

Tutorial 1: Basics of optimization

Optimization for machine learning, M2 MIAAGE ID Apprentissage

September 21, 2023



Exercise 1: Linear least squares

We consider a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, wherein $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$ for every $i = 1, \dots, n$. We seek a linear model that best fits the data, which we formulate as the following optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) := \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{w} - y_i)^2, \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$ are given by

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}.$$

This problem is among the most classical in data analysis. Its objective function is \mathcal{C}^2 , and the problem (1) always has at least one solution.

- a) Let $\mathbf{w}^* \in \mathbb{R}^d$ satisfy $\mathbf{X}\mathbf{w}^* = \mathbf{y}$ (hence \mathbf{w}^* is a solution of the linear system $\mathbf{X}\mathbf{w} = \mathbf{y}$). Justify then that \mathbf{w}^* is a global minimum of the objective function.
- b) The gradient of f at any $\mathbf{w} \in \mathbb{R}^d$ is given by $\nabla f(\mathbf{w}) = \mathbf{X}^T(\mathbf{X}\mathbf{w} - \mathbf{y})$. If \mathbf{w}^* is a local minimum of f , what is the value of $\nabla f(\mathbf{w}^*)$?
- c) The Hessian matrix of f at $\mathbf{w} \in \mathbb{R}^d$ is given by $\nabla^2 f(\mathbf{w}) = \mathbf{X}^T \mathbf{X}$. Note that it is constant with respect to \mathbf{w} , and that it only depends on the data matrix \mathbf{X} .
 - i) By construction, we have $\mathbf{X}^T \mathbf{X} \succeq \mathbf{0}$. What property on f does this imply?
 - ii) Suppose that $\mathbf{X}^T \mathbf{X} \succeq \mu \mathbf{I}_d$ with $\mu > 0$. Given $\mathbf{w} \in \mathbb{R}^d$, what can we say about $\nabla^2 f(\mathbf{w})$ in that case? What information does this provide about the set of solutions of problem (1)?

Exercise 2: Convex function

Let $q : \mathbb{R}^d \rightarrow \mathbb{R}$ be defined as $q(\mathbf{w}) = \frac{1}{4}\|\mathbf{w}\|^4$. This function is \mathcal{C}^2 , and for every $\mathbf{w} \in \mathbb{R}^d$, we have

$$\nabla q(\mathbf{w}) = \|\mathbf{w}\|^2 \mathbf{w}, \quad \nabla^2 q(\mathbf{w}) = 2\mathbf{w}\mathbf{w}^T + \|\mathbf{w}\|^2 \mathbf{I}_d.$$

- Using the expression of the Hessian matrix of q , show that the function q is convex. What does it imply on its local minima?
- Show that the zero vector $\mathbf{0}_{\mathbb{R}^d}$ is a local minimum of q . Does it satisfy the second-order sufficient condition?
- Given the answer to the previous question, can the function q be strongly convex?
- Justify that the function has a single global minimum.

Exercise 3: Quasiconvex functions

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is called **quasiconvex** if

$$\forall \mathbf{w}, \mathbf{v} \in \mathbb{R}^d, \forall t \in [0, 1], \quad f(t\mathbf{w} + (1-t)\mathbf{v}) \leq \max\{f(\mathbf{w}), f(\mathbf{v})\}. \quad (2)$$

Any convex function is quasiconvex, but the converse is not true.

Let f be a quasiconvex, \mathcal{C}^2 function. We consider:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}). \quad (3)$$

- Write the first- and second-order optimality conditions for problem (3).
- Since f is quasiconvex, it can be shown that

$$\forall \mathbf{w} \in \mathbb{R}^d, \forall \mathbf{v} \in \mathbb{R}^d, \quad \mathbf{v}^T \nabla f(\mathbf{w}) = 0 \Rightarrow \mathbf{v}^T \nabla^2 f(\mathbf{w}) \mathbf{v} \geq 0. \quad (4)$$

Let \mathbf{w}^* be a first-order stationary point. Justify that \mathbf{w}^* is also a second-order stationary point.

Solutions

Solutions for Exercise 1

Underlying goal: Introduce least-squares formulations. Apply the definitions of global minima/solutions and that of convexity.

a) If $\mathbf{X}\mathbf{w}^* = \mathbf{y}$, then

$$f(\mathbf{w}^*) = \frac{1}{2} \|\mathbf{X}\mathbf{w}^* - \mathbf{y}\|^2 = \frac{1}{2} \|\mathbf{0}\|^2 = 0.$$

Since f is always nonnegative (definition of a norm), we also have

$$\forall \mathbf{w} \in \mathbb{R}^d, f(\mathbf{w}) \geq 0 = f(\mathbf{w}^*).$$

The latter property corresponds to the definition of a global minimum for f , from which we conclude that \mathbf{w}^* is a global minimum of f or, equivalently, a solution of the unconstrained problem (1).

b) The function f is continuously differentiable (\mathcal{C}^2 , so \mathcal{C}^1). If \mathbf{w}^* is a local minimum of f , then $\nabla f(\mathbf{w}^*) = \mathbf{0}$ per the first-order optimality condition.

i) If $\mathbf{X}^T \mathbf{X} \succeq \mathbf{0}$, then $\nabla^2 f(\mathbf{w}) \succeq \mathbf{0}$ for any $\mathbf{w} \in \mathbb{R}^d$. This property is a characterization of convexity for a \mathcal{C}^2 function, from which we conclude that f is a convex function.

ii) Similarly to the previous question, the fact that $\mathbf{X}^T \mathbf{X} \succeq \mu \mathbf{I}_d$ means that $\nabla^2 f(\mathbf{w}) \succeq \mu \mathbf{I}_d$ for any $\mathbf{w} \in \mathbb{R}^d$. This is again a characterization of strong convexity for \mathcal{C}^2 functions, and therefore f is μ -strongly convex. As a result, there exists a unique solution for the optimization problem (or equivalently, f has a unique global minimum).

Solutions for Exercise 2

Goal: Introduce a bit more calculus to get students comfortable with scalar products and matrix-vector products. Give an example of global minimum that does not satisfy the sufficient optimality condition.

a) For any $\mathbf{w} \in \mathbb{R}^d$ and any $\mathbf{v} \in \mathbb{R}^d$, the linearity of both scalar products and matrix-vector products gives:

$$\begin{aligned} \mathbf{v}^T \nabla^2 q(\mathbf{w}) \mathbf{v} &= \mathbf{v}^T (2\mathbf{w}\mathbf{w}^T + \|\mathbf{w}\|^2 \mathbf{I}_d) \mathbf{v} \\ &= \mathbf{v}^T (2\mathbf{w}\mathbf{w}^T \mathbf{v} + \|\mathbf{w}\|^2 \mathbf{v}) \\ &= 2\mathbf{v}^T \mathbf{w}\mathbf{w}^T \mathbf{v} + \|\mathbf{w}\|^2 \mathbf{v}^T \mathbf{v} \\ &= 2(\mathbf{w}^T \mathbf{v})^2 + \|\mathbf{w}\|^2 \mathbf{v}^T \mathbf{v} \\ &= 2(\mathbf{w}^T \mathbf{v})^2 + \|\mathbf{w}\|^2 \|\mathbf{v}\|^2 \\ &\geq 0. \end{aligned}$$

Thus, for any $\mathbf{w} \in \mathbb{R}^d$, the Hessian matrix $\nabla^2 q(\mathbf{w})$ is positive semidefinite, i.e. $\nabla^2 q(\mathbf{w}) \succeq \mathbf{0}$. Consequently, the (\mathcal{C}^2) function q is convex, and all its local minima are global.

b) Since the function q is convex, every local minimum is global. Moreover, we have

$$q(\mathbf{w}) = \frac{1}{4}\|\mathbf{w}\|^4 \geq 0 = q(\mathbf{0}_{\mathbb{R}^d})$$

for any $\mathbf{w} \in \mathbb{R}^d$. The zero vector $\mathbf{0}_{\mathbb{R}^d}$ is thus a global minimum of q . If the zero vector were to satisfy the second-order sufficient optimality conditions, we would have $\nabla^2 q(\mathbf{0}_{\mathbb{R}^d}) \succ \mathbf{0}$. However, the expression for $\nabla^2 q$ gives

$$\nabla^2 q(\mathbf{0}_{\mathbb{R}^d}) = \mathbf{0},$$

and the zero matrix is only positive semidefinite (instead of positive definite). As a result, the zero vector does not satisfy the second-order sufficient optimality conditions. *Note: This does not contradict the fact that this vector is a global minimum, as the condition is sufficient but not necessary.*

c) If the function were strongly convex, there would exist $\mu > 0$ such that $\nabla^2 q(\mathbf{w}) \succeq \mu \mathbf{I}_d \succ \mathbf{0}$ for any \mathbf{w} , including the zero vector. Since the Hessian is zero at the zero vector, this cannot be true, from which we conclude that q is not strongly convex.

d) For every $\mathbf{w} \in \mathbb{R}^d$, we have $q(\mathbf{w}) \geq q(\mathbf{0}_{\mathbb{R}^d}) = 0$, hence the zero vector is a global minimum. Moreover, $q(\mathbf{w}) = 0$ if and only if $\mathbf{w} = \mathbf{0}_{\mathbb{R}^d}$, and thus the zero vector is the only global minimum of q .

Note: Classical argument in this last question, typical first question of an exam.

Solutions for Exercise 3

a) *The result is expected to be known.* The first-order necessary optimality conditions can be stated as follows. If a vector $\mathbf{w}^* \in \mathbb{R}^d$ is a local minimum of a \mathcal{C}^1 function f , then $\nabla f(\mathbf{w}^*) = \mathbf{0}$. The second-order necessary optimality conditions are a stronger characterization. If $\mathbf{w}^* \in \mathbb{R}^d$ is a local minimum of f , then

$$\nabla f(\mathbf{w}^*) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}.$$

b) Since \mathbf{w}^* is a first-order stationary point, it satisfies the first-order necessary conditions, hence $\nabla f(\mathbf{w}^*) = \mathbf{0}$ and

$$\forall \mathbf{v} \in \mathbb{R}^d, \quad \mathbf{v}^T \nabla f(\mathbf{w}^*) = \mathbf{v}^T \mathbf{0} = 0.$$

The left-hand side of the implication (4) thus holds for \mathbf{w}^* and any vector \mathbf{v} . Thus the right-hand also holds, i.e.

$$\mathbf{v}^T \nabla^2 f(\mathbf{w}^*) \mathbf{v} \geq 0 \quad \forall \mathbf{v} \in \mathbb{R}^d,$$

which is equivalent to $\nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}$. Therefore, the vector \mathbf{w}^* satisfies the second-order necessary optimality conditions, and it is a second-order stationary point.