

# OPTIMIZATION FOR MACHINE LEARNING

IASO/MASH - October 17, 2024

Today (Session 9/16): Subgradient methods

Homework: Currently in discussion

⇒ You will hear from us by November 4!

# INTRO

→ So far we have seen (a lot of) gradient descent and its applications

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

↑  
stepsize  $> 0$   
(aka learning rate)

GD iteration for minimizing a  $C^1$  function  $f$

↑  
can be computed using automatic differentiation

→ The rest of the course: alternatives to GD

→ Today: what if the function is not  $C^1$ ?  
what if the function is not differentiable?

Examples:

• Hinge loss / ReLU activation:  $h: t \mapsto \max(t, 0)$

$t < 0$ :  $h$  is differentiable at  $t$  and  $h'(t) = 0$

$t > 0$ : \_\_\_\_\_ at  $t$  and  $h'(t) = 1$

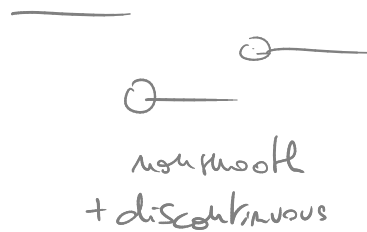
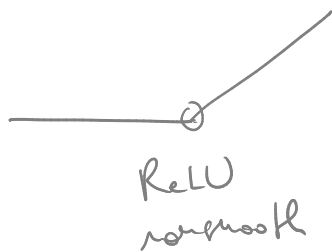
At  $t = 0$ , should we pick  $h'(0) = 0$  or  $h'(0) = 1$ ?

⇒ Actually, you can use any value between 0 and 1

•  $x \mapsto \|x\|$  where  $\|\cdot\|$  is any norm on  $\mathbb{R}^d$  is not differentiable at  $0_{\mathbb{R}^d}$  (that explains why we use  $\|\cdot\|_2^2$  instead of  $\|\cdot\|_2$  in linear regression  $\|x\|_2 = \sqrt{\sum_{i=1}^d x_i^2}$ )

Terminology:

We say that a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is nonsmooth if there exists a point at which it is not differentiable.



Nonsmooth functions are a broad class of functions

→ Focus: nonsmooth convex and continuous functions.

Key reference: Convex analysis (1970)  
R.T. Rockafellar

## ① Subgradients and subdifferentials

Setup:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  convex

Definition: Let  $x \in \mathbb{R}^d$ . A vector  $g \in \mathbb{R}^d$  is called a subgradient of  $f$  at  $x$  if

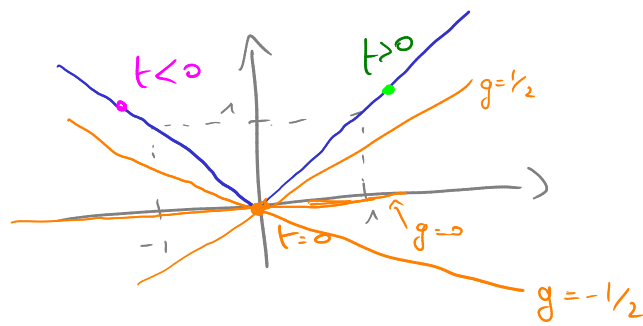
$$\forall y \in \mathbb{R}^d, \quad f(y) \geq f(x) + g^T(y-x)$$

The set of all subgradients of  $f$  at  $x$  is called the subdifferential of  $f$  at  $x$ , and denoted by  $\partial f(x) \subseteq \mathbb{R}^d$

Ex) Consider  $f: t \mapsto |t|$  in  $\mathbb{R}$  ( $d=1$ )

$f$  is not differentiable at 0.

$$f(t) = \begin{cases} t & \text{if } t \geq 0 \\ -t & \text{if } t \leq 0 \end{cases}$$



If  $t > 0$ :

$$g \in \underbrace{\partial f(t)}_{\in \mathbb{R}} \Leftrightarrow \forall u \in \mathbb{R}, f(u) \geq f(t) + g(u-t)$$

$$\Leftrightarrow \forall u \in \mathbb{R}, |u| \geq t + g(u-t)$$

$$\Leftrightarrow \forall u \in \mathbb{R}, |u-t| \geq g(u-t)$$

$$\left. \begin{array}{l} u=t \quad 0 \geq 0 \\ -2t \geq g(-2t) \\ g \geq 1 \\ u=2t \quad t \geq gt \Leftrightarrow 1 \geq g \end{array} \right\} \Rightarrow \begin{cases} \Leftrightarrow \forall u \in \mathbb{R}, \begin{cases} u-t \geq g(u-t) & \text{if } u \geq 0 \\ -u-t \geq g(u-t) & \text{if } u < 0 \end{cases} \\ \Leftrightarrow g=1 \end{cases}$$

$$\partial f(t) = \{1\} = \{f'(t)\}$$

Similarly, for any  $t < 0$ ,  $\partial f(t) = \{-1\} = \{f'(t)\}$

$$\text{If } t=0, \quad g \in \partial f(0) \Leftrightarrow \forall u \in \mathbb{R}, f(u) \geq \overset{0}{f(0)} + g(u-0)$$

$$\Leftrightarrow \forall u \in \mathbb{R}, |u| \geq gu$$

$$\Leftrightarrow \forall u \neq 0, \begin{cases} u \geq gu & \text{if } u > 0 \\ -u \geq gu & \text{if } u < 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} 1 \geq g \\ -1 \leq g \end{cases}$$

thus  $\partial f(0) = [-1, 1]$

Properties of the subdifferential  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  convex

(i) If  $f$  is differentiable at  $x \in \mathbb{R}^d$ , then  $\partial f(x) = \{ \nabla f(x) \}$

(ii) By convexity,  $\partial f(x)$  is always nonempty.

(iii) The converse of (i) is true: If  $\partial f(x) = \{ g \}$ , then  $f$  is differentiable at  $x$ , and  $g = \nabla f(x)$

NB: For  $C^1$  convex  $f$ ,  $f(y) \geq f(x) + \nabla f(x)^T (y-x) \quad \forall y \in \mathbb{R}^d$

$\hookrightarrow$  There exist calculus rules for the subdifferential

Subdifferential calculus

a) For any  $f_1, f_2: \mathbb{R}^d \rightarrow \mathbb{R}$  convex,

$$\forall x \in \mathbb{R}^d, \quad \partial \underbrace{(f_1 + f_2)}_{\text{convex}}(x) = \underbrace{\partial f_1(x) + \partial f_2(x)}_{\{g \in \mathbb{R}^d \mid g = g_1 + g_2, g_1 \in \partial f_1(x), g_2 \in \partial f_2(x)\}}$$

b) For any  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  convex,  $\forall \alpha > 0$

$$\partial \underbrace{(\alpha f)}_{\text{convex}}(x) = \underbrace{\alpha \partial f(x)}_{\{\alpha g \mid g \in \partial f(x)\}}$$

c) For any  $h: \mathbb{R}^m \rightarrow \mathbb{R}$  convex,  $\forall A \in \mathbb{R}^{m \times d}$ ,  $\forall b \in \mathbb{R}^m$ ,

let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $f: x \mapsto h(Ax + b)$

Then  $\forall x \in \mathbb{R}^d$ ,  $\partial f(x) = \underbrace{A^T \partial h(Ax+b)}_{\{A^T g \mid g \in \partial h(Ax+b)\}}$

NB: a) b) c) are generalizations of calculus rules for gradients/derivatives

d) Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  where  $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex

$$f: x \mapsto \max_{1 \leq i \leq m} f_i(x)$$

Then  $f$  is convex and

$$\partial f(x) = \text{conv} \{ \partial f_i(x) : f_i(x) = f(x) \}$$

$$\begin{aligned} \text{conv}(A) &= \text{conv hull of } A \\ &= \{ \alpha x + (1-\alpha)y \mid x \in A, y \in A, \alpha \in [0,1] \} \end{aligned}$$

NB: Property d) has no equivalent for differentiable functions

Example:  $f: x \mapsto \max(a_1^T x + b_1, a_2^T x + b_2)$

$a_1 \in \mathbb{R}^d$   
 $a_2 \in \mathbb{R}^d$   
 $b_1 \in \mathbb{R}, b_2 \in \mathbb{R}$

NB: If  $a_1 = a_2$ ,  $f(x) = a_1^T x + \max(b_1, b_2)$   $a_1 \neq a_2$

$a_1 \neq a_2$

$\hookrightarrow f$  is not differentiable at any  $x$  for which

$$\underbrace{a_1^T x + b_1}_{f_1(x)} = \underbrace{a_2^T x + b_2}_{f_2(x)}$$

By the calculus rules,

$$\partial f(x) = \text{conv} \{ \partial f_1(x), \partial f_2(x) \}$$

$$= \text{conv} \{ a_1, a_2 \} = \{ \alpha a_1 + (1-\alpha)a_2 : \alpha \in [0,1] \}$$

(If  $a_1 = 1, a_2 = -1, b_1 = b_2 = 0$ ,  $d=1$ , this is  $|x|$  !)

(If  $a_1 = 1, b_1 = 0, a_2 = 0, b_2 = 0$ ,  $d=1$ , this is  $\text{ReLU}(x) = \max(x, 0)$ )

## Subgradients and automatic differentiation (AD)

↳ Much more tricky than computing gradients!

→ Subdifferentials are sets, but AD is designed to output one subgradient given an input  $x$ : which subgradient do you get?

→ The subgradient you get depends on the way  $f(x)$  is encoded.

↳ Nice example (Bolte & Pauwels 2020)

"A mathematical model for automatic differentiation in machine learning"

$$\text{ReLU}(t) = \max(t, 0)$$

$$\partial \text{ReLU}(0) = [0, 1]$$

If  $\text{ReLU}(t)$  is encoded as  $\max(-t, 0) + t$ , then AD applied at 0 gives the subgradient  $g=1$

If  $\text{ReLU}(t)$  is encoded as  $\frac{1}{2}(\max(t, 0) + \max(-t, 0)) + t$   
then AD applied at 0 gives the subgradient  $g=0$

### Remark

For nonconvex functions:

→ The subdifferential can be empty ( $t \mapsto -|t|$  at 0)

→ Another subdifferential is used: the Clarke subdifferential, that coincides with the subdifferential above when the function is convex

Exercise: Compute the subdifferential of

$$\|\cdot\|_\infty: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$x \mapsto \|x\|_\infty = \max_{1 \leq i \leq d} |x_i|$$

and  $\|\cdot\|_1: \mathbb{R}^d \rightarrow \mathbb{R}$

$$x \mapsto \|x\|_1 = \sum_{i=1}^d |x_i|$$

## (2) Subgradient methods

Setup: (P) minimize  $f(x)$   $f: \mathbb{R}^d \rightarrow \mathbb{R}$  convex  
 $x \in \mathbb{R}^d$

Recall: when  $f \in C^1$ , we know that  $x^* \in \arg\min f(x)$   
 $\Leftrightarrow \nabla f(x^*) = 0_{\mathbb{R}^d}$

$\rightarrow$  Basis for GD:

If  $\nabla f(x) \neq 0_{\mathbb{R}^d}$ ,  $f$  decreases in the direction of  $-\nabla f(x)$

Q) How can we certify optimality or move towards a better point without gradients?

Theorem:

Consider problem (P) and  $x^* \in \mathbb{R}^d$ .

$$[x^* \in \arg\min_{x \in \mathbb{R}^d} f(x)] \Leftrightarrow [0_{\mathbb{R}^d} \in \partial f(x^*)]$$

Optimality condition for convex not smooth problems



Remark: If  $f \in C^1$ , this optimality condition reduces to

$$\nabla f(x^*) = 0_{\mathbb{R}^d}$$

$$0_{\mathbb{R}^d} \in \nabla f(x^*)$$

$$\Leftrightarrow 0_{\mathbb{R}^d} \in \{\nabla f(x^*)\}$$

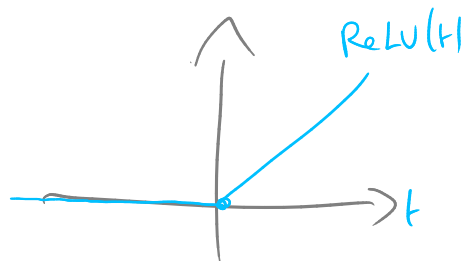
$$\Leftrightarrow 0_{\mathbb{R}^d} = \nabla f(x^*)$$

Ex)  $\text{ReLU}(t) = \max(t, 0)$

$\forall t < 0, \partial \text{ReLU}(t) = \{0\}$  contains 0

$\forall t > 0, \partial \text{ReLU}(t) = \{1\}$  does not contain 0

$\partial \text{ReLU}(0) = [0, 1]$  contains 0



$\text{argmin}_{t \in \mathbb{R}} \text{ReLU}(t) = \{t \leq 0\}$

$\hookrightarrow$  The optimality condition suggests that we can use the subdifferential in a way similar to the gradient in smooth optimization

**Subgradient method (general form)**

Initialization:  $x_0 \in \mathbb{R}^d$

Iteration k:  $x_{k+1} = x_k - \alpha_k g_k$ , where  $\alpha_k > 0$   
and  $g_k \in \partial f(x_k)$

Key challenge: Choosing the subgradient  $g_k$ !

$\rightarrow$  If  $f$  is differentiable at  $x_k$ , necessarily  $g_k = \nabla f(x_k)$

$\Rightarrow$  GD iteration

$\rightarrow$  Otherwise, the method needs a way to choose  $g_k$

• there are good theoretical choices

• but in practice  $g_k$  is fixed by code of  $f$  + AD tool

Ex)  $f: \mathbb{R} \rightarrow \mathbb{R}$   
 $x \mapsto |x|$

$x_0 = 0$

Any subgradient other than 0 ( $g \in [-1, 1], g \neq 0$ ) is an ascent direction for  $f$

$\forall \alpha > 0, f(x_0 - \alpha g) > f(x_0) = 0$

$\forall g \in \partial f(x_0), g \neq 0$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $x \mapsto \|x\|_1 = \sum_{i=1}^d |x_i|$

$\partial f\left(\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}\right) = \left\{ e_1 + \sum_{i=2}^d t_i e_i, t_i \in [-1, 1] \right\}$

$e_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$   $i^{\text{th}}$  row

$f(e_1) = 1$

Any vector  $g \in \partial f(e_1)$  with  $g = e_1 + \sum_{i=2}^d t_i e_i$

$\alpha > 0$   
 $f(e_1 - \alpha g) = |1 - \alpha| + \sum_{i=2}^d \alpha |t_i|$   
 $> |1 - \alpha| + \alpha$

$\sum_{i=2}^d |t_i| > 1$

is an ascent direction

if  $\alpha \in (0, 1]$   $f(e_1 - \alpha g) > 1 = f(e_1)$   
 $\alpha > 1$   $f(e_1 - \alpha g) > \alpha - 1 + \alpha = 2\alpha - 1 > \alpha > 1 = f(e_1)$

Proposition: Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  convex and let  $x \in \mathbb{R}^d$  such that  $0 \notin \partial f(x)$ .

Then, for any  $g_{\min} \in \underset{g \in \mathbb{R}^d}{\operatorname{argmin}} \{ \|g\|^2 \mid g \in \partial f(x) \}$  "minimum norm subgradient"

$g_{\min}$  defines a descent direction, i.e. that

$f(x - \alpha g_{\min}) < f(x)$  for sufficiently small  $\alpha$ .

↳ The proposition suggests to use  $g_k \in \underset{g \in \mathbb{R}^d}{\operatorname{argmin}} \{ \|g\|^2 \mid g \in \partial f(x_k) \}$  in the subgradient method...

↳ ... but this involves solving an auxiliary optimization problem, which is not tractable in general, especially when there is no explicit formula for the subdifferential

↳ What can we prove when the subgradient method uses an arbitrary subgradient at every iteration?

→  $f$  is not guaranteed to decrease at every iteration (unlike in GD, where we can guarantee that with a good stepsize)

→ The quantity of interest for convergence of subgradient methods cannot be  $f(x_k) - \min_{x \in \mathbb{R}^d} f(x)$

### Theorem

Consider the subgradient method applied to (P) and suppose that  $x_k$  ( $k$ th iterate) is not a minimum ( $\Leftrightarrow 0_{\mathbb{R}^d} \notin \partial f(x_k)$ )

Let  $x^* \in \arg \min_{x \in \mathbb{R}^d} f(x)$ . Then for any  $g_k \in \partial f(x_k)$ , there

exists  $\alpha > 0$  such that

$\alpha$  depends on  $g_k, x_k$  and  $x^*$

$$\|x_k - \alpha g_k - x^*\|^2 < \|x_k - x^*\|^2$$

distance between  $x_k - \alpha g_k$  and  $x^*$ 
distance between  $x_k$  and  $x^*$

Proof For any  $\alpha > 0$ ,

$$\|x_k - \alpha g_k - x^*\|^2 = \|(x_k - x^*) - \alpha g_k\|^2 = \|x_k - x^*\|^2 - 2\alpha g_k^T (x_k - x^*) + \alpha^2 \|g_k\|^2$$

Since  $g_k \in \partial f(x_k)$ , we have

$$f(y) \geq f(x_k) + g_k^T (y - x_k) \quad \forall y \in \mathbb{R}^d$$

hence

$$f(x^*) \geq f(x_k) + g_k^T (x^* - x_k)$$

$$\Leftrightarrow -g_k^T (x_k - x^*) \leq f(x^*) - f(x_k)$$

Therefore,

$$\|x_k - \alpha g_k - x^*\|^2 \leq \|x_k - x^*\|^2 + 2\alpha \underbrace{(f(x^*) - f(x_k))}_{< 0 \text{ because } x_k \text{ is not a minimum}} + \alpha^2 \|g_k\|^2$$

$$= \|x_k - x^*\|^2 - 2\alpha (f(x_k) - f(x^*)) + \alpha^2 \|g_k\|^2$$

For any  $\alpha \in (0, \frac{2(f(x_k) - f(x^*)))}{\|g_k\|^2})$ , we have

$$\begin{aligned} -2\alpha (f(x_k) - f(x^*)) + \alpha^2 \|g_k\|^2 &< -2\alpha (f(x_k) - f(x^*)) \\ &+ \alpha \|g_k\|^2 \times \left( \frac{2(f(x_k) - f(x^*))}{\|g_k\|^2} \right) \\ &= -2\alpha (f(x_k) - f(x^*)) \\ &+ 2\alpha (f(x_k) - f(x^*)) \\ &= 0 \end{aligned}$$

Thus, for any  $\alpha \in (0, \frac{2(f(x_k) - f(x^*)))}{\|g_k\|^2})$ ,  
 $0 \notin \partial f(x_k)$

$$\|x_k - \alpha g_k - x^*\|^2 < \|x_k - x^*\|^2$$

In practice, one can use:

- predefined constant stepsizes during a fixed number of iterations
- predefined sequence of decreasing stepsizes

For both stepsize strategies, it is possible to obtain convergence rates for the subgradient method.

### Theorem

Suppose that  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex, has a minimum and (for simplicity) that  $f$  has bounded subgradients, i.e.

$$\forall x \in \text{int} \mathcal{D}, \forall g \in \partial f(x), \|g\| \leq M < \infty$$

Let  $x^* \in \text{argmin}_{x \in \mathbb{R}^d} f(x)$

(true if  $f$  is  $M$ -Lipschitz continuous)

## 1) Fixed stepsize result

Run the subgradient method for  $K \geq 1$  iterations with

$$\alpha_k = \frac{\|x_0 - x^*\|}{M\sqrt{k}}$$

) theoretical value that works but you can also prove results for  $\alpha_k = \Theta/\sqrt{k}$

Then

$$f(\bar{x}_K) - f(x^*) \leq \frac{\|x_0 - x^*\| M}{\sqrt{K}}$$

where  $\bar{x}_K = \frac{1}{K} \sum_{k=0}^{K-1} x_k$  (average iterate)

## 2) Decreasing stepsize

Run the subgradient method with  $\alpha_k = \frac{\Theta}{\sqrt{k+1}}$   $k \geq 0$   
with  $\Theta > 0$

Then, after  $K \geq 1$  iterations, we have

$$f(\bar{x}_K) - f(x^0) \leq O\left(\frac{\|x_0 - x^*\|^2}{\sqrt{K}} + \frac{M^2 \log K}{\sqrt{K}}\right)$$

where  $\bar{x}_K = \frac{1}{\sum_{h=0}^{K-1} \alpha_h} \sum_{k=0}^{K-1} \alpha_k x_k$  (weighted average of the iterates)

### Observations:

• The convergence rates apply to the average of the iterates, for which they show that this average behaves more smoothly than the actual iterates

$\Rightarrow$  Similar results for  $\min_{0 \leq k \leq K-1} (f(x_k) - f(x^0))$

"Best-iterate result"

(NB: Computing the best iterate requires to evaluate  $f$ , whereas the average can be computed directly from the iterates)

- The rates  $O\left(\frac{1}{\sqrt{k}}\right)$  for fixed stepsize and  $O\left(\frac{\log k}{\sqrt{k}}\right)$  for decreasing stepsize

are worse than the rate of GD on a convex problem

$O(1/k) \Rightarrow$  nonsmooth convex minimization is harder than smooth convex minimization

- $O\left(\frac{1}{\sqrt{k^2}}\right)$  better than  $O\left(\frac{\log k}{\sqrt{k^2}}\right)$

but the first rate is obtained by fixing  $k$  and running the method with a stepsize that depends on  $k$