

## Basics of optimization theory



Cong Ma

University of Chicago, Autumn 2021

# Unconstrained optimization

---

Consider an unconstrained optimization problem

$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x})$$

- For simplicity, we assume  $f(\boldsymbol{x})$  is twice differentiable
- We assume the minimizer  $\boldsymbol{x}_{\text{opt}}$  exists, i.e.,

$$\boldsymbol{x}_{\text{opt}} := \arg \min_{\boldsymbol{x}} f(\boldsymbol{x})$$

## (Local) strong convexity and smoothness

---

### Definition 7.1

A twice differentiable function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is said to be  $\alpha$ -strongly convex in a set  $\mathcal{B}$  if for all  $\mathbf{x} \in \mathcal{B}$

$$\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}_n.$$

### Definition 7.2

A twice differentiable function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  is said to be  $\beta$ -smooth in a set  $\mathcal{B}$  if for all  $\mathbf{x} \in \mathcal{B}$

$$\|\nabla^2 f(\mathbf{x})\| \leq \beta.$$

# Gradient descent theory revisited

---

Gradient descent method with step size  $\eta > 0$

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t)$$

## Lemma 7.3

*Suppose  $f$  is  $\alpha$ -strongly convex and  $\beta$ -smooth in the local ball  $\mathcal{B}_\delta(\mathbf{x}_{\text{opt}}) := \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2 \leq \delta\}$ . Running gradient descent from  $\mathbf{x}^0 \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$  with  $\eta = 1/\beta$  achieves linear convergence*

$$\|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right)^t \|\mathbf{x}^0 - \mathbf{x}_{\text{opt}}\|_2, \quad t = 0, 1, 2, \dots$$

# Implications

---

- Condition number  $\beta/\alpha$  determines rate of convergence
- Attains  $\varepsilon$ -accuracy (i.e.,  $\|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2 \leq \varepsilon \|\mathbf{x}_{\text{opt}}\|_2$ ) within

$$O\left(\frac{\beta}{\alpha} \log \frac{1}{\varepsilon}\right)$$

iterations

- Needs initialization  $\mathbf{x}^0 \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$

## Proof of Lemma 7.3

---

Since  $\nabla f(\mathbf{x}_{\text{opt}}) = \mathbf{0}$ , we can rewrite GD as

$$\begin{aligned}\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}} &= \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) - [\mathbf{x}_{\text{opt}} - \eta \nabla f(\mathbf{x}_{\text{opt}})] \\ &= \left[ \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) d\tau \right] (\mathbf{x}^t - \mathbf{x}_{\text{opt}}),\end{aligned}$$

where  $\mathbf{x}(\tau) := \mathbf{x}_{\text{opt}} + \tau(\mathbf{x}^t - \mathbf{x}_{\text{opt}})$ . By local strong convexity and smoothness, one has

$$\alpha \mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}(\tau)) \preceq \beta \mathbf{I}_n, \quad \text{for all } 0 \leq \tau \leq 1$$

Therefore  $\eta = 1/\beta$  yields

$$\mathbf{0} \preceq \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) d\tau \preceq \left(1 - \frac{\alpha}{\beta}\right) \mathbf{I}_n,$$

which further implies

$$\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2$$

# Regularity condition

---

More generally, for update rule

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t),$$

where  $\mathbf{g}(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}^n$

## Definition 7.4

$\mathbf{g}(\cdot)$  is said to obey  $\text{RC}(\mu, \lambda, \delta)$  for some  $\mu, \lambda, \delta > 0$  if

$$2\langle \mathbf{g}(\mathbf{x}), \mathbf{x} - \mathbf{x}_{\text{opt}} \rangle \geq \mu \|\mathbf{g}(\mathbf{x})\|_2^2 + \lambda \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 \quad \forall \mathbf{x} \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$$

- Negative search direction  $\mathbf{g}$  is positively correlated with error  $\mathbf{x} - \mathbf{x}_{\text{opt}} \implies$  one-step improvement
- $\mu\lambda \leq 1$  by Cauchy-Schwarz

## RC = one-point strong convexity + smoothness

---

- One-point  $\alpha$ -strong convexity:

$$f(\mathbf{x}_{\text{opt}}) - f(\mathbf{x}) \geq \langle \nabla f(\mathbf{x}), \mathbf{x}_{\text{opt}} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 \quad (7.1)$$

- $\beta$ -smoothness:

$$\begin{aligned} f(\mathbf{x}_{\text{opt}}) - f(\mathbf{x}) &\leq f\left(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})\right) - f(\mathbf{x}) \\ &\leq \left\langle \nabla f(\mathbf{x}), -\frac{1}{\beta} \nabla f(\mathbf{x}) \right\rangle + \frac{\beta}{2} \left\| \frac{1}{\beta} \nabla f(\mathbf{x}) \right\|_2^2 \\ &= -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned} \quad (7.2)$$



# RC = one-point strong convexity + smoothness

Combining (7.1) and (7.2) yields

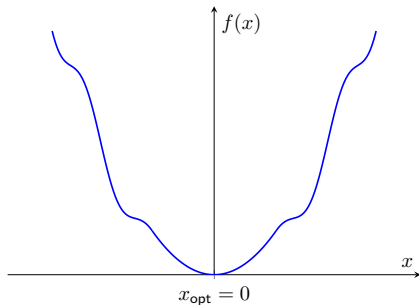
$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_{\text{opt}} \rangle \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 + \frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

— *RC holds with  $\mu = 1/\beta$  and  $\lambda = \alpha$*

# Extension of convex functions

---

When  $\mathbf{g}(x) = \nabla f(x)$ ,  $f$  is not necessarily convex



$$f(x) = \begin{cases} x^2, & |x| \leq 6, \\ x^2 + 1.5|x|(\cos(|x| - 6) - 1), & |x| > 6 \end{cases}$$

# Convergence under RC

## Lemma 7.5

Suppose  $g(\cdot)$  obeys  $\text{RC}(\mu, \lambda, \delta)$ . The update rule  $(x^{t+1} = x^t - \eta g(x^t))$  with  $\eta = \mu$  and  $x^0 \in \mathcal{B}_\delta(x_{\text{opt}})$  obeys

$$\|x^t - x_{\text{opt}}\|_2^2 \leq (1 - \mu\lambda)^t \|x^0 - x_{\text{opt}}\|_2^2$$

- $g(\cdot)$ : more general search directions
  - example: in vanilla GD,  $g(x) = \nabla f(x)$
- The product  $\mu\lambda$  determines the rate of convergence
- Attains  $\varepsilon$ -accuracy within  $O(\frac{1}{\mu\lambda} \log \frac{1}{\varepsilon})$  iterations

## Proof of Lemma 7.5

---

By definition, one has

$$\begin{aligned}\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2^2 &= \|\mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t) - \mathbf{x}_{\text{opt}}\|_2^2 \\&= \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \eta^2 \|\mathbf{g}(\mathbf{x}^t)\|_2^2 - 2\eta \langle \mathbf{g}(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}_{\text{opt}} \rangle \\&\leq \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \eta^2 \|\mathbf{g}(\mathbf{x}^t)\|_2^2 - \eta \left( \lambda \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \mu \|\mathbf{g}(\mathbf{x}^t)\|_2^2 \right) \\&= (1 - \eta\lambda) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \eta(\eta - \mu) \|\mathbf{g}(\mathbf{x}^t)\|_2^2 \\&\leq (1 - \eta\mu) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2\end{aligned}$$