

## Lecture 5

## Unconstrained Optimization for Smooth Strongly-Convex Functions, Part V

Lecturer: Bin Hu, Date:09/13/2018

In last lecture, we talked about momentum methods and a general model for first-order optimization methods.

In this lecture, we will present a general analysis framework for the general optimization model.

## 5.1 Remarks on the general optimization model

We have shown that gradient method, Heavy-ball method, and Nesterov's method are all special cases of the following general model.

$$\xi_{k+1} = A\xi_k + Bu_k \quad (5.1)$$

$$v_k = C\xi_k \quad (5.2)$$

$$u_k = \nabla f(v_k) \quad (5.3)$$

A natural question is what type of  $(A, B, C)$  will lead to an optimization method. One condition is that  $A$  should have eigenvalue at 1. You can verify this for the existing choice of  $A$  for the gradient method, Heavy-ball method, and Nesterov's method. The basic idea is that the iteration should not move if starting from the optimal point. The general model is equivalent to  $\xi_{k+1} = A\xi_k + B\nabla f(C\xi_k)$ . If starting from a point  $\xi^*$  such that  $\nabla f(C\xi^*) = 0$ , we want the iteration to stay at  $\xi^*$ , i.e.  $\xi^* = \xi_1 = A\xi_0 = A\xi^*$ . Therefore we have  $\xi^* = A\xi^*$  and 1 is an eigenvalue of  $A$ .

## 5.2 Dissipation Inequality and Physical Interpretation

An analysis routine for the general optimization model is provided by the dissipation inequality approach. We want to analyze the behavior of  $(\xi_k, u_k, v_k)$  satisfying (5.1) (5.2) and (5.3). For now let us first focus on (5.1).

Specifically, we look at  $\xi_{k+1} = A\xi_k + Bu_k$ . The sequence of  $\{u_k : k = 0, 1, \dots\}$  is called "input". Dissipation inequality just describes how the input  $u_k$  changes the energy of the "state"  $\xi_k$ .

**Definition 1.** *The system  $\xi_{k+1} = A\xi_k + Bu_k$  is dissipative with respect to the supply rate  $S(\xi, u)$  if there exists  $V : \mathbb{R}^{n_\xi} \mapsto \mathbb{R}^+$  such that*

$$V(\xi_{k+1}) - V(\xi_k) \leq S(\xi_k, u_k) \quad (5.4)$$

for all  $k$ . The function  $V$  is called a storage function, which quantifies the energy stored in the state  $\xi$ . The supply rate  $S$  is a function that quantifies the energy supplied to the state  $\xi_k$  by the input  $u_k$ . In addition, (5.4) is called the dissipation inequality.

The dissipation inequality (5.4) states that the internal energy increase is equal to the sum of the supplied energy and the energy dissipation. Since there will always be some energy dissipating from the system, hence the internal energy increase (which is exactly  $V(\xi_{k+1}) - V(\xi_k)$ ) is always bounded above by the energy supplied to the system (which is exactly  $S(\xi_k, u_k)$ ).

One important variant of the original dissipation inequality is the so-called exponential dissipation inequality:

$$V(\xi_{k+1}) - \rho^2 V(\xi_k) \leq S(\xi_k, u_k) \quad (5.5)$$

where  $0 < \rho^2 < 1$ . The dissipation inequality (5.5) just states that at least a  $(1 - \rho^2)$  fraction of the internal energy will dissipate at every step, and hence the internal energy at step  $k + 1$  is bounded above by the sum of the remaining energy  $\rho^2 V(\xi_k)$  and the supply energy  $S$ .

### 5.3 How to Use Dissipation Inequality?

Everything looks very abstract at this moment. Let's first ask if we can construct the dissipation inequality (5.5), what are we going to do about it? The answer is that the dissipation inequality (5.5) can be used to prove convergence rate bounds for the general optimization model (5.1) (5.2) and (5.3).

Notice by definition  $V_k \geq 0$  (the internal energy should be non-negative). Typically  $V$  is chosen to be a distance metric between  $\xi_k$  and the equilibrium point  $\xi^*$ . For example, for gradient method,  $V$  is chosen as  $V = \|x - x^*\|^2$ . The dissipation inequality can be typically used to prove two types of bounds.

1. If one already knows  $S \leq 0$ , then the dissipation inequality (5.5) states  $V(\xi_{k+1}) - \rho^2 V(\xi_k) \leq S(\xi_k, u_k) \leq 0$ . This gives a bound  $V(\xi_{k+1}) \leq \rho^2 V(\xi_k)$ . This proves a linear convergence rate  $\rho$  when  $V$  is used as a distance metric. For example, our analysis for the gradient method can be viewed as a special case of this type of analysis when  $V$  is chosen as  $\|x_k - x^*\|^2$ .
2. If one already knows  $S \leq \rho^2(f(x_k) - f(x^*)) - (f(x_{k+1}) - f(x^*))$ , then the dissipation inequality (5.5) states  $V(\xi_{k+1}) - \rho^2 V(\xi_k) \leq S(\xi_k, u_k) \leq \rho^2(f(x_k) - f(x^*)) - (f(x_{k+1}) - f(x^*))$ . This gives a bound  $V(\xi_{k+1}) + f(x_{k+1}) - f(x^*) \leq \rho^2 (V(\xi_k) + f(x_k) - f(x^*))$ . This proves a linear convergence rate  $\rho$  when  $V(\xi_k) + f(x_k) - f(x^*)$  is used as a distance metric. As we will see later, the convergence rate for Nesterov's accelerated method can be obtained using such an argument.

Typically it is much more difficult to construct a supply rate satisfying  $S \leq \rho^2(f(x_k) - f(x^*)) - (f(x_{k+1}) - f(x^*))$ . Consequently, the analysis of Nesterov's method is much more involved than the gradient method. Now we will talk about what  $S$  really looks like in convergence rate analysis.

## 5.4 How to Choose Supply Rate?

The supply rate  $S$  typically takes a form of a quadratic function:

$$S(\xi, u) = \begin{bmatrix} \xi - \xi^* \\ u \end{bmatrix}^T X \begin{bmatrix} \xi - \xi^* \\ u \end{bmatrix} \quad (5.6)$$

where  $X$  is some given matrix. The key question is how to choose  $X$ .

Recall that the general optimization model includes three parts: (5.1) (5.2) and (5.3). If we want to choose  $X$  to guarantee the supply rate  $S$  satisfying some inequality, e.g.  $S \leq 0$ , we have to use (5.2) and (5.3). Notice (5.2) and (5.3) just states  $u_k = \nabla f(C\xi_k)$ . Therefore, we could use the property of  $f$  to pick up  $X$  that can enforce  $S \leq 0$ . For example, if  $f$  is  $L$ -smooth and  $m$ -strongly convex, we know the following inequality holds for any  $u_k = \nabla f(C\xi_k)$

$$\begin{bmatrix} C\xi_k - C\xi^* \\ u_k \end{bmatrix} \begin{bmatrix} -2mLI & (m+L)I \\ (m+L)I & -2I \end{bmatrix} \begin{bmatrix} C\xi_k - C\xi^* \\ u_k \end{bmatrix} \geq 0. \quad (5.7)$$

The above inequality is just a restatement of (1b) in Homework 1. So we simply choose  $X = \begin{bmatrix} C^T & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} 2mLI & -(m+L)I \\ -(m+L)I & 2I \end{bmatrix} \begin{bmatrix} C & 0 \\ 0 & I \end{bmatrix}$  and then the supply rate (5.6) satisfies  $S \leq 0$  due to the fact  $u_k = \nabla f(C\xi_k)$ . (Verify this!)

How to construct a supply rate satisfying  $S \leq \rho^2(f(x_k) - f(x^*)) - (f(x_{k+1}) - f(x^*))$  for Nesterov's method? Basically we still use the properties of  $L$ -smoothness and  $m$ -strong convexity. We will cover the details of this construction in next lecture.

## 5.5 How to Construct the Dissipation Inequality?

Now suppose we have already constructed the supply rate (5.6) with desired properties. How can we construct the dissipation inequality (5.5) for such a supply rate? We can use the following approach.

**Theorem 2.** Suppose  $\xi_{k+1} = A\xi_k + Bu_k$  and  $\xi^* = A\xi^*$ . Consider a quadratic supply rate (5.6). If there exists a positive semidefinite matrix  $P \in \mathbb{R}^{n_\xi \times n_\xi}$  s.t.

$$\begin{bmatrix} A^T P A - \rho^2 P & A^T P B \\ B^T P A & B^T P B \end{bmatrix} - X \leq 0 \quad (5.8)$$

then we have  $V(\xi_{k+1}) - \rho^2 V(\xi_k) \leq S(\xi_k, u_k)$  with  $V(\xi) = (\xi - \xi^*)^T P (\xi - \xi^*)$ .

The proof will be presented in the next lecture. Now we apply the above theorem to recover the analysis condition we have obtained for the gradient method as a special case. Recall that we have obtained an analysis condition for gradient method. In 2b of Homework 1, we have used that condition to obtain a convergence rate formula for the gradient method. That condition can be viewed as a corollary of the above theorem. For gradient method, we have  $A = I$ ,  $B = -\alpha I$ , and  $C = I$ . As discussed in the last section, we can choose the following  $X$  to guarantee  $S \leq 0$ :

$$X = \begin{bmatrix} C^T & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} 2mLI & -(m+L)I \\ -(m+L)I & 2I \end{bmatrix} \begin{bmatrix} C & 0 \\ 0 & I \end{bmatrix} = \begin{bmatrix} 2mLI & -(m+L)I \\ -(m+L)I & 2I \end{bmatrix}$$

Now it is straightforward to verify that the condition (5.8) leads to the following condition

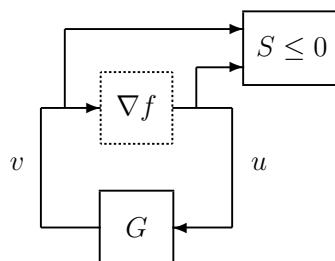
$$\begin{bmatrix} 1 - \rho^2 & -\alpha \\ -\alpha & \alpha^2 \end{bmatrix} + \begin{bmatrix} -2mL & m+L \\ m+L & -2 \end{bmatrix} \leq 0 \quad (5.9)$$

if we choose  $P = \frac{1}{\lambda}I$ . (The above condition looks familiar, right? 2b in Homework 1!)

The key message is that to apply the dissipation inequality for linear convergence rate analysis, one typically follows two steps:

1. Choose a proper quadratic supply rate function  $S$  satisfying certain desired properties, e.q.  $S(\xi_k, u_k) \leq 0$ .
2. Find a positive semidefinite matrix  $P$  satisfying (5.8) and obtain a quadratic storage function  $V$  which is then used to construct the dissipation inequality.

## 5.6 Graphical Interpretation



**Figure 5.1.** Removing  $\nabla f$  by Enforcing the Supply Rate Condition  $S(\xi_k, u_k) \leq 0$

Finally, we give a graphical interpretation for the dissipation inequality approach. The general optimization model (5.1) (5.2) and (5.3) can be viewed as a feedback loop of  $G$  and  $\nabla f$  where  $G$  maps  $u$  to  $v$  by the state space model  $\xi_{k+1} = A\xi_k + Bu_k$  and  $v_k = C\xi_k$  and  $\nabla f$  maps  $v$  to  $u$  as  $u_k = \nabla f(v_k)$ . When trying to analyze such a feedback interconnection, we

aim to draw conclusions on the pair  $(v, u)$  in the set  $\{(v, u) : v = G(u), u = \nabla f(v)\}$ . If for any  $u_k = \nabla f(v_k) = \nabla f(G\xi_k)$ , we have  $S(\xi_k, u_k) \leq 0$ , then we have

$$\{(v, u) : v = G(u), u = \nabla f(v)\} \subset \{(v, u) : v = G(u), S(\xi_k, u_k) \leq 0\} \quad (5.10)$$

If we can prove  $\xi_k$  converges at a certain linear rate for any pair  $(v, u)$  in the set  $\{(v, u) : v = G(u), S(\xi_k, u_k) \leq 0\}$ , then we guarantee that  $\xi_k$  converges at the same linear rate for any pair  $(v, u)$  satisfying  $v = G(u)$  and  $u = \nabla f(v)$  simultaneously. Hence we can completely remove  $\nabla f$  from our analysis by enforcing the condition  $S(\xi_k, u_k) \leq 0$ . A graphical interpretation for this idea is shown in Figure 5.1. We still have  $v = G(u)$ . But we remove  $\nabla f$  by enforcing the inequality  $S \leq 0$ .