

Tutorial 4: Regularization

Optimization for data science, M2 MIAGE ID/ID Apprentissage

January 21, 2026



Exercise 1: Elastic net

We consider a dataset formed by $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$. Given this dataset, we form a linear regression problem with so-called *elastic net* regularization :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda_2}{2} \|\mathbf{w}\|^2 + \lambda_1 \|\mathbf{w}\|_1, \quad (1)$$

where $\lambda_2 \geq 0$ and $\lambda_1 \geq 0$.

- What is the role of a regularization term in general?
- What is the purpose of the regularization term when $\lambda_1 = 0$ and $\lambda_2 > 0$?
- What is the purpose of the regularization when $\lambda_2 = 0$ and $\lambda_1 > 0$?
- Recall that the gradient of the function $\phi : \mathbf{w} \mapsto \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$ is given by

$$\nabla \phi(\mathbf{w}) = \frac{1}{n} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}).$$

Using this formula, write down the iteration of proximal gradient for problem (1).

- When $\lambda_2 = 0$ and $\lambda_1 > 0$, which algorithm is proximal gradient equivalent to?
- When $\lambda_1 > 0$ and $\lambda_2 > 0$, there does not exist an explicit formula for the proximal gradient iterates, and the proximal subproblem has to be solved approximately at every iteration. Propose an algorithm among those seen in class that could be employed to compute such an approximate solution, and justify your choice of that particular method.

Exercise 2: Reversed Huber loss

In this exercise, we consider the reverse philosophy of the Huber loss, that is, we propose to use a loss function that looks like the absolute value on $[-1, 1]$ and like a quadratic everywhere else.

The *reversed Huber loss* is thus defined as:

$$\begin{aligned} r : \mathbb{R} &\rightarrow \mathbb{R} \\ t &\mapsto r(t) := \begin{cases} |t| & \text{if } |t| < 1 \\ \frac{t^2+1}{2} & \text{otherwise.} \end{cases} \end{aligned} \quad (2)$$

This function is convex but nonsmooth, since it is not differentiable at 0.

As in Exercise 1, we consider linear models $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$ and a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$.

a) We first consider the convex, nonsmooth problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i^T \mathbf{w} - y_i). \quad (3)$$

- i) What mathematical tool can we use to design algorithms applicable to problem (3)?
- ii) Using this tool, how can the solutions of (3) be characterized?

b) We now study the family of problems:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \lambda \sum_{i=1}^d r([\mathbf{w}]_i), \quad (4)$$

where $f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ with every $f_i \in \mathcal{C}^1$ and depending on the data point (\mathbf{x}_i, y_i) , and $\lambda > 0$.

- i) How is this type of problem called? What is the purpose of the second term?
- ii) Write the generic proximal gradient iteration for this problem.
- iii) When is this algorithm worthy of consideration in practice?

Solutions

Solutions for Exercise 1

- a) A regularization term enforces a desired structure on the optimization variables.
- b) When $\lambda_1 = 0$ and $\lambda_2 > 0$, the regularization term is an ℓ_2 regularization term, that aims at reducing the variance of the solution with respect to the data.
- c) When $\lambda_2 = 0$ and $\lambda_1 > 0$, the regularization term is an ℓ_1 regularization term, that aims at promoting sparse solutions.
- d) The k th iteration of proximal gradient applied to problem (1) is

$$\mathbf{w}_{k+1} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ \phi(\mathbf{w}_k) + \frac{1}{n} (\mathbf{X}\mathbf{w}_k - \mathbf{y})^T \mathbf{X} (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|^2 + \frac{\lambda_2}{2} \|\mathbf{w}\|^2 + \lambda_1 \|\mathbf{w}\|_1 \right\}$$

where $\alpha_k > 0$.

- e) A subgradient algorithm can be used to solve the subproblem, since it is defined even in the presence of a nonsmooth function such as the ℓ_1 norm. The iteration cost would involve computing a subgradient. Alternatively, one may want to apply proximal gradient to this subproblem while treating the ℓ_1 norm as a regularization term. This would correspond to the ISTA method (with a quadratic cost function), and would still be tractable given that the iterations of ISTA are explicitly defined.

Solutions for Exercise 2

a)

- i) Since $\mathbf{w} \mapsto \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i^T \mathbf{w} - y_i)$ is convex, it is possible to define the subdifferential of v at any point: the elements of the subdifferential, called the subgradients, can be used in lieu of the gradient to construct optimization methods for solving problem (3).
- ii) Let $\phi_r : \mathbf{w} \mapsto \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i^T \mathbf{w} - y_i)$. A point $\bar{\mathbf{w}} \in \mathbb{R}^d$ is a global minimum of ϕ_r if and only if

$$\mathbf{0} \in \partial \phi_r(\bar{\mathbf{w}}),$$

where $\partial \phi_r(\cdot)$ denotes the subdifferential of ϕ_r .

b)

- i) Problem (4) is a regularized optimization problem. The goal of the second term, that does not depend on data, is to enforce desired properties for the solution.
- ii) At a point \mathbf{w}_k , the generic proximal gradient iteration (with a generic stepsize α_k) for this problem is:

$$\mathbf{w}_{k+1} \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|_2^2 + \lambda \sum_{i=1}^d r([\mathbf{w}]_i) \right\}.$$

- iii) The proximal gradient algorithm is only interesting when the cost of solving the subproblem is cheaper than that of solving the original problem.