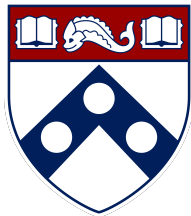


## Smoothing for nonsmooth optimization



Yuxin Chen

Wharton Statistics & Data Science, Fall 2023

# Outline

---

- Smoothing
- Smooth approximation
- Algorithm and convergence analysis

# Nonsmooth optimization

---

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

where  $f$  is convex but not always differentiable

- subgradient methods yield  $\varepsilon$ -accuracy in

$$O\left(\frac{1}{\varepsilon^2}\right) \text{ iterations}$$

- in contrast, if  $f$  is smooth, then accelerated GD yields  $\varepsilon$ -accuracy in

$$O\left(\frac{1}{\sqrt{\varepsilon}}\right) \text{ iterations}$$

— *significantly better than the nonsmooth case*

# Lower bound

---

— Nemirovski & Yudin '83

If one only has access to the first-order oracle (which takes as inputs a point  $x$  and outputs a subgradient of  $f$  at  $x$ ), then one cannot improve upon  $O(\frac{1}{\epsilon^2})$  in general

black box model

# Nesterov's smoothing idea

---

Practically, we rarely meet pure black box models; rather, we know something about the structure of the underlying problems

One possible strategy is:

1. approximate the nonsmooth objective by a smooth function
2. optimize the smooth approximation instead (using, e.g., Nesterov's accelerated method)

## **Smooth approximation**

# Smooth approximation

---

A convex function  $f$  is called  $(\alpha, \beta)$ -smoothable if, for any  $\mu > 0$ ,  $\exists$  convex function  $f_\mu$  s.t.

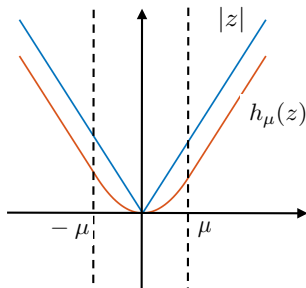
- $f_\mu(\mathbf{x}) \leq f(\mathbf{x}) \leq f_\mu(\mathbf{x}) + \beta\mu, \forall \mathbf{x}$  (approximation accuracy)
- $f_\mu$  is  $\frac{\alpha}{\mu}$ -smooth (smoothness)

—  $\mu$ : tradeoff between approximation accuracy and smoothness

Here,  $f_\mu$  is called a  $\frac{1}{\mu}$ -smooth approximation of  $f$  with parameters  $(\alpha, \beta)$

## Example: $\ell_1$ norm

---



Consider the Huber function

$$h_\mu(z) = \begin{cases} z^2/2\mu, & \text{if } |z| \leq \mu \\ |z| - \mu/2, & \text{else} \end{cases}$$

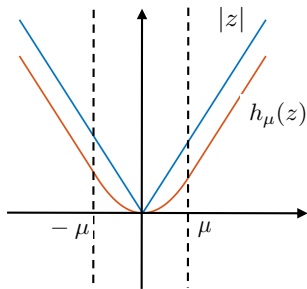
which satisfies

$$h_\mu(z) \leq |z| \leq h_\mu(z) + \mu/2 \quad \text{and} \quad h_\mu(z) \text{ is } \frac{1}{\mu}\text{-smooth}$$



## Example: $\ell_1$ norm

---



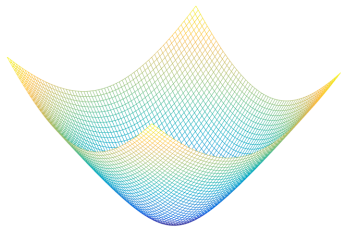
Therefore,  $f_\mu(\mathbf{x}) := \sum_{i=1}^n h_\mu(x_i)$  is  $\frac{1}{\mu}$ -smooth and obeys

$$f_\mu(\mathbf{x}) \leq \|\mathbf{x}\|_1 \leq f_\mu(\mathbf{x}) + \frac{n\mu}{2}$$

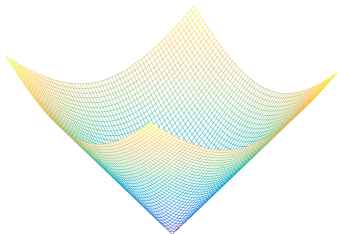
$$\implies \|\cdot\|_1 \text{ is } (1, n/2)\text{-smoothable}$$

## Example: $\ell_2$ norm

---



$f_\mu(\mathbf{x})$



$\|\mathbf{x}\|_2$

Consider  $f_\mu(\mathbf{x}) := \sqrt{\|\mathbf{x}\|_2^2 + \mu^2} - \mu$ , then for any  $\mu > 0$  and any  $\mathbf{x} \in \mathbb{R}^n$ ,

$$f_\mu(\mathbf{x}) \leq (\|\mathbf{x}\|_2 + \mu) - \mu = \|\mathbf{x}\|_2$$

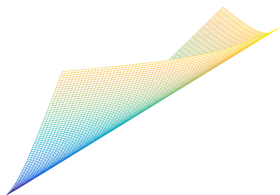
$$\|\mathbf{x}\|_2 \leq \sqrt{\|\mathbf{x}\|_2^2 + \mu^2} = f_\mu(\mathbf{x}) + \mu$$

In addition,  $f_\mu(\mathbf{x})$  is  $\frac{1}{\mu}$ -smooth (exercise)

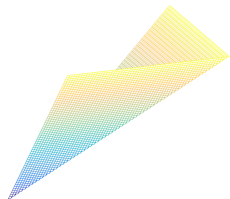
Therefore,  $\|\cdot\|_2$  is (1,1)-smoothable

## Example: max function

---



$$f_\mu(\mathbf{x})$$



$$\max\{x_1, x_2\}$$

Consider  $f_\mu(\mathbf{x}) := \mu \log \left( \sum_{i=1}^n e^{x_i/\mu} \right) - \mu \log n$ , then  $\forall \mu > 0$  and  $\forall \mathbf{x} \in \mathbb{R}^n$ ,

$$f_\mu(\mathbf{x}) \leq \mu \log \left( n \max_i e^{x_i/\mu} \right) - \mu \log n = \max_i x_i$$

$$\max_i x_i \leq \mu \log \left( \sum_{i=1}^n e^{x_i/\mu} \right) = f_\mu(\mathbf{x}) + \mu \log n$$

In addition,  $f_\mu(\mathbf{x})$  is  $\frac{1}{\mu}$ -smooth (exercise). Therefore,  $\max_{1 \leq i \leq n} x_i$  is  $(1, \log n)$ -smoothable

## Basic rules: addition

---

- $f_{\mu,1}$  is a  $\frac{1}{\mu}$ -smooth approximation of  $f_1$  with parameters  $(\alpha_1, \beta_1)$
- $f_{\mu,2}$  is a  $\frac{1}{\mu}$ -smooth approximation of  $f_2$  with parameters  $(\alpha_2, \beta_2)$

$\implies \lambda_1 f_{\mu,1} + \lambda_2 f_{\mu,2}$  ( $\lambda_1, \lambda_2 > 0$ ) is a  $\frac{1}{\mu}$ -smooth approximation of  $\lambda_1 f_1 + \lambda_2 f_2$  with parameters  $(\lambda_1 \alpha_1 + \lambda_2 \alpha_2, \lambda_1 \beta_1 + \lambda_2 \beta_2)$

## Basic rules: affine transformation

---

- $h_\mu$  is a  $\frac{1}{\mu}$ -smooth approximation of  $h$  with parameters  $(\alpha, \beta)$
- $f(\mathbf{x}) := h(\mathbf{A}\mathbf{x} + \mathbf{b})$

$\implies h_\mu(\mathbf{A}\mathbf{x} + \mathbf{b})$  is a  $\frac{1}{\mu}$ -smooth approximation of  $f$  with parameters  $(\alpha\|\mathbf{A}\|^2, \beta)$

## Example: $\|\mathbf{Ax} + \mathbf{b}\|_2$

---

Recall that  $\sqrt{\|\mathbf{x}\|_2^2 + \mu^2} - \mu$  is a  $\frac{1}{\mu}$ -smooth approximation of  $\|\mathbf{x}\|_2$  with parameters  $(1, 1)$

One can use the basic rule to show that

$$f_\mu(\mathbf{x}) = \sqrt{\|\mathbf{Ax} + \mathbf{b}\|_2^2 + \mu^2} - \mu$$

is a  $\frac{1}{\mu}$ -smooth approximation of  $\|\mathbf{Ax} + \mathbf{b}\|_2$  with parameters  $(\|\mathbf{A}\|^2, 1)$

## Example: $|x|$

---

Rewrite  $|x| = \max\{x, -x\}$ , or equivalently,

$$|x| = \max \{ \mathbf{A}x \} \quad \text{with } \mathbf{A} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

Recall that  $\mu \log(e^{x_1/\mu} + e^{x_2/\mu}) - \mu \log 2$  is a  $\frac{1}{\mu}$ -smooth approximation of  $\max\{x_1, x_2\}$  with parameters  $(1, \log 2)$

One can then invoke the basic rule to show that

$$f_\mu(x) := \mu \log(e^{x/\mu} + e^{-x/\mu}) - \mu \log 2$$

is  $\frac{1}{\mu}$ -smooth approximation of  $|x|$  with parameters  $(\|\mathbf{A}\|^2, \log 2) = (2, \log 2)$

# Smoothing via the Moreau envelope

---

The Moreau envelope (or Moreau-Yosida regularization) of a convex function  $f$  with parameter  $\mu > 0$  is defined as

$$M_{\mu f}(\mathbf{x}) := \inf_{\mathbf{z}} \left\{ f(\mathbf{z}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{z}\|_2^2 \right\}$$

- $M_{\mu f}$  is a smoothed or regularized form of  $f$
- minimizers of  $f =$  minimizers of  $M_f$   
 $\implies$  minimizing  $f$  and minimizing  $M_f$  are equivalent



## Connection with the proximal operator

---

- $\text{prox}_f(\mathbf{x})$  is the unique point that achieves the infimum that defines  $M_f$ , i.e.

$$M_f(\mathbf{x}) = f(\text{prox}_f(\mathbf{x})) + \frac{1}{2}\|\mathbf{x} - \text{prox}_f(\mathbf{x})\|_2^2$$

- $M_f$  is continuously differentiable with gradients (homework)

$$\nabla M_{\mu f}(\mathbf{x}) = \frac{1}{\mu}(\mathbf{x} - \text{prox}_{\mu f}(\mathbf{x}))$$

This means

$$\underbrace{\text{prox}_{\mu f}(\mathbf{x}) = \mathbf{x} - \mu \nabla M_{\mu f}(\mathbf{x})}_{\text{prox}_{\mu f}(\mathbf{x}) \text{ is the gradient step for minimizing } M_{\mu f}}$$

# Properties of the Moreau envelope

---

$$M_{\mu f}(\mathbf{x}) := \inf_z \left\{ f(z) + \frac{1}{2\mu} \|\mathbf{x} - z\|_2^2 \right\}$$

- $M_{\mu f}$  is convex (homework)
- $M_{\mu f}$  is  $\frac{1}{\mu}$ -smooth (homework)
- If  $f$  is  $L_f$ -Lipschitz, then  $M_{\mu f}$  is a  $\frac{1}{\mu}$ -smooth approximation of  $f$  with parameters  $(1, L_f^2/2)$

## Proof of smoothability

---

To begin with,

$$M_{\mu f}(\mathbf{x}) \leq f(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{x} - \mathbf{x}\|_2^2 = f(\mathbf{x})$$

In addition, let  $\mathbf{g}_x \in \partial f(\mathbf{x})$ , which obeys  $\|\mathbf{g}_x\|_2 \leq L_f$ . Hence,

$$\begin{aligned} M_{\mu f}(\mathbf{x}) - f(\mathbf{x}) &= \inf_z \left\{ f(\mathbf{z}) - f(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{x}\|_2^2 \right\} \\ &\geq \inf_z \left\{ \langle \mathbf{g}_x, \mathbf{z} - \mathbf{x} \rangle + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{x}\|_2^2 \right\} \\ &= -\frac{\mu}{2} \|\mathbf{g}_x\|_2^2 \geq -\frac{L_f^2}{2} \mu \end{aligned}$$

These together with the smoothness condition of  $M_f$  demonstrate that  $M_f$  is a  $\frac{1}{\mu}$ -smooth approximation of  $f$  with parameters  $(1, L_f^2/2)$

# Smoothing via conjugation

---

Suppose  $f = g^*$ , namely,

$$f(\mathbf{x}) = \sup_z \{\langle \mathbf{z}, \mathbf{x} \rangle - g(\mathbf{z})\}$$

One can build a smooth approximation of  $f$  by adding a strongly convex component to its dual, namely,

$$f_\mu(\mathbf{x}) = \sup_z \{\langle \mathbf{z}, \mathbf{x} \rangle - g(\mathbf{z}) - \mu d(\mathbf{z})\} = (g + \mu d)^*(\mathbf{x})$$

for some 1-strongly convex and continuous function  $d \geq 0$  (called **proximity function**)

# Smoothing via conjugation

---

2 properties:

- $g + \mu d$  is  $\mu$ -strongly convex  $\implies f_\mu$  is  $\frac{1}{\mu}$ -smooth
- $f_\mu(\mathbf{x}) \leq f(\mathbf{x}) \leq f_\mu(\mathbf{x}) + \mu D$  with  $D := \sup_{\mathbf{x}} d(\mathbf{x})$

$\implies f_\mu$  is a  $\frac{1}{\mu}$ -smooth approximation of  $f$  with parameters  $(1, D)$

## Example: $|x|$

---

Recall that

$$|x| = \sup_{|z| \leq 1} zx$$

If we take  $d(z) = \frac{1}{2}z^2$ , then smoothing via conjugation gives

$$f_{\mu}(x) = \sup_{|z| \leq 1} \left\{ zx - \frac{\mu}{2}z^2 \right\} = \begin{cases} x^2/2\mu, & |x| \leq \mu \\ |x| - \mu/2, & \text{else} \end{cases}$$

which is exactly the Huber function

## Example: $|x|$

---

Another way of conjugation:

$$|x| = \sup_{z_1, z_2 \geq 0, z_1 + z_2 = 1} (z_1 - z_2)x$$

If we take  $d(z) = z_1 \log z_1 + z_2 \log z_2 + \log 2$ , then smoothing via conjugation gives

$$f_\mu(x) = \mu \log (\cosh(x/\mu))$$

where  $\cosh x = \frac{e^x + e^{-x}}{2}$

## Example: norm

---

Consider  $\|\mathbf{x}\| = \sup_{\|z\|_* \leq 1} \langle z, \mathbf{x} \rangle$ , then smoothing via conjugation gives

$$f_\mu(\mathbf{x}) = \sup_{\|z\|_* \leq 1} \{\langle z, \mathbf{x} \rangle - \mu d(z)\}$$



# **Algorithm and convergence analysis**

# Algorithm

---

$$\text{minimize}_{\mathbf{x}} \quad F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$$

- $f$  is convex and  $(\alpha, \beta)$ -smoothable
- $h$  is convex but may not be differentiable

# Algorithm

---

Build  $f_\mu$  —  $\frac{1}{\mu}$ -smooth approximation of  $f$  with parameters  $(\alpha, \beta)$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta_t h}(\mathbf{y}^t - \eta_t \nabla f_\mu(\mathbf{y}^t))$$

$$\mathbf{y}^{t+1} = \mathbf{x}^{t+1} + \frac{\theta_t - 1}{\theta_{t+1}}(\mathbf{x}^{t+1} - \mathbf{x}^t)$$

where  $\mathbf{y}^0 = \mathbf{x}^0$ ,  $\theta_0 = 1$  and  $\theta_{t+1} = \frac{1 + \sqrt{1 + 4\theta_t^2}}{2}$

# Convergence

---

## Theorem 8.1 (informal)

Take  $\mu = \frac{\varepsilon}{2\beta}$ . Then one has  $F(\mathbf{x}^t) - F^{\text{opt}} \leq \varepsilon$  for any

$$t \gtrsim \frac{\sqrt{\alpha\beta}}{\varepsilon}$$

- iteration complexity:  $O(1/\varepsilon)$ , which improves upon that of subgradient methods  $O(1/\varepsilon^2)$

## Proof sketch

---

- **convergence rate for smooth problem:** to attain  $\frac{\varepsilon}{2}$ -accuracy for minimizing  $F_\mu(\mathbf{x}) := f_\mu(\mathbf{x}) + h(\mathbf{x})$ , one needs  $O\left(\sqrt{\frac{\alpha}{\mu}} \cdot \frac{1}{\sqrt{\varepsilon}}\right)$  iterations
- **approximation error:** set  $\beta\mu = \frac{\varepsilon}{2}$  to ensure  $|f(\mathbf{x}) - f_\mu(\mathbf{x})| \leq \frac{\varepsilon}{2}$
- since  $F(\mathbf{x}^t) - F(\mathbf{x}^{\text{opt}}) \leq \underbrace{|f(\mathbf{x}^t) - f_\mu(\mathbf{x}^t)|}_{\leq \varepsilon/2} + \underbrace{(F_\mu(\mathbf{x}^t) - F_\mu^{\text{opt}})}_{\leq \varepsilon/2}$ ,

the iteration complexity is

$$O\left(\sqrt{\frac{\alpha}{\mu}} \cdot \frac{1}{\sqrt{\varepsilon}}\right) = O\left(\sqrt{\frac{\alpha\beta}{\varepsilon}} \cdot \frac{1}{\sqrt{\varepsilon}}\right) = O\left(\frac{\sqrt{\alpha\beta}}{\varepsilon}\right)$$

# Reference

---

- "*Smooth minimization of non-smooth functions*," Y. Nesterov, *Mathematical programming*, 2005.
- "*First-order methods in optimization*," A. Beck, Vol. 25, SIAM, 2017.
- "*Optimization methods for large-scale systems, EE236C lecture notes*," L. Vandenberghe, UCLA.
- "*Mathematical optimization, MATH301 lecture notes*," E. Candes, Stanford.
- "*Smoothing and first order methods: A unified framework*," A. Beck, M. Teboulle, *SIAM Journal on Optimization*, 2012.
- "*Problem complexity and method efficiency in optimization*," A. Nemirovski, D. Yudin, Wiley, 1983.