

Linear Models and SVMs

Optimization for Machine Learning — Exercise #2

Monday 24th April, 2023

Part I: Theory

I.1. Linear model example

Exercise I.1 (Polynomial Curve Fitting). Given a set of points and their targets $\{x_i, t_i\}_{i=1}^N$ so that for $i \in [N]$, $x_i \in \mathbb{R}$ and $t_i \in \mathbb{R}$, the *curve fitting problem* is loosely defined as finding a function $f: \mathbb{R} \rightarrow \mathbb{R}$ such that $f(x_i) \approx t_i$ for all $i \in [N]$.

In order to find such a function, we restrict ourselves to a set of parametrized functions \mathcal{F} : each function can be parametrized with a vector $\mathbf{w} \in \mathbb{R}^{d+1}$.

To quantify the problem further, in this exercise, we limit ourselves to *polynomial functions* of degree d for the set \mathcal{F} , and can therefore write

$$f(x, \mathbf{w}) = w_0 + w_1x + \dots + w_dx^d = \sum_{k=0}^d w_k x^k \quad (1)$$

Notice how f is *linear* in \mathbf{w} , the parameter. Such model is called a *linear model*.

With N samples, we defined the loss (or error, or energy) of our parameter as the point-wise square distance between its estimation and the target:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - t_i)^2 \quad (2)$$

1. Is the function E convex in \mathbf{w} ? How to find the optimal parameter \mathbf{w}^* at which the loss is minimum?

2. Compute the gradient $\nabla_{\mathbf{w}} E(\mathbf{w}) = \begin{pmatrix} \partial_{w_0} E(\mathbf{w}) \\ \vdots \\ \partial_{w_d} E(\mathbf{w}) \end{pmatrix} \in \mathbb{R}^{d+1}$.

3. Show that the optimal parameter \mathbf{w}^* satisfies the following system of equation:

$$\forall k \in [d+1], \quad \sum_{j=0}^d A_{kj} w_j^* = T_k,$$

where

$$A_{kj} = \sum_{i=1}^N (x_i)^{k+j}, \quad T_k = \sum_{i=1}^N (x_i)^k t_i. \quad (3)$$

4. Is such a system of equation solvable? When / not?

Answers

1. Each term in E is convex in \mathbf{w} as the composition of a linear function $\mathbf{w} \mapsto f(x, \mathbf{w})$ and the convex functions $z \mapsto (z - t)^2$. The function E is then convex as the sum of convex functions in \mathbf{w} .

Therefore, the optimum \mathbf{w}^* can be found with the first-order optimality condition: \mathbf{w}^* satisfies $\nabla E(\mathbf{w}^*) = 0$.

2. For $k \in [d + 1]$, we write $\partial_k := \partial_{w_k}$. The gradient can then be computed as

$$\partial_k E(\mathbf{w}) = \partial_k \left(\frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - t_i)^2 \right) \quad (4)$$

$$= \sum_{i=1}^N (f(x_i, \mathbf{w}) - t_i) \partial_k f(x_i, \mathbf{w}) \quad (5)$$

$$= \sum_{i=1}^N (f(x_i, \mathbf{w}) - t_i) (x_i)^k, \quad (6)$$

$$= \sum_{i=1}^N \left(\sum_{j=0}^d w_j (x_i)^j - t_i \right) (x_i)^k, \quad (7)$$

and $\nabla E(\mathbf{w}) = (\partial_k E(\mathbf{w}))_{k \in [d+1]}$.

3. From 1, the minimizer \mathbf{w}^* is found by solving the equation $\nabla E(\mathbf{w}^*) = 0$, i.e., for all $k \in [d + 1]$,

$$\partial_k E(\mathbf{w}^*) \stackrel{!}{=} 0 \implies \sum_{i=1}^N \left(\sum_{j=0}^d w_j^* (x_i)^j (x_i)^k - t_i (x_i)^k \right) = 0 \quad (8)$$

$$\implies \sum_{j=0}^d \underbrace{\sum_{i=1}^N (x_i)^{j+k}}_{A_{kj}} w_j^* = \underbrace{\sum_{i=1}^N (x_i)^k t_i}_{T_k}. \quad (9)$$

4. This linear system can be written as $A\mathbf{w} = T$, where $A = (A_{kj})_{k,j \in [d+1] \times [d+1]} \in \mathbb{R}^{(d+1) \times (d+1)}$, and $T = (T_k)_{k \in [d+1]}$. Therefore, there are several cases possible.

- First case: $\text{rank } A = d + 1$ (i.e. A is invertible). In this case, there is a unique solution $\mathbf{w}^* = A^{-1}T$.

- Second case: $\text{rank } A < d + 1$. In this case,
 - if $T \in \text{span } A$, there is an infinite number of solutions $\mathbf{w}^* = \mathbf{w}_p + \mathbf{w}_k$, with \mathbf{w}_p any particular solution in $(\ker A)^\perp$ and $\mathbf{w}_k \in \ker A$.
 - if $T \notin \text{span } A$, there are 0 solutions.

Remark: Even when A is invertible, it might be ill-conditioned (meaning its inverse is unstable to compute).

It is usual to add a *regularizer* to the objective, penalizing “complex” models. This also can help selecting a model when several models are solutions to the optimization problem.

One of the most common regularizer is the parameter squared-norm: with a *penalizer weight* $\lambda \in \mathbb{R}_+$, the Equation (2) is modified to give

$$E_\lambda(\mathbf{w}) = E(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \frac{1}{2} \sum_{i=1}^N (f(x_i, \mathbf{w}) - t_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2. \quad (10)$$

5. What is the role of λ ?
6. Is E_λ convex?
7. Show that each component of the optimal weight \mathbf{w}_i^* is now found by solving

$$\forall k \in [d + 1], \quad \sum_{j=0}^d A_{kj} w_j + \lambda w_k = T_k,$$

with A_{kj} and T_k defined as in Equation (3).

Answers

5. λ is the regularizer weight, and make the error favour weights with small ℓ_2 norm. The higher the λ , the stronger weights with high ℓ_2 norm are avoided.
6. E_λ is still convex (as long as $\lambda \geq 0$), as a sum of convex functions in \mathbf{w} .
7. The first order optimality condition still applies, and we find, similarly to 3,

$$\begin{aligned} \partial_k E_\lambda(\mathbf{w}^*) \stackrel{!}{=} 0 &\implies \sum_{i=1}^N \sum_{j=0}^d w_j (x_i)^{j+k} - \sum_{i=1}^N t_i (x_i)^k + \frac{\lambda}{2} \partial_k \|\mathbf{w}\|^2 = 0 \\ &\implies \underbrace{\sum_{j=0}^d \sum_{i=1}^N (x_i)^{j+k} w_j}_{A_{kj}} + \lambda w_k = \underbrace{\sum_{i=1}^N t_i (x_i)^k}_{T_k}. \end{aligned}$$

Remark: Now, the system can be written $(A + \lambda I_{d+1})\mathbf{w}^* = T$. Even if A is singular, $(A + \lambda I_{d+1})$ will be invertible, making the problem well-posed.

Matrix expression It is sometimes preferable to deal with vector and matrices, rather than scalar expressions. When the model is linear in \mathbf{w} , it is possible to express it as a *linear product* between a matrix and a vector. The expression in Equation (1) can be thought as a dot product

between w and the vector of powers of x , that define as $\phi(x) := \begin{pmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^d \end{pmatrix}$, so that

$$f(x, \mathbf{w}) = \mathbf{w}^\top \phi(x).$$

Stacking all the N examples in a matrix, and denoting $\phi_i := \phi(x_i)$, we define

$$\Phi := \begin{pmatrix} | & | & \cdots & | \\ \phi_1 & \phi_2 & \cdots & \phi_N \\ | & | & \cdots & | \end{pmatrix} \in \mathbb{R}^{(d+1) \times N}$$

and can therefore compute the model *on the whole dataset* in one expression: $\mathbf{y}(\mathbf{w}) = \Phi^\top \mathbf{w} \in \mathbb{R}^N$. Each entry i of \mathbf{y} corresponds to a different sample x_i . Then, stacking the targets into a vector $\mathbf{t} \in \mathbb{R}^N$, the error function (2) can equivalently written as

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(\mathbf{w}) - \mathbf{t}\|^2,$$

and the regularized error as

$$E_\lambda(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(\mathbf{w}) - \mathbf{t}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

8. Show that $\nabla E_\lambda(\mathbf{w}) = \Phi(\Phi^\top \mathbf{w} - \mathbf{t}) + \lambda \mathbf{w}$, so that \mathbf{w}^* solves the linear equation

$$(\Phi\Phi^\top + \lambda I_{d+1})\mathbf{w}^* = \Phi\mathbf{t}.$$

Answer

8. Since the unregularized error E is obtained with $\lambda = 0$, we can focus on the regularized objective E_λ , for $\lambda \geq 0$.

We compute the gradient using the product rule from Exercise Sheet #1:

$$\begin{aligned} \nabla_{\mathbf{w}} E_\lambda(\mathbf{w}) &= \nabla_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{y}(\mathbf{w}) - \mathbf{t}\|^2 \right) + \nabla_{\mathbf{w}} \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 \right) \\ &= J_{\mathbf{y}}(\mathbf{w})^\top (\mathbf{y}(\mathbf{w}) - \mathbf{t}) + \lambda \mathbf{w} \quad (\mathbf{y}(\mathbf{w}) = \Phi^\top \mathbf{w} \implies J_{\mathbf{y}}(\mathbf{w}) = \Phi^\top) \\ &= \Phi(\Phi^\top \mathbf{w} - \mathbf{t}) + \lambda \mathbf{w}. \end{aligned}$$

Therefore, the first-order optimality condition can be written

$$\nabla E_\lambda(\mathbf{w}^*) = 0 \implies (\Phi\Phi^\top + \lambda I_{d+1})\mathbf{w}^* = \Phi\mathbf{t}.$$

I.2. Subgradients

When a convex loss function $E: \mathbb{R}^d \rightarrow \mathbb{R}$ is not differentiable, its *subgradient* can be used. It is defined as the set, for $x \in \mathbb{R}^d$,

$$\partial E(x) = \{g \in \mathbb{R}^d \mid \forall y \in \mathbb{R}^d, E(y) \geq E(x) + \langle g, y - x \rangle\}.$$

If E is differentiable at x , then $\partial E(x) = \{\nabla E(x)\}$.

For instance, for $E: \mathbb{R} \rightarrow \mathbb{R}$, $x \mapsto E(x) = |x|$, E is differentiable at any $x \neq 0$, with gradient -1 on $(-\infty, 0)$ and 1 on $(0, +\infty)$.

At $x = 0$, we compute, for any $y \in \mathbb{R}$ and $g \in \mathbb{R}$:

$$\begin{aligned} E(y) \geq E(0) + \langle g, y - 0 \rangle &\iff |y| \geq \langle g, y \rangle \\ &\iff |y| \geq gy \end{aligned}$$

This condition has to be true for *any* $y \in \mathbb{R}$. This is only true when $g \in [-1, 1]$. Therefore,

$$\partial E(x) = \begin{cases} \{-1\} & \text{if } x < 0 \\ [-1, 1] & \text{if } x = 0 \\ \{1\} & \text{if } x > 0 \end{cases}$$

Geometrically, this can be interpreted as having, for the absolute value at the origin, any lines with slope between -1 and 1 lower-bounding the graph of the function.

Exercise I.2. Let $E: \mathbb{R} \rightarrow \mathbb{R}$, $x \mapsto E(x) = \max(0, 1 - x)$.

1. Where is E differentiable?

2. Show that $\partial E(x) = \begin{cases} \{-1\} & \text{if } x < 1 \\ [-1, 0] & \text{if } x = 1 \\ \{0\} & \text{if } x > 1 \end{cases}$

Answer

1. E is differentiable on $(-\infty, 1)$ and on $(1, +\infty)$, but not at 1 .

2. Since E is differentiable on $(-\infty, 1)$, $\partial E(x) = \{\nabla E(x)\} = \{-1\}$ for $x \in (-\infty, 1)$. Likewise, $\partial E(x) = \{0\}$ for $x \in (1, +\infty)$.

For $x = 1$, we are looking for $g \in \mathbb{R}$ such that, for all $y \in \mathbb{R}$,

$$E(y) \geq E(1) + g \cdot (y - 1).$$

Assuming such g exists, it necessarily satisfies, for all $y \in \mathbb{R}$,

$$\max(0, 1 - y) \geq g \cdot (y - 1) \tag{1}$$

- For $y = 1$, g can be arbitrary (since we get $0 = g \cdot 0$ in this case, which is true for all g).
- For $y > 1$, the necessary condition become

$$0 \geq g(y - 1) \implies 0 \geq g \quad \text{since } y - 1 > 0$$

- For $y < 1$, the necessary condition becomes

$$1 - y \geq g(y - 1) \implies -1 \leq g$$

Since the condition (1) has to be true for *all* $y \in \mathbb{R}$, g has to satisfy $-1 \leq g \leq 0$.

One verifies that taking $g \in [-1, 0]$ always satisfies (1).

$$\text{Therefore, } \partial E(x) = \begin{cases} \{-1\} & \text{if } x < 1, \\ [-1, 0] & \text{if } x = 1, \\ \{0\} & \text{if } x > 1. \end{cases}$$

Part II: Programming

Exercise II.1 (Model fitting). This exercise implements some results found in Exercise I.1.

1. **Generation of the target.** In this toy example, we generate the N points ourselves. The true target t_i will be sinusoidal, with some noise, i.e. $t_i = \sin(2\pi x_i) + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$. The different scales (for σ, x_i) are given as $\sigma = 0.1$, and $x_i \sim \mathcal{U}([0, 1])$, uniform distribution on the segment $[0, 1]$.

The generation of the data is performed by the function `gen_sin_data` in the file `ex02/Utils.py`.

2. Implement the parametrization function (1) as $\hat{f}(\mathbf{x}, \mathbf{w})$, where the dimensions d is implied by the size of \mathbf{w} .
3. Implement the error function E defined in (2), and its gradient $\nabla E(\mathbf{w})$.
4. Find \mathbf{w}^* , either by
 - a) gradient descent; or
 - b) solving the linear system of equations (3).