

OPTIMIZATION FOR MACHINE LEARNING

October 5, 2023

Today: Last lecture of the first block
Second block: November 8 - November 16

IASD
App

NON SMOOTH OPTIMIZATION AND REGULARIZATION

What can we do when we don't have gradients?

How can we guarantee that a solution has a specific structure?

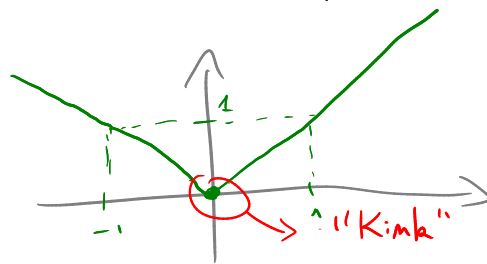
① Optimizing without gradients

minimize $f(w)$ $f: \mathbb{R}^d \rightarrow \mathbb{R}$
 $w \in \mathbb{R}^d$

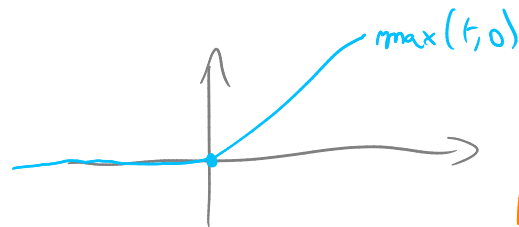
→ So far we have assumed that f was C^1

→ But many problems involve functions that are not C^1 and that do not even have a gradient at certain points

Ex) $t \mapsto |t|$ in \mathbb{R}
 does not have a derivative (i.e. a gradient) at $t=0$



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



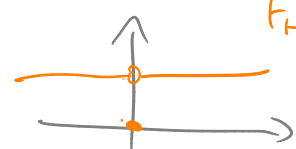
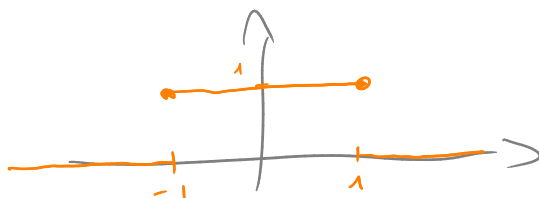
L -smooth = $C_{L,1}^{1,1}$
 smooth $\equiv C^1$

Def (for this course) A function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is **nonsmooth** if there exists a vector $w \in \mathbb{R}^d$ at which f does not have a gradient



A nonsmooth function need not be discontinuous (Ex: $t \mapsto |t|$)

$t \mapsto \begin{cases} 1 & \text{if } |t| \leq 1 \\ 0 & \text{otherwise} \end{cases}$



$t \mapsto \begin{cases} 0 & \text{if } t=0 \\ 1 & \text{otherwise} \end{cases}$

↳ We will focus on convex nonsmooth functions, that are necessarily continuous

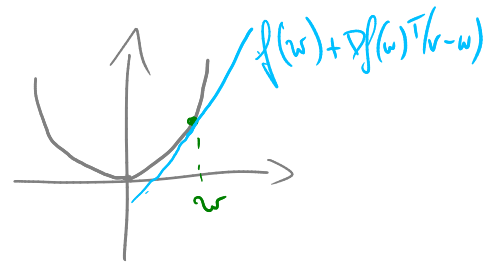
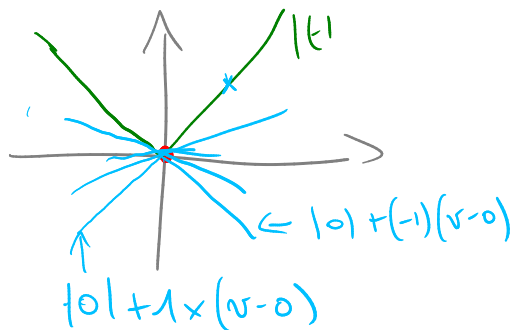
Def. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be convex (hence continuous) function. and let $w \in \mathbb{R}^d$.

A vector $g \in \mathbb{R}^d$ is called a subgradient of f at w

if $\forall v \in \mathbb{R}^d, f(v) \geq f(w) + g^T(v-w)$.

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

↳ For C^1 convex f , we have $f(v) \geq f(w) + \nabla f(w)^T(v-w) \quad \forall v \in \mathbb{R}^d$
In that case, $\partial f(w) = \{ \nabla f(w) \}$.



$$\partial | \cdot | (0) = [-1, 1]$$

$$\partial | \cdot | (1) = \{1\} \Rightarrow \text{If } f \text{ has a gradient at } w, \text{ then } \partial f(w) = \{ \nabla f(w) \}$$

Th. If $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a nonsmooth convex function, then

$$[\bar{w} \in \mathbb{R}^d \text{ is a global minimum of } f] \Leftrightarrow [0_{\mathbb{R}^d} \in \partial f(\bar{w})]$$

↳ Generalization of the C^1 setting, equivalent when f is C^1
($\partial f(\bar{w}) = \{ \nabla f(\bar{w}) \}$ when f is C^1)

↳ Using the notion of subgradient, we can build a class of algorithms for minimizing convex nonsmooth functions. \Rightarrow Subgradient methods

Basic subgradient method iteration:

$$\forall k \geq 0, \quad w_{k+1} = w_k - \alpha_k g_k \quad \text{where } g_k \in \partial f(w_k)$$

$$\alpha_k > 0$$

\Rightarrow This is the GD iteration when $\partial f(w_k) = \{\nabla f(w_k)\}$!

Some results about subgradient methods

↳ Does not always work: depends on the choice of subgradient g_k !

Ex) $f(w) = |w|$ in dimension 1 $w_k = 0 \quad \alpha_k > 0$
 $\partial f(w_k) = [-1, 1]$

If $g_k = 1 \in \partial f(w_k)$, then $w_{k+1} < 0$
 $|w_{k+1}| > |w_k| = 0$

↳ A good strategy for choosing g_k :

$$g_k \in \underset{g \in \mathbb{R}^d}{\operatorname{argmin}} \{ \|g\| \mid g \in \partial f(w_k) \}$$

(+) If $0 \in \partial f(w_k)$, then $g_k = 0$

(-) This is an optimization problem to be solved and it may require to compute the entire subdifferential \Rightarrow possibly expensive

↳ Possible to derive convergence rates for subgradient methods

Recall: For $f \in C^1$ convex $f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{k}\right)$

after K iterations of GD

Ex) For f convex nonsmooth, after $K \gg 1$ iterations of the subgradient method with well-chosen subgradients, then

$$f\left(\frac{1}{\sum_{k=0}^K \alpha_k} \sum_{k=0}^K \alpha_k w_k\right) - \min_{w \in \mathbb{R}^d} f(w) \leq O\left(\frac{1}{\sqrt{K}}\right)$$

slower than \nearrow the GD rate

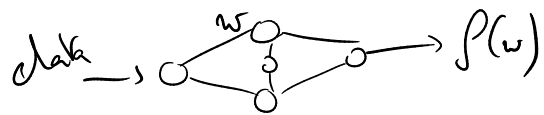
↓
weighted average of the iterates (if $\alpha_k = \alpha > 0$, that is $\frac{1}{K+1} \sum_{k=0}^K w_k$)

Computing gradients (or subgradients)

• Most deep learning libraries rely on automatic differentiation to compute gradients (or subgradients)

• Typical architecture:

→ Forward pass to evaluate $f(w)$



→ Reverse/Backward pass to compute $\nabla f(w)$ without re-accessing the data

Automatic differentiation (AD)
If you have a code that computes $f(w)$, you can apply AD to compute $\nabla f(w)$ at the cost of at most 5 times that of $f(w)$

$$w \mapsto b \mapsto a \mapsto f$$

Key concept: Chain rule

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial b} \frac{\partial b}{\partial w}$$

Example: $d=1$, $f(w) = (1 + w^2 e^{-w})^{1/3}$ Goal: Compute $\nabla f(0) = f'(0)$

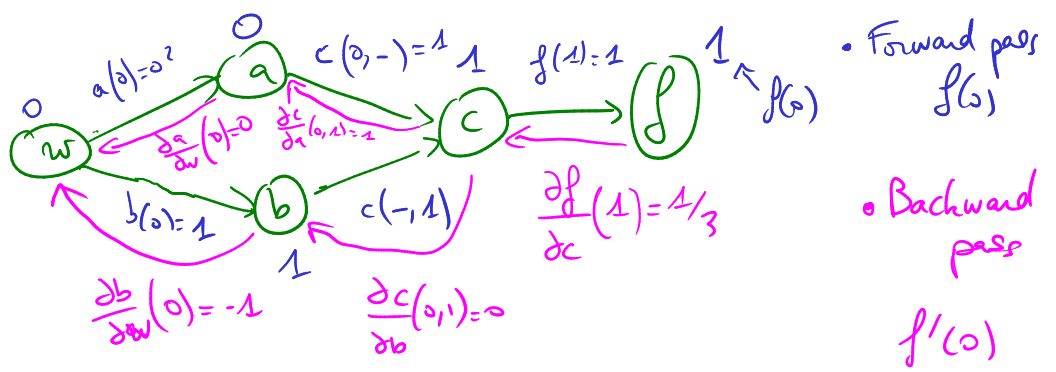
1) Build a computational graph/network to compute $f(w)$

$$a = w^2$$

$$b = e^w$$

$$c = 1 + ab$$

$$f = c^{1/3}$$



→ what you code to evaluate f
 ← generated automatically by automatic differentiation

$$\frac{df}{dw}(0) = \frac{\partial f}{\partial c}(\cdot) \left(\frac{\partial c}{\partial b}(\cdot) \frac{db}{dw}(\cdot) + \frac{\partial c}{\partial a}(\cdot) \frac{da}{dw}(\cdot) \right)$$

Takeaway:

→ there exists software that, given a code to compute $f(w)$, constructs automatically a code to compute $\nabla f(w)$

→ But recently (2021!), researchers found that the way this is done in neural networks does not lead to useful subgradient values for optimization

② Regularization

↳ minimize $f(w)$ where $f: \mathbb{R}^d \rightarrow \mathbb{R}$ depends on some data
 $w \in \mathbb{R}^d$

Ex) $f(w) = \frac{1}{2m} \|Xw - y\|^2$

\uparrow data $X \in \mathbb{R}^{m \times d}$
 \uparrow $y \in \mathbb{R}^m$

↳ So far we have seen examples of data-fitting problems

↳ In ML, the goal is generally not limited to "fit the data"

- Other objectives:
- generalization: model should apply to unseen data
 - simplicity: model should have as few parameters as possible
 - robustness: model should be robust to outliers

⇒ In optimization, one way to incorporate these objectives is through regularization

Regularized optimization problem

minimize $w \in \mathbb{R}^d$ $f(w) + \lambda \Omega(w)$

regularization parameter $\lambda > 0$
 controls the balance between data-fitting and regularization

↑
 data-fitting term, typically a finite sum over data points

(e.g. $\frac{1}{m} \sum_{i=1}^m (x_i^T w - y_i)^2$)

↑
 $\Omega: \mathbb{R}^d \rightarrow \mathbb{R}$
 regularization term penalizes vectors that do not satisfy desirable properties

↳ Ω represents some "soft" constraints on w

minimize $f(w)$
 s.t. $w = 0$
 hard constraints

⇒

minimize w $f(w) + \frac{\lambda}{2} \|w\|^2$
 soft constraint/penalty

Important example

$$\mathcal{L}(w) = \frac{1}{2} \|w\|^2 = \frac{1}{2} \sum_{j=1}^d w_j^2$$

↳ regularization / ridge regularization / Tychonov regularization / etc

Aka batch normalization in the context of stochastic gradient

minimize $f(w) + \frac{\lambda}{2} \|w\|^2$ for $\lambda > 0$
 $w \in \mathbb{R}^d$

- If f convex, then the regularized objective is λ -strongly convex
⇒ 1 unique global minimum
(often helpful for generalization purposes)
- As $\lambda \rightarrow \infty$, the problem becomes closer to minimize $\frac{\lambda}{2} \|w\|^2$
and $\text{argmin}(\) \rightarrow \{0\}$
- More broadly, this regularization penalizes solutions of the un-regularized problem that are large in norm
- It reduces the variance of the solution(s) with respect to the data

Sparsity and regularization

→ Sparse models ($w \in \mathbb{R}^d$ with many zero components) are often desirable because they are simple and even more interpretable

→ One model for sparsity

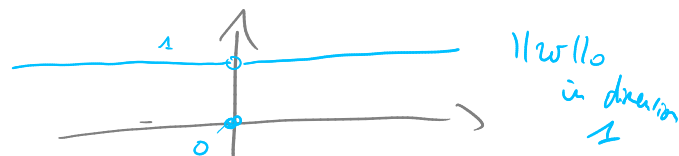
minimize $f(w) + \lambda \|w\|_0$
 $w \in \mathbb{R}^d$

↳ regularizer

$$\|w\|_0 = \left| \{j \in \{1, \dots, d\} \mid w_j \neq 0\} \right|$$

↳ cardinality (number of elements)

$\|w\|_0$ is nonconvex
discontinuous (not smooth)
combinatorial



⇒ Hard to optimize

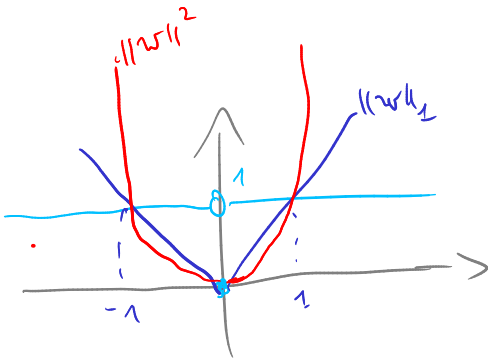
↳ A more tractable formulation uses a convex approximation of $\|w\|_0$ called the l_1 norm

minimize $f(w) + \lambda \|w\|_1$
 $w \in \mathbb{R}^d$

↳ l_1 regularization
 LASSO regularization

continuous
 convex
 non-smooth

$$\|w\|_1 = \sum_{j=1}^d |w_j|$$



• As $\lambda \rightarrow \infty$, solutions to the problem become sparser

• Unlike l_2 regularization, in which all components become small but nonzero as $\lambda \rightarrow \infty$, for l_1 regularization the number of zero components increases

③ Proximal gradient methods for regularized optimization

Setup : minimize $f(w) + \lambda \Omega(w)$ $\lambda > 0$
 $w \in \mathbb{R}^d$

where $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is C^1 (smooth)
 and $\Omega: \mathbb{R}^d \rightarrow \mathbb{R}$ is convex, possibly non-smooth

PG Proximal gradient iteration

$w_{k+1} \in \underset{w \in \mathbb{R}^d}{\text{argmin}}$

$$\left\{ \underbrace{f(w_k) + \nabla f(w_k)^T (w - w_k)}_{\approx f(w) \text{ around } w_k} + \frac{1}{2\alpha_k} \|w - w_k\|^2 + \underbrace{\lambda \Omega(w)}_{\text{left untouched}} \right\}$$

Every iteration of PG solves an auxiliary optimization problem (a subproblem)

Proximal term: penalizes vectors that are far from w_k

→ If $\Omega(w) = 0 \forall w \in \mathbb{R}^d$, then you can show that w_{k+1} is uniquely defined and that $w_{k+1} = w_k - \alpha_k \nabla f(w_k)$ GD

→ The PG subproblem can be solved explicitly for classical choices of Ω regardless of the value of f

Ex) $\Omega(w) = \|w\|_1$

PG subproblem $w_{k+1} \in \arg \min_w \left\{ f(w_k) + \nabla f(w_k)^T (w - w_k) + \frac{1}{2\alpha_k} \|w - w_k\|^2 + \lambda \|w\|_1 \right\}$

has a unique solution defined coordinatewise by

$$\forall j=1..d, \quad [w_{k+1}]_j = \begin{cases} [w_k - \alpha_k \nabla f(w_k)]_j - \lambda \alpha_k & \text{if } [w_k - \alpha_k \nabla f(w_k)]_j > \lambda \alpha_k \\ [w_k - \alpha_k \nabla f(w_k)]_j + \lambda \alpha_k & \text{if } [w_k - \alpha_k \nabla f(w_k)]_j < -\lambda \alpha_k \\ 0 & \text{otherwise} \end{cases}$$

⇒ In that particular case, proximal gradient iterations have a closed form and the method is called ISTA in image and signal processing (Iterative Soft-Thresholding Algorithm)

NB: ISTA + acceleration on convex problems ⇒ Fast ISTA (FISTA)