

Error Bound Can Give Near Optimal Convergence Rate for Inexact Accelerated Proximal Gradient Method

Author 1 Name, Author 2 Name ^{*}

November 14, 2025

This paper is currently in draft mode. Check source to change options.

Abstract

This is still a draft. [6].

2010 Mathematics Subject Classification: Primary 47H05, 52A41, 90C25; Secondary 15A09, 26A51, 26B25, 26E60, 47H09, 47A63. **Keywords:**

1 Introduction

Notations. Let $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$, we denote g^* to be the Fenchel conjugate. $I : \mathbb{R}^n \rightarrow \mathbb{R}^n$ denotes the identity operator. For a multivalued mapping $T : \mathbb{R}^n \rightarrow 2^{\mathbb{R}^n}$, $\text{gra } T$ denotes the graph of the operator, defined as $\{(x, y) \in \mathbb{R}^n \times \mathbb{R}^n : y \in Tx\}$. Π_S denotes the projection onto a set $S \subseteq \mathbb{R}^n$.

{def:esp-subgrad} 1.1 Preliminaries

Definition 1.1 (ϵ -subgradient) Let $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be proper, lsc. Let $\epsilon \geq 0$. Then the ϵ -subgradient of g at some $\bar{x} \in \text{dom } g$ is given by:

$$\partial g_\epsilon(\bar{x}) := \{v \in \mathbb{R}^n \mid \langle v, x - \bar{x} \rangle \leq g(x) - g(\bar{x}) + \epsilon \forall x \in \mathbb{R}^n\}.$$

*Subject type, Some Department of Some University, Location of the University, Country. E-mail: author.namee@university.edu.

When $\bar{x} \notin \text{dom } g$, it has $\partial g_\epsilon(\bar{x}) = \emptyset$.

Remark 1.2 $\partial_\epsilon g$ is a multivalued operator and, it's not monotone, unless $\epsilon = 0$, which makes it equivalent to Fenchel subgradient ∂g .

If we assume lsc, proper and convex g , we will now introduce results in the literatures that we will use.

Fact 1.3 (ϵ -Fenchel inequality) Let $\epsilon \geq 0$, then:

$$x^* \in \partial_\epsilon f(\bar{x}) \iff f^*(x^*) + f(\bar{x}) \leq \langle x^*, \bar{x} \rangle + \epsilon \implies \bar{x} \in \partial_\epsilon f^*(x^*).$$

They are all equivalent if $f^{**}(\bar{x}) = f(\bar{x})$.

Remark 1.4 The above fact is taken from Zalinascu [5, Theorem 2.4.2].

We will now define inexact proximal point based on ϵ -subgradient

Definition 1.5 (inexact proximal point) For all $x \in \mathbb{R}^n, \epsilon \geq 0, \lambda > 0$, \tilde{x} is an inexact evaluation of proximal point at x , if and only if it satisfies:

$$\lambda^{-1}(x - \tilde{x}) \in \partial_\epsilon g(\tilde{x}).$$

We denote it by $\tilde{x} \approx_\epsilon \text{prox}_{\lambda g}(x)$.

Remark 1.6 This definition is nothing new, for example see Villa et al. [4, Definition 2.1]

Fact 1.7 (the resolvant identity) Let $T : \mathbb{R}^n \rightarrow 2^{\mathbb{R}^n}$, then it has:

$$(I + T)^{-1} = (I - (I + T^{-1})^{-1}).$$

Theorem 1.8 (inexact Moreau decomposition) Let $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be a closed, convex and proper function. It has the equivalence

$$\tilde{y} \approx_\epsilon \text{prox}_{\lambda^{-1}g^*}(\lambda^{-1}y) \iff y - \lambda\tilde{y} \approx_\epsilon \text{prox}_{\lambda g}(y).$$

Proof. Consider $\tilde{y} \approx_{\epsilon} \text{prox}_{\lambda^{-1}g^*}(\lambda^{-1}y)$, then it has:

$$\begin{aligned}
& \tilde{y} \in (I + \lambda^{-1}\partial_{\epsilon}g^*)^{-1}(\lambda^{-1}y) \\
\iff & (\lambda^{-1}y, \tilde{y}) \in \text{gra}(I + \lambda^{-1}\partial_{\epsilon}g^*)^{-1} \\
\stackrel{(1)}{\iff} & (\lambda^{-1}y, \tilde{y}) \in \text{gra}(I - (I + \partial_{\epsilon}g \circ (\lambda I))^{-1}) \\
\iff & (\lambda^{-1}y, \lambda^{-1}y - \tilde{y}) \in \text{gra}(I + \partial_{\epsilon}g \circ (\lambda I))^{-1} \\
\iff & (\lambda^{-1}y - \tilde{y}, \lambda^{-1}y) \in \text{gra}(I + \partial_{\epsilon}g \circ (\lambda I)) \\
\iff & (y - \lambda\tilde{y}, \lambda^{-1}y) \in \text{gra}(\lambda^{-1}I + \partial_{\epsilon}g) \\
\iff & (y - \lambda\tilde{y}, y) \in \text{gra}(I + \lambda\partial_{\epsilon}g) \\
\iff & y - \lambda\tilde{y} \in (I + \lambda\partial_{\epsilon}g)^{-1}y \\
\iff & y - \lambda\tilde{y} \approx_{\epsilon} \text{prox}_{\lambda g}(y).
\end{aligned}$$

At (1) we can use Fact 1.7, and it has $(\lambda^{-1}\partial_{\epsilon}g^*)^{-1} = \partial_{\epsilon}g \circ (\lambda I)$ by Fact 1.3 and the assumption that g is closed, convex and proper. \blacksquare

1.2 Inexact proximal gradient inequality

{ass:for-inxt-pg-ineq}

Assumption 1.9 (for inexact proximal gradient) The assumption is about (F, f, g, L) . We assume that

- (i) $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex, L Lipschitz smooth function (i.e: ∇f is L a Lipschitz continuous mapping).
- (ii) $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is a convex, proper, and lsc function which we do not have its exact proximal operator.
- (iii) The over all objective is $F = f + g$.

No, we develop the theory based on the use of epsilon subgradient as in Definition 1.1. Let $\rho > 0$, the exact proximal gradient operator defined for (f, g, L) satisfying Assumption 1.9 has

$$\begin{aligned}
T_{\rho}(x) &= \underset{z \in \mathbb{R}^n}{\text{argmin}} \left\{ g(z) + \langle \nabla f(x), z \rangle + \frac{\rho}{2} \|z - x\|^2 \right\} \\
&= \text{prox}_{\rho^{-1}g}(x - \rho^{-1}\nabla f(x)).
\end{aligned}$$

The following definition extends the proximal gradient operator to the inexact case using the concept of ϵ -subgradient as given by Definition 1.1.

Definition 1.10 (inexact proximal gradient) Let (f, g, L) satisfies Assumption 1.9. Let $\epsilon \geq 0, \rho > 0$. Then, $\tilde{x} \approx_{\epsilon} T_{\rho}(x)$ is an inexact proximal gradient if it satisfies variational

inequality:

$$\mathbf{0} \in \nabla f(x) + \rho(x - \tilde{x}) + \partial_\epsilon g(\tilde{x}).$$

Remark 1.11 We assumed that we can get exact evaluation of ∇f at any points $x \in \mathbb{R}^n$.

{lemma:other-repr-inxt-pg}

Lemma 1.12 (other representations of inexact proximal gradient)

Let (f, g, L) satisfies Assumption 1.9, $\epsilon \geq 0, \rho > 0$, then for all $\tilde{x} \approx_\epsilon T_\rho(x)$, it has the following equivalent representations:

$$\begin{aligned} (x - \rho^{-1}\nabla f(x)) - \tilde{x} &\in \rho^{-1}\partial_\epsilon g(\tilde{x}) \\ \iff \tilde{x} &\in (I + \rho^{-1}\partial_\epsilon g(\tilde{x}))^{-1}(x - \rho^{-1}\nabla f(x)) \\ \iff x &\approx_\epsilon \text{prox}_{\rho^{-1}g}(x - \rho^{-1}\nabla f(x)) \end{aligned}$$

Proof. It's direct. ■

{thm:inxt-pg-ineq}

Theorem 1.13 (inexact over-regularized proximal gradient inequality)

Let (F, f, g, L) satisfies Assumption 1.9, $\epsilon \geq 0, B \geq 0, \rho > 0$. Consider $\tilde{x} \approx_\epsilon T_{B+\rho}(x)$. Denote $F = f + g$. If in addition, \tilde{x}, B satisfies the line search condition $D_f(\tilde{x}, x) \leq B/2\|x - \tilde{x}\|^2$, then it has $\forall z \in \mathbb{R}^n$:

$$-\epsilon \leq F(z) - F(\tilde{x}) + \frac{B + \rho}{2}\|x - z\|^2 - \frac{B + \rho}{2}\|z - \tilde{x}\|^2 - \frac{\rho}{2}\|\tilde{x} - x\|^2.$$

Proof. By Definition 1.10 write the variational inequality that describes $\tilde{x} \approx_\epsilon T_B(x)$, and the definition of epsilon subgradient (Definition 1.1) it has for all $z \in \mathbb{R}^n$:

$$\begin{aligned} -\epsilon &\leq g(z) - g(\tilde{x}) - \langle (B + \rho)(\tilde{x} - x) - \nabla f(x), z - \tilde{x} \rangle \\ &= g(z) - g(\tilde{x}) - (B + \rho)\langle \tilde{x} - x, z - \tilde{x} \rangle + \langle \nabla f(x), z - \tilde{x} \rangle \\ &\stackrel{(1)}{\leq} g(z) + f(z) - g(\tilde{x}) - f(\tilde{x}) - (B + \rho)\langle \tilde{x} - x, z - \tilde{x} \rangle - D_f(z, x) + D_f(\tilde{x}, x) \\ &\stackrel{(2)}{\leq} F(z) - F(\tilde{x}) - (B + \rho)\langle \tilde{x} - x, z - \tilde{x} \rangle + \frac{B}{2}\|\tilde{x} - x\|^2 \\ &= F(z) - F(\tilde{x}) + \frac{B + \rho}{2}(\|x - z\|^2 - \|\tilde{x} - x\|^2 - \|z - \tilde{x}\|^2) + \frac{B}{2}\|\tilde{x} - x\|^2 \\ &= F(z) - F(\tilde{x}) + \frac{B + \rho}{2}\|x - z\|^2 - \frac{B + \rho}{2}\|z - \tilde{x}\|^2 - \frac{\rho}{2}\|\tilde{x} - x\|^2. \end{aligned}$$

At (1), we used considered the following:

$$\begin{aligned} \langle \nabla f(x), z - x \rangle &= \langle \nabla f(x), z - x + x - \tilde{x} \rangle \\ &= \langle \nabla f(x), z - x \rangle + \langle \nabla f(x), x - \tilde{x} \rangle \\ &= -D_f(z, x) + f(z) - f(x) + D_f(\tilde{x}, x) - f(\tilde{x}) + f(x) \\ &= -D_f(z, x) + f(z) + D_f(\tilde{x}, x) - f(\tilde{x}). \end{aligned}$$

At (2), we used the fact that f is convex hence $-D_f(z, x) \leq 0$ always, and in the statement hypothesis we assumed that B has $D_f(\tilde{x}, x) \leq B/2\|\tilde{x} - x\|^2$. We also used $F = f + g$. ■

Remark 1.14 When $\epsilon = 0, \rho = 0$, this reduces to proximal gradient inequality in the exact case. In this inequality, observe that the parameter ϵ controls the inexactness of the proximal gradient evaluation. More specifically, ϵ_k controls the absolute perturbations of the proximal gradient inequality compared to its exact counterpart. ρ on the other hand, it is the over-relaxation of proximal gradient operator, and it compensates the perturbations caused by \approx_ϵ relative to the term $\|\tilde{x} - x\|^2$.

1.3 Optimizing an inexact proximal point problem

{sec:optz-inxt-pp-problem}

In this section we will present the optimization problem that obtains a \tilde{x} such that $\tilde{x} \approx_{\epsilon} \text{prox}_{\lambda g}(z)$. Eventually we want to evaluate $T_\rho(x)$ of some $F = f + g$ inexactly using Lemma 1.12. To do that one would need to evaluate $\text{prox}_{\rho^{-1}g}$ inexactly which is defined in Definition 1.5.

Most of these results that will follow are from the literature. To start, we must assume the following about a function $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$, with g closed, convex and proper.

{ass:for-inxt-prox}

Assumption 1.15 (linear composite of convex nonsmooth function)

This assumption is about $(g, \omega, A, K_\omega, K_{\omega^*})$. Let $m \in \mathbb{N}, n \in \mathbb{R}^n$, we assume that

- (i) $A \in \mathbb{R}^{m \times n}$ is a matrix.
- (ii) $\omega : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is a closed and convex function such that it has exact proximal operator $\text{prox}_{\lambda \omega}$ and, with known conjugate ω^* is known. In addition assume that $\text{dom } \omega = \mathbb{R}^n$ so that
- (iii) $g : \mathbb{R}^m \rightarrow \mathbb{R} := \omega(Ax)$ and by previous item it satisfies constraint qualification $\text{rng } A \cap \text{ri dom } \omega \neq \emptyset$.
- (iv) ω is K_ω Lipschitz continuous, and ω^* is K_{ω^*} Lipschitz continuous and in addition, $\text{dom } \omega^*$ is a compact set. As an immediate consequence, the operator $\partial\omega$ is a bounded on \mathbb{R}^n .

Now, we are ready to discuss how to choose $\tilde{x} \approx_{\epsilon} \text{prox}_{\lambda g}(x)$. Fix $y \in \mathbb{R}^n, \lambda > 0$, we are ultimately interested in minimizing:

{eqn:primal-pp}

$$\Phi_\lambda(u) := \omega(Au) + \frac{1}{2\lambda}\|u - y\|^2 \quad (1.1)$$

Observe that $\text{rng } A \cap \text{ri dom } g \neq \emptyset$ in Assumption 1.15 shows g is also closed convex and proper. The function Φ_λ is coersive due to its quadratic term and hence it must admit a

minimizer [3, Theorem 1.9]. Recall the following famous theoretical result in the convex programming literature that we had adapted into our context.

{fact:fn-rck-duality}

Fact 1.16 (Fenchel Rockafellar Duality [1, Proposition 15.22])

Let $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be closed convex and proper, $g : \mathbb{R}^m \rightarrow \overline{\mathbb{R}}$, $A \in \mathbb{R}^{m \times n}$. If $\mathbf{0} \in \text{int}(\text{dom } g - A \text{ dom } f)$, then

$$\inf_{u \in \mathbb{R}^n} \{f(u) + g(Au)\} + \min_{v \in \mathbb{R}^m} \{f^* \circ (-A^\top) + g^*(v)\} = 0.$$

Remark 1.17 The theorem is not exactly the same as what is claimed in the original text, since we are in a finite dimensional setting. To see the original theorem cited in the finite dimension, we need the space $\mathcal{H} = \mathbb{R}^n$ and then use [1, Proposition 6.12].

In our context, we are interested in the dual of the proximal problem Φ_λ which makes $f = u \mapsto \frac{1}{2\lambda}\|u - y\|^2$, and $g = \omega$, and it has $f^*(v) = \frac{1}{2\lambda}\|\lambda v + y\|^2 - \frac{1}{2\lambda}\|y\|^2$ (see Appendix A.1). Consequently, $[f^* \circ (-A^\top)](v) = \frac{1}{2\lambda}\|-\lambda A^\top v + y\|^2 - \frac{1}{2\lambda}\|y\|^2$. And therefore, Φ_λ admits Fenchel Rockafellar dual (or simply the dual) objective in \mathbb{R}^m :

{eqn:dual-pp}

$$\Psi_\lambda(v) := f^* \circ (-A^\top) + g^*(v) = \frac{1}{2\lambda}\|\lambda A^\top v - y\|^2 + \omega^*(v) - \frac{1}{2\lambda}\|y\|^2. \quad (1.2)$$

We define the duality gap

{eqn:duality-gap-pp}

$$\mathbf{G}_\lambda(u, v) := \Phi_\lambda(u) + \Psi_\lambda(v). \quad (1.3)$$

Note that in this case the smooth part is quadratic and, $\text{dom } f = \mathbb{R}^n$, hence it translates to $\mathbf{0} \in \text{int}(\text{dom } g - A \text{ dom } f) = \text{int}(\text{dom } g - \text{rng } A)$. This will hold because of $\text{rng } A \cap \text{ri dom } g \neq \emptyset$ in Assumption 1.15. Therefore strong duality holds and it exists (\hat{u}, \hat{v}) such that we have the following:

$$\mathbf{G}_\lambda(\hat{u}, \hat{v}) = 0 = \min_u \Phi_\lambda(u) + \min_v \Psi_\lambda(v)$$

{thm:primal-dual-trans}

The following theorem quantifies a sufficient conditions for $\tilde{x} \approx_\epsilon \text{prox}_{\lambda g}(x)$. The theorem below is from [4, Proposition 2.2].

Theorem 1.18 (primal translate to dual [4, Proposition 2.2]) Let (g, ω, A) satisfies assumption 1.15, $\epsilon \geq 0$, then

$$(\forall z \approx_\epsilon \text{prox}_{\lambda g}(y)) (\exists v \in \text{dom } \omega^*) : z = y - \lambda A^\top v.$$

This theorem that follows is from Villa et al. [4, Proposition 2.3].

{thm:dlty-gap-inxt-pp}

Theorem 1.19 (duality gap of inexact proximal problem [4, Proposition 2.3])

Let (g, ω, A) satisfies Assumption 1.15, for all $\epsilon \geq 0$, $v \in \mathbb{R}^n$ consider the following conditions:

- (i) $\mathbf{G}_\lambda(y - \lambda A^\top v, v) \leq \epsilon.$
- (ii) $A^\top v \approx_\epsilon \text{prox}_{\lambda^{-1}g^*}(\lambda^{-1}y).$
- (iii) $y - \lambda A^\top v \approx_\epsilon \text{prox}_{\lambda g}(y).$

They have (a) \implies (b) \iff (c). If in addition $\omega^*(v) = g^*(A^\top v)$, then all three conditions are equivalent. \blacksquare

Proof. The proof of (a) \implies (b), and the case when (a) \iff (b), we refer readers to Villa et al. [4, Proposition 2.3], and to show (b) \iff (c) use Theorem 1.8. \blacksquare

The following theorem is enhanced from Villa et al. [4, Theorem 5.1].

Un nouvel candidat ce que remplace le denière théorem suivant.

{thm:minimizing-dual-pp}

Theorem 1.20 (minimizing the dual of the proximal problem)

Assume that we have (g, ω, A) given by Assumption 1.15. Let the Φ_λ be given by (1.1), and dual Ψ_λ by (1.2). Let \bar{v} be a minimizer of Ψ_λ and suppose that sequence $(v_j)_{j \geq 0}$ minimizes dual Ψ_λ . Then, the followings are true:

- (i) If \bar{v} is an minimizer of dual Ψ_λ , then $\bar{z} = y - \lambda A^\top \bar{v}$ is a minimizer of primal Φ_λ .
- (ii) Let $z_j = y - \lambda A^\top v_j$, it has $\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}) \geq \frac{1}{2\lambda} \|z_j - \bar{z}\|^2$, and therefore it has $z_j \rightarrow \bar{z}$.
- (iii) And, by our assumption, the primal optimality gap is bounded by dual for all $v \in \mathbb{R}^m$ by the inequality:

$$\begin{aligned} & \Phi_\lambda(z_j) - \Phi_\lambda(\bar{z}) \\ & \leq \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + \lambda^{-1} \|z_j - y\| + \frac{1}{2\lambda} \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \right). \end{aligned}$$

Proof. There are two intermediate results to prior to proving the items in the theorem. In preparations, for all $v \in \mathbb{R}^m$, it's obvious that we have:

$$\begin{aligned} & \frac{1}{2\lambda} \|\lambda A^\top v - y\|^2 - \frac{1}{2\lambda} \|\lambda A^\top \bar{v} - y\|^2 + \langle A\bar{z}, v - \bar{v} \rangle \\ &= \frac{1}{2\lambda} \|\lambda A^\top v - \lambda A^\top \bar{v} + \lambda A^\top \bar{v} - y\|^2 - \frac{1}{2\lambda} \|\lambda A^\top \bar{v} - y\|^2 + \langle A\bar{z}, v - \bar{v} \rangle \\ &= \frac{1}{2\lambda} \|\lambda A^\top (v - \bar{v})\|^2 + \frac{1}{\lambda} \langle \lambda A^\top (v - \bar{v}), \lambda A^\top \bar{v} - y \rangle + \langle \bar{z}, A^\top (v - \bar{v}) \rangle \quad (1.4) \\ &\stackrel{(1)}{=} \frac{1}{2\lambda} \|\lambda A^\top (v - \bar{v})\|^2 - \frac{1}{\lambda} \langle \lambda A^\top (v - \bar{v}), \bar{z} \rangle + \langle \bar{z}, A^\top (v - \bar{v}) \rangle \\ &= \frac{1}{2\lambda} \|\lambda A^\top (v - \bar{v})\|^2. \end{aligned}$$

At (1), we assumed for the theorem statement that $\bar{z} = y - \lambda A^\top \bar{v}$. By the optimality of \bar{v} on the dual problem, we have the following:

$$\begin{aligned} & \mathbf{0} \in A(\lambda A^\top \bar{v} - y) + \partial\omega^*(\bar{v}) \\ \{eqn:thm:minimizing-dual-pp:pitem2\} \quad & \iff A\bar{z} \in \partial\omega^*(\bar{v}) \\ & \iff \omega^*(v) - \omega^*(\bar{v}) \geq \langle A\bar{z}, v_n - \bar{v} \rangle. \end{aligned} \tag{1.5}$$

We will now prove (i), because $A\bar{z} \in \partial\omega^*(\bar{v}) \iff \partial\omega(A\bar{z}) \ni \bar{v}$, we can derive that $\bar{z} = y - \lambda A^\top \bar{v}$ is primal optimal because

$$y - \bar{z} = \lambda A^\top \bar{v} \in \lambda A^\top \partial\omega(A\bar{z}).$$

Re-arranging yields: $\mathbf{0} \in \bar{z} - y + \lambda A^\top \partial\omega(A\bar{z}) = \partial\Phi_\lambda(\bar{z})$.

Now we will prove (ii). Recall that $z_j = y - \lambda A^\top v_j$, therefore it has:

$$\begin{aligned} \Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}) &= \frac{1}{2\lambda} \|\lambda A^\top v_j - y\|^2 - \frac{1}{2\lambda} \|\lambda A^\top \bar{v} - y\|^2 + \omega^*(v_j) - \omega^*(\bar{v}) \\ &\stackrel{(2)}{\geq} \frac{1}{2\lambda} \|\lambda A^\top v_j - y\|^2 - \frac{1}{2\lambda} \|\lambda A^\top \bar{v} - y\|^2 + \langle A\bar{z}, v_j - \bar{v} \rangle \\ &\stackrel{(3)}{=} \frac{1}{2\lambda} \|\lambda A^\top (v_j - \bar{v})\|^2 \\ &\stackrel{(4)}{=} \frac{1}{2\lambda} \|z_j - \bar{z}\|^2. \end{aligned}$$

At (2) we used (1.5), at (3) we used (1.4). At (4), we again $\bar{z} = y - \lambda A^\top \bar{v}$. We have \bar{z} being the optimal solution of primal Φ_λ . Now, if we set $v = v_j$, and then:

$$0 = \lim_{j \rightarrow \infty} \Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}) \geq \lim_{j \rightarrow \infty} \|z_j - \bar{z}\|^2.$$

Implying that $z_j \rightarrow \bar{z}$ also.

We will now prove (iii). We have for all $z \in \mathbb{R}^n$, the following:

$$\begin{aligned} & \Phi_\lambda(z_j) - \Phi_\lambda(\bar{z}) \\ &= \omega(Az_j) - \omega(A\bar{z}) + \frac{1}{2\lambda} (\|z_j - y\|^2 - \|\bar{z} - y\|^2) \\ &\leq K_\omega \|A\| \|z_j - \bar{z}\| + \frac{1}{2\lambda} (\|z_j - y\| + \|\bar{z} - y\|) (\|z_j - y\| - \|\bar{z} - y\|) \\ &\leq K_\omega \|A\| \|z_j - \bar{z}\| + \frac{1}{2\lambda} (\|z_j - y\| + \|\bar{z} - y\|) \|z_j - \bar{z}\| \\ &\leq K_\omega \|A\| \|z_j - \bar{z}\| + \frac{1}{2\lambda} (\|z_j - \bar{z}\| + 2\|\bar{z} - y\|) \|z_j - \bar{z}\| \end{aligned}$$

$$\begin{aligned}
&= \|z_j - \bar{z}\| \left(K_\omega \|A\| + \lambda^{-1} \|\bar{z} - y\| + \frac{\|z_j - \bar{z}\|}{2\lambda} \right) \\
&\stackrel{(ii)}{\leq} \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + \lambda^{-1} \|\bar{z} - y\| + \frac{\sqrt{2\lambda}}{2\lambda} \sqrt{\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})} \right) \\
&\stackrel{(5)}{\leq} \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + \max_{z \in \partial g(\bar{z})} \|z\| + \frac{\sqrt{2\lambda}}{2\lambda} \sqrt{\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})} \right) \\
&\stackrel{(6)}{=} \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + K_\omega + \frac{\sqrt{2\lambda}}{2\lambda} \sqrt{\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})} \right).
\end{aligned}$$

At (5), recall the primal minimizer $\mathbf{0} \in \partial\Phi_\lambda(\bar{z}) = (I + \lambda\partial g)^{-1}(y)$ which means $\lambda^{-1}(y - z) \in \partial g(\bar{z})$, which implies that $\lambda^{-1}\|y - z\| \leq \max_{z \in \partial g} \|z\|$. At (6) we used the assumption that ω is K_ω Lipschitz continuous stated in Assumption 1.15, and we invoke Lemma A.2. ■

1.4 Literature reviews

1.5 Our contributions

2 The inexact accelerated proximal gradient with controlled errors

In this section, we present an accelerated algorithm with controlled error using Definition 1.10, and show that it can have a convergence rate under certain error conditions.

{def:inxt-apg}

Definition 2.1 (our inexact accelerated proximal gradient)

Suppose that (F, f, g, L) and, sequences $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$ satisfies the following

- (i) $(\alpha_k)_{k \geq 0}$ is a sequence such that $\alpha \in (0, 1]$ for all $k \geq 0$.
- (ii) $(B_k)_{k \geq 0}$ has $B_k > 0 \forall k$, it characterizes any potential line search, back tracking routine.
- (iii) $(\rho_k)_{k \geq 0}$ be a sequence such that $\rho_k \geq 0$, characterizing the over-relaxation of the proximal gradient operator.
- (iv) $(\epsilon_k)_{k \geq 0}$ has $\epsilon_k > 0$ for all $k \geq 0$, it characterizes the errors of inexact proximal evaluation.
- (v) (f, g, L) satisfies Assumption 1.9, and let $F = f + g$.

Denote $L_k = B_k + \rho_k$ for short. Given any initial condition $v_{-1}, x_{-1} \in \mathbb{R}^n$, the algorithm

generates the sequences $(y_k, x_k, v_k)_{k \geq 0}$ such that they satisfy for all $k \geq 0$:

{def:inxt-apg:yk}

$$y_k = \alpha_k v_{k-1} + (1 - \alpha_k) x_{k-1}, \quad (2.1)$$

{def:inxt-apg:xk}

$$x_k \approx_{\epsilon_k} T_{L_k}(y_k), \quad (2.2)$$

$$D_f(x_k, y_k) \leq \frac{B_k}{2} \|x_k - y_k\|^2, \quad (2.3)$$

{def:inxt-apg:vk}

$$v_k = x_{k-1} + \alpha_k^{-1}(x_k - x_{k-1}). \quad (2.4)$$

{lemma:inxt-apg-cnvg-prep1}

Lemma 2.2 (inexact accelerated proximal gradient preparation stage I)

Let (F, f, g, L) , and $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$, be given by Definition 2.1. Denote $L_k = B_k + \rho_k$. Then, for any $\bar{x} \in \mathbb{R}^n$, the sequences $(y_k, x_k, v_k)_{k \geq 0}$ generated satisfy for all $k \geq 1$ the inequality:

$$\begin{aligned} & \frac{\rho_k}{2} \|x_k - y_k\|^2 - \epsilon_k \\ & \leq (1 - \alpha_k)(F(x_{k-1}) - F(\bar{x})) + F(\bar{x}) - F(x_k) \\ & + \max \left(1 - \alpha_k, \frac{\alpha_k^2 L_k}{\alpha_{k-1}^2 L_{k-1}} \right) \frac{\alpha_{k-1}^2 L_{k-1}}{2} \|\bar{x} - v_{k-1}\|^2 - \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2. \end{aligned}$$

When, $k = 1$ it instead has:

$$\begin{aligned} & \frac{\rho_0}{2} \|x_0 - y_0\|^2 - \epsilon_0 \\ & \leq (1 - \alpha_0)(F(x_{-1}) - F(\bar{x})) + F(\bar{x}) - F(x_0) + \frac{\alpha_0^2 L_0}{2} \|\bar{x} - v_{-1}\|^2 - \frac{\alpha_0^2 L_0}{2} \|\bar{x} - v_0\|^2. \end{aligned}$$

Proof. Two intermediate results are in order before we can prove the inequality. Define $z_k := \alpha_k \bar{x} + (1 - \alpha_k) x_{k-1}$ for short. It has for all $k \geq 1$ the equality:

{eqn:inxt-apg-cnvg-prep1-a}

$$\begin{aligned} z_k - x_k &= \alpha_k \bar{x} + (1 - \alpha_k) x_{k-1} - x_k \\ &= \alpha_k x^+ + (x_{k-1} - x_k) - \alpha_k x_{k-1} \\ &= \alpha_k \bar{x} - \alpha_k v_k. \end{aligned} \quad (a)$$

It also has for all $k \geq 1$ the equality:

{eqn:inxt-apg-cnvg-prep1-b}

$$\begin{aligned} z_k - y_k &= \alpha_k \bar{x} + (1 - \alpha_k) x_{k-1} - y_k \\ &= \alpha_k \bar{x} - \alpha_k v_{k-1}. \end{aligned} \quad (b)$$

Let's denote $L_k = B_k + \rho_k$ for short. Recall that (f, g, L) satisfies Assumption 1.9, if we choose $x = y_k$ so $\tilde{x} = x_k \approx_{\epsilon_k} T_{L_k}(y_k)$, and set $z = z_k, \epsilon = \epsilon_k$ then Theorem 1.13 has:

$$\begin{aligned}
& \frac{\rho_k}{2} \|x_k - y_k\|^2 - \epsilon_k \\
& \leq F(z_k) - F(x_k) + \frac{L_k}{2} \|y_k - z_k\|^2 - \frac{L_k}{2} \|z_k - x_k\|^2 \\
& \stackrel{(1)}{\leq} \alpha_k F(\bar{x}) + (1 - \alpha_k) F(x_{k-1}) - F(x_k) + \frac{L_k}{2} \|y_k - z_k\|^2 - \frac{L_k}{2} \|z_k - x_k\|^2 \\
& \stackrel{(2)}{=} (1 - \alpha_k)(F(x_{k-1}) - F(\bar{x})) + F(\bar{x}) - F(x_k) + \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_{k-1}\|^2 - \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 \\
& \leq (1 - \alpha_k)(F(x_{k-1}) - F(\bar{x})) + F(\bar{x}) - F(x_k) \\
& + \max \left(1 - \alpha_k, \frac{\alpha_k^2 L_k}{\alpha_{k-1}^2 L_{k-1}} \right) \frac{\alpha_{k-1}^2 L_{k-1}}{2} \|\bar{x} - v_{k-1}\|^2 - \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2.
\end{aligned}$$

At (1) we used the fact that $F = f + g$ hence F is convex. At (2) we used (a), (b). Finally, if $k = 0$, then take the RHS of $\stackrel{(1)}{=}$ then:

$$\begin{aligned}
& \frac{\rho_0}{2} \|x_0 - y_0\|^2 - \epsilon_0 \\
& \leq (1 - \alpha_0)(F(x_{-1}) - F(\bar{x})) + F(\bar{x}) - F(x_0) + \frac{\alpha_0^2 L_0}{2} \|\bar{x} - v_{-1}\|^2 - \frac{\alpha_0^2 L_0}{2} \|\bar{x} - v_0\|^2.
\end{aligned}$$

■

The following assumption encapsulates assumptions on the errors such that a near optimal convergence rate is still attainable by an algorithm that satisfies Definition 2.1.

Assumption 2.3 (valid error schedule) The following assumption is about an algorithm satisfying Definition 2.1, its parameters $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$ in relation to its iterates $(y_k, x_k, v_k)_{k \geq 0}$ and, some additional parameters $(\beta_k)_{k \geq 0}, \mathcal{E}_0$ and p . Let

- (i) $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}, (F, f, g, L)$ and $(y_k, x_k, v_k)_{k \geq 0}$ be given by Definition 2.1.
- (ii) $\mathcal{E}_0 \geq 0$ be arbitrary;
- (iii) the sequence $(\beta_k)_{k \geq 0}$ be defined as $\beta_k := \prod_{i=1}^k \max \left(1 - \alpha_i, \frac{\alpha_i^2 L_i}{\alpha_{i-1}^2 L_{i-1}} \right)$ for all $k \geq 1$, with the base case being $\beta_0 = 1$;
- (iv) $p \geq 1$ is some constant which will bound the error ϵ_k relative to $\rho_k \|x_k - y_k\|^2, \beta_k$ and, k .

In addition, we assume that the error parameter $\epsilon_k \geq 0$ and over-relaxation parameter ρ_k , iterates x_k, y_k and β_k together satisfies for all $k \geq 0$ the relations:

$$\frac{-\mathcal{E}_0 \beta_k}{k^p} \leq \frac{\rho_k}{2} \|x_k - y_k\|^2 - \epsilon_k.$$

{prop:inxt-apg-cnvg-generic}

The following proposition is a prototype of the convergence rate together with the error schedule that delivers convergence of algorithms satisfying Definition 2.1.

Proposition 2.4 (convergence with valid error schedule)

Let (F, f, g, L) , $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$, $(\beta_k)_{k \geq 0}$, \mathcal{E}_0, p as assumed in Assumption 2.3. Fix any $\bar{x} \in \mathbb{R}^n$ for all $k \geq 0$ and assume that $\alpha_0 = 1$. Then for the iterates generated $(y_k, x_k, v_k)_{k \geq 0}$ by the algorithm, for all $k \geq 0$ they will satisfy:

$$F(x_k) - F(\bar{x}) + \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 \leq \beta_k \left(\frac{L_0}{2} \|\bar{x} - v_{-1}\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} \right).$$

Proof. Consider results from Lemma 2.2 has $\forall k \geq 1$:

$$\begin{aligned} & \frac{\rho_k}{2} \|x_k - y_k\|^2 - \epsilon_k \\ & \leq (1 - \alpha_k)(F(x_{k-1}) - F(\bar{x})) + F(\bar{x}) - F(x_k) \\ & \quad + \max \left(1 - \alpha_k, \frac{\alpha_k^2 L_k}{\alpha_{k-1}^2 L_{k-1}} \right) \frac{\alpha_{k-1}^2 L_{k-1}}{2} \|\bar{x} - v_{k-1}\|^2 - \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2. \\ & \leq \max \left(1 - \alpha_k, \frac{\alpha_k^2 L_k}{\alpha_{k-1}^2 L_{k-1}} \right) \left(F(x_{k-1}) - F(\bar{x}) + \frac{\alpha_{k-1}^2 L_{k-1}}{2} \|\bar{x} - v_{k-1}\|^2 \right) \\ & \quad + F(\bar{x}) - F(x_k) - \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 \end{aligned}$$

For notation brevity, we introduce β_k, Λ_k :

$$\begin{aligned} \beta_0 &= 1, \\ \beta_k &:= \prod_{i=1}^k \max \left(1 - \alpha_i, \frac{\alpha_i^2 L_i}{\alpha_{i-1}^2 L_{i-1}} \right), \\ \Lambda_k &:= -F(\bar{x}) + F(x_k) + \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2. \end{aligned}$$

Now, suppose that in addition there is a non-negative sequence $(\mathcal{E}_k)_{k \geq 0}$ such that

- (i) For all $k \geq 0$, it has $\frac{-\mathcal{E}_k}{k^p} \leq \frac{\rho_k}{2} \|x_k - y_k\|^2 - \epsilon_k$ where $p \geq 1$,
- (ii) For all $k \geq 1$, it has $\mathcal{E}_k = \frac{\beta_k}{\beta_{k-1}} \mathcal{E}_{k-1}$, with $\mathcal{E}_0 \geq 0$.

These conditions are equivalent to the assumption that $\frac{-\mathcal{E}_0 \beta_k}{k^p} \leq \frac{\rho_k}{2} \|x_k - y_k\|^2 - \epsilon_k$ (which was stated in Assumption 2.3). One can show that by unrolling recurrence on \mathcal{E}_k . Then (2.5) implies $\forall k \geq 1$:

$$\frac{-\mathcal{E}_k}{k^p} \leq \frac{\beta_k}{\beta_{k-1}} \Lambda_{k-1} - \Lambda_k \iff \Lambda_k \leq \frac{\beta_k}{\beta_{k-1}} \Lambda_{k-1} + \frac{\mathcal{E}_k}{k^p}. \quad (2.5)$$

Now, we show the convergence of Λ_k , using the relations of $\mathcal{E}_k, \Lambda_k, \beta_k$ above.

$$\begin{aligned}
\Lambda_k &\leq \frac{\beta_k}{\beta_{k-1}} \Lambda_{k-1} + \frac{\mathcal{E}_k}{k^p} \\
&\leq \frac{\beta_k}{\beta_{k-1}} \Lambda_{k-1} + \frac{\beta_k}{\beta_{k-1}} \frac{\mathcal{E}_{k-1}}{k^p} \\
&= \frac{\beta_k}{\beta_{k-1}} \left(\Lambda_{k-1} + \frac{\mathcal{E}_{k-1}}{k^p} \right) \\
&\leq \frac{\beta_k}{\beta_{k-1}} \left(\frac{\beta_{k-1}}{\beta_{k-2}} \Lambda_{k-2} + \frac{\mathcal{E}_{k-1}}{(k-1)^p} + \frac{\mathcal{E}_{k-1}}{k^p} \right) \\
&= \frac{\beta_k}{\beta_{k-2}} \left(\Lambda_{k-2} + \frac{\mathcal{E}_{k-2}}{(k-1)^p} + \frac{\mathcal{E}_{k-2}}{k^p} \right) \\
&\dots \\
&\leq \frac{\beta_k}{\beta_1} \left(\Lambda_1 + \mathcal{E}_1 \sum_{n=2}^k \frac{1}{n^p} \right) \\
&\leq \frac{\beta_k}{\beta_1} \left(\frac{\beta_1}{\beta_0} \Lambda_0 + \mathcal{E}_1 \sum_{n=1}^k \frac{1}{n^p} \right) \\
&= \frac{\beta_k}{\beta_0} \left(\Lambda_0 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} \right).
\end{aligned}$$

Therefore, it points to the following inequality:

$$\begin{aligned}
F(x_k) - F(\bar{x}) + \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 \\
\leq \beta_k \left(F(x_0) - F(\bar{x}) + \frac{\alpha_0^2 L_0}{2} \|\bar{x} - v_0\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} \right).
\end{aligned}$$

Finally, when $\alpha_0 = 1$, then the results from 2.2 with $k = 0$ simplifies the above inequality and give:

$$F(x_k) - F(\bar{x}) + \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 \leq \beta_k \left(\frac{L_0}{2} \|\bar{x} - v_{-1}\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} \right).$$

■

Now, it only remains to determine the sequence α_k to derive a type of convergence rate for the algorithm because from the above theorem, we have the convergence rate β_k and, the error parameters ϵ_k, ρ_k both controlled by the sequence $(\alpha_k)_{k \geq 0}$.

2.1 Convergence results of the outer loop

This section will give specific instances of the error control sequence $(\epsilon_k)_{k \geq 0}$, $(\rho_k)_{k \geq 0}$ and, momentum sequence $(\alpha_k)_{k \geq 0}$ such that an optimal convergence rate of $\mathcal{O}(1/k^2)$ can be achieved.

{ass:opt-mmntm-seq}

Assumption 2.5 (the optimal momentum sequence)

Keeping everything assumed in Assumption 2.3 about $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$, (F, f, g, L) , $(y_k, x_k, v_k)_{k \geq 0}$, $(\beta_k)_{k \geq 0}$, \mathcal{E}_0 and p . We assume in addition that the sequence $(\alpha_k)_{k \geq 0}$ satisfies for all $k \geq 0$ the equality: $(1 - \alpha_k) = \alpha_k^2 L_k \alpha_{k-1}^{-2} L_{k-1}^{-1}$ and, $p > 1$.

{lemma:opt-mmntm-seq}

Lemma 2.6 (the optimal momentum sequence is indeed valid and optimal)

Let $(\alpha_k)_{k \geq 0}$, $(\beta_k)_{k \geq 0}$ be given by Assumption 2.5. If we choose $\alpha_0 \in (0, 1]$ then for all $k \geq 1$ it has:

$$\alpha_k = \frac{L_{k-1}}{2L_k} \left(-\alpha_{k-1}^2 + \left(\alpha_{k-1}^4 + 4\alpha_{k-1} \frac{L_k}{L_{k-1}} \right)^{1/2} \right) \in (0, 1) \quad (2.6)$$

For the sequence $(\beta_k)_{k \geq 0}$ it has $\forall k \geq 1$

$$\left(1 + \alpha_0 \sqrt{L_0} \sum_{i=1}^k \sqrt{L_i^{-1}} \right)^{-2} \leq \beta_k = \frac{\alpha_k^2 L_k}{\alpha_0^2 L_0} \leq \left(1 + \frac{\alpha_0 \sqrt{L_0}}{2} \sum_{i=1}^k \sqrt{L_i^{-1}} \right)^{-2}. \quad (2.7)$$

Proof. Firstly, we will show (2.6). We will prove using induction. Fix any $k \geq 1$. Assume inductively that $\alpha_{k-1} \in (0, 1]$. We can solve for α_k using the recursive equality $(1 - \alpha_k) = \alpha_k^2 L_k \alpha_{k-1}^{-2} L_{k-1}^{-1}$ from Assumption 2.5. Writing α_{k-1} as α , and L_k/L_{k-1} as q . Then, we solve for α_k , the quadratic equation always admits one root that is strictly positive which is given as:

$$\begin{aligned} \alpha_k &= \frac{1}{2} \left(-\frac{\alpha^2}{q} + \sqrt{\frac{\alpha^4}{q^2} + \frac{4\alpha^2}{q}} \right) \\ &= \frac{\alpha^2}{2q} \left(-1 + \sqrt{1 + \frac{4q}{\alpha^2}} \right) \\ &\stackrel{(1)}{<} \frac{\alpha^2}{2q} \left(-1 + 1 + \frac{2q}{\alpha^2} \right) \\ &= 1 \end{aligned}$$

At (1) we completed a square in the radical, and we used the assumption $\alpha_k > 0$ and, $L_k > 0, L_{k-1} > 0$ because we had $B_k > 0$, therefore the following chain of inequality holds:

$$\begin{aligned} 1 + \frac{4q}{\alpha^2} &= 1 + \frac{4q}{\alpha^2} + \frac{4q^2}{\alpha^4} - \frac{4\alpha^2}{\alpha^2} \\ &= \left(1 + \frac{2q}{\alpha^2}\right)^2 - \frac{4q^2}{\alpha^4} \\ &< \left(1 + \frac{2q}{\alpha^2}\right)^2. \end{aligned}$$

To see that $\alpha_k > 0$, recall the same fact that $L_k > 0$, and the inductive hypothesis $\alpha_{k-1} \in (0, 1]$ then $4q/\alpha^2 > 0$ so obviously $\alpha_k = \frac{\alpha^2}{2q} \left(-1 + \sqrt{1 + 4q/\alpha^2}\right) > 0$ because the quantity inside the radical is strictly larger than 1. Therefore, inductively it holds that $\alpha_k \in (0, 1)$ too.

We will now show (2.7). From the assumption that $(\alpha_k)_{k \geq 0}$ has $(1 - \alpha_k) = \alpha_k^2 L_k \alpha_{k-1}^{-2} L_{k-1}^{-1}$ for all $k \geq 0$ and, the definition of β_k , it yields the following equalities:

$$\beta_k = \prod_{i=1}^k \max\left(1 - \alpha_i, \frac{\alpha_i^2 L_i}{\alpha_{i-1}^2 L_{i-1}}\right) = \prod_{i=1}^k (1 - \alpha_i) = \prod_{i=1}^k \frac{\alpha_i^2 L_i}{\alpha_0^2 L_0} = \frac{\alpha_k^2 L_k}{\alpha_0^2 L_0}.$$

With the above relation, and the definitions of the sequences $(\alpha_k)_{k \geq 0}, (\beta_k)_{k \geq 0}$ it satisfies for all $k \geq 1$ the properties:

- (a) β_k is monotone decreasing and $\beta_k > 0$ for all $k \geq 0$ because $\beta_k = \prod_{i=1}^k (1 - \alpha_i)$ and, $\alpha_k \in (0, 1]$.
- (b) It has the equalities $\beta_k / \beta_{k-1} = (1 - \alpha_k) = \frac{\alpha_k^2 L_k}{\alpha_{k-1}^2 L_{k-1}}$ for all $k \geq 1$.

Using the above observations, we can show the chain of equalities $\alpha_k^2 = (1 - \beta_k / \beta_{k-1})^2 = \beta_k L_0 \alpha_0^2 L_k^{-1}$ for all $k \geq 0$. This is true by first considering the relations $\prod_{i=1}^k (1 - \alpha_i) = \beta_k$:

$$\begin{aligned} \{ \text{eqn:opt-mmntm-seq-pitem1} \} \quad (1 - \alpha_k) &= \beta_k / \beta_{k-1} \\ &\iff \alpha_k = 1 - \beta_k / \beta_{k-1} \\ &\implies \alpha_k^2 = (1 - \beta_k / \beta_{k-1})^2. \end{aligned} \tag{2.8}$$

Next, the recursive relation of $(\alpha_k)_{k \geq 0}$ gives

$$\begin{aligned} \{ \text{eqn:opt-mmntm-seq-pitem2} \} \quad \alpha_k^2 &= (1 - \alpha_k) \alpha_{k-1}^2 L_{k-1} L_k^{-1} \\ &= (1 - \alpha_k) \left(\frac{\alpha_{k-1}^2 L_{k-1}}{\alpha_0^2 L_0}\right) \frac{\alpha_0^2 L_0}{L_k} \\ &= (\beta_k \beta_{k-1}^{-1}) (\beta_{k-1}) L_0 \alpha_0^2 L_k^{-1} \\ &= \beta_k L_0 \alpha_0^2 L_k^{-1}. \end{aligned} \tag{2.9}$$

Combining (2.8), (2.9) and the fact that $\beta_k > 0 \forall k \geq 0$, it would mean for all $i \geq 1$ it has:

$$\begin{aligned} L_0 \alpha_0^2 L_i^{-1} &= \beta_i^{-1} \left(1 - \frac{\beta_k}{\beta_{k-1}}\right)^2 \\ &= \beta_i (\beta_i^{-1} - \beta_{i-1}^{-1})^2 \\ &= \beta_i \left(\beta_i^{-1/2} - \beta_{i-1}^{-1/2}\right)^2 \left(\beta_i^{-1/2} + \beta_{i-1}^{-1/2}\right)^2 \\ &= \left(\beta_i^{-1/2} - \beta_{i-1}^{-1/2}\right)^2 \left(1 + \beta_i^{1/2} \beta_{i-1}^{-1/2}\right)^2. \end{aligned}$$

Since β_i is monotone decreasing, it has $0 < \beta_i^{1/2} \beta_{i-1}^{-1/2} \leq 1$, this gives:

$$\beta_i^{-1/2} - \beta_{i-1}^{-1/2} \leq \alpha_0 \sqrt{\frac{L_0}{L_i}} \leq 2 \left(\beta_i^{-1/2} - \beta_{i-1}^{-1/2}\right).$$

Performing a telescoping sum for $i = 1, 2, \dots, k$, use the fact that $\beta_0 = 1$ will yield the desired results after some algebraic manipulations. \blacksquare

Remark 2.7

{prop:opt-cnvg-outr-loop}

Proposition 2.8 ($\mathcal{O}(1/k^2)$ outer loop function value convergence)

Let (f, g, L) , $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$, $(\beta_k)_{k \geq 0}$, \mathcal{E}_0, p be given by Assumption 2.5. Assume in addition that

- (i) there exists $\bar{x} \in \mathbb{R}^n$ that is a minimizer of $F = f + g$;
- (ii) the sequence $L_k := B_k + \rho_k$ is bounded, and there exists an L_{\max} such that for all $k \geq 0$ it has $L_{\max} \geq \max_{k \geq 0} L_i$.

Then, it has $\forall k \geq 0$:

$$F(x_k) - F(\bar{x}) + \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 \leq \left(1 + \frac{k\alpha_0 \sqrt{L_0}}{2\sqrt{L_{\max}}}\right)^{-2} \left(\frac{L_0}{2} \|\bar{x} - v_{-1}\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p}\right).$$

Since, $p > 0$ the sum is convergent and hence the above inequality claims an overall convergence rate $\mathcal{O}(1/k^2)$.

Proof. We use the results from Lemma 2.6 because of the same assumption on $(\alpha_k)_{k \geq 0}$, then using the fact that L_k is bounded:

$$\beta_k \leq \left(1 + \frac{\alpha_0 \sqrt{L_0}}{2} \sum_{i=1}^k \sqrt{L_i^{-1}}\right)^{-2} \leq \left(1 + \frac{k\alpha_0 \sqrt{L_0}}{2\sqrt{L_{\max}}}\right)^{-2}.$$

Then, apply Theorem 2.4 to obtain the desired results. \blacksquare

Remark 2.9 In this remark, we assert the fact that all assumptions made in the theorem can be satisfied on practice, and we will also bring attention to the current, and future roles played by some parameters used in the inexact algorithm.

Pay attention that α_k had been determined in the above theorem (as seemed in Lemma 2.6), and $(B_k)_{k \geq 0}$ is reserved for potential line search routine, the only parameter left undetermined in Definition 2.1 is the over-relaxation sequence $(\rho_k)_{k \geq 0}$. Since we only need the entire sequence $L_k = \rho_k + B_k$ to be bounded above, we give the freedom to the practitioners to choose $(\rho_k)_{k \geq 0}$. However, this sequence ρ_k is not useless because it counters ϵ_k in the proximal gradient inequality, this has huge impact in the earlier stage (the first few iterations) of the algorithm $\|x_k - y_k\|$ is large. Of course, the parameter \mathcal{E}_0 is also free to choose.

Finally, observe that from the above proof, in case when $p = 1$, the convergence rate would be $\mathcal{O}(\log(k)/k^2)$.

In the next subsection, we will show that under the assumption of the above theorem, there exists an error sequence ϵ_k such that it can never approach 0 at an arbitrarily fast rate.

2.2 The fastest rate of which the error schedule can shrink

To have overall convergence claim of the algorithm, it's necessary to prevent the error schedule $(\epsilon_k)_{k \geq 0}$ from crashing into zero too quickly. Following Assumption 2.3, in this section, we will provide the lower bound results for ϵ_k in Theorem 2.8 to show that in the worst case it cannot approach zero arbitrarily fast, if we choose the largest possible ϵ_k using β_k .

{lemma:err-schedule-lbnd}

Lemma 2.10 (error schedule lower bound)
Let, $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$, $(\beta_k)_{k \geq 0}$, \mathcal{E}_0, p be given by Assumption 2.5, $L_k := \rho_k + B_k$. Let $(\epsilon_k)_{k \geq 0}$ be given by $\epsilon_k := \frac{\mathcal{E}_0 \beta_k}{k^p} + \rho_k \frac{\|x_k - y_k\|^2}{2}$, then it will be a valid error sequence and so that it satisfies the assumption. If in addition, there exists L_{\min} such that for all $k \geq 0$ such that it has $L_{\min} \leq \min_{1 \leq i \leq k} L_i$ and, we assume $\mathcal{E} > 0$, then ϵ_k admits the non-trivial lower bound:

$$\epsilon_k \geq \frac{\mathcal{E}_0}{k^p} \left(1 + k\alpha_0 \sqrt{L_0} \sqrt{L_{\min}^{-1}} \right)^{-2} = \mathcal{O}(k^{-2-p}).$$

Proof. Recall Assumption 2.3, the largest valid error schedule is $\epsilon_k = \frac{\mathcal{E}_0\beta_k}{k^p} + \rho_k \frac{\|x_k - y_k\|^2}{2}$. Then it has

$$\begin{aligned}\epsilon_k &\geq \frac{\mathcal{E}_0\beta_k}{k^p} \\ &\stackrel{(1)}{\geq} \left(1 + \alpha_0\sqrt{L_0} \sum_{i=1}^k \sqrt{L_i^{-1}}\right)^{-2} \frac{\mathcal{E}_0}{k^p} \\ &\stackrel{(2)}{\geq} \frac{\mathcal{E}_0}{k^p} \left(1 + k\alpha_0\sqrt{L_0} \sqrt{L_{\min}^{-1}}\right)^{-2} \\ &= \mathcal{O}(k^{-2-p}).\end{aligned}$$

At (1), we used Lemma 2.6. At (2), we used that $L_{\min} \leq L_i$ for all $i = 0, 1, 2, \dots$. \blacksquare

2.3 The gradient mapping also shrinks

In this section, we show one extremely favorable properties and its assumptions such that, successive iterates v_k, v_{k-1} converges to zero as $k \rightarrow \infty$, and hence the norm of the gradient mapping converges at a rate of $\mathcal{O}(1/k)$. This property will be crucial for analyzing the total complexity for the algorithm because it is the termination criteria.

{prop:vk-gm}

Proposition 2.11 (iterates v_k and gradient mapping)

Let (F, f, g, L) , $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$, $(\beta_k)_{k \geq 0}$, \mathcal{E}_0, p as assumed in Assumption 2.3. Assume there exists $\bar{x} \in \mathbb{R}^n$ which is the minimizer of F and, we set $\alpha_0 = 1$. Then for the iterates generated $(y_k, x_k, v_k)_{k \geq 0}$ by the algorithm, for all $k \geq 0$ all followings are true.

{prop:vk-gm:result1}
{prop:vk-gm:result2}

- (i) it has $v_k - v_{k-1} = \alpha_k^{-1}(x_k - y_k) = \alpha_k^{-1}L_k^{-1}\mathcal{G}_{L_k}(y_k)$.
- (ii) as a consequence of the former, we have the following bound:

$$\|\bar{x} - v_k\| \leq \|\bar{x} - v_{-1}\| + \left(\frac{2\mathcal{E}_0}{L_0} \sum_{n=1}^k \frac{1}{n^p}\right)^{1/2}.$$

Proof. Our first result follows. It has $y_k = \alpha_k v_{k-1} + (1 - \alpha_k)x_{k-1}$ from (2.1), and $v_k = x_{k-1} + \alpha_k^{-1}(x_k - x_{k-1})$ from (2.4). Using this information we have

$$\begin{aligned}v_k - v_{k-1} &= x_{k-1} + \alpha_k^{-1}(x_k - x_{k-1}) - \alpha_k^{-1}(y_k - (1 - \alpha_k)x_{k-1}) \\ &= (1 - \alpha_k^{-1})x_{k-1} + \alpha_k^{-1}x_k - \alpha_k^{-1}y_k + (\alpha_k^{-1} - 1)x_{k-1} \\ &= \alpha_k^{-1}(x_k - y_k) = \alpha_k^{-1}L_k^{-1}L_k(x_k - y_k) = \alpha_k^{-1}L_k^{-1}\mathcal{G}_{L_k}(y_k).\end{aligned}$$

Our second result follows. The assumption for Proposition 2.4 are satisfied here hence:

$$\begin{aligned}
0 &\leq \beta_k \left(\frac{L_0}{2} \|\bar{x} - v_{k-1}\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} \right) - \frac{\alpha_k^2 L_k}{2} \|\bar{x} - v_k\|^2 - (F(x_k) - F(\bar{x})) \\
&\stackrel{(1)}{\implies} 0 \leq \frac{L_0}{2} \|\bar{x} - v_{k-1}\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} - \frac{\alpha_k^2 L_k}{2\beta_k} \|\bar{x} - v_k\|^2 \\
&\stackrel{(2)}{=} \frac{L_0}{2} \|\bar{x} - v_{k-1}\|^2 + \mathcal{E}_0 \sum_{n=1}^k \frac{1}{n^p} - \frac{L_0}{2} \|\bar{x} - v_k\|^2 \\
&\iff 0 \leq \|\bar{x} - v_{k-1}\| - \|\bar{x} - v_k\| + \left(\frac{2\mathcal{E}_0}{L_0} \sum_{n=1}^k \frac{1}{n^p} \right)^{1/2}.
\end{aligned}$$

At (1), we assumed \bar{x} is the minimizer so $-F(x_k) + F(\bar{x}) < 0$ and so we replaced it with a zero, and since $\beta_k > 0$ always we can divided on both side of the inequality without changing the sign. At (2), we used $\beta_k = \frac{\alpha_k^2 L_k}{\alpha_0^2 L_0}$ from (2.7), since we assumed $\alpha_0 = 1$ here it has $\beta_k = \alpha_k^2 L_k / L_0$ hence the coefficient $\frac{\alpha_k^2 L_k}{2\beta_k} = \frac{L_0}{2}$. ■

Now, we will show the convergence rate of the normed gradient mapping is $\mathcal{O}(1/k)$ for the best chosen sequence. The theorem that follows will accomplish that

Theorem 2.12 ($\mathcal{O}(1/k)$ convergence of the gradient mapping)

N'oubliez pas de finir cette partie.

Proof.

■

3 Linear convergence for the proximal problem in the inner loop

In this section, we continue the discussion from Section 1.3. As an important reminder, we will fix the vector $y \in \mathbb{R}^n$, which is in the inexact proximal problem as a constant in this entire section.

The inner loop is another algorithm that evaluates $x_k \approx_\epsilon T_{(B_k + \rho_k)}(y_k)$ for a given value of $\epsilon, B + \rho$ and at the point y_k . Let's assume that the outer loop iteration k is fixed throughout

this entire section to simplify our discussion since, in this section we will only focus on the convergence rate of the inner loop for one specific iteration of the outer loop.

Let $\lambda = (B_k + \rho_k)^{-1}$, the algorithm needs to resolve the following equivalent inexact proximal point problem:

$$x_k \approx_{\epsilon} \text{prox}_{\lambda g}(y_k - \lambda \nabla f(y_k)).$$

Unfortunately recall that $g = \omega \circ A$ in the context of the outer loop hence it's impossible to directly evaluate the proximal operator of g and hence we optimize the function Φ_λ as given by (1.1).

We will show that there exists an algorithm generating the sequences z_j, v_j such that $\mathbf{G}_\lambda(z_j, v_j)$ converges linearly if Ψ_λ satisfies the error bound conditions. Using results available in the literature, we will characterize the exact scenarios of $\omega \circ A$ when this is possible to achieve. To start, the following assumption is the general error bound condition of a convex with additive composite structure.

{def:pg-and-gm} **Definition 3.1 (the exact proximal gradient and gradient mapping)**

Let $F = f + g$ where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a differentiable function and, $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ is closed convex and proper. For all $\tau > 0$, we define the exact proximal gradient mapping and, gradient mapping:

- (i) *The proximal gradient operator $T_\tau(x) := \text{prox}_{\tau^{-1}g}(x - \tau^{-1}\nabla f(x))$,*
- (ii) *and the gradient mapping operator $\mathcal{G}_\tau(x) := \tau(x - T_\tau(x))$.*

{def:err-bnd} **Definition 3.2 (the error bound condition)**

Under the same assumptions in Definition 3.1, the function $F = f + g$ satisfies the error bound condition if there exists $\gamma > 0$ such that

$$\|\mathcal{G}_\tau(x)\| \geq \gamma \text{dist}(x|S). \quad (3.1)$$

{ass:pg-eb} **Assumption 3.3 (problem with proximal gradient and error bound)**

The following assumption is about (F, f, g, L, S, γ) . Assume that

- (i) (f, g, L) satisfies Assumption 1.9,
- (ii) let $\tau > 0$ be the step size inverse, let T_τ be the proximal operator of $f + g$ as given by $T_\tau(x) := \text{prox}_{\tau^{-1}g}(x - \tau^{-1}\nabla f(x))$,
- (iii) $S = \argmin_x f(x) + g(x) \neq \emptyset$,
- (iv) the objective function is given by $F = f + g$ and, it satisfies error bound (Definition 3.2).

{def:ista} **Definition 3.4 (the proximal gradient method)**

Suppose that (f, g, L) satisfies Assumption 1.9. Let $\tau \geq L$, and $x_0 \in \mathbb{R}^n$. Then an algorithm

is a proximal gradient method if it generates iterates $(x_k)_{k \geq 0}$ such that they satisfy for all $k \geq 1$:

$$x_{k+1} = \text{prox}_{\tau^{-1}g}(x_k + \tau^{-1}\nabla f(x_k)).$$

3.1 Error bound and linear convergence of ISTA

The following theorem characterizes linear convergence of the proximal gradient method under gradient mapping error bound condition.

{thm:lin-cnvg-ista-eb}

Theorem 3.5 (linear convergence under gradient mapping error bound)

Assume that (F, f, g, L, S, γ) is given by Assumption 3.3. Under this assumption, the iterates $(x_k)_{k \geq 0}$ given by Definition 3.4 satisfies for all $k \geq 0, \bar{x} \in S$ and $\tau \geq L$:

$$F(x_{k+1}) - F(\bar{x}) \leq \left(1 - \frac{\gamma}{2\tau}\right)(F(x_k) - F(\bar{x})), \text{ where } \frac{\gamma}{2\tau} \in (0, 1)$$

Hence, the algorithm generates $F(x_k) - F(\bar{x}) \leq \mathcal{O}((1 - \gamma/(2\tau))^k)$.

Proof. Two important immediate results will be presented first. Consider the proximal gradient inequality from Theorem 1.13, but with $\rho = 0, \epsilon = 0, B = \tau$, then for all x such that $\|\mathcal{G}_\tau(x)\| > 0$ it has for $\tilde{x} = T_\tau(x), z \in \mathbb{R}^n$ the inequality

$$\begin{aligned} F(\tilde{x}) - F(z) &\leq \frac{\tau}{2}\|x - z\|^2 - \frac{\tau}{2}\|z - \tilde{x}\|^2 \\ &= -\frac{\tau}{2}\|x - \tilde{x}\|^2 + \tau\langle x - z, x - \tilde{x} \rangle \\ &= -\frac{1}{2\tau}\|\mathcal{G}_\tau(x)\|^2 + \langle x - z, \mathcal{G}_\tau(x) \rangle \\ &\leq -\frac{1}{2\tau}\|\mathcal{G}_\tau(x)\|^2 + \|x - z\|\|\mathcal{G}_\tau(x)\| \\ &= \|\mathcal{G}_\tau(x)\|^2 \left(\frac{\|x - z\|}{\|\mathcal{G}_\tau(x)\|} - \frac{1}{2\tau} \right). \end{aligned}$$

Now, for all $z = \bar{x} \in S$, from Assumption 3.6 $\exists \gamma > 0$ such that:

$$\frac{\|x - z\|}{\|\mathcal{G}_\tau(x)\|} \leq \frac{\|x - z\|}{\gamma \text{dist}(x|S)} \leq \frac{1}{\gamma}.$$

Hence, for all $\bar{x} \in S$ it has

$$0 \leq F(\tilde{x}) - F(\bar{x}) \leq \|\mathcal{G}_\tau(x)\|^2 \left(\frac{1}{\gamma} - \frac{1}{2\tau} \right). \quad (3.2)$$

Obviously it has $\gamma^{-1} - (1/2)\tau^{-1} > 0$. When $z = x$, we have the inequality:

$$\{ineq:lin-cnvg-ista-eb-pitem2\} \quad F(\tilde{x}) - F(x) \leq -\frac{1}{2\tau} \|\mathcal{G}_\tau(x)\|^2. \quad (3.3)$$

To derive the linear convergence, we use (3.2) with $x = x_k, \tilde{x} = x_{k+1}$:

$$\begin{aligned} 0 &\leq \|\mathcal{G}_\tau(x_k)\|^2 \left(\frac{1}{\gamma} - \frac{1}{2\tau} \right) - F(x_{k+1}) + F(\bar{x}) \\ &= \frac{1}{2\tau} \|\mathcal{G}_\tau(x_k)\|^2 \left(\frac{2\tau}{\gamma} - 1 \right) - F(x_{k+1}) + F(\bar{x}) \\ &\stackrel{(1)}{\leq} \left(\frac{2\tau}{\gamma} - 1 \right) (F(x_k) - F(x_{k+1})) - F(x_{k+1}) + F(\bar{x}) \\ &= \left(\frac{2\tau}{\gamma} - 1 \right) (F(x_k) - F(\bar{x}) + F(\bar{x}) - F(x_{k+1})) - F(x_{k+1}) + F(\bar{x}) \\ &= \frac{2\tau}{\gamma} (F(\bar{x}) - F(x_{k+1})) + \left(\frac{2\tau}{\gamma} - 1 \right) (F(x_k) - F(\bar{x})). \end{aligned}$$

At (1) we used (3.3). Multiple both side by $\frac{\gamma}{2\tau}$ then we are done. \blacksquare

3.2 Conditions for linear convergence of the proximal problem

{sec:conds-lin-cnvg-pp}

In this section, we will focus on the sufficient characterization of the proximal problem which allows proximal gradient method to achieve linear convergence rate. The following assumption characterizes a set of sufficient conditions of the proximal problem such that linear convergence rate of applying ISTA to dual proximal objective Ψ_λ can be achieved.

{ass:lin-cnvg-for-pp}

Assumption 3.6 (conditions for linear convergence of proximal problem)

This assumption is about $(g, \omega, A, y, \Phi_\lambda, \Psi_\lambda, \gamma_\lambda)$. Here are the assumptions

{ass:lin-cnvg-for-pp:item1}

(i) $y \in \mathbb{R}^n$ is the vector of which the proximal problem is anchored at, fixed it to be an arbitrary vector.

{ass:lin-cnvg-for-pp:item2}

(ii) Assume (g, ω, A) satisfies Assumption 1.15. The primal objective Φ_λ is given by (1.1), and dual Ψ_λ by (1.2). This means if we let $h(v) := \frac{1}{2\lambda} \|\lambda A^\top v - y\|^2 - \frac{1}{2\lambda} \|y\|^2$, then $\Psi_\lambda(v) = h(v) + \omega^*(v)$.

{ass:lin-cnvg-for-pp:item3}

(iii) Next, assume $\Psi_\lambda = h + \omega^*$ satisfies error bound condition (Assumption 3.3) with $\Psi_\lambda = F$ and, $\gamma = \gamma_\lambda$ and $f = h$. Note that we can do this because h is quadratic hence obviously Lipschitz continuous and Lipschitz smooth.

The following definition specifies the algorithm that can achieve linear convergence rate with the assumptions above.

{def:ista-inner-lp}

Definition 3.7 (the ISTA inner loop algorithm)

Let $\lambda > 0$, $\epsilon > 0$, and $(g, \omega, A, y, \Phi_\lambda, \Psi_\lambda, \gamma_\lambda)$ satisfies Assumption 3.6.

- (i) Let $v_0 \in \text{dom } \omega^*$ be a feasible initial guess of Ψ_λ , and let $\tau \geq \lambda \|AA^\top\|$ be the inverse step size.
- (ii) Define $z_0 = y - \lambda A^\top v_0$ and smooth part of Ψ_λ as $h := v \mapsto \frac{1}{2\lambda} \|\lambda A^\top v - y\|^2$.

The algorithm that solves the proximal problem generates the primal dual sequences (z_j, v_j) such that for all $j = 0, 1, 2, \dots$, they satisfy:

$$\begin{aligned} v_{j+1} &= \text{prox}_{\tau^{-1}\omega^*}(v_j - \tau^{-1} \nabla h(v_j)), \\ z_{j+1} &= y - \lambda A^\top v_{j+1}. \end{aligned}$$

Terminates if $\mathbf{G}_\lambda(z_j, v_j) \leq \epsilon$ where \mathbf{G}_λ is given by (1.3), then the result we want is z_j .

Remark 3.8 The value of $\mathbf{G}_\lambda(z_j, v_j)$ is easy to compute, it only requires access to matrix A, A^\top , and the function ω . In case when the proximal operator for ω^* is nontrivial, we can use the Moreau identity and the proximal operator of ω instead. The gradient $\nabla f(v)$ is easy to compute, and it is: $AA^\top v - Ay$.

The following propositions precisely show that the linear convergence is achievable when Assumption 3.6 holds.

{prop:inn-loop-lin-cnvg}

Proposition 3.9 (sufficient conditions of linear convergence of the inner loop)

Let the parameters $(g, \omega, A, y, \Phi_\lambda, \Psi_\lambda, \gamma_\lambda)$ of a proximal problem satisfy Assumption 3.6. Let $\tau, v_0, \epsilon > 0$ be given by Definition 3.7 along with iterates $(z_j, v_j)_{j \geq 0}$. Let \bar{v} be a minimizer of Ψ_λ , then the followings are true:

- (i) Let \bar{v} be minimizer of Ψ_λ , the sequence $\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})$ converges linearly to zero.
- (ii) There exists constants K_Ψ^λ such that Ψ_λ is K_Ψ^λ Lipschitz continuous on $\text{dom } \Psi_\lambda$, as a consequence there exists constant C_Ψ^λ such that, $\sup_{v \in \mathbb{R}^m} \Phi_\lambda(v) - \Phi_\lambda(\bar{v}) \leq K_\Psi^\lambda \sup_{v \in \text{dom } \Psi_\lambda} \|v - \bar{v}\| = C_\Psi^\lambda < \infty$.
- (iii) The duality gap has a linear convergence rate by the inequality:

$$\mathbf{G}_\lambda(z_j, v_j) \leq \left(1 - \frac{\gamma_\lambda}{2\tau}\right)^{j/2} C_\lambda \text{ where: } C_\lambda = \sqrt{2\lambda C_\Psi^\lambda} \left(K_\omega \|A\| + K_\omega + \frac{\sqrt{2\lambda C_\Psi^\lambda}}{\lambda} \right).$$

{prop:inn-loop-lin-cnvg:item4}

- (iv) If $\mathbf{G}_\lambda(z_j, v_j) \leq \epsilon$, then $z_j \approx_\epsilon \text{prox}_{\lambda g}(y)$ and the number of iterations sufficient to achieve the accuracy would be

$$j \geq \left\lceil \frac{2 \ln \left(\frac{\epsilon}{C_\lambda} \right)}{\ln \left(1 - \frac{\gamma_\lambda}{2\tau} \right)} \right\rceil.$$

Proof. To start, we prove (i). Recall that it has $\Psi_\lambda = h + \omega^*$ with h being Lipschitz smooth and, ω^* closed convex proper from item (ii) in Assumption 3.6. In addition Ψ_λ also satisfies error bound in item (iii) of Assumption 3.6 with $\gamma = \gamma_\lambda$. Finally, $(z_j, v_j)_{j \geq 0}$ given by Definition 3.7 has $\tau \geq \lambda \|A^\top A\|$ so it's essentially Definition 3.4 with $x_j = v_j$. Therefore results from Theorem 3.5 applies with $F = \Psi_\lambda, x_j = v_j$ and it has:

$$\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}) \leq \left(1 - \frac{\gamma_\lambda}{2\tau}\right)^j (\Psi_\lambda(v_0) - \Psi_\lambda(\bar{v})). \quad (3.4)$$

Next, we show (ii). The smooth part of Φ_λ is the quadratic $h(v) = \frac{1}{2\lambda} \|A^\top v - y\|^2 - \frac{1}{2\lambda} \|y\|^2$, it's locally Lipschitz on the compact domain of $\text{dom } \Psi_\lambda = \text{dom } \omega^*$ is compact by Assumption 1.15, hence the smooth part is Lipschitz on $\text{dom } \omega^*$. In addition, ω^* is assumed to be Lipschitz on $\text{dom } \omega^*$, therefore the entire function Ψ is Lipschitz on ω^* , we denote the Lipschitz constant by K_Ψ^λ . Because ω^* is bounded, therefore there exists constant C_Ψ^λ such that.

$$\sup_{v \in \mathbb{R}^m} \Phi_\lambda(v) - \Phi_\lambda(\bar{v}) \leq K_\Psi^\lambda \sup_{v \in \text{dom } \Psi_\lambda} \|v - \bar{v}\| = C_\Psi^\lambda < \infty.$$

To see the duality gap, Theorem 1.20(iii) applies because v_j by (3.4) provides us the sequence $(v_j)_{j \geq 0}$ such $\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})$ converges.

$$\begin{aligned} & \mathbf{G}_\lambda(z_j, v_j) \\ &= \Phi_\lambda(z_j) - \Phi_\lambda(\bar{z}) + 0 \\ &= \Phi_\lambda(z_j) - \Phi_\lambda(\bar{z}) + \Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}) \\ &\leq \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + K_\omega + \frac{\sqrt{2\lambda}}{2\lambda} \sqrt{\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})} \right) \\ &\quad + \Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}) \\ &= \sqrt{2\lambda(\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + K_\omega + \frac{\sqrt{2\lambda}}{\lambda} \sqrt{\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})} \right) \\ &\stackrel{(1)}{\leq} \left(1 - \frac{\gamma_\lambda}{2\tau}\right)^{j/2} \sqrt{2\lambda(\Psi_\lambda(v_0) - \Psi_\lambda(\bar{v}))} \left(K_\omega \|A\| + K_\omega + \frac{\sqrt{2\lambda}}{\lambda} \sqrt{\Psi_\lambda(v_j) - \Psi_\lambda(\bar{v})} \right) \\ &\stackrel{(2)}{\leq} \left(1 - \frac{\gamma_\lambda}{2\tau}\right)^{j/2} \sqrt{2\lambda C_\Psi^\lambda} \left(K_\omega \|A\| + K_\omega + \frac{\sqrt{2\lambda C_\Psi^\lambda}}{\lambda} \right). \end{aligned}$$

At (1) we used (3.4). At (2) we used (ii). To see (iv), use Theorem 1.19 and the rest is just some algebra because $\gamma/(2\tau) \leq 1$ from 3.5. ■

Remark 3.10 In practice, use the proximal operator of ω^* to choose a feasible v_0 , equivalently we can translate any primal feasible solution into a dual feasible solution using Theorem 1.18.

As a prelude, to have a global convergence rate it requires an upper bound of the initial optimality gap of Ψ_λ of for all iterations of the outer loop. Suppose for each outer iteration k , the inner loop is initialized to start on $v_0^{(k)}$, then there must be an upper bound for all $\Psi_\lambda(v_0^{(k)}) - \Psi_\lambda(\bar{v})$.

3.3 Concrete examples where we can have linear convergence

Continuing our discussion from the previous section, the goal of this section is to show that there exists specific type of ω that appears in applications and satisfies the error bound condition for the proximal problem Φ_λ , and hence the results from the previous section are applicable and relevant to practical applications.

Our major results are stated in Proposition 3.19.

{def:q-scnvx}

Definition 3.11 (smooth quasi strongly convex [2, Definition 1])

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex, differentiable and, L lipschitz smooth. Let $X \subseteq \mathbb{R}^n$ and suppose that the set of minimizers $X^+ = \operatorname{argmin}_{x \in X} f(x) \neq \emptyset$ exists and, denote f^+ to be the minimum of f on X . It is smooth quasi strongly convex (SQSC) on the set $X \subseteq \mathbb{R}^n$ if $\exists \kappa > 0$ such that $\forall x \in X$, let $\bar{x} = \Pi_{X^+}x$, it satisfies

$$0 \leq f^+ - f(x) - \langle \nabla f(x), \bar{x} - x \rangle - \frac{\kappa}{2} \|x - \bar{x}\|^2 \quad (\forall x \in X).$$

{def:q-grwth}

Definition 3.12 (quadratic growth condition)

Let $F = f + g$ where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex differentiable and, $g : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be closed convex and, proper. Assume in addition that the set of minimizer $S := \operatorname{argmin}_{x \in \mathbb{R}^n} F \neq \emptyset$, denote F^+ to be the minimum, then it satisfies the quadratic growth conditions if there exists $\kappa > 0$ such that

$$(\forall x \in \mathbb{R}^n) F(x) \geq F^+ + \frac{\kappa}{2} \operatorname{dist}^2(x|S).$$

Obviously, The following theorem shows that under the standard framework of smooth nonsmooth additive compsite objective function, the quadratic growth condition implies error bound.

{thm:qfg-means-eb}

Theorem 3.13 (quadratic growth implies error bound)

Suppose (f, g, L) satisfies Assumption 1.9, let $F = f + g$. In addition, let's assume that

- (i) the set of minimizers $S = \operatorname{argmin}_x F \neq \emptyset$,
- (ii) the function F satisfies quadratic growth (Defintion 3.12) with some parameter $\kappa > 0$.

Then, by choosing any $\tau \geq L$, the function F also satisfies error bound (Definition 3.2) with:

$$\gamma = \frac{\tau\kappa}{\kappa + \tau + \sqrt{\tau(\kappa + \tau)}}.$$

Proof. Denote $x^+ = T_\tau(x)$, $\bar{x} = \Pi_S x$, and $\bar{x}^+ = \Pi_S x^+$. We use the exact proximal gradient inequality by applying Theorem 1.13 with $\epsilon = 0, \rho = 0$ and, $z = \bar{x}^+$ (we can apply this because $\tau \geq L$) which gives:

$$\begin{aligned} 0 &\geq F(x^+) - F(\bar{x}^+) - \tau \langle \bar{x}^+ - x^+, x^+ - x \rangle - \frac{\tau}{2} \|x - x^+\|^2 \\ &\stackrel{(1)}{\geq} \frac{\kappa}{2} \|x^+ - \bar{x}^+\|^2 - \|\tau(\bar{x}^+ - x^+)\| \|x^+ - x\| - \frac{1}{2\tau} \|\tau(x - x^+)\|^2 \\ &= \frac{\kappa}{2} \|x^+ - \bar{x}^+\|^2 - \frac{\tau}{2} (2 \|\bar{x}^+ - x^+\| \|x^+ - x\| + \|x - x^+\|^2) \\ &\stackrel{(2)}{=} \frac{\kappa + \tau}{2} \|x^+ - \bar{x}^+\|^2 - \frac{\tau}{2} (\|\bar{x}^+ - x^+\| + \|x^+ - x\|)^2 \\ &\stackrel{(3)}{\geq} \frac{\kappa + \tau}{2} (\|x - \bar{x}\| - \|x - x^+\|)^2 - \frac{\tau}{2} (\|x - \bar{x}\|)^2 \\ \iff 0 &\geq \sqrt{\kappa + \tau} (\|x - \bar{x}\| - \|x - x^+\|) - \sqrt{\tau} (\|x - \bar{x}\|) \\ &= (\sqrt{\kappa + \tau} - \sqrt{\tau}) \|x - \bar{x}\| - \sqrt{\kappa + \tau} \|x - x^+\|. \end{aligned}$$

At (1), we used quadratic growth assumption and substitute the inequality of Definition 3.12, and the Cauchy inequality. At (2), consider substituting $a = \bar{x}^+ - x^+, b = x^+ - x$, so the second terms with the parenthesis has inside $2\|a\|\|b\| + \|b\|^2 = (\|a\| + \|b\|)^2 - \|a\|^2$. At (3), we used the properties of projecting onto the set of minimizer S , which gives:

$$\begin{aligned} \|x - \bar{x}\| &\stackrel{(4)}{\leq} \|x - \bar{x}^+\| \leq \|x - x^+\| + \|x^+ - \bar{x}^+\| \\ &\implies \|x - \bar{x}\| - \|x - x^+\| \leq \|x^+ - \bar{x}^+\| \end{aligned}$$

At (4) we used the fact that $\bar{x} = \Pi_S(x)$ is closer to S than the point \bar{x}^+ . Finally, re-arranging the results it should yield the following inequality:

$$\begin{aligned} \|\mathcal{G}_\tau(x)\| &= \|\tau(x - x^+)\| \\ &\geq \tau \left(\frac{\sqrt{\kappa + \tau} - \sqrt{\tau}}{\sqrt{\kappa + \tau}} \right) \|x - \bar{x}\| \\ &= \frac{\tau\kappa}{\kappa + \tau + \sqrt{\tau(\kappa + \tau)}} \|x - \bar{x}\|. \end{aligned}$$

The above is error bound conditions as defined in Definition 3.2, with $\gamma = \frac{\tau\kappa}{\kappa + \tau + \sqrt{\tau(\kappa + \tau)}}$. \blacksquare

Remark 3.14 This theorem is not a new result, and the converse of the statement remains true and is not hard to show. We adapted claim and the proof here so the notations stays consistent and to make it self contained.

{fact:qscnvx-q-growth}

Fact 3.15 (quasi strongly convex implies quadratic growth [2, Theorem 4])

Let f, X, κ be given by Definition 3.11. Then the function $F = f + \delta_X$ satisfies Definition 3.12 (quadratic growth) with the same κ .

{fact:hoffm-eb}

Paraphrased in Necoara et al. [2] is the following classic fact on Hoffman error bound.

Fact 3.16 (Hoffman error bound) Consider a nonempty polyhedral set $P = \{x \in \mathbb{R}^n : Ax = b, Cx \leq d\}$ defined via some $A \in \mathbb{R}^{p \times n}, C \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m, d \in \mathbb{R}^m$. Then there exists a constant $\theta > 0$ which only depends on A, C such that:

$$(\forall x \in \mathbb{R}^n) \quad \text{dist}(x|P) \leq \theta \text{dist}((Ax - b, Cx - d) | \{\mathbf{0}\} \times \mathbb{R}_-^m).$$

Remark 3.17 In the literature, the exact, the smallest value of θ is a big topic. It has something do with the angle and geometric between a cone, or an affine space of two convex sets. Here, we will only state its existence without giving it a precise expression.

{fact:polyhedral-qscnvx-fxn}

Fact 3.18 (a type of smooth quasi strongly convex function [2, Theorem 8])

Consider any $C \in \mathbb{R}^{m \times n}, d \in \mathbb{R}^m$ such that it defines nonempty polyhedral set $X = \{x : Cx \leq d\}$. Let h be $\sigma > 0$ strongly convex and L_h Lipschitz smooth, consider $f = h \circ A$ where $A \in \mathbb{R}^{p \times n}$. Then

- (i) The set of minimizer X^+ is nonempty and it's a polyhedral which its definition involves matrices C, A .
- (ii) The function f is quasi strongly convex with $\kappa = \sigma/\theta^2, L_f = L_h \|A\|^2$ where θ is the Hoffman constant as given by Fact 3.16 of the polyhedral set X^+ .

Now, we give our first major results of this section. The proposition that follows characterizes a precise case where Assumption 3.6 is true, and it's a case widely available in applications.

{prop:1nrm-prox-problem}

Proposition 3.19 (1-norm problem has a linear convergence rate)

Let (g, ω, A) satisfies Assumption 1.15. In addition, if $\omega = \|\cdot\|_1$, then the function $\Psi_\lambda, \Phi_\lambda$ as given by (1.1), (1.2) satisfies Assumption 3.6 and progressively it can be breakdown into:

{prop:1nrm-prox-problem:result1}

(i) Ψ is quasi strongly convex, with $\kappa = \lambda/\theta^2$ where θ is a Hoffman constant of polyhedral set $X^+ = \underset{x \in \mathbb{R}^n}{\text{argmin}} \Psi_\lambda$, and Lipschitz smoothness constant $L_\Psi = \lambda \|A^\top\|^2$.

{prop:1nrm-prox-problem:result2}

(ii) Ψ satisfies quadratic growth condition (Definition 3.12), with the same constant κ from its quasi strong convexity.

{prop:lnrm-prox-problem:result3}

- (iii) Ψ satisfies the error bound condition (Definition 3.2) for all step size $\tau \geq L_\Psi$. And the error bound constant has

$$\gamma = \frac{\frac{\lambda}{\theta^2}}{\frac{\lambda}{\tau\theta^2} + 1 + \sqrt{\frac{\lambda}{\theta^2\tau} + 1}}.$$

{prop:lnrm-prox-problem:result4}

- (iv) The primal, and dual satisfies Assumption 3.6, hence Proposition 3.9 applies and, a linear convergence rate is possible.

And they have the relation (i) \Rightarrow (ii) \Rightarrow (iii).

Proof. We first show (i). Recall $\Psi_\lambda(v) = \frac{1}{2\lambda}\|\lambda A^\top v - y\|^2 + \omega^*(v) - \frac{1}{2\lambda}\|y\|^2$ where $\omega^* = \delta(z|\{x : \|x\|_1 \leq 1\})$ so it's a quadratic function with a polyhedral constraint. Observe that that $\Psi_\lambda(v) = h(A^\top v) + \delta_X$, and here it has $h = z \mapsto (\lambda/2)\|z - \lambda^{-1}y\|^2$ so h is a λ strongly convex and smooth function, and $X = \{x : \|x\|_1 \leq 1\}$ which can be represented by a polyhedral set characterized by linear inequalities. Therefore, Fact 3.18 applies, and Ψ_λ is quasi strongly convex with $\kappa = \lambda/\theta^2$, $L_\Phi = \lambda\|A^\top\|^2$.

To see (ii), use Fact 3.15. To see (iii), use Theorem 3.13. To see (iv), we already have most of the results from the assumption and (iii), and it remains to verify that Φ_λ is Lipschitz continuous on $\text{dom}(\omega \circ A)$ with is true because $\Phi_\lambda = \|Au\|_1 + \frac{1}{2\lambda}\|u - y\|^2$, it's the sum of a quadratic and a Lipschitz continuous function $\|Au\|_1$. ■

Remark 3.20 It is very difficult to obtain a lower estimate for γ in practice.

Definition 3.21 (piecewise linear-quadratic function [3, Definition 10.20])

A function $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ is piecewise linear-quadratic when $\text{dom } f$ is the union of finitely many piecewise polyhedral set, such that on each polyhedral partition of f there exists $\alpha \in \mathbb{R}$, $a \in \mathbb{R}^n$, $C \in \mathbb{R}^{n \times n}$ as a symmetric matrix is where f has the representation $x \rightarrow \langle x, Cx \rangle + \langle a, x \rangle + c$.

The following proposition characterizes another case where error bound conditions of problem

{example:inn-lp-lin-cnvg}

holds.

Example 3.22 (inner loop linear convergence rate examples)

4 Total complexity of the algorithm

This section puts results regarding the total complexity of the proposed inexact proximal gradient algorithm.

4.1 Inner loop complexity

We remind that the parameters in the inner loop will change depending on the outer loop's iteration. To discuss the convergence we unfortunately have to re-introduce the proximal problem in the context of the accelerated proximal gradient method of the outer loop.

Let (F, f, g, L) and sequence $(\alpha_k, B_k, \rho_k, \epsilon_k)_{k \geq 0}$ satisfies Definition 2.1. Fix any $k \geq 0$ to be the iteration counter of the outer loop. Let (g, ω, A) satisfies Assumption 1.15. The inner loop is an algorithm that solves the inexact proximal gradient problem $x \approx_{\epsilon_k} T_{B_k + \rho_k}(y_k)$ which is equivalent to evaluating:

$$x_k \approx_{\epsilon_k} \text{prox}_{L_k^{-1}g}(y_k - L_k^{-1}\nabla f(y_k)).$$

Where, $L_k := B_k + \rho_k$. Let $\lambda^{(k)} := L_k^{-1}$, $\tilde{y}_k := y_k - L_k^{-1}\nabla f(y_k)$, the proximal problem boils down to optimizing the following function $\Phi_\lambda^{(k)}$ as defined by:

$$\{\text{eqn:primal-pp-k}\} \quad \Phi_\lambda^{(k)}(u) := \frac{1}{2\lambda^{(k)}} \|u - \tilde{y}_k\|^2 + \omega(Au). \quad (4.1)$$

And its dual which is

$$\{\text{eqn:dual-pp-k}\} \quad \Psi_\lambda^{(k)}(v) := \frac{1}{2\lambda^{(k)}} \|\lambda^{(k)}A^\top v - \tilde{y}_k\|^2 + \omega^*(v) - \frac{1}{2\lambda^{(k)}} \|\tilde{y}_k\|^2. \quad (4.2)$$

The primal dual gap is $\mathbf{G}_\lambda^{(k)}(u, v) = \Phi_\lambda^{(k)}(u) + \Psi_\lambda^{(k)}(v)$. The above are identical to proximal problem defined in Section 1.3 except for the introduction of iteration counter k from the outer loop, and we had specified \tilde{y}_k in relation to the outer loop. Finally, the inner loop algorithm is responsible for optimizing the optimality gap, it satisfies that $\mathbf{G}_\lambda^{(k)}(u, v) \leq \epsilon_k$ and all results from Section 3.2 applies.

The following assumption specifies conditions where the complexity of the inner using ISTA can be globally bounded for all iteration of the accelerated proximal gradient algorithm of the outer loop independent of the initial guess $v_0^{(k)}, \lambda^{(k)}$.

{ass:inn-cmplx}

Assumption 4.1 (globally bounded inner loop complexity)

For any integer $k \geq 0$, this assumption is about parameters of the proximal problem $(g, \omega, A, \tilde{y}_k, \Phi_\lambda^{(k)}, \Psi_\lambda^{(k)}, L_\Phi^{\lambda, (k)}, \gamma_\lambda^{(k)})$ introduced at the beginning of this section, iterates $(z_j^{(k)}, v_j^{(k)})_{j \geq 0}$, the primal dual solutions $(\bar{z}^{(k)}, \bar{v}^{(k)})$ and additional constants $\lambda^{\min}, \lambda^{\max}, \gamma^{\min}$. The assumptions now follow.

{ass:inn-cmplx:item1}

- (i) As specified in this section, the parameter $\lambda^{(k)} = (B_k + \rho_k)^{-1}$ are from the outer loop. We denote $L_k = B_k + \rho_k$ for short. We assume in addition, there exists $\lambda^{\min}, \lambda^{\max}$ such that:

$$-\infty < \lambda^{\min} \leq \inf_{k \in \mathbb{N} \cup \{0\}} \lambda^{(k)} \leq \sup_{k \in \mathbb{N} \cup \{0\}} \lambda^{(k)} \leq \lambda^{\max} < \infty.$$

- {ass:inn-cmplx:item2} (ii) There exists the smallest error bound constant γ^{\min} such that it lower bound all the constants for the proximal problem, i.e: $\exists \gamma_{\min} : 0 < \gamma^{\min} \leq \inf_{k \in \mathbb{N} \cup \{0\}} \gamma_{\lambda}^{(k)}$.
- {ass:inn-cmplx:item3} (iii) $\tilde{y}_k, \Phi_{\lambda}^{(k)}, \Psi_{\lambda}^{(k)}, L_{\Phi}^{\lambda, (k)}, \gamma_{\lambda}^{(k)}$ all satisfies error bound conditions (Assumption 3.6) with $y = \tilde{y}_k, \Phi_{\lambda} = \Phi_{\lambda}^{(k)}, \Psi_{\lambda} = \Psi_{\lambda}^{(k)}, L_{\Phi}^{\lambda} = L_{\Phi}^{\lambda, (k)}$ and $\gamma_{\lambda} = \gamma_{\lambda}^{(k)}$.
- {ass:inn-cmplx:item4} (iv) The iterates of the inner loop $(z_j^{(k)}, v_j^{(k)})_{j \geq 0}$ are produced by an ISTA algorithm satisfying Definition 3.7. We choose the same stepsize $\tau = 1/\|A^T A\|$ for all iteration of ISTA.

Here, we discuss shorter which assumption are easy to satisfies, and which are more demanding. Item (i) is easy to satisfies since popular line search and backtracking method ensures that B_k will be bounded above and blow, and ρ_k is entirely to the choice of the practitioner so any bounded sequence works. Item (ii) requires some extra argument but one example of ω remains feasible. Observe that it is trivially true if the domain of $\Psi_{\lambda}^{(k)}$ is always bounded because initial guess $v_0^{(k)} \in \text{dom } \Psi_{\lambda}^{(k)}$, which is equivalent to ω^* has a bounded domain which can be satisfied if $\omega = \|\cdot\|_1$. In addition, the results of Proposition 3.19(iii), 3.19(iv), item (iii) can be satisfied because $\gamma_{\lambda}^{(k)}$ only depends $\lambda^{(k)}$ it would be bounded from (i).

N'oubliez pas se faire ici.

The proposition that follows characterize the bare minimum requirements so that the inner loop's complexity depends only on required accuracies ϵ , and parameter λ for all initial guess of the outer loop.

{prop:inner-lp-cmplx}

Proposition 4.2 (inner loop complexity can be bounded globally)

Let $(g, \omega, A, \tilde{y}_k, \Phi_{\lambda}^{(k)}, \Psi_{\lambda}^{(k)}, L_{\Phi}^{\lambda, (k)}, \gamma_{\lambda}^{(k)})$, $(z_j^{(k)}, v_j^{(k)})_{j \geq 0}$, $(\bar{z}^{(k)}, \bar{v}^{(k)})$ and additional parameters $\lambda^{\min}, \lambda^{\max}, \gamma^{\min}, C_0$ such that they satisfy Assumption 4.1. Then, the total number of inner loop iteration sufficient to achieve $\mathbf{G}_{\lambda}^{(k)}(z_j^{(k)}, v_j^{(k)}) \leq \epsilon_k$, at each iteration k denoted by $J_{\epsilon}^{(k)}$ can be upper bounded by:

Proof. To prove the above theorem, we use the results from Proposition 3.9 into the context of the iterations in the outer loop.

■

4.2 Overall complexity

We derive the overall complexity in this section. Before we start, it's necessary to discuss the best initial guess vector $v_0^{(k)}$ of the inner loop so that item (ii) in Assumption 4.1 can be satisfied for as many objective functions as possible.

A Super boring chores

Super boring stuff that I just want to skip but somehow necessary.

{lemma:chore1} **Lemma A.1 (That conjugate for the dual of proximal problem)**
Let $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}} := \frac{1}{2\lambda} \|u - v\|^2$, then its conjugate is given by

$$f^*(v) = \frac{1}{2\lambda} \|\lambda v + y\|^2 - \frac{1}{2\lambda} \|y\|^2.$$

Proof. Recall the following properties for any closed, proper convex function $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$. Let $a \in \mathbb{R}^n$ be any vector, let $\alpha \in \mathbb{R}, c \in \mathbb{R}$ be some scalar, then we have these three properties of conjugating convex funtion.

- (i) $(\alpha f)^* = \alpha f^* \circ (\alpha^{-1}I)$.
- (ii) $(f + c)^*(y) = f^*(y) - c$.
- (iii) $(x \mapsto f(x) + \langle x, y \rangle)^*(y) = f^*(y - a)$.

From here we have:

$$\begin{aligned} f^*(v) &= \left(u \mapsto \lambda^{-1} \left(\frac{1}{2} \|u\|^2 - \langle u, y \rangle \right) + \frac{1}{2\lambda} \|y\|^2 \right)^*(v) \\ &= \left(u \mapsto \lambda^{-1} \left(\frac{1}{2} \|u\|^2 - \langle u, y \rangle \right) \right)^*(v) - \frac{1}{2\lambda} \|y\|^2 \\ &= \left[\lambda^{-1} \left(u \mapsto \left(\frac{1}{2} \|u\|^2 - \langle u, y \rangle \right) \right)^* \circ (\lambda I) \right] (v) - \frac{1}{2\lambda} \|y\|^2 \\ &= \left[\lambda^{-1} \left(u \mapsto \left(\frac{\|\cdot\|^2}{2} \right)^* (u + y) \right) \circ (\lambda I) \right] (v) - \frac{1}{2\lambda} \|y\|^2 \\ &= \left[\lambda^{-1} \left(u \mapsto \frac{\|u + y\|^2}{2} \right) \circ (\lambda I) \right] (v) - \frac{1}{2\lambda} \|y\|^2 \\ &= \lambda^{-1} \left(\frac{1}{2} \|\lambda v + y\|^2 \right) - \frac{1}{2\lambda} \|y\|^2. \end{aligned}$$

■

{lemma:lipz-cnvx-fxn} **Lemma A.2 (Lipschitz constant of convex function)** *Let $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be a closed, convex, proper function. Let ∂f be its convex subgradient. Then,*

- (i) *for all $x \in \mathbb{R}^n, y \in \mathbb{R}^n$ it has: $|f(x) - f(y)| \leq \sup_{x \in \text{dom } \partial f} (\partial f(x) | \mathbf{0}) \|y - x\|$.*
- (ii) *If in addition, the function is K_f Lipschitz continuous globally on \mathbb{R}^n , then: $(\forall y \in \mathbb{R}^n)(\forall v_y \in \partial f(y)) : K_f \geq \|v_y\|$.*

Proof. The proof is direct. Let $x, y \in \mathbb{R}^n$ be arbitrary. Choose $v_x \in \partial f(y)$ and $v_y \in \partial f(y)$ such that $\|v_x\| = \text{dist}(\partial f(x) | \mathbf{0}), \|v_y\| = \text{dist}(\partial f(y) | y)$ which is possible because $\partial f(x)$ is closed for all $x \in \text{dom } \partial f$. Therefore, we have the following.

$$\begin{aligned} |f(x) - f(y)| &\leq \max(f(x) - f(y), f(y) - f(x)) \\ &\stackrel{(1)}{\leq} \max(-\langle v_x, y - x \rangle, -\langle v_y, x - y \rangle) \\ &\leq \max(\|v_x\|, \|v_y\|) \|y - x\| \\ &\leq \left(\sup_{x \in \text{dom } \partial f} \text{dist}(\partial f(x) | \mathbf{0}) \right) \|y - x\|. \end{aligned}$$

At (1), we used the fact that $f(x) - f(y) \leq -\langle v_x, y - x \rangle$ and, $f(y) - f(x) \leq -\langle v_y, x - y \rangle$ by convex subgradient inequality.

For the second results it's direct, choose any $y \in \text{dom } f = \mathbb{R}^n, v_y \in \partial f(y)$, then there exists some $x \in \mathbb{R}^n$:

$$\begin{aligned} 0 &\leq \frac{f(x) - f(y)}{\|x - y\|} - \frac{\langle v_y, x - y \rangle}{\|x - y\|} \\ &\leq K_f - \left\langle v_y, \frac{x - y}{\|x - y\|} \right\rangle \\ &= K_f - \|v_y\|. \end{aligned}$$

AT (1), consider choosing any x such that $\frac{x-y}{\|x-y\|} = \frac{v_y}{\|v_y\|}$, this is possible because the domain of f is \mathbb{R}^n . ■

References

- [1] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, CMS Books in Mathematics, Springer International Publishing, Cham, 2017.
- [2] I. NECOARA, Y. NESTEROV, AND F. GLINEUR, *Linear convergence of first order methods for non-strongly convex optimization*, Mathematical Programming, 175 (2019), pp. 69–107.

- [3] R. T. ROCKAFELLAR AND R. J. B. WETS, *Variational Analysis*, vol. 317 of Grundlehren der mathematischen Wissenschaften, Springer, Berlin, Heidelberg, 1998.
- [4] S. VILLA, S. SALZO, L. BALDASSARRE, AND A. VERRI, *Accelerated and inexact forward-backward algorithms*, SIAM Journal on Optimization, 23 (2013), pp. 1607–1633.
- [5] C. ZALINESCU, *Convex analysis in general vector spaces*, World Scientific, River Edge, N.J. ; London, 2002.
- [6] M. ZHANG, M. ZHANG, F. ZHANG, A. CHADDAD, AND A. EVANS, *Robust brain MR image compressive sensing via re-weighted total variation and sparse regression*, Magnetic Resonance Imaging, 85 (2022), pp. 271–286.