

# MM Optimization Algorithms

**Chathuranga Weeraddana**

March 2022

## LECTURE 4: KEY INEQUALITIES FOR MM (PART III)

# MAJORIZATION AND PARTIAL OPTIMIZATION

# Partial Minimization

- ▶ variety of functions can be represented as partial minima <sup>1</sup>

$$f(x) = \min_{y \in \mathcal{Y}} g(x, y) \quad (1)$$

- ▶ such a function  $f$  can readily be majorized at  $x^{(n)} \in \mathbb{R}^p$ , i.e.,

$$\begin{aligned} f(x) &= \min_{y \in \mathcal{Y}} g(x, y) & (2) \\ &\leq g(x, y^{(n)}) = h(x | x^{(n)}) \end{aligned}$$

where  $y^{(n)} = \arg \min_{y \in \mathcal{Y}} g(x^{(n)}, y)$

- ▶  $h(x | x^{(n)})$  is the **restriction** of  $g$  to set  $\{(x, y^{(n)}) \mid x \in \mathbb{R}^p\}$

---

<sup>1</sup>It is assumed that the minimum over  $y \in \mathcal{Y}$  is attained for each  $x$ .

## EXAMPLES

# Block Descent

- ▶ suppose you are given the problem

$$\begin{aligned} & \text{minimize} && g(x, y) \\ & \text{subject to} && x \in \mathcal{X} \\ & && y \in \mathcal{Y} \end{aligned}$$

- ▶ the problem is equivalent to

$$\begin{aligned} & \text{minimize} && f(x) = \min_{y \in \mathcal{Y}} g(x, y) \\ & \text{subject to} && x \in \mathcal{X} \end{aligned}$$

- ▶ MM principle:  $f(x) \leq h(x|x^{(n)}) = g(x, y^{(n)})$
- ▶ can the constraint be of the form  $(x, y) \in \mathcal{Z}$ ?

## A SIMPLE PROBLEM

- ▶ consider the following problem

$$\begin{array}{ll} \text{minimize} & \|Ax - y\| \\ \text{subject to} & y \in \mathcal{Y} \end{array}$$

where the decision variables are  $x, y$

- ▶ the iterative algorithm reduces to:

$$x^{(n+1)} = (A^T A)^{-1} A^T P_{\mathcal{Y}}(Ax^{(n)})$$

where  $P_{\mathcal{Y}}(\cdot)$  is the projection onto  $\mathcal{Y}$

## DISTANCE BETWEEN TWO SETS

- ▶  $\mathcal{X}, \mathcal{Y} \rightarrow$  two disjoint closed sets
- ▶ compute  $\text{dist}(\mathcal{X}, \mathcal{Y})$  the optimal value of

$$\begin{aligned} & \text{minimize} && \|x - y\| \\ & \text{subject to} && x \in \mathcal{X}, y \in \mathcal{Y} \end{aligned}$$

where the decision variables are  $x, y$

- ▶ the iterative algorithm reduces to:

$$x^{(n+1)} = P_{\mathcal{X}}(P_{\mathcal{Y}}(x^{(n)}))$$

- ▶ optimality if  $\mathcal{X}$  (or both sets) is nonconvex?

# Proximal Minimization Algorithm

- ▶ suppose you are given the problem

$$\underset{x}{\text{minimize}} \quad f(x)$$

- ▶ trivial to see:  $f(x) = \min_y [f(x) + (1/2\mu)\|x - y\|^2]$ ,  $\mu > 0$

- ▶ thus a majorization function of  $f$  is given by  $h(x|x^{(n)})$  where

$$h(x|x^{(n)}) = f(x) + (1/2\mu)\|x - x^{(n)}\|^2$$

- ▶ output of  $h(\cdot |x^{(n)})$  minimization compromises between
  - ▶ minimizing  $f$  and being near to  $x^{(n)}$  (controlled by  $\mu$ )

- ▶ the algorithm if MM principle is applied

$$x^{(n+1)} = \text{prox}_{\mu f}(x^{(n)})$$

where  $\text{prox}_{\mu f}$  is called the proximal operator of  $\mu f$ ,

$$\text{prox}_{\mu f}(v) = \arg \min_x f(x) + (1/2\mu)\|x - v\|^2$$

- ▶ the resulting algorithm is a **proximal minimization algorithm**
  - ▶ also called: proximal iteration or the proximal point algorithm <sup>2</sup>

---

<sup>2</sup>See § 4.1 of *Proximal Algorithms* by N. Parikh and S. Boyd, now Foundations and Trends in Optimization 2013.

- ▶ why compute a sequence of proximal operators?
  - ▶ subproblems usually admits easy closed-form solutions
  - ▶ can be solved sufficiently quickly
  - ▶ minimizing of  $(f + \text{quadratic})$  is *easier* than minimizing  $f$ 
    - ▶ handle ill-conditioned situations  $\rightarrow$  higher reliability
    - ▶ fewer iterations or faster convergence
  - ▶ amenable to distributed optimization
- ▶ an application: iterative refinement  $\rightarrow$  a homework exercise

# Schur Compliment Majorization

- ▶ more specifically
  - ▶ we are given a vector  $x$
  - ▶ we have a parameter matrix  $D$  given by

$$D = \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}$$

- ▶ now it is useful to **bound**

$$1) \quad x^\top A^{-1} x \quad 2) \quad \ln |A|$$

# Bounding $x^\top A^{-1}x$

► we have <sup>3</sup>

$$\begin{aligned} f(x) &= \inf_z g(x, z) \\ &= \inf_x \begin{bmatrix} x \\ z \end{bmatrix}^\top \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}^{-1} \begin{bmatrix} x \\ z \end{bmatrix} \\ &= x^\top A^{-1}x \\ &= \begin{bmatrix} x \\ B^\top A^{-1}x \end{bmatrix}^\top \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}^{-1} \begin{bmatrix} x \\ B^\top A^{-1}x \end{bmatrix} \end{aligned}$$

---

<sup>3</sup>See § A.5.5, *Convex Optimization* by S. Boyd and L. Vandenberghe, 2004.

- ▶ from the last two equations we get for all  $A, B$  and  $C$

$$\begin{aligned} x^T A^{-1} x &= \begin{bmatrix} x \\ B^T A^{-1} x \end{bmatrix}^T \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}^{-1} \begin{bmatrix} x \\ B^T A^{-1} x \end{bmatrix} \\ &\leq \begin{bmatrix} x \\ z^{(n)} \end{bmatrix}^T \begin{bmatrix} A & B \\ B^T & C \end{bmatrix}^{-1} \begin{bmatrix} x \\ z^{(n)} \end{bmatrix} \end{aligned} \quad (3)$$

with  $z^{(n)} = B^{(n)T} (A^{(n)})^{-1} x$

- ▶ note that the inequality holds with equality when

$$A = A^{(n)}, \quad B = B^{(n)}, \quad C = C^{(n)}$$

# Bounding $\ln |A|$

- ▶ the Schur complement of  $A$  in the matrix  $D$  is given by

$$C - B^T A^{-1} B$$

- ▶ for clarity let us define  $G$  as

$$G = (C - B^T A^{-1} B)^{-1}$$

- ▶ we have the following determinant identity

$$|D| = |A| \times |C - B^T A^{-1} B| = |A|/|G|$$

- ▶ now we can bound  $\ln |A|$  as follows:

$$\begin{aligned}\ln |A| &= \ln |D| + \ln |G| \\ &\leq \ln |D| + \ln |G^{(n)}| + \text{Tr}[(G^{(n)})^{-1}(G - G^{(n)})] \\ &= \ln |D| + \ln |G^{(n)}| - \text{Tr}(I) + \text{Tr}[(G^{(n)})^{-1}G] \\ &= \ln |D| + \ln |G^{(n)}| - \text{Tr}(I) + \text{Tr}[D^{-1}F^{(n)}] \quad (4)\end{aligned}$$

- ▶ the last equality follows from that  $F^{(n)} \triangleq \begin{bmatrix} 0 & 0 \\ 0 & (G^{(n)})^{-1} \end{bmatrix}$  and

$$D^{-1} = \begin{bmatrix} A^{-1} + A^{-1}BGB^{\top}A^{-1} & -A^{-1}BG \\ -GB^{\top}A^{-1} & G \end{bmatrix}$$

- ▶ note that the inequality holds with equality when

$$A = A^{(n)}, \quad B = B^{(n)}, \quad C = C^{(n)}$$

# Fenchel Conjugate

- ▶ Fenchel conjugate <sup>4</sup> of a function  $f$

$$f^*(x) = \sup_y \{x^\top y - f(y)\} \quad (5)$$

- ▶ in general we have for

$$f^*(x) = \sup_y \{x^\top y - f(y)\} \quad (6)$$

$$\geq \bar{y}^\top x - f(\bar{y}) \quad (7)$$

$$= g(x|x^{(n)}) \quad (8)$$

where  $\bar{y} \in \partial f^*(x^{(n)}) = \arg \max_y \{x^{(n)\top} y - f(y)\}$

---

<sup>4</sup>For more details see pages 15-17.

## APPENDICES

# Legendre-Fenchel Transform

- ▶ for any function  $f : \mathbb{R}^N \rightarrow \bar{\mathbb{R}}$  define <sup>5</sup>

$$f^*(x) = \sup_y \{x^\top y - f(y)\} \quad (9)$$

- ▶  $f^*$  is called the conjugate to  $f$
- ▶ biconjugate to  $f$  is given by  $f^{**} = (f^*)^*$ , where

$$f^{**}(y) = \sup_x \{y^\top x - f^*(x)\} \quad (10)$$

---

<sup>5</sup> $\bar{\mathbb{R}} = \mathbb{R} \cup \{\infty\}$ .

- ▶ the mapping  $f \rightarrow f^*$  from  $\text{fcns}(\mathbb{R}^N)$ <sup>6</sup> into  $\text{fcns}(\mathbb{R}^N)$ 
  - ▶ is called the **Legendre-Fenchel Transform**<sup>7</sup>
- ▶ if  $f$  is proper, lsc, and convex, so is  $f^*$  and  $f^{**} = f$

---

<sup>6</sup> $\text{fcns}(\mathbb{R}^N)$ : the collection of all extended-real-valued functions on  $\mathbb{R}^N$

<sup>7</sup>See pp. 473-476 *Variational Analysis* by R. T. Rockafellar and R. J-B Wets, 3rd printing 2009.

- ▶ for any proper, lsc, convex function  $f$

$$\bar{x} \in \partial f(\bar{y}) \iff \bar{y} \in \partial f^*(\bar{x}) \iff f(\bar{y}) + f^*(\bar{x}) = \bar{x}^\top \bar{y}$$

where

$$\partial f(\bar{y}) = \arg \max_x \{\bar{y}^\top x - f^*(x)\} \quad \partial f^*(\bar{x}) = \arg \max_y \{\bar{x}^\top y - f(y)\}$$

- ▶ in general,

$$f(y) + f^*(x) \geq x^\top y \quad \text{for all } x, y$$

- ▶ see Proposition 11.3, R. T. Rockafellar