_+}(y)$ and $\lambda_t(y)$ can be respectively expressed as

$$\begin{aligned} \lambda_{t_+}(y) &:= \left\| \tilde{\nabla}^2 d_{t_+}(y) \left(y - \text{prox}_{h_{t_+}}^{\tilde{\nabla}^2 d_{t_+}(y)} (y - \tilde{\nabla}^2 d_{t_+}(y)^{-1} \tilde{\nabla} d_{t_+}(y)) \right) \right\|_{y, t_+}^* \\ &= \| \bar{y}_+ - y \|_{y, t_+}, \\ \lambda_t(y) &:= \left\| \tilde{\nabla}^2 d_t(y) \left(y - \text{prox}_{h_t}^{\tilde{\nabla}^2 d_t(y)} (y - \tilde{\nabla}^2 d_t(y)^{-1} \tilde{\nabla} d_t(y)) \right) \right\|_{y, t}^* \\ &= \| \bar{u} - y \|_{y, t}. \end{aligned} \quad (70)$$

If we denote by $\bar{h}(y) := \phi^*(-y)$, then $h_t(y) = \frac{M_t^2}{4} \bar{h}(y)$. By the definition of \bar{u} and \bar{y}_+ , we can write

$$\begin{cases} -\frac{4}{M_t^2} [\tilde{\nabla}^2 d_t(y)(\bar{u} - y) + \tilde{\nabla} d_t(y)] &\in \partial \bar{h}(\bar{u}), \\ -\frac{4}{M_{t_+}^2} [\tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y) + \tilde{\nabla} d_{t_+}(y)] &\in \partial \bar{h}(\bar{y}_+). \end{cases}$$

Using the monotonicity of $\partial \bar{h}(\cdot)$, we can show that

$$\left\langle \frac{4}{M_t^2} \tilde{\nabla}^2 d_t(y)(\bar{u} - y) - \frac{4}{M_{t_+}^2} \tilde{\nabla}^2 d_{t_+}(y)(\bar{y}_+ - y) + \frac{4}{M_t^2} \tilde{\nabla} d_t(y) - \frac{4}{M_{t_+}^2} \tilde{\nabla} d_{t_+}(y), \bar{y}_+ - \bar{u} \right\rangle \geq 0.$$

Rearranging this inequality, we obtain

$$\begin{aligned} & \left\langle \frac{4}{M_t^2} \tilde{\nabla}^2 d_t(y)(\bar{u} - y) - \frac{4}{M_{t+}^2} \tilde{\nabla}^2 d_{t+}(y)(\bar{u} - y) + \frac{4}{M_t^2} \tilde{\nabla} d_t(y) - \frac{4}{M_{t+}^2} \tilde{\nabla} d_{t+}(y), \bar{y}_+ - \bar{u} \right\rangle \\ & \geq \frac{4}{M_{t+}^2} \|\bar{y}_+ - \bar{u}\|_{y, t+}^2. \end{aligned}$$

By the Cauchy-Schwarz inequality, we can derive that

$$\begin{aligned} & \left\| \overbrace{\frac{4}{M_t^2} \tilde{\nabla}^2 d_t(y)(\bar{u} - y) - \frac{4}{M_{t+}^2} \tilde{\nabla}^2 d_{t+}(y)(\bar{u} - y)}^{\mathcal{T}_1} + \overbrace{\frac{4}{M_t^2} \tilde{\nabla} d_t(y) - \frac{4}{M_{t+}^2} \tilde{\nabla} d_{t+}(y)}^{\mathcal{T}_2} \right\|_{y, t+}^* \\ & \geq \frac{4}{M_{t+}^2} \|\bar{y}_+ - \bar{u}\|_{y, t+}. \end{aligned} \quad (71)$$

To estimate \mathcal{T}_1 , we first show the relationship between $\nabla^2 \psi_t(\tilde{x}_t^*(y))$ and $\nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y))$. Then, we use it to get the relationship between $\tilde{\nabla}^2 d_t(y)$ and $\tilde{\nabla}^2 d_{t+}(y)$. Recall that $\nabla^2 \psi_t(x) := \frac{M_t^2}{4} [\nabla^2 g(x) + t \nabla^2 f(x)]$. Moreover, if $t \in [0, 1]$, then $\frac{M_t^2}{4} := \max\{1, \frac{1}{t}\} = \frac{1}{t}$. Therefore, we can write

$$\nabla^2 \psi_t(x) = \frac{1}{t} \nabla^2 g(x) + \nabla^2 f(x) \quad \text{and} \quad \nabla^2 \psi_{t+}(x) = \frac{1}{t+} \nabla^2 g(x) + \nabla^2 f(x).$$

For any $0 \leq t_+ \leq t \leq 1$, we have

$$\nabla^2 \psi_t(x) \preceq \nabla^2 \psi_{t+}(x) \preceq \frac{1}{t+} \nabla^2 g(x) + \frac{t}{t+} \nabla^2 f(x) = \frac{t}{t+} \nabla^2 \psi_t(x). \quad (72)$$

In addition, using the self-concordance of ψ_t and (72), we also have

$$\begin{aligned} \nabla^2 \psi_t(\tilde{x}_t^*(y)) & \preceq \frac{1}{(1 - |\tilde{x}_{t+}^*(y) - \tilde{x}_t^*(y)|_{\tilde{x}_t^*(y), t})^2} \nabla^2 \psi_t(\tilde{x}_{t+}^*(y)) \\ & \stackrel{(72)}{\preceq} \frac{1}{(1-\Delta)^2} \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y)), \\ \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y)) & \stackrel{(72)}{\preceq} \frac{t}{t+} \nabla^2 \psi_t(\tilde{x}_{t+}^*(y)) \preceq \frac{t}{t+(1 - |\tilde{x}_{t+}^*(y) - \tilde{x}_t^*(y)|_{\tilde{x}_t^*(y), t})^2} \nabla^2 \psi_t(\tilde{x}_t^*(y)) \\ & = \frac{t}{t+(1-\Delta)^2} \nabla^2 \psi_t(\tilde{x}_t^*(y)). \end{aligned} \quad (73)$$

If we take the inverse of both sides of (73), then we get

$$\begin{cases} \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y))^{-1} \preceq \frac{1}{(1-\Delta)^2} \nabla^2 \psi_t(\tilde{x}_t^*(y))^{-1}, \\ \nabla^2 \psi_t(\tilde{x}_t^*(y))^{-1} \preceq \frac{t}{t+(1-\Delta)^2} \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y))^{-1}. \end{cases}$$

Since $\tilde{\nabla}^2 d_t(y) = \frac{M_t^4}{16} A \nabla^2 \psi_t(\tilde{x}_t^*(y))^{-1} A^\top$ and $\tilde{\nabla}^2 d_{t+}(y) = \frac{M_{t+}^4}{16} A \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y))^{-1} A^\top$, the last inequalities imply

$$\frac{16}{M_{t+}^4} \tilde{\nabla}^2 d_{t+}(y) \preceq \frac{16}{M_t^4 (1-\Delta)^2} \tilde{\nabla}^2 d_t(y) \quad \text{and} \quad \frac{16}{M_t^4} \tilde{\nabla}^2 d_t(y) \preceq \frac{16t}{M_{t+}^4 t+(1-\Delta)^2} \tilde{\nabla}^2 d_{t+}(y),$$

which are respectively equivalent to

$$\tilde{\nabla}^2 d_{t+}(y) \preceq \frac{M_{t+}^4}{M_t^4 (1-\Delta)^2} \tilde{\nabla}^2 d_t(y) \quad \text{and} \quad \tilde{\nabla}^2 d_t(y) \preceq \frac{M_t^4 t}{M_{t+}^4 t+(1-\Delta)^2} \tilde{\nabla}^2 d_{t+}(y).$$

Since $\frac{M_t^2}{4} = \frac{1}{t}$ and $\frac{M_{t+}^2}{4} = \frac{1}{t+}$, we obtain from the above inequalities that

$$\tilde{\nabla}^2 d_{t+}(y) \preceq \frac{t^2}{t_+^2(1-\tilde{\Delta})^2} \tilde{\nabla}^2 d_t(y) \quad \text{and} \quad \tilde{\nabla}^2 d_t(y) \preceq \frac{t_+}{t(1-\tilde{\Delta})^2} \tilde{\nabla}^2 d_{t+}(y). \quad (74)$$

Now we can estimate the first term \mathcal{T}_1 in (71) as

$$\begin{aligned} [\|\mathcal{T}_1\|_{y,t+}^*]^2 &= \left[\left\| \frac{4}{M_t^2} \tilde{\nabla}^2 d_t(y)(\bar{u} - y) - \frac{4}{M_{t+}^2} \tilde{\nabla}^2 d_{t+}(y)(\bar{u} - y) \right\|_{y,t+}^* \right]^2 \\ &= \left[\left\| t \tilde{\nabla}^2 d_t(y)(\bar{u} - y) - t_+ \tilde{\nabla}^2 d_{t+}(y)(\bar{u} - y) \right\|_{y,t+}^* \right]^2 \\ &= (\bar{u} - y)^\top \left([t \tilde{\nabla}^2 d_t(y) - t_+ \tilde{\nabla}^2 d_{t+}(y)] \tilde{\nabla}^2 d_{t+}(y)^{-1} [t \tilde{\nabla}^2 d_t(y) - t_+ \tilde{\nabla}^2 d_{t+}(y)] \right) (\bar{u} - y) \\ &= (\bar{u} - y)^\top \left(t^2 \tilde{\nabla}^2 d_t(y) \tilde{\nabla}^2 d_{t+}(y)^{-1} \tilde{\nabla}^2 d_t(y) - 2tt_+ \tilde{\nabla}^2 d_t(y) + t_+^2 \tilde{\nabla}^2 d_{t+}(y) \right) (\bar{u} - y) \\ &\stackrel{(74)}{\leq} (\bar{u} - y)^\top \left(\frac{t_+ t}{(1-\tilde{\Delta})^2} - 2tt_+ + \frac{t^2}{(1-\tilde{\Delta})^2} \tilde{\nabla}^2 d_t(y) \right) (\bar{u} - y). \\ &= \frac{t^2 - 2t_+ t(1-\tilde{\Delta})^2 + tt_+}{(1-\tilde{\Delta})^2} \|\bar{u} - y\|_{y,t}^2. \end{aligned} \quad (75)$$

To estimate the second term \mathcal{T}_2 of (71), by the definition of $\tilde{\nabla} d_t$, we have

$$\begin{aligned} [\|\mathcal{T}_2\|_{y,t+}^*]^2 &= \left[\left\| \frac{4}{M_t^2} \tilde{\nabla} d_t(y) - \frac{4}{M_{t+}^2} \tilde{\nabla} d_{t+}(y) \right\|_{y,t+}^* \right]^2 \\ &= \left[\left\| A \tilde{x}_t^*(y) - A \tilde{x}_{t+}^*(y) \right\|_{y,t+}^* \right]^2 \\ &= (\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y))^\top A^\top \tilde{\nabla}^2 d_{t+}(y)^{-1} A (\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y)) \\ &\stackrel{(21)}{=} \frac{16}{M_{t+}^4} (\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y))^\top A^\top \left(A \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y))^{-1} A^\top \right)^{-1} A (\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y)) \\ &\leq \frac{16}{M_{t+}^4} (\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y))^\top \nabla^2 \psi_{t+}(\tilde{x}_{t+}^*(y)) (\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y)) \\ &= \frac{16}{M_{t+}^4} |\tilde{x}_t^*(y) - \tilde{x}_{t+}^*(y)|_{\tilde{x}_{t+}^*(y), t+}^2 = \frac{16}{M_{t+}^4} \tilde{\Delta}_+^2. \end{aligned} \quad (76)$$

Here, we use the fact that $A^\top (AQ^{-1}A^\top)^{-1}A \preceq Q$ for any symmetric positive definite matrix Q and any full-row rank matrix A .

Plugging (75) and (76) into (71), we get

$$\frac{4}{M_{t+}^2} \|\bar{y}_+ - \bar{u}\|_{y,t+} \leq \frac{4}{M_{t+}^2} \tilde{\Delta}_+ + \frac{\sqrt{t^2 - 2t_+ t(1-\tilde{\Delta})^2 + tt_+}}{1-\tilde{\Delta}} \|\bar{u} - y\|_{y,t}.$$

This inequality is equivalent to

$$\|\bar{y}_+ - \bar{u}\|_{y,t+} \leq \tilde{\Delta}_+ + \frac{\sqrt{(\frac{t}{t_+})^2 - 2\frac{t}{t_+}(1-\tilde{\Delta})^2 + \frac{t}{t_+}}}{1-\tilde{\Delta}} \|\bar{u} - y\|_{y,t}. \quad (77)$$

Finally, we can derive

$$\begin{aligned} \lambda_{t+}(y) &\stackrel{(70)}{=} \|\bar{y}_+ - y\|_{y,t+} \\ &\leq \|\bar{y}_+ - \bar{u}\|_{y,t+} + \|\bar{u} - y\|_{y,t+} \\ &\stackrel{(77)}{\leq} \tilde{\Delta}_+ + \frac{\sqrt{(\frac{t}{t_+})^2 - 2\frac{t}{t_+}(1-\tilde{\Delta})^2 + \frac{t}{t_+}}}{1-\tilde{\Delta}} \lambda_t(y) + \|\bar{u} - y\|_{y,t+} \\ &\stackrel{(74)}{\leq} \tilde{\Delta}_+ + \frac{\sqrt{(\frac{t}{t_+})^2 - 2\frac{t}{t_+}(1-\tilde{\Delta})^2 + \frac{t}{t_+}}}{1-\tilde{\Delta}} \lambda_t(y) + \frac{t}{t_+(1-\tilde{\Delta})} \|\bar{u} - y\|_{y,t} \\ &\stackrel{(70)}{\leq} \tilde{\Delta}_+ + \left[\frac{\sqrt{(\frac{t}{t_+})^2 - 2\frac{t}{t_+}(1-\tilde{\Delta})^2 + \frac{t}{t_+}}}{1-\tilde{\Delta}} + \frac{t}{t_+} \right] \lambda_t(y), \end{aligned} \quad (78)$$

which is exactly (40) due to the update $t_+ := \sigma t$. \square

In order to prove Lemma 5.3 we need the following auxiliary result.

Lemma A.1 *Let $a \in (0, 1)$ and $b \in (0, 1)$ be two positive numbers such that $a + b < 1$. Let*

$$\mathcal{N}(a, b) := \left\{ (u, v) \in \mathbb{R}_+^2 \mid \frac{u^2}{1+u} \leq au + bv, \frac{v^2}{1+v} \leq av + bu \right\}.$$

Then, $\mathcal{N}(a, b) \subseteq \left\{ (u, v) \in \mathbb{R}_+^2 \mid u \leq \frac{a+b}{1-a-b}, v \leq \frac{a+b}{1-a-b} \right\}$.

Proof Suppose $(u, v) \in \mathcal{N}(a, b)$ and $u > \frac{a+b}{1-a-b}$. Then, according to $\frac{u^2}{1+u} \leq au + bv$, we have

$$v \geq \frac{1}{b} \left(\frac{u^2}{1+u} - au \right) = \frac{u}{b} \left(\frac{u}{1+u} - a \right) > \frac{u}{b} \left(\frac{\frac{a+b}{1-a-b}}{1 + \frac{a+b}{1-a-b}} - a \right) = u > \frac{a+b}{1-a-b}. \quad (79)$$

Therefore, we can show that

$$\frac{1}{b} \left(\frac{v^2}{1+v} - av \right) = \frac{v}{b} \left(\frac{v}{1+v} - a \right) \stackrel{(79)}{>} \frac{v}{b} \left(\frac{\frac{a+b}{1-a-b}}{1 + \frac{a+b}{1-a-b}} - a \right) = v. \quad (80)$$

However, because $\frac{v^2}{1+v} \leq av + bu$, one can show that

$$u \geq \frac{1}{b} \left(\frac{v^2}{1+v} - av \right) \stackrel{(80)}{>} v.$$

This contradicts (79). Consequently, we must have $u \leq \frac{a+b}{1-a-b}$. Use the symmetry between u and v , we also have $v \leq \frac{a+b}{1-a-b}$. \square

A.2.3 The proof of Lemma 5.3: Upper bound on the solution difference $\tilde{\Delta}$.

First, by the self-concordance of ψ_t , we have

$$\begin{aligned} \frac{\tilde{\Delta}^2}{1+\tilde{\Delta}} &\leq \langle \tilde{x}_{t_+}^*(y) - \tilde{x}_t^*(y), \nabla \psi_t(\tilde{x}_{t_+}^*(y)) - \nabla \psi_t(\tilde{x}_t^*(y)) \rangle \\ &\leq \tilde{\Delta} |\nabla \psi_t(\tilde{x}_t^*(y))|_{\tilde{x}_t^*(y), t}^* + \tilde{\Delta}_+ |\nabla \psi_t(\tilde{x}_{t_+}^*(y))|_{\tilde{x}_{t_+}^*(y), t_+}^*. \end{aligned} \quad (81)$$

Next, since $\nabla \psi_t(x) = \frac{1}{t} \nabla g(x) + \nabla f(x)$ and $\nabla \psi_{t_+}(x) = \frac{1}{t_+} \nabla g(x) + \nabla f(x)$, we have

$$\nabla \psi_t(\tilde{x}_{t_+}^*(y)) = \frac{t_+}{t} \nabla \psi_{t_+}(\tilde{x}_{t_+}^*(y)) + \frac{(t-t_+)}{t} \nabla f(\tilde{x}_{t_+}^*(y)).$$

Therefore, we can bound

$$\begin{aligned} |\nabla \psi_t(\tilde{x}_{t_+}^*(y))|_{\tilde{x}_{t_+}^*(y), t_+}^* &\leq \frac{t_+}{t} |\nabla \psi_{t_+}(\tilde{x}_{t_+}^*(y))|_{\tilde{x}_{t_+}^*(y), t_+}^* \\ &\quad + \frac{(t-t_+)}{t} |\nabla f(\tilde{x}_{t_+}^*(y))|_{\tilde{x}_{t_+}^*(y), t_+}^*. \end{aligned} \quad (82)$$

Now, since $\nabla^2 \psi_{t_+}(x) \succeq \nabla^2 f(x)$, we can show that

$$|\nabla f(\tilde{x}_{t_+}^*(y))|_{\tilde{x}_{t_+}^*(y), t_+}^* \leq \left[\nabla f(\tilde{x}_{t_+}^*(y))^\top \nabla^2 f(\tilde{x}_{t_+}^*(y))^{-1} \nabla f(\tilde{x}_{t_+}^*(y)) \right]^{1/2} \leq \sqrt{\nu_f}.$$

Substituting this and (82) into (81), and using the definition of $\hat{\delta}_+$ and $\hat{\delta}$, we get

$$\frac{\tilde{\Delta}^2}{1 + \tilde{\Delta}} \leq \tilde{\Delta}\hat{\delta} + \left(\frac{\hat{\delta}_+ t_+}{t} + \frac{(t - t_+)}{t} \sqrt{\nu_f} \right) \tilde{\Delta}_+.$$

Finally, using $t_+ = \sigma t$, we obtain the first estimate of (42) from the last inequality.

Similarly, by following the same argument as in the proof of the first estimate in (42), we can show that

$$\frac{\tilde{\Delta}_+^2}{1 + \tilde{\Delta}_+} \leq \tilde{\Delta}_+ \hat{\delta}_+ + \left(\frac{\hat{\delta} t}{t_+} + \frac{(t - t_+)}{t_+} \sqrt{\nu_f} \right) \tilde{\Delta},$$

which is the second estimate of (42).

Assume that we choose $\hat{\delta} \leq \delta$ and $\hat{\delta}_+ \leq \delta$ for some $\delta \in (0, 1)$. Since $t_+ = \sigma t$, if we denote by $c_\nu(\sigma) := \frac{\delta}{\sigma} + \frac{(1-\sigma)}{\sigma} \sqrt{\nu_f} \in (0, 1)$. Assume further that $\delta + c_\nu(\sigma) < 1$. Then, it is clear that $\frac{\hat{\delta}_+ t_+}{t} + \frac{(t - t_+)}{t} \sqrt{\nu_f} \leq c_\nu(\sigma)$ and $\frac{\hat{\delta} t}{t_+} + \frac{(t - t_+)}{t_+} \sqrt{\nu_f} \leq c_\nu(\sigma)$. Applying Lemma A.1, we can see that $(\tilde{\Delta}, \tilde{\Delta}_+) \in \mathcal{N}(\delta, c_\nu(\sigma))$. Hence, we have

$$\tilde{\Delta} \leq \frac{\delta + c_\nu(\sigma)}{1 - \delta - c_\nu(\sigma)} \quad \text{and} \quad \tilde{\Delta}_+ \leq \frac{\delta + c_\nu(\sigma)}{1 - \delta - c_\nu(\sigma)},$$

which proves (44). \square

A.3 The proof of result in Subsection 5.5: Finding an initial point.

The proof of Theorem 5.3 requires the following key lemma.

Lemma A.2 *Let $\{\hat{y}_j\}$ be the sequence generated by Algorithm 2, where the step-size α_j is chosen such that $\alpha_j \in (0, 1]$ and $\frac{\alpha_j \hat{\lambda}_j}{1 - \delta_j} < 1$. Then, the following estimate holds*

$$D_{t_0}(\hat{y}^{j+1}) \leq D_{t_0}(\hat{y}^j) - \alpha_j [\hat{\lambda}_j^2 - (\epsilon_j + \delta_j) \hat{\lambda}_j] + \omega_* \left(\frac{\alpha_j \hat{\lambda}_j}{1 - \delta_j} \right), \quad (83)$$

where D_t is defined in (14) and $\omega_*(\tau) := -\tau - \ln(1 - \tau)$. The optimal step-size α_j that minimizes the right-hand side of (83) is

$$\alpha_j := \frac{(\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)^2}{\left[1 + (1 - \delta_j)(\hat{\lambda}_j - \epsilon_j - \delta_j) \right] \hat{\lambda}_j} \in (0, 1). \quad (84)$$

The corresponding estimate from (83) with this step-size is

$$D_{t_0}(\hat{y}^{j+1}) \leq D_{t_0}(\hat{y}^j) - \omega \left((\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j) \right). \quad (85)$$

In particular, if we set $\delta_j = \epsilon_j = 0$, then we get the original damped-step proximal-Newton step-size $\alpha_j = \frac{1}{1 + \lambda_j}$ and the estimate $D_{t_0}(\hat{y}^{j+1}) \leq D_{t_0}(\hat{y}^j) - \omega(\lambda_j)$ for $\omega(\tau) := \tau - \ln(1 + \tau)$.

Proof Firstly, from the self concordance of d_{t_0} defined in (14) and $\hat{y}_{j+1} = (1 - \alpha)\hat{y}_j + \alpha\hat{s}_j$, we can show that

$$\begin{aligned} d_{t_0}(\hat{y}^{j+1}) + h_{t_0}(\hat{y}^{j+1}) &\leq d_{t_0}(\hat{y}^j) + \langle \nabla d_{t_0}(\hat{y}^j), \hat{y}^{j+1} - \hat{y}^j \rangle + \omega_*(\|\hat{y}^{j+1} - \hat{y}^j\|_{\hat{y}^j, t_0}) \\ &\quad + (1 - \alpha)h_{t_0}(\hat{y}^j) + \alpha h_{t_0}(\hat{s}^j) \\ &= d_{t_0}(\hat{y}^j) + \alpha \langle \nabla d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle + \omega_*(\alpha \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0}) \\ &\quad + (1 - \alpha)h_{t_0}(\hat{y}^j) + \alpha h_{t_0}(\hat{s}^j) \\ &= (1 - \alpha)(d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{y}^j)) + \omega_*(\alpha \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0}) \\ &\quad + \alpha(d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{s}^j) + \langle \nabla d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle). \end{aligned} \quad (86)$$

Next, we will prove that

$$d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{s}^j) + \langle \nabla d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle \leq d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{y}^j) - \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2} + \delta_j \hat{\lambda}_j, \quad (87)$$

where $\lambda_j := \|\hat{y}^j - s^j\|_{\hat{y}^j, t_0}$.

Indeed, by the Cauchy-Schwarz inequality, we have

$$\begin{aligned} \langle \nabla d_{t_0}(\hat{y}^j) - \tilde{\nabla} d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle &\leq \|\nabla d_{t_0}(\hat{y}^j) - \tilde{\nabla} d_{t_0}(\hat{y}^j)\|_{\hat{y}^j, t_0}^* \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0} \\ &\leq \delta_j \hat{\lambda}_j. \end{aligned} \quad (88)$$

Since

$$\hat{s}^j \approx s^j := \text{prox}_{h_{t_0}}^{\tilde{\nabla}^2 d_{t_0}(\hat{y}^j)} \left(\hat{y}^j - \tilde{\nabla}^2 d_{t_0}(\hat{y}^j)^{-1} \tilde{\nabla} d_{t_0}(\hat{y}^j) \right),$$

we have

$$\begin{aligned} \langle \tilde{\nabla} d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle + h_{t_0}(\hat{s}^j) &\leq \langle \tilde{\nabla} d_{t_0}(\hat{y}^j), s^j - \hat{y}^j \rangle + h_{t_0}(s^j) \\ &\quad + \frac{1}{2} \|s^j - \hat{y}^j\|_{\hat{y}^j, t_0}^2 - \frac{1}{2} \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0}^2 + \frac{\epsilon_j^2}{2} \\ &= \langle \tilde{\nabla} d_{t_0}(\hat{y}^j), s^j - \hat{y}^j \rangle + h_{t_0}(s^j) + \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2}. \end{aligned} \quad (89)$$

Using $0 \in \tilde{\nabla} d_{t_0}(\hat{y}^j) + \tilde{\nabla}^2 d_{t_0}(s^j - \hat{y}^j) + \partial h_{t_0}(s^j)$, we can further estimate

$$\begin{aligned} \langle \tilde{\nabla} d_{t_0}(\hat{y}^j), s^j - \hat{y}^j \rangle + h_{t_0}(s^j) &+ \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2} \\ &= \langle -\tilde{\nabla}^2 d_{t_0}(s^j - \hat{y}^j) - \nabla h_{t_0}(s^j), s^j - \hat{y}^j \rangle + h_{t_0}(s^j) + \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2} \\ &= \langle \nabla h_{t_0}(s^j), \hat{y}^j - s^j \rangle + h_{t_0}(s^j) - \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2} \\ &\leq h_{t_0}(\hat{y}^j) - \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2}, \end{aligned} \quad (90)$$

where $\nabla h_{t_0}(s^j) \in \partial h_{t_0}(s^j)$. Combining (89) and (90), we get

$$\langle \tilde{\nabla} d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle + h_{t_0}(\hat{s}^j) \leq h_{t_0}(\hat{y}^j) - \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2}. \quad (91)$$

Now, we can prove (87) as follows:

$$\begin{aligned} d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{s}^j) + \langle \nabla d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle &= d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{s}^j) + \langle \tilde{\nabla} d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle + \langle \nabla d_{t_0}(\hat{y}^j) - \tilde{\nabla} d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle \\ &\stackrel{(88)}{\leq} d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{s}^j) + \langle \tilde{\nabla} d_{t_0}(\hat{y}^j), \hat{s}^j - \hat{y}^j \rangle + \delta_j \hat{\lambda}_j \\ &\stackrel{(91)}{\leq} d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{y}^j) - \frac{1}{2}\lambda_j^2 - \frac{1}{2}\hat{\lambda}_j^2 + \frac{\epsilon_j^2}{2} + \delta_j \hat{\lambda}_j. \end{aligned}$$

Combining (86) and (87), and notice that $\omega_*(\alpha \|\hat{s}^j - \hat{y}^j\|_{\hat{y}^j, t_0}) \leq \omega_*(\frac{\alpha \hat{\lambda}_j}{1 - \delta_j})$ we can deduce

$$d_{t_0}(\hat{y}^{j+1}) + h_{t_0}(\hat{y}^{j+1}) \leq d_{t_0}(\hat{y}^j) + h_{t_0}(\hat{y}^j) - \alpha \left(\frac{\lambda_j^2}{2} + \frac{\hat{\lambda}_j^2}{2} - \frac{\epsilon_j^2}{2} - \delta_j \hat{\lambda}_j \right) + \omega_* \left(\frac{\alpha \hat{\lambda}_j}{1 - \delta_j} \right).$$

Using the fact that $\lambda_j \geq \hat{\lambda}_j - \epsilon_j$ and the definition $D_{t_0} := d_{t_0} + h_{t_0}$, we obtain (83).

Next, if we maximize $\zeta(\alpha) := \alpha[\hat{\lambda}_j^2 - (\epsilon_j + \delta_j)\hat{\lambda}_j] - \omega_* \left(\frac{\alpha \hat{\lambda}_j}{1 - \delta_j} \right)$, we have $\alpha^* := \frac{(\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)^2}{[1 + (1 - \delta_j)(\hat{\lambda}_j - \epsilon_j - \delta_j)]\hat{\lambda}_j}$ as defined by (84). Plugging α^* into $\zeta(\alpha)$, we get

$$\zeta(\alpha^*) = \frac{(\hat{\lambda}_j - \epsilon_j - \delta_j)^2(1 - \delta_j)^2}{1 + (1 - \delta_j)(\hat{\lambda}_j - \epsilon_j - \delta_j)} - \omega_* \left(\frac{(\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)}{1 + (\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)} \right).$$

Since $\frac{x^2}{1+x} - \omega_*\left(\frac{x}{1+x}\right) = \omega(x)$, we finally have $\zeta(\alpha^*) = \omega\left((\hat{\lambda}_j - \epsilon_j - \delta_j)(1 - \delta_j)\right)$, which proves (85). If $\delta_j = \epsilon_j = 0$, then α_j reduces to $\frac{1}{1+\lambda_j}$ and we obtain $D_{t_0}(\hat{y}^{j+1}) \leq D_{t_0}(\hat{y}^j) - \omega(\lambda_j)$ from (85). \square

References

1. A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sci.*, 2(1):183–202, 2009.
2. A. Ben-Tal and A. Nemirovski. *Lectures on modern convex optimization: Analysis, algorithms, and engineering applications*, volume 3 of *MPS/SIAM Series on Optimization*. SIAM, 2001.
3. D.P. Bertsekas and J. N. Tsitsiklis. *Parallel and distributed computation: Numerical methods*. Prentice Hall, 1989.
4. J.R. Birge. Decomposition and Partitioning Methods for Multistage Stochastic Linear Programs. *Operations Research*, 33(5):989–1007, 1985.
5. A. Bitlislioglu, I. Pejic, and C. Jones. Interior-point decomposition for multi-agent optimization. In *20th IFAC World Congress*, number EPFL-CONF-228343, 2017.
6. Marián Boguñá, Romualdo Pastor-Satorras, Albert Díaz-Guilera, and Alex Arenas. Models of social networks based on social distance attachment. *Physical Review E*, 70(5):056122, 2004.
7. S. Boyd and L. Vandenberghe. *Convex Optimization*. University Press, Cambridge, 2004.
8. A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *J. Math. Imaging Vis.*, 40(1):120–145, 2011.
9. A.J. Conejo, R. Mínguez, E. Castillo, and R. García-Bertrand. *Decomposition Techniques in Mathematical Programming: Engineering and Science Applications*. Springer-Verlag, 2006.
10. G. B. Dantzig. *Linear Programming and Extensions*. Princeton University Press, 1963.
11. O. Devolder, F. Glineur, and Y. Nesterov. First-order methods of smooth convex optimization with inexact oracle. *Math. Program.*, 146(1–2):37–75, 2014.
12. M. Fukuda and M. Kojima. *Interior-point methods for Lagrangian duals of semidefinite programs*. Inst. of Technology, 2000.
13. M. Fukuda, M. Kojima, and M. Shida. Lagrangian dual interior-point methods for semidefinite programs. *SIAM J. Optim.*, 12:1007–1031, 2002.
14. S. Gros. A newton algorithm for distributed semi-definite programs using the primal-dual interior-point method. In *53rd IEEE Conference on Decision and Control*, pages 3222–3227. IEEE, 2014.
15. Bjarni V Halldórsson and Reha H Tütüncü. An interior-point method for a class of saddle-point problems. *Journal of Optimization Theory and Applications*, 116(3):559–590, 2003.
16. M. Kojima, N. Megiddo, S. Mizuno, and et al. Horizontal and vertical decomposition in interior point methods for linear programs. Technical report., Information Sciences, Tokyo Institute of Technology, Tokyo, 1993.
17. X. Li, D. Sun, and K.-C. Toh. A highly efficient semismooth Newton augmented Lagrangian method for solving Lasso problems. *SIAM J. Optim.*, 28(1):433–458, 2018.
18. I. Necoara and J.A.K. Suykens. Interior-point Lagrangian decomposition method for separable convex optimization. *J. Optim. Theory and Appl.*, 143(3):567–588, 2009.
19. Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, 2004.
20. Y. Nesterov. Barrier subgradient method. *Math. Program., Ser. B*, 127:31–56, 2011.
21. Y. Nesterov. Gradient methods for minimizing composite objective function. *Math. Program.*, 140(1):125–161, 2013.
22. Y. Nesterov and A. Nemirovski. *Interior-point Polynomial Algorithms in Convex Programming*. Society for Industrial Mathematics, 1994.
23. Y. Nesterov and J.-Ph. Vial. Augmented self-concordant barriers and nonlinear optimization problems with finite complexity. *Math. Program.*, 99:149–174, 2004.

24. S. K. Pakazad, A. Hansson, and M. S. Andersen. Distributed primal-dual interior-point methods for solving loosely coupled problems using message passing. *Optim. Method Softw.*, 32(3):401–435, 2017.
25. D.P. Palomar and M. Chiang. A Tutorial on Decomposition Methods for Network Utility Maximization. *IEEE J. Selected Areas in Communications*, 24(8):1439–1451, 2006.
26. R.T. Rockafellar. *Convexity and Duality in Optimization*, chapter Monotropic Programming: A generalization of linear programming and network programming., pages 10–036. Springer-Verlag, 1985.
27. M. Shida. An interior-point smoothing technique for Lagrangian relaxation in large-scale convex programming. *Optimization*, 57(1):183–200, 2008.
28. N. Spring, R. Mahajan, and D. Wetherall. Measuring ISP topologies with rocketfuel. In *SIGCOMM*, volume 32, pages 133–145, 2002.
29. T. Sun and Q. Tran-Dinh. Generalized Self-Concordant Functions: A Recipe for Newton-Type Methods. *Math. Program. (online first)*, pages 1–63, 2018.
30. K.-Ch. Toh, M.J. Todd, and R.H. Tütüncü. On the implementation and usage of SDPT3 – a Matlab software package for semidefinite-quadratic-linear programming. Tech. Report 4, NUS Singapore, 2010.
31. Q. Tran-Dinh, A. Kyrillidis, and V. Cevher. An inexact proximal path-following algorithm for constrained convex minimization. *SIAM J. Optim.*, 24(4):1718–1745, 2014.
32. Q. Tran-Dinh, A. Kyrillidis, and V. Cevher. Composite self-concordant minimization. *J. Mach. Learn. Res.*, 15:374–416, 2015.
33. Q. Tran-Dinh, A. Kyrillidis, and V. Cevher. A single phase proximal path-following framework. *Math. Oper. Res.*, 43(4):1326–1347, 2018.
34. Q. Tran-Dinh, I. Necoara, C. Savorgnan, and M. Diehl. An inexact perturbed path-following method for Lagrangian decomposition in large-scale separable convex optimization. *SIAM J. Optim.*, 23(1):95–125, 2013.
35. Q. Tran-Dinh, T. Sun, and S. Lu. Self-concordant inclusions: A unified framework for path-following generalized Newton-type algorithms. *Math. Program. (online first)*, pages 1–51, 2018.
36. P. Tsiaflakis, M. Diehl, and M. Moonen. Distributed spectrum management algorithms for multi-user DSL networks. *IEEE Transactions on Signal Processing*, 56(10):4825–4843, 2008.
37. M. Yamashita, K. Fujisawa, M. Fukuda, K. Kobayashi, K. Nakta, and M. Nakata. *Handbook on Semidefinite, Cone and Polynomial Optimization: Theory, Algorithms, Software and Applications*, chapter Latest developments in the SDPA Family for solving large-scale SDPs, pages 687–714. Springer-Verlag, New York, USA, 2011.
38. L. Yang, D. Sun, and K.-C. Toh. SDPNAL+: a majorized semismooth Newton-CG augmented Lagrangian method for semidefinite programming with nonnegative constraints. *Math. Program. Comput.*, 7(3):331–366, 2015.
39. G. Zhao. A Lagrangian dual method with self-concordant barriers for multistage stochastic convex programming. *Math. Program.*, 102:1–24, 2005.
40. X.-Y. Zhao, D. Sun, and K.-C. Toh. A Newton-CG augmented Lagrangian method for semidefinite programming. *SIAM J. Optim.*, 20(4):1737–1765, 2010.