

OPTIMIZATION FOR MACHINE LEARNING

November 6, 2023

Today: Last Lecture

Thursday: Last year's exam (with Florian Le Biennec)

↳ Course project is online! Deadline: January 15

REGULARIZATION

Motivation: A typical learning problem

minimize $w \in \mathbb{R}^d$

$$\frac{1}{n} \sum_{i=1}^n f_i(w) + \lambda \Omega(w)$$

$\lambda > 0$

$\Omega: \mathbb{R}^d \rightarrow \mathbb{R}$

"Finite sum"
Every f_i depends on a data point in a dataset of size n

"Regularization term" that does not depend on the data

- 1) We have seen how to minimize the finite sum using stochastic gradient when $f_i \in C^1 \forall i$:
 \Rightarrow But what if the functions are not C^1 ?

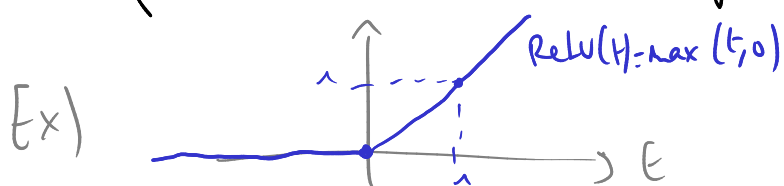
Ex) Neural network with ReLU activations
 $t \mapsto \max(t, 0)$ not C^1

- 2) What is the function Ω used for?
Do we need to take into account in the optimization process?

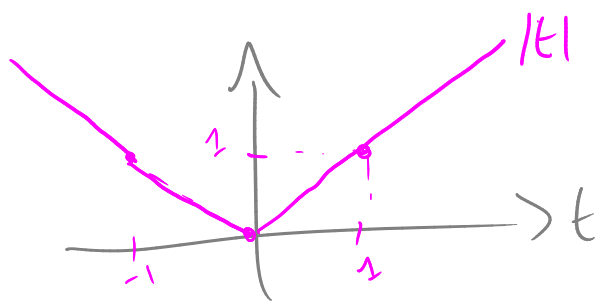
① Optimization without gradients

\hookrightarrow So far all the methods we have seen rely on gradients

\hookrightarrow But many functions used in ML are nonsmooth, i.e. there are points at which these functions do not have a gradient



does not have a gradient at 0



$t \mapsto |t|$ does not have a gradient at 0

\Rightarrow Not having a gradient is an issue for both designing an optimization algorithm and for checking that a point is a minimum

\Rightarrow For convex functions, we can build a tool that plays the role of the gradient in optimization

Def. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex function and $w \in \mathbb{R}^d$. A vector $g \in \mathbb{R}^d$ is called a subgradient of f at w if $\forall z \in \mathbb{R}^d, f(z) \geq f(w) + g^T(z-w)$

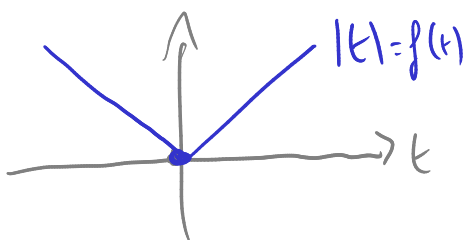
The set of all subgradients for f at w is called the subdifferential of f at w , and denoted by $\partial f(w) \subseteq \mathbb{R}^d$

\hookrightarrow why is that an appropriate definition?

- If f has a gradient at w , then $\partial f(w) = \{ \nabla f(w) \}$
 \Rightarrow More general concept than gradient

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

\Rightarrow Optimality condition based on the subdifferential (can check whether a point is a minimum)



$$\forall t \in \mathbb{R}, \partial f(t) = \begin{cases} \{1\} & \text{if } t > 0 \\ \{-1\} & \text{if } t < 0 \\ [-1, 1] & \text{if } t = 0 \end{cases}$$

Remarks

• Using subgradients instead of gradients, we can generalize all gradient-based algorithms (gradient descent, stochastic gradient, ...) so that apply to convex nonsmooth functions

$$\text{Ex) } \Leftrightarrow w_{k+1} = w_k - \alpha_k \nabla f(w_k) \quad \text{for } f \in C^1$$

$$\text{Subgradient method } w_{k+1} = w_k - \alpha_k g_k \quad g_k \in \partial f(w_k)$$

⚠ The subgradient method is more difficult to implement and analyze

• In practice, gradients and subgradients are computed using automatic differentiation (backward pass in a neural network) and it works

② Regularization terms

↳ In general there are many ways to fit some data, especially with lots of parameters

⇒ minimize $\frac{1}{n} \sum_{i=1}^n f_i(w)$ with d large enough
with $w \in \mathbb{R}^d$ ⇒ there exist many solutions to the problem

↳ Typically we have preferences on the models/parameters we would like to obtain (simplicity/as few parameters as possible, robustness/not too sensitive to data, generalization/model also fits unseen data)

⇒ We can encode these preferences into the objective of the optimization problem using regularization

Two main examples

• l_2 regularization / ridge regularization

$$\text{minimize}_{w \in \mathbb{R}^{d \times d}} \frac{1}{n} \sum_{i=1}^n f_i(w) + \frac{\lambda}{2} \|w\|^2$$

$$\|w\|^2 = \sum_{j=1}^d w_j^2$$

$\lambda > 0$

($\lambda = 0 \Rightarrow$ No regularization)

- As $\lambda \rightarrow +\infty$, the problem gets closer to

$$\text{minimize}_{w \in \mathbb{R}^d} \frac{\lambda}{2} \|w\|^2, \text{ that has solution } \bar{w} = 0$$

- With $\lambda > 0$, the solution of the regularized problem is less sensitive to changes in the data than the solution of the problem without regularization

\Rightarrow We say that l_2 regularization reduces the variance of the solution with respect to the data

- Often implemented in ML for Stochastic Gradient under the name "Weight decay"

(increasing λ reduces the value of the w_i^2)

\Rightarrow Typically improves generalization

• l_1 regularization / LASSO

$$\text{minimize}_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(w) + \lambda \|w\|_1$$

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

$\lambda > 0$

Nonsmooth function
(but convex!)

\hookrightarrow For $\lambda > 0$, the solution of the regularized problem has more zero coordinates than the solution without regularization

sparse
lots of zero coefficients

⇒ We say that l_1 regularization promotes sparse solutions

↳ For linear models, l_1 regularization automatically performs feature selection
⇒ More generally, helps in identifying the most important parameters

Optimizing a regularized problem

$$\underset{w \in \mathbb{R}^d}{\text{minimize}} \quad \underbrace{\frac{1}{n} \sum_{i=1}^n f_i(w)}_{f(w)} + \lambda \Omega(w)$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R} \quad C^1$$

$$\Omega: \mathbb{R}^d \rightarrow \mathbb{R} \quad \text{convex}$$

↳ Proximal gradient techniques exploit the structure of regularized problems

Iteration k:

$$w_{k+1} \in \underset{w \in \mathbb{R}^d}{\text{argmin}} \left\{ \underbrace{f(w) + \nabla f(w_k)^T (w - w_k)}_{\approx f(w)} + \underbrace{\frac{1}{2\alpha_k} \|w - w_k\|^2}_{\alpha_k > 0} + \underbrace{\lambda \Omega(w)}_{\text{regularization (unchanged)}} \right\}$$

"proximal term" (\approx stepsize) ↓

• Every iteration requires to solve an optimization problem ⇒ Only worth doing if the "subproblem" (problem at iteration k) is easier to solve than the original one

• Special cases

• $\alpha \Omega(w) = 0 \quad \forall w \in \mathbb{R}^d$: the iteration becomes

$w_{k+1} = w_k - \alpha_k \nabla f(w_k)$, which is the GD iteration

* $\Omega(w) = \frac{\lambda}{2} \|w\|^2$ (l_2 regularization)

$$w_{k+1} = w_k - \frac{\alpha_k}{1 + \lambda \alpha_k} \nabla f(w_k)$$

"GD with a modified stepsize"

* $\Omega(w) = \|w\|_1$ (l_1 regularization)

→ Explicit formula for w_{k+1}

→ Proximal gradient in that setting corresponds to a method called ISTA from signal processing

NB

Extensions of proximal gradient:

- acceleration
- stochastic
- subgradients