

## Smoothing for non-smooth optimization, lecture 2

Last time:

- Basics of convex analysis

Today:

- Standard gradient scheme (smooth):  $O(1/\epsilon)$ 
  - unconstrained
  - simply-constrained
- Optimal gradient scheme (simply-constrained):  $O(1/\sqrt{\epsilon})$
- Optimal subgradient scheme (non-smooth):  $O(1/\epsilon^2)$

1

## Unconstrained optimization

$$\min \{f(x) : x \in E\}$$

### Notation

- $\bar{f} := \min \{f(x) : x \in E\}$
- $\bar{x} := \operatorname{argmin} \{f(x) : x \in E\}$

### Assume (until we state otherwise)

- $E$ : finite-dimensional Euclidean space
- $f : E \rightarrow \mathbb{R}$  is a smooth convex function
- $\nabla f$  is Lipschitz with constant  $L$

2

### Algorithm 1

- Pick  $x_0 \in E$
  - For  $k = 0, 1, \dots$
- $$x_{k+1} := x_k - h_k \nabla f(x_k)$$

**Theorem 2** If  $h_j = \frac{1}{L}$  for  $j = 1, 2, \dots$ , then

$$f(x_k) - \bar{f} \leq \frac{2L(f(x_0) - \bar{f})\|x_0 - \bar{x}\|^2}{2L\|x_0 - \bar{x}\|^2 + k(f(x_0) - \bar{f})} \sim \frac{1}{k}$$

3

### Immediate (but important) observations

- Speed of convergence to optimal value is  $O(1/k)$
- Complexity to get within  $\epsilon$  of optimal value is  $O(1/\epsilon)$
- Gradient-step update at each iteration:

$$x - \frac{1}{L} \nabla f(x) = \operatorname{argmin}_y \left\{ \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2 \right\}.$$

$$x \rightarrow x - \frac{1}{L} \nabla f(x)$$

4

**Proof.** Since  $f$  is convex and  $\nabla f$  is Lipschitz with constant  $L$ ,

- (1) •  $0 \leq f(y) - f(x) - \langle \nabla f(x), y - x \rangle \leq \frac{L}{2} \|y - x\|^2$  for all  $x, y$
- (2) •  $\frac{1}{L} \|\nabla f(y) - \nabla f(x)\|^2 \leq \langle \nabla f(y) - \nabla f(x), y - x \rangle$  for all  $x, y \in C$

Hence for  $x \in E$

$$f\left(x - \frac{1}{L} \nabla f(x)\right) \leq f(x) + \langle \nabla f(x), -\frac{1}{L} \nabla f(x) \rangle + \frac{L}{2} \|\frac{1}{L} \nabla f(x)\|^2$$

$$f\left(x - \frac{1}{L} \nabla f(x)\right) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2,$$

$$f(x) - \bar{f} \leq \langle \nabla f(x), x - \bar{x} \rangle \leq \|\nabla f(x)\| \|x - \bar{x}\|,$$

and

$$\langle \nabla f(x), x - \bar{x} \rangle \geq \frac{1}{L} \|\nabla f(x)\|^2.$$

Thus

$$\|x - \frac{1}{L} \nabla f(x) - \bar{x}\|^2 \leq \|x - \bar{x}\|^2 - \frac{1}{L^2} \|\nabla f(x)\|^2 \leq \|x - \bar{x}\|^2.$$

5

Put  $\Delta_k := f(x_k) - \bar{f}$ ,  $r_k := \|x_k - \bar{x}\|$ . We get

$$\Delta_{k+1} \leq \Delta_k - \frac{1}{2L} \|\nabla f(x_k)\|^2 \leq \Delta_k - \frac{\Delta_k^2}{2Lr_k^2} \leq \Delta_k - \frac{\Delta_k^2}{2Lr_0^2}.$$

Therefore

$$\frac{1}{\Delta_{k+1}} \geq \frac{1}{\Delta_k} + \frac{\Delta_k}{2Lr_0^2 \Delta_{k+1}} \geq \frac{1}{\Delta_k} + \frac{1}{2Lr_0^2} \geq \frac{1}{\Delta_{k-1}} + \frac{2}{2Lr_0^2}$$

Consequently

$$\frac{1}{\Delta_k} \geq \frac{1}{\Delta_0} + \frac{k}{2Lr_0^2},$$

that is,

$$\Delta_k \leq \frac{2L\Delta_0 r_0^2}{2Lr_0^2 + k\Delta_0}$$

□

## Simply-constrained optimization

$$\min \{f(x) : x \in Q\}$$

Assume  $Q \subseteq E$  is a “simple” compact convex set.

Overwrite previous notation

- $\bar{f} := \min \{f(x) : x \in Q\}$
- $\bar{x} := \operatorname{argmin} \{f(x) : x \in Q\}$

**Key ingredient:** Consider the *gradient mapping*

$$T_Q(x) := \operatorname{argmin}_y \left\{ \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2 : y \in Q \right\}$$

## Algorithm 3

- Pick  $x_0 \in E$
- For  $k = 0, 1, \dots$   
 $x_{k+1} := T_Q(x_k)$

## Theorem 4

$$f(x_k) - \bar{f} \leq \frac{(f(x_0) - \bar{f})(2L\|x_0 - \bar{x}\|^2 + f(x_0) - \bar{f})}{2L\|x_0 - \bar{x}\|^2 + f(x_0) - \bar{f} + k}.$$

Similar speed of convergence to that in the unconstrained case.

7

8

**Proof.** (similar to that of Theorem 2.)

Main steps: Put  $g(x) := L(x - T_Q(x))$  (sort of  $\nabla f(x)$ ). Then

$$f(T_Q(x)) \leq f(x) - \frac{1}{2L} \|g(x)\|^2,$$

$$f(T_Q(x)) - \bar{f} \leq \langle g(x), x - \bar{x} \rangle \leq \|g(x)\| \|x - \bar{x}\|,$$

and

$$\langle g(x), x - \bar{x} \rangle \geq \frac{1}{2L} \|g(x)\|^2.$$

9

Hence

$$\Delta_{k+1} \leq \Delta_k - \frac{1}{2L} \|g(x_k)\|^2 \leq \Delta_k - \frac{\Delta_{k+1}^2}{2Lr_0^2},$$

and therefore,

$$\Delta_k \leq \frac{\Delta_0(2Lr_0^2 + \Delta_0)}{2Lr_0^2 + \Delta_0 + k}.$$

□

The main steps above are consequences of the following lemma.

**Lemma 5** Assume  $x \in E, T_Q(x), g(x)$  are as above. Then for any  $y \in Q$

$$f(y) \geq f(T_Q(x)) + \langle g(x), y - x \rangle + \frac{1}{2L} \|g(x)\|^2.$$

(for  $y = x \Rightarrow f(x) \geq f(T_Q(x)) + \frac{1}{2L} \|g(x)\|^2$ )

## Optimal gradient scheme

Assume  $d$  is a prox-function of  $Q$ :

- $d$  is strongly convex on  $Q$  with modulus  $\rho > 0$
- $\min \{d(x) : x \in Q\} = 0$
- $\min \{d(x) - \langle g, x \rangle : x \in Q\}$  is easily computable

**Idea:** construct  $\{x_k\}, \{y_k\} \subseteq Q$  such that

$$\frac{(k+1)(k+2)}{4} f(y_k) \leq \min_x \left\{ \frac{L}{\rho} d(x) + \sum_{i=0}^k \frac{i+1}{2} (f(x_i) + \langle \nabla f(x_i), x - x_i \rangle) : x \in Q \right\} \quad (1)$$

**Theorem 6** If  $\{x_k\}, \{y_k\}$  satisfy (1) then

$$f(y_k) - \bar{f} \leq \frac{4Ld(\bar{x})}{\rho(k+1)(k+2)}$$

**Proof.** Since  $f(x_i) + \langle \nabla f(x_i), x - x_i \rangle \leq f(x)$  for  $i = 1, \dots, k$  from (1) we get

$$\begin{aligned} \frac{(k+1)(k+2)}{4} f(y_k) &\leq \min \left\{ \frac{L}{\rho} d(x) + \frac{(k+1)(k+2)}{4} f(x) : x \in Q \right\} \\ &\leq \frac{L}{\rho} d(\bar{x}) + \frac{(k+1)(k+2)}{4} f(\bar{x}) \end{aligned}$$

□

## Remarks

- If we can generate  $\{x_k\}, \{y_k\}$  that satisfy (1), then we get an algorithm with speed of convergence  $O(1/k^2)$
- Complexity to get within  $\epsilon$ -approximate solution:  $O(1/\sqrt{\epsilon})$
- This would be an *optimal* complexity result: for any "black-box" gradient-based algorithm there are smooth problems that require  $\Omega(1/\sqrt{\epsilon})$  iterations to find an  $\epsilon$ -approximate solution.

13

## Algorithm 7

- Set  $x_0 := \operatorname{argmin}\{d(x) : x \in Q\}$ ,  $y_0 := T_Q(x_0)$
- For  $k = 0, 1, \dots$ 

$$z_k := \operatorname{argmin}_{x \in Q} \left\{ \frac{L}{\rho} d(x) + \sum_{i=0}^k \frac{i+1}{2} (f(x_i) + \langle \nabla f(x_i), x - x_i \rangle) : x \in Q \right\}$$

$$x_{k+1} := \frac{2}{k+3} z_k + \frac{k+1}{k+3} y_k$$

$$y_{k+1} := T_Q(x_{k+1})$$

**Theorem 8** The sequences  $\{x_k\}, \{y_k\}$  generated by Algorithm 7 satisfy (1).

14

**Proof.** Will use the following claim:

**Claim:** If  $d$  is strongly convex on  $Q$  with modulus  $\rho$  and  $x_0 = \operatorname{argmin}\{d(x) : x \in Q\}$  then

$$d(x) \geq d(x_0) + \frac{\rho}{2} \|x - x_0\|^2, \quad \text{for all } y \in Q.$$

$\langle \nabla f(x), T_\alpha(x) - x \rangle + \frac{L}{2} \|T_\alpha(x) - x\|^2$

By construction, for  $x \in Q$  we have

$$f(T_Q(x)) \leq f(x) + \underbrace{\min_y \left\{ \langle \nabla f(x), y - x \rangle + \frac{1}{2} L \|y - x\|^2 : y \in Q \right\}}_{\text{for } x = x_0} \leq \frac{L}{\rho} d(y)$$

Since  $d(y) \geq \frac{\rho}{2} \|y - x_0\|^2$ , we have

$$f(T_Q(x_0)) \leq \min_y \left\{ f(x_0) + \langle \nabla f(x_0), y - x_0 \rangle + \frac{L}{\rho} d(y) : y \in Q \right\}.$$

Thus (1) holds for  $k = 0$ . Next proceed by induction on  $k$ .

15

Put  $A_k := \frac{(k+1)(k+2)}{4}$ ,  $\alpha_k := \frac{k+1}{2}$ ,  $\tau_k := \frac{2}{k+3}$ , and

$$\psi_k := \min_x \left\{ \frac{L}{\rho} d(x) + \sum_{i=0}^k \frac{i+1}{2} (f(x_i) + \langle \nabla f(x_i), x - x_i \rangle) : x \in Q \right\}.$$

Want to show that if  $A_k f(y_k) \leq \psi_k$  then  $A_{k+1} f(y_{k+1}) \leq \psi_{k+1}$ .

16

Proceed in two main steps.

Step 1:

$$\psi_{k+1} \geq A_{k+1} \left( f(x_{k+1}) + \min_{x \in Q} \left\{ \frac{\tau_k^2}{2} L \|x - z_k\|^2 + \tau_k \langle \nabla f(x_{k+1}), x - z_k \rangle \right\} \right).$$

Step 2:

$$\begin{aligned} & \min_x \left\{ \frac{\tau_k^2}{2} L \|x - z_k\|^2 + \tau_k \langle \nabla f(x_{k+1}), x - z_k \rangle : x \in Q \right\} \\ & \geq \min_y \left\{ \frac{1}{2} L \|y - z_k\|^2 + \langle \nabla f(x_{k+1}), y - z_k \rangle : y \in Q \right\} \\ & \geq f(y_{k+1}) - f(x_{k+1}). \end{aligned}$$

17

Details of step 1:

$$\begin{aligned} \psi_{k+1} & \geq \min_x \left\{ \psi_k + \frac{L}{2} \|x - z_k\|^2 + \alpha_{k+1} (f(x_{k+1}) + \langle \nabla f(x_{k+1}), x - x_{k+1} \rangle) : x \in Q \right\} \\ & \geq A_{k+1} \left( f(x_{k+1}) + \min_x \left\{ \frac{\tau_k^2}{2} L \|x - z_k\|^2 + \tau_k \langle \nabla f(x_{k+1}), x - z_k \rangle \right\} : x \in Q \right). \end{aligned}$$

The first inequality holds because  $d$  is strongly convex on  $Q$  with modulus  $\rho$ .

The second inequality follows from  $A_{k+1}\tau_k = \alpha_{k+1} \leq 1/\tau_k$  and

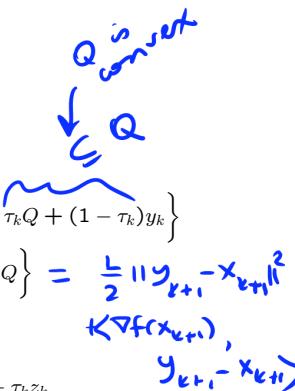
$$\begin{aligned} & \psi_k + \alpha_{k+1} (f(x_{k+1}) + \langle \nabla f(x_{k+1}), x - x_{k+1} \rangle) \\ & \quad \xrightarrow{\text{use}} \geq A_k f(y_k) + \alpha_{k+1} (f(x_{k+1}) + \langle \nabla f(x_{k+1}), x - x_{k+1} \rangle) \\ & \quad \geq A_k (f(x_{k+1}) + \langle \nabla f(x_{k+1}), y_k - x_{k+1} \rangle) \\ & \quad \quad \alpha_{k+1} (f(x_{k+1}) + \langle \nabla f(x_{k+1}), x - x_{k+1} \rangle) \\ & \quad = A_{k+1} f(x_{k+1}) + \alpha_{k+1} \langle \nabla f(x_{k+1}), x - z_k \rangle. \end{aligned}$$

$$\Psi_k \geq A_k f(y_k)$$

18

Details of step 2:

$$\begin{aligned} & \min_x \left\{ \frac{\tau_k^2}{2} L \|x - z_k\|^2 + \tau_k \langle \nabla f(x_{k+1}), x - z_k \rangle : x \in Q \right\} \\ & = \min_y \left\{ \frac{1}{2} L \|y - x_{k+1}\|^2 + \langle \nabla f(x_{k+1}), y - x_{k+1} \rangle : y \in \tau_k Q + (1 - \tau_k) y_k \right\} \\ & \geq \min_y \left\{ \frac{1}{2} L \|y - x_{k+1}\|^2 + \langle \nabla f(x_{k+1}), y - x_{k+1} \rangle : y \in Q \right\} = \frac{L}{2} \|y_{k+1} - x_{k+1}\|^2 \\ & \quad \left\langle \nabla f(x_{k+1}), y_{k+1} - x_{k+1} \right\rangle \\ & \geq f(y_{k+1}) - f(x_{k+1}). \end{aligned}$$



The first equality follows from

$$\begin{aligned} \tau_k(x - z_k) &= \tau_k x + (1 - \tau_k) y_k - (1 - \tau_k) y_k - \tau_k z_k \\ &= \tau_k x + (1 - \tau_k) y_k + x_{k+1}. \end{aligned}$$

□

19

### Simply-constrained non-smooth optimization

$$\min \{f(x) : x \in Q\}$$

Assume

- $Q \subseteq E$  is a “simple” compact convex set
- $f$  is convex and Lipschitz continuous with constant  $M$
- $f$  is not necessarily smooth
- At each  $x \in Q$  can compute some  $g \in \partial f(x)$

20

### Subgradient scheme

**Notation:** Let  $\pi_Q : E \rightarrow Q$  be the projection map  
 $x \mapsto \operatorname{argmin} \{\|y - x\| : y \in Q\}$

### Algorithm 9

- Pick  $x_0 \in Q$
- For  $k = 0, 1, \dots$   
 Compute  $g_k \in \partial f(x_k)$   
 $x_{k+1} := \pi_Q \left( x_k - h_k \frac{g_k}{\|g_k\|} \right)$

**Theorem 10** If  $h_k = h > 0$  then

$$f(x_k) - \bar{f} \leq \frac{M \left( \|x_0 - \bar{x}\|^2 + \sum_{j=0}^k h_j^2 \right)}{\sum_{j=0}^k h_j}.$$

Let  $D := \max \{\|y - x\| : x, y \in Q\}$ .

**Corollary 11** If  $h_k = \frac{D}{\sqrt{N+1}}$  for  $k = 0, 1, \dots, N$  then

$$f(x_N) - \bar{f} \leq \frac{MD}{\sqrt{N+1}}$$

21

22

### Remarks

- Complexity to get within  $\epsilon$  of optimal value is  $O(1/\epsilon^2)$
- This is an *optimal* complexity result: for any “black-box” subgradient-based algorithm there are non-smooth problems that require  $\Omega(1/\epsilon^2)$  iterations to find an  $\epsilon$ -approximate solution.

### References for today's material

- Y. Nesterov, “Introductory Lectures on Convex Optimization,” Kluwer Academic Publishers, 2004.
- Y. Nesterov, “Smooth minimization of non-smooth functions,” Math Program. 103 (2005), 127–152.

23

24