

# MM Optimization Algorithms

**Chathuranga Weeraddana**

April 2022

## LECTURE 3: KEY INEQUALITIES FOR MM (PART II)

## QUADRATIC UPPER BOUND PRINCIPLE

# Quadratic Upper Bound Principle

- ▶ key mechanisms
  - ▶ majorization via gradient Lipschitz continuity
  - ▶ majorization via bounded Hessian
  - ▶ minorization via strong convexity

# Gradient Lipschitz Continuity

▶ suppose

▶  $f$  is differentiable

▶ gradient Lipschitz continuous with constant  $L$ , i.e.,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \text{ for all } x, y \quad (1)$$

► from (44) <sup>1</sup>

$$f(x) = f(y) + \int_0^1 \nabla f(y + t(x - y))^\top (x - y) dt \quad (2)$$

$$\begin{aligned} &= f(y) + \nabla f(y)^\top (x - y) + \\ &\quad \int_0^1 [\nabla f(y + t(x - y)) - \nabla f(y)]^\top (x - y) dt \quad (3) \end{aligned}$$

$$\begin{aligned} &\leq f(y) + \nabla f(y)^\top (x - y) + \\ &\quad \int_0^1 \|\nabla f(y + t(x - y)) - \nabla f(y)\| \|x - y\| dt \quad (4) \end{aligned}$$

$$\leq f(y) + \nabla f(y)^\top (x - y) + L\|x - y\|^2 \int_0^1 t dt \quad (5)$$

$$= f(y) + \nabla f(y)^\top (x - y) + (L/2)\|x - y\|^2 = g(x|y) \quad (6)$$

---

<sup>1</sup>See page 34 and substitute  $a = y$  and  $b = x - y$ .

# Bounded Hessian

- ▶ suppose
  - ▶  $f$  is twice differentiable
  - ▶  $f$  has bounded Hessians, i.e.,

$$\exists B \succ 0 \text{ s.t. } B - \nabla^2 f(x) \succeq 0 \text{ for all } x \quad (7)$$

► from (46) <sup>2</sup>

$$f(x) = f(y) + \nabla f(y)^\top (x - y) + \quad (8)$$

$$\int_0^1 \int_0^t (x - y)^\top \nabla^2 f(y + \tau(x - y))(x - y) d\tau dt \quad (9)$$

$$\leq f(y) + \nabla f(y)^\top (x - y) + \int_0^1 \int_0^t (x - y)^\top B(x - y) d\tau dt \quad (10)$$

$$= f(y) + \nabla f(y)^\top (x - y) + (1/2)(x - y)^\top B(x - y) \quad (11)$$

$$= g(x|y) \quad (12)$$

---

<sup>2</sup>See page 35 and substitute  $a = y$  and  $b = x - y$ .



# Strong Convexity

▶ suppose

▶  $f$  is twice differentiable and strongly convex

▶ thus,

$$\exists m > 0 \text{ s.t. } \nabla^2 f(x) - mI \succeq 0 \text{ for all } x \quad (13)$$

► from (46)

$$f(x) = f(y) + \nabla f(y)^\top (x - y) + \quad (14)$$

$$\int_0^1 \int_0^t (x - y)^\top \nabla^2 f(y + \tau(x - y)) (x - y) d\tau dt \quad (15)$$

$$\geq f(y) + \nabla f(y)^\top (x - y) +$$

$$m \int_0^1 \int_0^t (x - y)^\top (x - y) d\tau dt \quad (16)$$

$$= f(y) + \nabla f(y)^\top (x - y) + (m/2) \|x - y\|^2 \quad (17)$$

$$= g(x|y) \quad (18)$$

## EXAMPLES

# Landweber's Method

- ▶ consider a **positive definite** matrix  $A \in \mathbb{S}^n$
- ▶ we have to find the **solution**  $x^*$  for  $Ax = -b$
- ▶ the task requires order of  $n^3$  flops
  - ▶  $A^{-1}$  is to be computed
  - ▶ if  $n$  is very large the task may be computationally challenging
- ▶ but we may rely on MM principle  $\rightarrow$  avoid matrix inversion

- ▶ note that

$$x^* = \arg \min_x f(x) = (1/2) x^T A x + b^T x \quad (19)$$

- ▶ moreover, we have <sup>3</sup>

$$\|\nabla f(x) - \nabla f(y)\| \leq \lambda_{\max}(A) \|x - y\| \quad (20)$$

$$\leq \|A\|_* \|x - y\| \quad (21)$$

where  $\|\cdot\|_*$  is any matrix norm

---

<sup>3</sup>See p. 497, *Matrix Analysis and Applied Linear Algebra* by C. D. Meyer, 2000.

- ▶ thus, with any  $\rho > \|A\|_*$ , we have

$$f(x) \leq f(x^{(n)}) + \nabla f(x^{(n)})^\top (x - x^{(n)}) + \frac{\rho}{2} \|x - x^{(n)}\|^2 \quad (22)$$

$$= g(x|x^{(n)}) \quad (23)$$

- ▶  $x$  can be updated as

$$x^{(n+1)} = x^{(n)} - (1/\rho)\nabla f(x^{(n)}) \quad (24)$$

$$= x^{(n)} - (1/\rho)(Ax^{(n)} + b) \quad (25)$$

- ▶ what if  $A$  is positive semidefinite? indefinite? not symmetric?

# Jacobi Iterations

- ▶ **symmetric** matrix  $D$  and  $N$  with positive definite  $D + N$
- ▶ we have to find the **solution**  $x^*$  for  $(D + N)x = -b$
- ▶ note that

$$x^* = \arg \min_x f(x) = (1/2) x^T (D + N)x + b^T x \quad (26)$$

- ▶ e.g., computing the Newton step for

$$f(x) = \sum_{i=1}^N f_i(x_i) + r(Ax - b) \quad (27)$$

$x \in \mathbb{R}^M$ ,  $A \in \mathbb{R}^{p \times M}$ ,  $b \in \mathbb{R}^p$ ,  $r$  is a regularization function

- ▶ Newton's step  $\Delta x_{\text{nt}}$  for  $f$  at  $x$  is given by

$$\nabla^2 f(x) \Delta x_{\text{nt}} = -\nabla f(x) \quad (28)$$

- ▶ the Hessian  $\nabla^2 f(x)$  is of the form

$$\nabla^2 f(x) = D + N \quad (29)$$

- ▶ identify  $D$  and  $N \rightarrow$  structured matrix  $+N$



- ▶ e.g., continues ..
  - ▶ if  $N$  is low rank  $\rightarrow$  things can be handled efficiently
  - ▶ if  $N$  is not 'sufficiently' low rank  $\rightarrow$  apply MM principle
- ▶ find  $L > \|N\|_*$  and majorize  $(1/2) x^T N x$
- ▶ thus  $f$  is majorized

- more specifically, we have

$$f(x) = (1/2) x^T (D + N)x + b^T x \quad (30)$$

$$= (1/2) x^T D x + b^T x + (1/2) x^T N x \quad (31)$$

$$\leq (1/2) x^T D x + b^T x + g(x|x^{(n)}) \quad (32)$$

since  $(1/2) x^T N x \leq g(x|x^{(n)})$ , where

$$g(x|x^{(n)}) = (\rho/2) \|x - x^{(n)}\|^2 + x^{(n)T} N x - (1/2) x^{(n)T} N x^T$$

- ▶  $x$  can be updated as

$$x^{(n+1)} = (D + \rho I)^{-1} \left( \rho x^{(n)} - Nx^{(n)} - b \right) \quad (33)$$

- ▶  $(D + \rho I)$  is efficiently invertible  $\rightarrow$  has a rich structure
  - ▶ e.g., block diagonal

# NONNEGATIVE QUADRATIC PROGRAMMING

# The Related Problem

- ▶ consider the following problem

$$\begin{aligned} & \text{minimize} && (1/2)x^T R x + s^T x \\ & \text{subject to} && Cx \preceq d \end{aligned}$$

- ▶  $x \in \mathbb{R}^N$ , **positive definite**  $R$ , suppose it has some structure
- ▶  $C \in \mathbb{R}^{p \times N}$ ,  $d \in \mathbb{R}^p$
- ▶ **Lagrangian**  $L$  is given by

$$L(x, \mu) = (1/2)x^T R x + s^T x + \mu^T (Cx - d) \quad (34)$$

$$\mu \in \mathbb{R}^p$$

- ▶ **minimizer**  $x(\mu)$  of the Lagrangian is given by

$$x(\mu) = -R^{-1}(s + C^T \mu) \quad (35)$$

- ▶ **dual function**  $h$  is given by

$$h(\mu) = -(1/2)\mu^T P \mu + q^T \mu + r \quad (36)$$

where

$$P = CR^{-1}C^T, \quad q = -(d + CR^{-1}s), \quad r = -\frac{1}{2}s^T R^{-1}s$$

- ▶  $P$  **lacks** any structure even if  $R$  does

## DUAL PROBLEM

- ▶ dual problem is given by

$$\begin{aligned} & \text{maximize} && h(\mu) = -(1/2)\mu^\top P\mu + q^\top\mu + r \\ & \text{subject to} && \mu \succeq 0 \end{aligned}$$

- ▶ how to solve the dual problem?
  - ▶ interior-point methods <sup>4</sup> → recall  $P$  lacks any structure
    - ▶ not easily implemented for large  $P$
  - ▶ coordinate descent
  - ▶ MM principle

---

<sup>4</sup>See § 11.3.1 of *Convex Optimization* by S. Boyd and L. Vandenberghe, 2004.

## COORDINATE DESCENT TO SOLVE THE DUAL

- ▶ restrict  $h$  to a line  $\mathcal{L}_i = \{\mu^{(n)} + te_i \mid t \in \mathbb{R}\}$
- ▶ minimize  $h$  over the restriction  $\mathcal{L}_i$  <sup>5</sup>

$$\underset{t \in \mathbb{R}}{\text{minimize}} \quad h_i(t) = h(\mu^{(n)} + te_i)$$

- ▶ we can unfold  $h(\mu^{(n)} + te_i)$  in a straightforward manner, i.e., <sup>6</sup>

$$\begin{aligned} h_i(t) &= -(1/2)(\mu^{(n)} + te_i)^\top P(\mu^{(n)} + te_i) + q^\top(\mu^{(n)} + te_i) \\ &= -(P_{ii}/2) t^2 + \left( q_i - \sum_{k=1}^p P_{ik} \mu_i^{(n)} \right) t + \text{irrelevant const.} \end{aligned}$$

---

<sup>5</sup>Let us first ignore the constrain  $\mu \succeq 0$  and assimilate it later.

<sup>6</sup>The constant  $r$  is dropped since it is irrelevant.



- ▶ compute the derivative  $h'_i$  of  $h$  to determine  $t^*$ , i.e.,

$$-P_{ii}t + \left( q_i - \sum_{k=1}^p P_{ik}\mu_i^{(n)} \right) = 0 \implies t^* = \frac{q_i - \sum_{k=1}^p P_{ik}\mu_i^{(n)}}{p_{ii}}$$

- ▶ so the  $i$ th coordinate of current  $\mu^{(n)}$  is updated as

$$\begin{aligned} \mu_i^{(n)} &:= \mu_i^{(n)} + t^* \\ &= \mu_i^{(n)} + \frac{1}{p_{ii}} \left( q_i - \sum_{k=1}^p P_{ik}\mu_i^{(n)} \right) \end{aligned}$$

- ▶  $\mu \succeq 0$  can be assimilated as?

$$\mu_i^{(n)} := \max \left\{ 0, \mu_i^{(n)} + \frac{1}{p_{ii}} \left( q_i - \sum_{k=1}^p P_{ik}\mu_i^{(n)} \right) \right\}$$

- ▶ iterate from  $i = 1$  to  $i = p$  and cycles back to  $i = 1$

---

**Algorithm 1** Coordinate Decent

---

**Input:**  $\mu^{(0)} \succeq 0, n = 0$

```
1: while a stopping criterion true do
2:   for  $i \leftarrow 1$  to  $p$  do
3:      $\mu_i^{(n)} \leftarrow \max \left\{ 0, \mu_i^{(n)} + \frac{1}{p_{ii}} \left( q_i - \sum_{k=1}^p P_{ik} \mu_k^{(n)} \right) \right\}$ 
4:   end for
5:    $\mu^{(n+1)} = \mu^{(n)}$  and  $n \leftarrow n + 1$ 
6: end while
7: return  $\mu^{(n+1)}$ 
```

---

- ▶ then the solution is given by (35) with  $\mu = \mu^{(n+1)}$

## MM PRINCIPLE TO SOLVE THE DUAL

- ▶ recall the objective function <sup>7</sup>

$$\begin{aligned}
 h(\mu) &= -\frac{1}{2} \sum_{i=1}^N P_{ii} \mu_i^2 + \sum_{i=1}^N q_i \mu_i \\
 &\quad - \frac{1}{2} \sum_{\{i,j|i \neq j, P_{ij} \geq 0\}} P_{ij} \mu_i \mu_j - \frac{1}{2} \sum_{\{i,j|i \neq j, P_{ij} < 0\}} P_{ij} \mu_i \mu_j \\
 &= -\frac{1}{2} \sum_{i=1}^N P_{ii} \mu_i^2 + \sum_{i=1}^N q_i \mu_i \\
 &\quad - \frac{1}{2} \sum_{\{i,j|i \neq j, P_{ij} \geq 0\}} P_{ij} \mu_i \mu_j + \frac{1}{2} \sum_{\{i,j|i \neq j, P_{ij} < 0\}} |P_{ij}| \mu_i \mu_j \\
 &\geq -\frac{1}{2} \sum_{i=1}^N P_{ii} \mu_i^2 + \sum_{i=1}^N q_i \mu_i \\
 &\quad - \frac{1}{2} \sum_{\{i,j|i \neq j, P_{ij} \geq 0\}} P_{ij} \left[ \frac{\mu_i^{(n)}}{2\mu_j^{(n)}} \mu_j^2 + \frac{\mu_j^{(n)}}{2\mu_i^{(n)}} \mu_i^2 \right] \\
 &\quad + \frac{1}{2} \sum_{\{i,j|i \neq j, P_{ij} < 0\}} |P_{ij}| \mu_i^{(n)} \mu_j^{(n)} \left[ 1 + \ln \left( \frac{\mu_i}{\mu_i^{(n)}} \right) + \ln \left( \frac{\mu_j}{\mu_j^{(n)}} \right) \right] \\
 &= g(\mu | \mu^{(n)})
 \end{aligned}$$

---

<sup>7</sup>The constant  $r$  is dropped since it is irrelevant.

- here the last inequality follows from

$$\mu_i \mu_j \leq \frac{\mu_i^{(n)}}{2\mu_j^{(n)}} \mu_j^2 + \frac{\mu_j^{(n)}}{2\mu_i^{(n)}} \mu_i^2$$

and

$$-\mu_i \mu_j \leq -\mu_i^{(n)} \mu_j^{(n)} \left[ 1 + \ln \left( \frac{\mu_i}{\mu_i^{(n)}} \right) + \ln \left( \frac{\mu_j}{\mu_j^{(n)}} \right) \right]$$

- compute the derivative  $g'(\cdot | \mu^{(n)})$  of  $g(\cdot | \mu^{(n)})$  to yield

$$\left[ \sum_{\{i|P_{ki}>0\}} (\mu_i^{(n)} / \mu_k^{(n)}) P_{ki} \right] \mu_k - \left[ \sum_{\{i|P_{ki}<0\}} \mu_i^{(n)} \mu_k^{(n)} |P_{ki}| \right] \frac{1}{\mu_k} - q_k = 0$$

$$\implies \alpha \mu_k^2 + \beta \mu_k + \gamma \quad \text{form}$$

$$\implies \text{take the positive root as } \mu_k^{(n+1)}$$

► in particular we get

$$\mu_k^{(n+1)} = \frac{q_k + \sqrt{q_k^2 + 4 \left[ \sum_{\{i|P_{ki}>0\}} \mu_i^{(n)} P_{ki} \right] \left[ \sum_{\{i|P_{ki}<0\}} \mu_i^{(n)} |P_{ki}| \right]}}{\left[ \sum_{\{i|P_{ki}>0\}} \mu_i^{(n)} P_{ki} \right]} \quad (37)$$

- ▶ iterates can be performed in parallel

---

## Algorithm 2 MM Principle

---

**Input:**  $\mu^{(0)} \succeq 0, n = 0$

- 1: **while** a stopping criterion true **do**
  - 2:      $\forall k, \mu_k^{(n+1)}$  is computed from (37) and  $n \leftarrow n + 1$
  - 3: **end while**
  - 4: **return**  $\mu^{(n+1)}$
- 

- ▶ then the solution is given by (35) with  $\mu = \mu^{(n+1)}$
- ▶ main differences between MM based algorithm and the coordinate descent?

# ARITHMETIC-GEOMETRIC MEAN INEQUALITY



# A Majorization to Monomials

- ▶ weighted arithmetic-geometric mean inequality

$$\prod_{i=1}^p x_i^{\alpha_i} \leq \sum_{i=1}^p \alpha_i x_i \quad \text{for all } x \succeq 0 \quad (38)$$

- ▶  $\alpha_i$  are given,  $\alpha_i > 0$ <sup>8</sup> and  $\sum_i \alpha_i = 1$
- ▶ (38): a majorization to  $\prod_{i=1}^p x_i^{\alpha_i}$  at  $\{\gamma 1 \in \mathbb{R}^p \mid \gamma \in \mathbb{R}_+\}$
- ▶ a majorization function to  $\prod_{i=1}^p x_i^{\beta_i}$  at
  - ▶ arbitrary  $y \succeq 0$  when  $\beta \succ 0$ ?

---

<sup>8</sup>If  $\alpha_i = 0$ , the corresponding  $x_i$  is irrelevant.

## A GENERAL MAJORIZATION FUNCTION

- ▶ let  $\beta_{\text{sum}} = \sum_i \beta_i$
- ▶ substitute  $x_i \leftarrow (x_i/y_i)^{\beta_{\text{sum}}}$  and  $\alpha_i \leftarrow \beta_i/\beta_{\text{sum}}$  in (38)
- ▶ thus, we get

$$\prod_{i=1}^p x_i^{\beta_i} \leq \left[ \prod_{i=1}^p y_i^{\beta_i} \right] \left[ \sum_{i=1}^p \frac{\beta_i}{\beta_{\text{sum}}} \left( \frac{x_i}{y_i} \right)^{\beta_{\text{sum}}} \right] \quad \text{for all } x \succeq 0$$
$$= g(x|y) \quad \text{for all } x \succeq 0$$

# A Minorization to Monomials

- ▶ we rely on the supporting hyperplane inequality

$$\log z \leq z - 1 \quad \text{for all } z \in \mathbb{R}_{++} \quad (39)$$

- ▶ suppose  $\beta \succcurlyeq 0$  is given,  $x_i > 0$
- ▶ substitute  $z = \prod_{i=1}^p (x_i/y_i)^{\beta_i}$  in (39), i.e.,

$$\begin{aligned} \prod_{i=1}^p x_i^{\beta_i} &\geq \prod_{i=1}^p y_i^{\beta_i} [1 + \sum_{i=1}^p \beta_i \ln (x_i/y_i)] \quad \text{for all } x \succcurlyeq 0 \\ &= \prod_{i=1}^p y_i^{\beta_i} [1 + \sum_{i=1}^p \beta_i \ln x_i - \sum_{i=1}^p \beta_i \ln y_i] \\ &= g(x|y) \end{aligned}$$

## EXAMPLES

# A Majorization to A Signomial

- ▶ consider the signomial  $f$

$$f(x) = \frac{1}{x_1^3} + \frac{3}{x_1 x_2^2} + x_1 x_2 - \sqrt{x_1 x_2} \quad (40)$$

- ▶ majorization function to  $f$  at  $y$ ?

- ▶  $1/(x_1 x_2^2) \leq y_1^2/(3y_2^2 x_1^3) + (2y_2)/(3y_1 x_2^3)$
- ▶  $x_1 x_2 \leq (y_2 x_1^2)/(2y_1) + (y_1 x_2^2)/(2y_2)$
- ▶  $\sqrt{x_1 x_2} \geq (1/2)\sqrt{y_1 y_2}(2 + \ln x_1 + \ln x_2 - \ln y_1 - \ln y_2)$

## APPENDICES

# Composition with Affine Function

- ▶ suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is differentiable
- ▶ then define  $\tilde{f} : \mathbb{R} \rightarrow \mathbb{R}$  by

$$\tilde{f}(\tau) = f(a + \tau b)$$

is differentiable and <sup>9</sup>

$$\tilde{f}'(\tau) = \frac{d\tilde{f}(\tau)}{d\tau} = \nabla f(a + \tau b)^\top b \quad (41)$$

---

<sup>9</sup>see § A.4.2 of Convex Optimization by S. Boyd and L. Vandenberghe, 2004.

# Composition with Affine Function

- ▶ suppose  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is twice differentiable
- ▶ then  $\tilde{f}$  [cf. (41)] is twice differentiable and <sup>10</sup>

$$\tilde{f}''(\tau) = \frac{d^2 \tilde{f}(\tau)}{d\tau^2} = b^\top \nabla^2 f(a + \tau b) b \quad (42)$$

---

<sup>10</sup>see § A.4.4 of Convex Optimization by S. Boyd and L. Vandenberghe, 2004.



# Newton-Leibniz Formula

- ▶ recall  $\tilde{f}(\tau) = f(a + \tau b)$
- ▶ let us apply Newton-Leibniz formula <sup>11</sup>

$$\tilde{f}(1) = \tilde{f}(0) + \int_0^1 \tilde{f}'(t) dt \quad (43)$$

- ▶ thus from (41), (43) becomes

$$f(a + b) = f(a) + \nabla f(a)^\top b + \int_0^1 \nabla f(a + tb)^\top b dt \quad (44)$$

---

<sup>11</sup>Based on elementary classical analysis.

# Taylor with the Integral Remainder

- ▶ Taylor formula with the integral remainder
- ▶ recall  $\tilde{f}(\tau) = f(a + \tau b)$
- ▶ we have <sup>12</sup>

$$\tilde{f}(1) = \tilde{f}(0) + \tilde{f}'(0) + \int_0^1 \int_0^t \tilde{f}''(\tau) d\tau dt \quad (45)$$

- ▶ thus from (41) and (42), (45) becomes

$$f(a+b) = f(a) + \nabla f(a)^\top b + \int_0^1 \int_0^t b^\top \nabla^2 f(a+\tau b) b d\tau dt \quad (46)$$

---

<sup>12</sup>Based on elementary classical analysis.