

Proximal quasi-Newton method

Andersen Ang

Department of Combinatorics and Optimization, University of Waterloo, Canada

msxang@uwaterloo.ca, where $\mathbf{x} = \lfloor \pi \rfloor$ Homepage: angms.science

First draft: September 7, 2022 Last update: October 17, 2022

Table of Contents

- 1 Proximal Newton's method
- 2 Convergence theory of proximal Newton's method: part 1
 - What happens at convergence: Fixed-point property
 - What happens at not convergence: descent property
 - Step size $\alpha_k \in (0, 1]$
- 3 Convergence theory of proximal Newton's method: part 2
 - Convergence rate with no line search, part 1
 - Convergence rate with no line search, part 2
 - Convergence rate with no line search, part 3
 - A lemma on decreasing sequence
 - Convergence rate $\frac{1}{k}$
- 4 Summary

Minimizing composite function

- ▶ A standard class of convex optimization problem

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

- ▶ $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and L -smooth
 - ▶ $g : \mathbb{R}^n \rightarrow \bar{\mathbb{R}} \cup \{+\infty\}$ is convex and possibly nonsmooth
-
- ▶ Example of (\mathcal{P}) : Basis Pursuit Denoising

$$(\mathcal{P}_1) : \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

- ▶ $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \ll n$ (standard setting)
- ▶ $\mathbf{b} \in \mathbb{R}^m$ is observation
- ▶ $\mathbf{x} \in \mathbb{R}^n$ is variable

Gradient descent and proximal gradient method

Gradient descent

- ▶ Problem class: $\min_{\mathbf{x}} f(\mathbf{x})$, f convex and differentiable
- ▶ Method: gradient descent: $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$.

Proximal gradient method

- ▶ Problem class: $\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x})$, f, g convex, f differentiable
- ▶ Method: proximal gradient $\mathbf{x}_{k+1} = \text{prox}_{\alpha_k g}(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k))$
- ▶ Proximal operator: $\text{prox}_{\gamma g}(\mathbf{x}) := \underset{\boldsymbol{\xi}}{\text{argmin}} \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}\|_2^2 + \gamma g(\boldsymbol{\xi})$.

Proximal gradient method = gradient descent plus g

Both gradient descent and proximal gradient method are based on a local quadratic model of f :

Gradient descent

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2\alpha_k} \|\boldsymbol{\xi} - \mathbf{x}_k\|_2^2 \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2} \left\langle \frac{1}{\alpha_k} \mathbf{I}(\boldsymbol{\xi} - \mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \right\rangle \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} \frac{1}{2\alpha_k} \left\| \boldsymbol{\xi} - \left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \right) \right\|_2^2\end{aligned}$$

Proximal gradient method

$$\begin{aligned}\mathbf{x}_{k+1} &= \operatorname{prox}_{\alpha_k g} \left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \right) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} \frac{1}{2} \left\| \boldsymbol{\xi} - \left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \right) \right\|_2^2 + \alpha_k g(\boldsymbol{\xi}) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} \frac{1}{2\alpha_k} \left\| \boldsymbol{\xi} - \left(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \right) \right\|_2^2 + g(\boldsymbol{\xi}) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2} \left\langle \frac{1}{\alpha_k} \mathbf{I}(\boldsymbol{\xi} - \mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \right\rangle + g(\boldsymbol{\xi})\end{aligned}$$

Generalizing proximal gradient method

- ▶ Proximal gradient method

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2} \langle \frac{1}{\alpha_k} \mathbf{I}(\boldsymbol{\xi} - \mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + g(\boldsymbol{\xi})$$

- ▶ Generalize it by replacing $\frac{1}{\alpha_k} \mathbf{I}$ with a positive definite matrix $\mathbf{H}_k \succ \mathbf{0}$

$$\begin{aligned} \mathbf{x}_{k+1} &= \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{H}_k(\boldsymbol{\xi} - \mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + g(\boldsymbol{\xi}) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{\mathbf{H}_k}^2 + g(\boldsymbol{\xi}) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} \frac{1}{2} \left\| \boldsymbol{\xi} - \left(\mathbf{x}_k - \mathbf{H}_k^{-1} \nabla f(\mathbf{x}_k) \right) \right\|_{\mathbf{H}_k}^2 + g(\boldsymbol{\xi}) \end{aligned}$$

where $\|\mathbf{x}\|_{\mathbf{H}} := \sqrt{\langle \mathbf{H}\mathbf{x}, \mathbf{x} \rangle}$.

Proximal Newton's update

$$\begin{aligned}\mathbf{x}_{k+1} &= \operatorname{argmin}_{\boldsymbol{\xi}} f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x} \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{\mathbf{H}_k}^2 + g(\boldsymbol{\xi}) \\ &= \operatorname{argmin}_{\boldsymbol{\xi}} \frac{1}{2} \left\| \boldsymbol{\xi} - \left(\mathbf{x}_k - \mathbf{H}_k^{-1} \nabla f(\mathbf{x}_k) \right) \right\|_{\mathbf{H}_k}^2 + g(\boldsymbol{\xi}) \\ &=: \operatorname{prox}_g^{\mathbf{H}_k} \left(\mathbf{x}_k - \mathbf{H}_k^{-1} \nabla f(\mathbf{x}_k) \right)\end{aligned}$$

- ▶ If $\mathbf{H}_k = \frac{1}{\alpha_k} \mathbf{I}$, Proximal Newton's method reduces to the classical proximal gradient method.
- ▶ If $\mathbf{H}_k =$ the Hessian of f at \mathbf{x}_k , such update is called Proximal Newton's method.
 - ▶ Advantage: \mathbf{H}_k brings the curvature information at \mathbf{x}_k , improves the local quadratic model.
- ▶ If \mathbf{H}_k is not the Hessian of f at \mathbf{x}_k but an approximate, such update is called Proximal quasi-Newton's method.
- ▶ $\operatorname{prox}_g^{\mathbf{H}_k}$ itself is an optimization problem: can be hard/expensive to solve in general.
 - ▶ A trade-off: you sacrifice for having high computational complexity on solving $\operatorname{prox}_g^{\mathbf{H}_k}$ to gain speed up on the convergence of $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$.

Proximal Newton's update with a line search

$$\mathbf{x}_{k+1} = \text{prox}_g^{\mathbf{H}_k} \left(\mathbf{x}_k - \mathbf{H}_k^{-1} \nabla f(\mathbf{x}_k) \right) \quad (\text{Proximal Newton's update})$$

- ▶ Adding a linear combination with the previous point \mathbf{x}_k

$$\begin{aligned} \mathbf{x}_{k+1} &= (1 - \alpha) \mathbf{x}_k + \alpha \underbrace{\text{prox}_g^{\mathbf{H}_k} \left(\mathbf{x}_k - \mathbf{H}_k^{-1} \nabla f(\mathbf{x}_k) \right)}_{=: \mathbf{x}_{\mathbf{H}_k}^*} \\ &= (1 - \alpha) \mathbf{x}_k + \alpha \mathbf{x}_{\mathbf{H}_k}^* \end{aligned}$$

where $\alpha > 0$ is a stepsize that $\alpha \in (0, 1]$, to be shown later.

- ▶ Why adding a line search: sometimes it improves the convergence.
- ▶ No line search: it means $\alpha_k = 1$ for all k .

Proximal Newton's (with a line search) algorithm

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

- ▶ Starting with an initial point \mathbf{x}_0 , let

$$m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{\mathbf{H}_k}^2,$$

Algorithm 1: Proximal Newton's algorithm

Result: An approximate solution to (\mathcal{P})

- 1 **for** $k = 1, 2, \dots$ **do**
 - 2 $\mathbf{x}_{\mathbf{H}_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi})$ solving the proximal Newton
 subproblem;
 - 3 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{\mathbf{H}_k}^* - \mathbf{x}_k)$ convex combination (line search on α_k)
 - 4 **end**
-

- ▶ If $\alpha_k = 1$ then the algorithm reduces to

$$\mathbf{x}_{k+1} = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}) = \operatorname{prox}_g^{\mathbf{H}_k} \left(\mathbf{x}_k - \mathbf{H}_k^{-1} \nabla f(\mathbf{x}_k) \right).$$

Last slide before mindfuck: summary of the theory

- ▶ The sequence $F(\mathbf{x}_k)$ converges to F^* with with rate $\mathcal{O}(\frac{1}{k})$.
- ▶ The next 20+ pages are just to explain this.
- ▶ **Warning: attempt to understand them will drain your stamina and cause brain damage.**
- ▶ I warned you, if you are ok then let's go.

Table of Contents

- 1 Proximal Newton's method
- 2 Convergence theory of proximal Newton's method: part 1
 - What happens at convergence: Fixed-point property
 - What happens at not convergence: descent property
 - Step size $\alpha_k \in (0, 1]$
- 3 Convergence theory of proximal Newton's method: part 2
 - Convergence rate with no line search, part 1
 - Convergence rate with no line search, part 2
 - Convergence rate with no line search, part 3
 - A lemma on decreasing sequence
 - Convergence rate $\frac{1}{k}$
- 4 Summary

Fixed-point property

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

$$\text{algo} : \begin{cases} \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), & m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{H_k}^2 \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k) \end{cases}$$

Lemma 0. $H_k(\mathbf{x}_{H_k}^* - \mathbf{x}_k) = \mathbf{0}$ if and only if \mathbf{x}_k solves (\mathcal{P}) .

Proof. $H_k(\mathbf{x}_{H_k}^* - \mathbf{x}_k) = \mathbf{0} \stackrel{H_k \succ \mathbf{0}}{\iff} \mathbf{x}_{H_k}^* - \mathbf{x}_k = \mathbf{0} \iff \mathbf{x}_{H_k}^* = \mathbf{x}_k$.

Subgradient 1st-order optimality on $\mathbf{x}_{H_k}^*$: $\mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}) \iff \nabla m_{H_k}(\mathbf{x}_{H_k}^*, \mathbf{x}_k) + \partial g(\mathbf{x}_{H_k}^*) \ni \mathbf{0}$.

Recall

- ▶ the quadratic term in m_{H_k} is $\|\cdot\|_{H_k}^2$, so when we take ∇ there will be a H_k matrix,
- ▶ $H_k \succ \mathbf{0}$ so m_{H_k} is strongly convex so the solution set to $\underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi})$ is singleton, which is $\mathbf{x}_{H_k}^*$

i.e, we have

$$\nabla f(\mathbf{x}_k) + H_k(\mathbf{x}_{H_k}^* - \mathbf{x}_k) + \partial g(\mathbf{x}_{H_k}^*) \ni \mathbf{0},$$

which gives $\nabla f(\mathbf{x}_k) + \partial g(\mathbf{x}_k) \ni \mathbf{0}$ if $\mathbf{x}_{H_k}^* = \mathbf{x}_k$, which is when \mathbf{x}_k is a solution of (\mathcal{P}) . □

Descent property

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

$$\text{algo} : \begin{cases} \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), & m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{H_k}^2 \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k) \end{cases}$$

Lemma 1. While $\mathbf{x}_k \notin \operatorname{argmin} F$, for sufficiently small stepsize $\alpha_k > 0$, we have $F(\mathbf{x}_{k+1}) < F(\mathbf{x}_k)$.

- **Proof.** By definition, $\mathbf{x}_{H_k}^*$ is the minimizer of $\phi := m_{H_k} + g$, i.e., therefore $\phi(\mathbf{x}_{H_k}^*) \leq \phi(\mathbf{x}_{k+1})$. Furthermore, ϕ is strongly convex since $H_k \succ \mathbf{0}$, thus we have strict inequality

$$\begin{aligned} \phi(\mathbf{x}_{H_k}^*) &< \phi(\mathbf{x}_{k+1}) \\ f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{H_k}^* - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{H_k}^* - \mathbf{x}_k\|_{H_k}^2 + g(\mathbf{x}_{H_k}^*) &< f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{H_k}^2 + g(\mathbf{x}_{k+1}). \end{aligned}$$

- Remove $f(\mathbf{x}_k)$, let $\mathbf{d}_k := \mathbf{x}_{H_k}^* - \mathbf{x}_k$ and by definition $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ thus $\mathbf{x}_{k+1} - \mathbf{x}_k = \alpha_k \mathbf{d}_k$

$$\langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{1}{2} \|\mathbf{d}_k\|_{H_k}^2 + g(\mathbf{x}_{H_k}^*) < \alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{\alpha_k^2}{2} \|\mathbf{d}_k\|_{H_k}^2 + \underline{\underline{g(\mathbf{x}_{k+1})}}. \quad (\Delta)$$

- By assumption g is convex so it satisfies the Jensen inequality of convexity

$$g(\mathbf{x}_{k+1}) = g(\mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k)) = g((1 - \alpha_k)\mathbf{x}_k + \alpha_k \mathbf{x}_{H_k}^*) \leq (1 - \alpha_k)g(\mathbf{x}_k) + \alpha_k g(\mathbf{x}_{H_k}^*). \quad (\nabla)$$

Put (∇) into (Δ) gives

$$\langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{1}{2} \|\mathbf{d}_k\|_{H_k}^2 + g(\mathbf{x}_{H_k}^*) < \alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{\alpha_k^2}{2} \|\mathbf{d}_k\|_{H_k}^2 + \underline{\underline{(1 - \alpha_k)g(\mathbf{x}_k) + \alpha_k g(\mathbf{x}_{H_k}^*)}}. \quad (\square)$$

- Simply (□)

$$\begin{aligned}
 (1 - \alpha_k) \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{1 - \alpha_k^2}{2} \|\mathbf{d}_k\|_{\mathbf{H}_k}^2 + (1 - \alpha_k)g(\mathbf{x}_{\mathbf{H}_k}^*) &< (1 - \alpha_k)g(\mathbf{x}_k) && \text{cancel and group terms} \\
 \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \frac{1 + \alpha_k}{2} \|\mathbf{d}_k\|_{\mathbf{H}_k}^2 + g(\mathbf{x}_{\mathbf{H}_k}^*) &< g(\mathbf{x}_k) && \text{divide by } 1 - \alpha_k \quad (*) \\
 \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + g(\mathbf{x}_{\mathbf{H}_k}^*) - g(\mathbf{x}_k) &< -\frac{1 + \alpha_k}{2} \|\mathbf{d}_k\|_{\mathbf{H}_k}^2 && \text{rearrange}
 \end{aligned}$$

- Consider $F(\mathbf{x}_{k+1}) - F(\mathbf{x}_k)$ and apply the Jansen inequality on g again

$$\begin{aligned}
 F(\mathbf{x}_{k+1}) - F(\mathbf{x}_k) = f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) + g(\mathbf{x}_{k+1}) - g(\mathbf{x}_k) &\leq f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) + (1 - \alpha_k)g(\mathbf{x}_k) + \alpha_k g(\mathbf{x}_{\mathbf{H}_k}^*) - g(\mathbf{x}_k) \\
 &= f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) - \underline{\underline{\alpha_k g(\mathbf{x}_k) + \alpha_k g(\mathbf{x}_{\mathbf{H}_k}^*)}} \quad (\#)
 \end{aligned}$$

- Recall the Taylor expansion for f at $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ is

$$f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) = f(\mathbf{x}_k) + \alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \mathcal{O}((\alpha_k)^2). \quad (\dagger)$$

- Put (†) into (#) gives

$$\begin{aligned}
 F(\mathbf{x}_{k+1}) - F(\mathbf{x}_k) &\leq \alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle + \mathcal{O}((\alpha_k)^2) - \alpha_k g(\mathbf{x}_k) + \alpha_k g(\mathbf{x}_{\mathbf{H}_k}^*) \\
 &= \alpha_k \left(\langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle - g(\mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*) \right) + \mathcal{O}((\alpha_k)^2) \\
 &\stackrel{(*)}{\leq} \alpha_k \left(-\frac{1 + \alpha_k}{2} \|\mathbf{d}_k\|_{\mathbf{H}_k}^2 \right) + \mathcal{O}((\alpha_k)^2) \\
 &= -\frac{\alpha_k}{2} \|\mathbf{d}_k\|_{\mathbf{H}_k}^2 + -\frac{\alpha_k^2}{2} \|\mathbf{d}_k\|_{\mathbf{H}_k}^2 + \mathcal{O}((\alpha_k)^2) \\
 &< 0
 \end{aligned}$$

where the last inequality is true for sufficiently small α_k . □

Stepsize range $\alpha_k \in (0, 1]$

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

$$\text{algo} : \begin{cases} \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), & m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{H_k}^2 \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k) \end{cases}$$

Lemma 2. If $F(\mathbf{x}_{H_k}^*) \leq m_{H_k}(\mathbf{x}_{H_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{H_k}^*)$ holds for $\mathbf{x}_{H_k}^*$, then $F(\mathbf{x}_{k+1}) < F(\mathbf{x}_k)$ for any $\alpha_k \in (0, 1]$.

► **Proof.** First we show that $F(\mathbf{x}_{H_k}^*) < F(\mathbf{x}_k)$

$$\begin{aligned} F(\mathbf{x}_{H_k}^*) &\leq m_{H_k}(\mathbf{x}_{H_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{H_k}^*) && \text{by assumption} \\ &< m_{H_k}(\mathbf{x}_k, \mathbf{x}_k) + g(\mathbf{x}_k) && \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), \quad m_{H_k} \text{ strongly convex} \\ &= F(\mathbf{x}_k) && F(\mathbf{x}_k) = m_{H_k}(\mathbf{x}_k, \mathbf{x}_k) + g(\mathbf{x}_k) \end{aligned}$$

► Apply Jansen inequality on F

$$F(\mathbf{x}_{k+1}) = F\left((1 - \alpha)\mathbf{x}_k + \alpha_k \mathbf{x}_{H_k}^*\right) \leq (1 - \alpha)F(\mathbf{x}_k) + \alpha_k F(\mathbf{x}_{H_k}^*) \stackrel{F(\mathbf{x}_{H_k}^*) < F(\mathbf{x}_k)}{<} F(\mathbf{x}_k).$$

Jansen inequality holds when $\alpha_k \in [0, 1]$, but when $\alpha_k = 0$ we have $\mathbf{x}_{k+1} = \mathbf{x}_k \implies F(\mathbf{x}_{k+1}) = F(\mathbf{x}_k)$, so we have $F(\mathbf{x}_{k+1}) < F(\mathbf{x}_k)$ for $\alpha_k \in (0, 1]$. □

► Problem and algorithm

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}.$$

$$\text{algo} : \begin{cases} m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{\mathbf{H}_k}^2. \\ \mathbf{x}_{\mathbf{H}_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}). \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{\mathbf{H}_k}^* - \mathbf{x}_k). \end{cases}$$

► What we have proved

- Fixed-point property: if $\mathbf{x}_k \in \left\{ \operatorname{argmin} F \right\} \iff \mathbf{H}_k(\mathbf{x}_{\mathbf{H}_k}^* - \mathbf{x}_k) = \mathbf{0}$.
- If $\mathbf{x}_k \notin \left\{ \operatorname{argmin} F \right\}$, then for sufficiently small $\alpha_k > 0$, we have $F(\mathbf{x}_{k+1}) < F(\mathbf{x}_k)$.
- If $F(\mathbf{x}_{\mathbf{H}_k}^*) \leq m_{\mathbf{H}_k}(\mathbf{x}_{\mathbf{H}_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*)$, then $F(\mathbf{x}_{k+1}) < F(\mathbf{x}_k)$ for any $\alpha_k \in (0, 1]$.
- The “weakness” of such result: only descent but not sufficient descent. Recall that for proximal gradient update we have sufficient descent condition.

► Problem and algorithm

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}.$$

$$\text{algo} : \begin{cases} m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{\mathbf{H}_k}^2. \\ \mathbf{x}_{\mathbf{H}_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}). \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{\mathbf{H}_k}^* - \mathbf{x}_k). \end{cases}$$

► Lemma 2. If $F(\mathbf{x}_{\mathbf{H}_k}^*) \leq m_{\mathbf{H}_k}(\mathbf{x}_{\mathbf{H}_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*)$, then $F(\mathbf{x}_{k+1}) < F(\mathbf{x}_k)$ for any $\alpha_k \in (0, 1]$.

► Lemma 2 requires $F(\mathbf{x}_{\mathbf{H}_k}^*) \leq m_{\mathbf{H}_k}(\mathbf{x}_{\mathbf{H}_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*)$, this condition restricts the scope of \mathbf{H}_k : we cannot use any arbitrary \mathbf{H}_k but one that *majorize* f at $\mathbf{x}_{\mathbf{H}_k}^*$.

$$\begin{aligned} F(\mathbf{x}_{\mathbf{H}_k}^*) &\leq m_{\mathbf{H}_k}(\mathbf{x}_{\mathbf{H}_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*) \\ \iff f(\mathbf{x}_{\mathbf{H}_k}^*) + g(\mathbf{x}_{\mathbf{H}_k}^*) &\leq m_{\mathbf{H}_k}(\mathbf{x}_{\mathbf{H}_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*) \end{aligned}$$

$$\iff f(\mathbf{x}_{\mathbf{H}_k}^*) \leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{\mathbf{H}_k}^* - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{\mathbf{H}_k}^* - \mathbf{x}_k\|_{\mathbf{H}_k}^2$$

Table of Contents

- 1 Proximal Newton's method
- 2 Convergence theory of proximal Newton's method: part 1
 - What happens at convergence: Fixed-point property
 - What happens at not convergence: descent property
 - Stepsize $\alpha_k \in (0, 1]$
- 3 Convergence theory of proximal Newton's method: part 2
 - Convergence rate with no line search, part 1
 - Convergence rate with no line search, part 2
 - Convergence rate with no line search, part 3
 - A lemma on decreasing sequence
 - Convergence rate $\frac{1}{k}$
- 4 Summary

Convergence with no line search ($\alpha_k = 1$), part 1

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

$$\text{algo} : \begin{cases} \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), & m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{H_k}^2 \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k) \end{cases}$$

Lemma 3. Suppose $F(\mathbf{x}_{k+1}) \leq m_{H_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) + g(\mathbf{x}_{k+1})$, then

$$F(\mathbf{x}) - F(\mathbf{x}_{k+1}) \geq \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{H_k}^2 + \langle H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*), \mathbf{x} - \mathbf{x}_k \rangle \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

► Before the proof, we use compact notation to shorten the proof: $F_k := F(\mathbf{x}_k)$, $\nabla f_k := \nabla f(\mathbf{x}_k)$, $g_k := g(\mathbf{x}_k)$, \dots

► **Proof.** First, based on the fact that f and g are convex,

$$\begin{aligned} f(\mathbf{x}) &\geq f_k + \langle \nabla f_k, \mathbf{x} - \mathbf{x}_k \rangle && f \text{ is convex} \\ g(\mathbf{x}) &\geq g_{k+1} + \langle \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle && g \text{ is convex} \end{aligned}$$

where $\partial g_{k+1} := \partial g(\mathbf{x}_{k+1})$ is a subgradient of g at \mathbf{x}_{k+1} .

Note: we consider f at \mathbf{x}_k and g at \mathbf{x}_{k+1} , therefore this step is actually tricky.

► Sum the two inequalities

$$F(\mathbf{x}) \geq f_k + \langle \nabla f_k, \mathbf{x} - \mathbf{x}_k \rangle + g_{k+1} + \langle \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle \tag{X}$$

- By hypothesis $F_{k+1} \leq m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) + g_{k+1}$

$$\begin{aligned}
 F(\mathbf{x}) - F_{k+1} &\geq F(\mathbf{x}) - m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) - g_{k+1} \\
 &\stackrel{(X)}{\geq} f_k + \langle \nabla f_k, \mathbf{x} - \mathbf{x}_k \rangle + g_{k+1} + \langle \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) - g_{k+1} \\
 &= f_k + \langle \nabla f_k, \mathbf{x} - \mathbf{x}_k \rangle + \langle \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k).
 \end{aligned}$$

The second line explains why we are considering f at \mathbf{x}_k and g at \mathbf{x}_{k+1} in step (X).

- By definition $m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) := f_k + \langle \nabla f_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2$

$$\begin{aligned}
 F(\mathbf{x}) - F_{k+1} &\geq \langle \nabla f_k, \mathbf{x} - \mathbf{x}_k \rangle + \langle \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \langle \nabla f_k, \mathbf{x}_{k+1} - \mathbf{x}_k \rangle - \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2 \\
 &= \langle \nabla f_k, \mathbf{x} - \mathbf{x}_{k+1} \rangle + \langle \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2 \\
 &= \langle \nabla f_k + \partial g_{k+1}, \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2
 \end{aligned} \tag{O}$$

- We now deal with $\nabla f_k + \partial g_{k+1}$. By definition \mathbf{x}_{k+1} is the minimizer of $m_{\mathbf{H}_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi})$, so

$$\mathbf{0} \in \nabla m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) + \partial g_{k+1} \iff \mathbf{0} \in \nabla f_k + \mathbf{H}_k(\mathbf{x}_{k+1} - \mathbf{x}_k) + \partial g_{k+1}$$

in other words $\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{k+1}) \in \nabla f(\mathbf{x}_k) + \partial g(\mathbf{x}_{k+1})$. Put this into (O) gives

$$F(\mathbf{x}) - F_{k+1} \geq \langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{k+1}), \mathbf{x} - \mathbf{x}_{k+1} \rangle - \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2.$$

The proof completes by considering $\langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{k+1}), \mathbf{x} - \mathbf{x}_k + \mathbf{x}_k - \mathbf{x}_{k+1} \rangle$. □

Convergence with no line search ($\alpha_k = 1$), part 2

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

$$\text{algo} : \begin{cases} \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), & m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{H_k}^2 \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k) \end{cases}$$

Lemma 4. Suppose $F(\mathbf{x}_{k+1}) \leq m_{H_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) + g(\mathbf{x}_{k+1})$, then

$$F_k - F_{k+1} \geq \frac{(F_{k+1} - F^*)^2}{2\bar{\lambda}\delta^2}$$

- ▶ $F^* = F(\mathbf{x}^*)$ for $\mathbf{x}^* \in \operatorname{argmin} F$
- ▶ $\bar{\lambda} = \max_{i=1, \dots, k} \lambda(H_i)$
- ▶ Define the sublevel set $L_{F(\mathbf{x}_0)} := \left\{ \mathbf{x} \in \mathbb{R}^n : F(\mathbf{x}) \leq F(\mathbf{x}_0) \right\}$.
- ▶ $\delta := \text{diameter of } L_{F(\mathbf{x}_0)}$.

- To prove lemma 4 we use lemma 3 we just proved

$$F(\mathbf{x}) - F(\mathbf{x}_{k+1}) \geq \langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{k+1}), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2.$$

with $\mathbf{x} = \mathbf{x}^*$.

$$\begin{aligned} F^* - F_{k+1} &\geq \langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{k+1}), \mathbf{x}^* - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2 \\ &= \frac{1}{2} \langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{k+1}), 2\mathbf{x}^* - 2\mathbf{x}_k + \mathbf{x}_k - \mathbf{x}_{k+1} \rangle \\ &= \frac{1}{2} \langle \mathbf{H}_k(\mathbf{x}^* - \mathbf{x}_k) - \mathbf{H}_k(\mathbf{x}^* - \mathbf{x}_{k+1}), (\mathbf{x}^* - \mathbf{x}_k) + (\mathbf{x}^* - \mathbf{x}_{k+1}) \rangle \\ &= \frac{1}{2} \left(\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k}^2 - \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}^2 \right) \\ &= \frac{1}{2} \left(\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} - \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \right) \left(\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \right) \end{aligned}$$

- For $\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} - \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}$, by triangle inequality

$$\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} \geq \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \implies \|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} - \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \geq -\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}$$

therefore

$$\begin{aligned} F^* - F_{k+1} &\geq \frac{-1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} \left(\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \right) \\ \implies \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} &\geq \frac{-2(F^* - F_{k+1})}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}} = 2 \frac{F_{k+1} - F^*}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}} \end{aligned} \tag{AA}$$

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} \geq 2 \frac{F_{k+1} - F^*}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}} \quad (\text{AA})$$

- Now focus on the norm $\|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}$. Based on norm inequality,

$$\|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \leq \|\mathbf{H}_k\|_2 \cdot \|\mathbf{x}^* - \mathbf{x}_k\|_2 = \lambda(\mathbf{H}_k) \|\mathbf{x}^* - \mathbf{x}_k\|_2 \leq \bar{\lambda} \|\mathbf{x}^* - \mathbf{x}_k\|_2 \quad (\text{U})$$

where the last inequality is due to definition that $\bar{\lambda} = \max_{i=1, \dots, k} \lambda(\mathbf{H}_i)$.

- (U) holds for any k , so

$$\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k} \leq \bar{\lambda} (\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2 + \|\mathbf{x}^* - \mathbf{x}_k\|_2)$$

Therefore

$$\frac{1}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_{\mathbf{H}_k} + \|\mathbf{x}^* - \mathbf{x}_k\|_{\mathbf{H}_k}} \geq \frac{\bar{\lambda}^{-1}}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2 + \|\mathbf{x}^* - \mathbf{x}_k\|_2} \quad (\text{BB})$$

- Put (BB) into (AA)

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} \geq \frac{2}{\bar{\lambda}} \frac{F_{k+1} - F^*}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2 + \|\mathbf{x}^* - \mathbf{x}_k\|_2} \quad (\text{CC})$$

- To proceed, we need to make use of sublevel set argument to bound the terms $\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2$ and $\|\mathbf{x}^* - \mathbf{x}_k\|_2$

1-page review of the diameter of a sublevel set

- ▶ Let $d(\mathbf{x}, \mathbf{y})$ be the distance between two points \mathbf{x}, \mathbf{y} .
- ▶ Let \mathcal{A} be a compact set.
- ▶ Diameter of \mathcal{A} , denoted as $\text{diam } \mathcal{A}$, is defined as the supremum (if it exists) of the distance $d(\mathbf{x}, \mathbf{y})$ across all pairs of points \mathbf{x}, \mathbf{y} in \mathcal{A}

$$\text{diam } \mathcal{A} := \sup \left\{ d(\mathbf{x}, \mathbf{y}) \mid \mathbf{x}, \mathbf{y} \in \mathcal{A} \right\}.$$

- ▶ Now our \mathcal{A} is the sublevel set $L_{F(\mathbf{x}_0)} := \left\{ \mathbf{x} \in \mathbb{R}^n : F(\mathbf{x}) \leq F(\mathbf{x}_0) \right\}$.
- ▶ Since F is a convex function, thus the sublevel set $L_{F(\mathbf{x}_0)}$ is a convex set and it is bounded.
- ▶ In Euclidean space \mathbb{R}^n , distance $d(\mathbf{x}, \mathbf{y})$ is the Euclidean norm $\|\mathbf{x} - \mathbf{y}\|_2$
- ▶ The diameter of $L_{F(\mathbf{x}_0)}$ is thus

$$\text{diam } L_{F(\mathbf{x}_0)} = \sup \left\{ \|\mathbf{x} - \mathbf{y}\|_2 : F(\mathbf{x}) \leq F(\mathbf{x}_0), F(\mathbf{y}) \leq F(\mathbf{x}_0) \right\}$$

i.e., $\text{diam } L_{F(\mathbf{x}_0)}$ is the upper bound of the distance $\|\mathbf{x} - \mathbf{y}\|_2$ for any \mathbf{x}, \mathbf{y} that $F(\mathbf{x}) \leq F(\mathbf{x}_0), F(\mathbf{y}) \leq F(\mathbf{x}_0)$.

- ▶ Such supremum exists for $\text{diam } L_{F(\mathbf{x}_0)}$ because F is convex so $L_{F(\mathbf{x}_0)}$ is a convex and bounded.

Using the sublevel set argument

- ▶ We want to proceed with

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} \geq \frac{2}{\bar{\lambda}} \frac{F_{k+1} - F^*}{\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2 + \|\mathbf{x}^* - \mathbf{x}_k\|_2} \quad (\text{CC})$$

- ▶ We have

- ▶ sublevel set $L_{F(\mathbf{x}_0)} := \{\mathbf{x} \in \mathbb{R}^n : F(\mathbf{x}) \leq F(\mathbf{x}_0)\}$.

- ▶ $\delta := \text{diam } L_{F(\mathbf{x}_0)} = \sup\{\|\mathbf{x} - \mathbf{y}\|_2 : F(\mathbf{x}) \leq F(\mathbf{x}_0), F(\mathbf{y}) \leq F(\mathbf{x}_0)\}$

- ▶ $F^* = F(\mathbf{x}^*) \leq F(\mathbf{x}_0)$ by definition \mathbf{x}^* is a (global) minimizer

- ▶ $F(\mathbf{x}_k) \leq F(\mathbf{x}_0)$ and $F(\mathbf{x}_{k+1}) \leq F(\mathbf{x}_0)$ by the descent property of the update (Lemma 2)

all these imply the three points \mathbf{x}^* , \mathbf{x}_{k+1} , \mathbf{x}_k are inside the sublevel set $L_{F(\mathbf{x}_0)}$ and thus the distances $\|\mathbf{x}^* - \mathbf{x}_k\|_2$ and $\|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2$ are upper bounded by δ

- ▶ Thus we have $\|\mathbf{x}^* - \mathbf{x}_k\|_2 + \|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2 \leq 2\delta \implies \frac{1}{\|\mathbf{x}^* - \mathbf{x}_k\|_2 + \|\mathbf{x}^* - \mathbf{x}_{k+1}\|_2} \geq \frac{1}{2\delta}$. Now (CC) becomes

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k} \geq \frac{F_{k+1} - F^*}{\bar{\lambda}\delta} \quad (\text{DD})$$

- ▶ Lastly, using Lemma 3 again on $\mathbf{x} = \mathbf{x}_k$ gives

$$F(\mathbf{x}_k) - F(\mathbf{x}_{k+1}) \geq \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{\mathbf{H}_k}^2. \quad (\text{EE})$$

Combine (DD) and (EE) finishes the proof. □

Convergence with no line search ($\alpha_k = 1$), part 3

$$(\mathcal{P}) : \min_{\mathbf{x}} \left\{ F(\mathbf{x}) := f(\mathbf{x}) + g(\mathbf{x}) \right\}$$

$$\text{algo} : \begin{cases} \mathbf{x}_{H_k}^* = \underset{\boldsymbol{\xi}}{\operatorname{argmin}} m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) + g(\boldsymbol{\xi}), & m_{H_k}(\boldsymbol{\xi}, \mathbf{x}_k) := f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \boldsymbol{\xi} - \mathbf{x}_k \rangle + \frac{1}{2} \|\boldsymbol{\xi} - \mathbf{x}_k\|_{H_k}^2 \\ \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k (\mathbf{x}_{H_k}^* - \mathbf{x}_k) \end{cases}$$

Lemma 5. For $\alpha_k \equiv 1$, $H_k \succ \mathbf{0}$, suppose $F(\mathbf{x}_{k+1}) \leq m_{H_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) + g(\mathbf{x}_{k+1})$, then for all $\mathbf{x} \in \mathbb{R}^n$,

$$\begin{aligned} F(\mathbf{x}) &\geq F(\mathbf{x}_{k+1}) + \langle H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2} \|H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*)\|_{H_k^{-1}}^2 \\ &\geq F(\mathbf{x}_{k+1}) + \langle H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\lambda_1(H_k)} \|H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*)\|_2^2 \end{aligned}$$

► **Proof.** First, starting from Lemma 3

$$\begin{aligned} F(\mathbf{x}) &\geq F(\mathbf{x}_{k+1}) + \langle H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_{H_k}^2 \\ &= F(\mathbf{x}_{k+1}) + \langle H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2} \|H_k(\mathbf{x}_{k+1} - \mathbf{x}_k)\|_{H_k^{-1}}^2 \\ &\geq F(\mathbf{x}_{k+1}) + \langle H_k(\mathbf{x}_k - \mathbf{x}_{H_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\lambda_1(H_k)} \|H_k(\mathbf{x}_{k+1} - \mathbf{x}_k)\|_2^2 \end{aligned}$$

► For the last inequality sign, see explanation next page.

Inequality on positive definite matrices

- ▶ Let $\lambda_1(\cdot)$ be the maximum eigenvalue of a matrix, then from the fact that $\lambda_1 \geq \lambda_{i>1}$

$$\lambda_1 \geq \lambda_2, \lambda_1 \geq \lambda_3, \dots, \lambda_1 \geq \lambda_n \underbrace{\geq 0}_{\mathbf{H}_k \succ \mathbf{0}}, \implies \frac{1}{\lambda_n} \geq \frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_2} \geq \frac{1}{\lambda_1}.$$

- ▶ Therefore, for matrices, we have the following inequality

$$\underbrace{\mathbf{H}_k}_{\lambda_1, \lambda_2, \lambda_3, \dots} \preceq \underbrace{\lambda_1(\mathbf{H}_k)\mathbf{I}}_{\lambda_1, \lambda_1, \lambda_1, \dots} \implies \underbrace{(\mathbf{H}_k)^{-1}}_{\frac{1}{\lambda_n}, \frac{1}{\lambda_{n-1}}, \dots, \frac{1}{\lambda_1}} \succeq \underbrace{\frac{1}{\lambda_1(\mathbf{H}_k)}\mathbf{I}}_{\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_1}}$$

- ▶ Therefore, for any $\boldsymbol{\xi} \in \mathbb{R}^n$,

$$\|\boldsymbol{\xi}\|_{(\mathbf{H}_k)^{-1}}^2 \geq \|\boldsymbol{\xi}\|_{\frac{1}{\lambda_1(\mathbf{H}_k)}\mathbf{I}}^2 = \frac{1}{\lambda_1(\mathbf{H}_k)} \|\boldsymbol{\xi}\|_2^2.$$

- ▶ Put $\boldsymbol{\xi} = \mathbf{H}_k(\mathbf{x}_{k+1} - \mathbf{x}_k)$ gives

$$\|\mathbf{H}_k(\mathbf{x}_{k+1} - \mathbf{x}_k)\|_{\mathbf{H}_k^{-1}}^2 \geq \frac{1}{\lambda_1(\mathbf{H}_k)} \|\mathbf{H}_k(\mathbf{x}_{k+1} - \mathbf{x}_k)\|_2^2.$$

Comment on Lemma 5

Lemma 5. For $\alpha_k \equiv 1$, $\mathbf{H}_k \succ \mathbf{0}$, suppose $F(\mathbf{x}_{k+1}) \leq m_{\mathbf{H}_k}(\mathbf{x}_{k+1}, \mathbf{x}_k) + g(\mathbf{x}_{k+1})$, then for all $\mathbf{x} \in \mathbb{R}^n$,

$$\begin{aligned} F(\mathbf{x}) &\geq F(\mathbf{x}_{k+1}) + \langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2} \|\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*)\|_{\mathbf{H}_k^{-1}}^2 \\ &\geq F(\mathbf{x}_{k+1}) + \langle \mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2\lambda_1(\mathbf{H}_k)} \|\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*)\|_2^2 \end{aligned}$$

- Put $\mathbf{x} = \mathbf{x}_k$, the inner product term is gone and we have

$$\begin{aligned} F_k &\geq F_{k+1} + \frac{1}{2} \|\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*)\|_{\mathbf{H}_k^{-1}}^2 \\ &\geq F_{k+1} + \frac{1}{2\lambda_1(\mathbf{H}_k)} \|\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*)\|_2^2 \end{aligned}$$

- Re-arrange

$$F_{k+1} \leq F_k - \frac{1}{2\lambda_1(\mathbf{H}_k)} \|\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*)\|_2^2$$

Recall the discussion on Lemma 0, the term $\mathbf{H}_k(\mathbf{x}_k - \mathbf{x}_{\mathbf{H}_k}^*)$ is related to the subgradient of F

- In fact, such express resembles the sufficient descent condition of gradient descent ([see here](#))

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2^2 = f(\mathbf{x}_k) - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2.$$

A lemma on decreasing sequence

Lemma 6. If $\{\omega_k\}_{k \in \mathbb{N}}$

▶ is a decreasing sequence

▶ $\omega_k - \omega_{k+1} \geq \frac{(\omega_{k+1} - \omega^*)^2}{\mu}$ for all k

▶ $\omega_1 - \omega^* \leq 4\mu$

then for all k

$$\omega_k - \omega^* \leq \frac{4\mu}{k}.$$

▶ **Proof.** By induction. Case $k = 1$ is true by assumption.

▶ Case $k + 1$

$$\omega_{k+1} - \omega^* = \omega_k - \omega^* + \omega_{k+1} - \omega_k \leq \omega_k - \omega^* - \frac{(\omega_{k+1} - \omega^*)^2}{\mu} \leq \frac{4\mu}{k} - \frac{(\omega_{k+1} - \omega^*)^2}{\mu}.$$

So we have

$$\omega_{k+1} - \omega^* \leq \frac{4\mu}{k} - \frac{(\omega_{k+1} - \omega^*)^2}{\mu} \iff \frac{(\omega_{k+1} - \omega^*)^2}{\mu} + \omega_{k+1} - \omega^* - \frac{4\mu}{k} \leq 0.$$

$$\frac{(\omega_{k+1} - \omega^*)^2}{\mu} + \omega_{k+1} - \omega^* - \frac{4\mu}{k} \leq 0$$

- Let $\nu = \omega_{k+1} - \omega^*$ and $p_k = \frac{4\mu}{k}$

$$\frac{\nu^2}{\mu} + \nu - p_k \leq 0$$

which has a nonnegative solution given by

$$\nu = \frac{-\mu + \sqrt{\mu^2 + 4p_k\mu}}{2} = \frac{2p_k}{1 + \sqrt{1 + \frac{4p_k}{\mu}}}$$

- $f(x) = \frac{1}{1 + \sqrt{1+x}}$ is convex and on any interval $[0, \alpha]$ it is bounded above by its secant interpolation.
- To finish the proof, consider two separate cases, when $k = 1$ to show that the lemma holds for $k + 1 = 2$, and when $k \geq 2$ to show that the lemma holds for $k + 1 \geq 3$.
- For the second case we need to make use of the fact that $\frac{1}{k} - \frac{1}{k^2} \leq \frac{1}{k+1}$.

Convergence rate $\frac{1}{k}$

Theorem 1. Suppose

► $\mathbf{x}_{k+1} = \mathbf{x}_{\mathbf{H}_k}^*$ ($\alpha_k \equiv 1$)

► $\delta := \text{diam}L_F(\mathbf{x}_0) = \sup \left\{ \|\mathbf{x} - \mathbf{y}\|_2 : F(\mathbf{x}) \leq F(\mathbf{x}_0), F(\mathbf{y}) \leq F(\mathbf{x}_0) \right\}$ (sublevel set for Lemma 4)

► $\bar{\lambda} = \max_{i=1, \dots, k} \lambda(\mathbf{H}_i)$ (sublevel set for Lemma 4)

► $F(\mathbf{x}_{\mathbf{H}_k}^*) \leq m_{\mathbf{H}_k}(\mathbf{x}_{\mathbf{H}_k}^*, \mathbf{x}_k) + g(\mathbf{x}_{\mathbf{H}_k}^*)$ (condition on \mathbf{H}_k for Lemma 2)

then for all k

$$F(\mathbf{x}_k) - F^* \leq \frac{\max \left\{ 8\bar{\lambda}^2 \delta^2, F(\mathbf{x}_0) - F^* \right\}}{k}.$$

► **Proof.** Combine Lemma 4 and Lemma 6. □

Table of Contents

- 1 Proximal Newton's method
- 2 Convergence theory of proximal Newton's method: part 1
 - What happens at convergence: Fixed-point property
 - What happens at not convergence: descent property
 - Stepsize $\alpha_k \in (0, 1]$
- 3 Convergence theory of proximal Newton's method: part 2
 - Convergence rate with no line search, part 1
 - Convergence rate with no line search, part 2
 - Convergence rate with no line search, part 3
 - A lemma on decreasing sequence
 - Convergence rate $\frac{1}{k}$
- 4 Summary

Last page - summary

Discussed

- ▶ Proximal Newton's method.
- ▶ Convergence of Proximal Newton's method.

Not discussed

- ▶ Proximal Newton's method with Nesterov's acceleration
- ▶ Problem-specific implementation of \mathbf{H}_k (Identity minus rank one in the paper)
- ▶ Other implementation of \mathbf{H}_k (Identity minus multiple rank matrix)

Reference

Sahar Karimi and Stephen Vavasis, "IMRO: A Proximal Quasi-Newton Method for Solving ℓ_1 -Regularized Least Squares Problems", SIAM J. Optimization, 2017.

End of document