

# Convex Optimization, Spring 2021

## Homework 5

Due on June 9, 2021

Read all the instructions below carefully before you start working on the assignment, and before you make a submission.

- You are required to write down all the major steps towards making your conclusions; otherwise you may obtain limited points ( 20%) of the problem.
- Write your homework in English; otherwise you will get no points of this homework.
- Do your homework by yourself. Any form of plagiarism will lead to 0 point of this homework. If more than one plagiarisms during the semester are identified, we will prosecute all violations to the fullest extent of the university regulations, including but not limited to failing this course, academic probation, or expulsion from the university.
- If you have any doubts regarding the grading, you need to contact the instructor or the TAs within two days since the grade is announced.

## I. Proximal Gradient Methods

A useful extension of Moreau decomposition can be stated as

**Fact 1** (refer to Lecture 11). *Suppose  $f$  is closed convex and  $\lambda > 0$ . Then*

$$\mathbf{x} = \text{prox}_{\lambda f}(\mathbf{x}) + \lambda \text{prox}_{\frac{1}{\lambda} f^*}(\mathbf{x}/\lambda). \quad (1)$$

Please verify this fact (the derivation process is required). (30 points)

Solution:

*Proof.* Using Moreau decomposition, for any  $\mathbf{x}$ ,

$$\text{prox}_{\lambda f}(\mathbf{x}) = \mathbf{x} - \text{prox}_{(\lambda f)^*}(\mathbf{x}) = \mathbf{x} - \text{prox}_{\lambda f^*}(\mathbf{x}/\lambda), \quad (2)$$

where the second equality follows by Theorem 4.14(a) in [1]. By Theorem 6.12 in [1],

$$\text{prox}_{\lambda f^*}(\mathbf{x}/\lambda) = \lambda \text{prox}_{\lambda^{-1} f^*}(\mathbf{x}/\lambda), \quad (3)$$

which, combined with (2), yields (1).  $\square$

## II. Accelerated Gradient Methods

The momentum coefficient of the fast iterative shrinkage-thresholding algorithm (FISTA) (refer to lecture 11) at  $t$ -th iterate is given as

$$\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2} \quad \text{with} \quad \theta_0 = 1.$$

(1) Show that

$$\theta_t \geq \frac{t+2}{2} \quad (4)$$

holds true for all  $t \geq 1$ . (20 points)

(2) Show that

$$\frac{\theta_t - 1}{\theta_{t+1}} = 1 - \frac{3}{t} + o\left(\frac{1}{t}\right). \quad (5)$$

(20 points)

Solution:

(1) *Proof.* The proof is by induction on  $t$ .

– Base case: Obviously, for  $t = 0$ ,  $\theta_0 = 1 \geq \frac{0+2}{2}$ .

– Inductive step: Suppose that the claim holds for  $t$ , meaning that  $\theta_t \geq \frac{t+2}{2}$ . We now show that  $\theta_{t+1} \geq \frac{t+3}{2}$ . By the recursive relation defining the sequence and the induction assumption,

$$\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2} \geq \frac{1 + \sqrt{1 + (t+2)^2}}{2} \geq \frac{1 + \sqrt{(t+2)^2}}{2} = \frac{t+3}{2} \quad (6)$$

This completes the proof.  $\square$

(2) *Proof.* Since  $\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2} = \frac{2\theta_t^2}{\sqrt{1 + 4\theta_t^2} - 1}$ , we have

$$\frac{\theta_t - 1}{\theta_{t+1}} = \frac{(\theta_t - 1)(\sqrt{1 + 4\theta_t^2} - 1)}{2\theta_t^2} \quad (7)$$

Due to the fact that  $\theta_t \geq \frac{t+2}{2}, \forall t \geq 1$  (we require you to prove in the next part), we have  $\frac{1}{\theta_t} \leq \frac{2}{t+2}$ . Thus

$$\frac{(\theta_t - 1)(\sqrt{1 + 4\theta_t^2} - 1)}{2\theta_t^2} = \frac{(\theta_t - 1)(2\theta_t - 1)}{2\theta_t^2} + o\left(\frac{1}{t}\right) \quad (8)$$

$$= 1 - \frac{3}{2\theta_t} + o\left(\frac{1}{t}\right) \quad (9)$$

$$= 1 - \frac{3}{t} + o\left(\frac{1}{t}\right). \quad (10)$$

This completes the proof.  $\square$

### III. Alternating Direction Method of Multipliers

Suppose we observe  $\mathbf{M}$  (refer to lecture 14), which is superposition of low-rank component  $\mathbf{L}$  and sparse outliers  $\mathbf{S}$ , provide the ADMM update (the derivation process is required) for each variable at the  $t$ -th iteration

$$\begin{aligned} & \underset{\mathbf{L}, \mathbf{S}}{\text{minimize}} \quad \|\mathbf{L}\|_* + \lambda \|\mathbf{S}\|_1 \\ & \text{subject to} \quad \mathbf{L} + \mathbf{S} = \mathbf{M}, \end{aligned} \quad (11)$$

where  $\|\mathbf{L}\|_* := \sum_{i=1}^n \sigma_i(\mathbf{L})$  is nuclear norm, and  $\|\mathbf{S}\|_1 = \sum_{i,j} |S_{i,j}|$  is entrywise  $\ell_1$ -norm. (30 points)

Solution:

By introducing dual variable  $\mathbf{\Lambda}$ , ADMM update rules for solving (11) are given by

$$\mathbf{L}^{t+1} = \arg \min_{\mathbf{L}} \left\{ \|\mathbf{L}\|_* + \frac{\rho}{2} \left\| \mathbf{L} + \mathbf{S}^t - \mathbf{M} + \frac{1}{\rho} \mathbf{\Lambda}^t \right\|_{\text{F}}^2 \right\} \quad (12)$$

$$\mathbf{S}^{t+1} = \arg \min_{\mathbf{S}} \left\{ \lambda \|\mathbf{S}\|_1 + \frac{\rho}{2} \left\| \mathbf{L}^{t+1} + \mathbf{S} - \mathbf{M} + \frac{1}{\rho} \mathbf{\Lambda}^t \right\|_{\text{F}}^2 \right\} \quad (13)$$

$$\mathbf{\Lambda}^{t+1} = \mathbf{\Lambda}^t + \rho (\mathbf{L}^{t+1} + \mathbf{S}^{t+1} - \mathbf{M}) \quad (14)$$

Define the function  $h_t(\mathbf{L}) = \|\mathbf{L}\|_* + \frac{\rho}{2} \left\| \mathbf{L} + \mathbf{S}^t - \mathbf{M} + \frac{1}{\rho} \mathbf{\Lambda}^t \right\|_{\text{F}}^2$ . Since  $h_t(\mathbf{L})$  is strictly convex, it is easy to see that there exists a unique minimizer. Now  $\hat{\mathbf{L}}$  minimizes  $h_t$  if and only if  $\mathbf{0}$  is a subgradient of the functional  $h_t$  at the point  $\hat{\mathbf{L}}$ , i.e.,

$$\mathbf{0} \in \hat{\mathbf{L}} - \left( \mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t \right) + \frac{1}{\rho} \partial \|\hat{\mathbf{L}}\|_*, \quad (15)$$

where  $\partial \|\hat{\mathbf{L}}\|_*$  is the set of subgradients of the nuclear norm. Let  $\mathbf{L}$  be an arbitrary matrix and  $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$  be its SVD. It is known [2] that

$$\partial \|\mathbf{L}\|_* = \{ \mathbf{U}\mathbf{V}^* + \mathbf{W} : \mathbf{U}^* \mathbf{W} = \mathbf{0}, \quad \mathbf{W}\mathbf{V} = \mathbf{0}, \quad \|\mathbf{W}\|_2 \leq 1 \}. \quad (16)$$

Set  $\hat{\mathbf{L}} = \text{SVT}_{\rho^{-1}} \left( \mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t \right)$ . Now we need to show that  $\hat{\mathbf{L}}$  obeys (15). Decompose the SVD of  $\mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t$  as  $\mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t = \mathbf{U}_0 \mathbf{\Sigma}_0 \mathbf{V}_0^* + \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^*$ , where  $\mathbf{U}_0, \mathbf{V}_0$  (resp.  $\mathbf{U}_1, \mathbf{V}_1$ ) are the singular vectors associated with singular values greater than  $1/\rho$  (resp. smaller than or equal to  $1/\rho$ ). With these notations, we have  $\hat{\mathbf{L}} = \mathbf{U}_0 (\mathbf{\Sigma}_0 - 1/\rho \mathbf{I}) \mathbf{V}_0^*$ , and therefore,

$$\mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t - \hat{\mathbf{L}} = 1/\rho (\mathbf{U}_0 \mathbf{V}_0^* + \mathbf{W}), \quad \mathbf{W} = \rho \mathbf{U}_1 \mathbf{\Sigma}_1 \mathbf{V}_1^* \quad (17)$$

By definition,  $\mathbf{U}_0^* \mathbf{W} = \mathbf{0}, \mathbf{W}\mathbf{V}_0 = \mathbf{0}$  and since the diagonal elements of  $\mathbf{\Sigma}_1$  have magnitudes bounded by  $1/\rho$ , we also have  $\|\mathbf{W}\|_2 \leq 1$ . Hence,  $\mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t - \hat{\mathbf{L}} \in 1/\rho \partial \|\hat{\mathbf{L}}\|_*$ , which implies

$$\mathbf{L}^{t+1} = \text{SVT}_{\rho^{-1}} \left( \mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t \right). \quad (18)$$

Define the function  $g_t(\mathbf{S}) = \lambda \|\mathbf{S}\|_1 + \frac{\rho}{2} \left\| \mathbf{L}^{t+1} + \mathbf{S} - \mathbf{M} + \frac{1}{\rho} \mathbf{\Lambda}^t \right\|_{\text{F}}^2$ . Similarly,  $\hat{\mathbf{S}}$  should satisfy

$$\mathbf{0} \in \hat{\mathbf{S}} - \left( \mathbf{M} - \mathbf{L}^{t+1} - \frac{1}{\rho} \mathbf{\Lambda}^t \right) + \frac{\lambda}{\rho} \partial \|\hat{\mathbf{S}}\|_1, \quad (19)$$

where  $\partial \|\hat{\mathbf{S}}\|_1$  is the set of subgradients of the  $\ell_1$ -norm. Since  $\|\mathbf{S}\|_1 = \sum_{i,j} |S_{i,j}|$ , we have

$$\sum_{i,j} \text{sign}(S_{i,j}) \mathbf{E}_{i,j} \in \partial \|\hat{\mathbf{S}}\|_1, \quad (20)$$

where  $\mathbf{E}_{i,j}$  is a matrix with only  $E_{i,j} = 1$ . That combines with (19) implies

$$\mathbf{S}^{t+1} = \hat{\mathbf{S}} = \text{ST}_{\lambda \rho^{-1}} \left( \mathbf{M} - \mathbf{L}^{t+1} - \frac{1}{\rho} \mathbf{\Lambda}^t \right). \quad (21)$$

Hence, the exact the ADMM rules are summarized as

$$\mathbf{L}^{t+1} = \text{SVT}_{\rho^{-1}} \left( \mathbf{M} - \mathbf{S}^t - \frac{1}{\rho} \mathbf{\Lambda}^t \right) \quad (22)$$

$$\mathbf{S}^{t+1} = \text{ST}_{\lambda\rho^{-1}} \left( \mathbf{M} - \mathbf{L}^{t+1} - \frac{1}{\rho} \mathbf{\Lambda}^t \right) \quad (23)$$

$$\mathbf{\Lambda}^{t+1} = \mathbf{\Lambda}^t + \rho (\mathbf{L}^{t+1} + \mathbf{S}^{t+1} - \mathbf{M}) . \quad (24)$$

## REFERENCES

- [1] A. Beck, *First-Order Methods in Optimization*. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2017.
- [2] A. S. Lewis, “The mathematics of eigenvalue optimization,” *Math. Program.*, vol. 97, no. 1-2, pp. 155–176, 2003.