

Convex optimization

February 9, 2021

Abstract

As part of my machine learning study, I write what I learned about basic optimization problem. Among the optimization problems, I describe about the convex optimization. The convex optimization is the problem of minimizing convex functions over convex sets. Due to the nature of the convex function, the convex optimization is relatively easy to solve the optimal solution of the objective function, and the solution method has been established to some extent. This time, I show about a classifier that makes use of convex optimization such as a log-linear classifier.

Convex optimization

Before writing about the convex optimization, briefly write what the optimization problem is. An optimization problem is a problem of finding valuable values and function value that optimizes (minimize or maximize) a function under certain constraints. And, the function to be optimized is called the objective function. For example, suppose there is an objective function $f(x_1, x_2) = x_1 x_2$ and a constraint $x_1 - x_2 - a = 0$ (a is a constant). Also, assume that the optimization problem minimizes $f(x_1, x_2)$. The formal writing of these is as follows,

$$\begin{cases} \min. & f(x_1, x_2) = x_1 x_2 \\ \text{s.t.} & x_1 - x_2 - a = 0 \end{cases}$$

and s.t. means subject to. Solving this gives $x_1 x_2 = (x_1 - \frac{a}{2})^2 - \frac{a^2}{4}$, and is obtained optimal solutions $-\frac{a^2}{4}$ and $(x_1, x_2) = (\frac{a}{2}, -\frac{a}{2})$ with the smallest $f(x_1, x_2)$. Next, I describe what a convex set and convex function is and those properties.

0.1 Convex set and convex function

First, I write the definition of the convexity of the set on \mathbb{R}^n .

Definition 0.1 (Convex set).

$X \subset \mathbb{R}^n$ is a convex set. $\stackrel{\text{def}}{\iff} \forall \mathbf{x}_1, \mathbf{x}_2 \in X \text{ and } \forall \lambda \in [0, 1] \rightarrow \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in X$

As you can see by drawing the convex set in the figure, the convex set means a set in which the line segment connecting arbitrary points in the convex set does not protrude from the set. Next, although I write about the definition of a convex function, there are two types of convex

functions, before that, I state the definition of epigraph, which is a concept that connects convex sets and convex downward functions.

Definition 0.2 (Epigraph).

Given the real-valued function f , the following set is called the epigraph of f

$$\text{epi } f := \{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R} : y \geq f(\mathbf{x})\}$$

Definition 0.3 (Convex function).

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function $\stackrel{\text{def}}{\iff}$ $\text{epi } f$ is a convex set.*¹

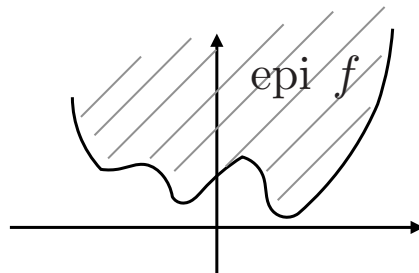


Figure1 epigraph

The epigraph in Figure1 is an epigraph that is not a convex set. Using these definitions (0.1 ~ 0.3), the following theorem holds.

Theorem 0.1.

The fact that the function f on \mathbb{R}^n is a convex function is equal to the following condition : for $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ and $\forall \lambda \in [0, 1]$,

$$\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2)$$

Proof. For $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$ and $\forall \lambda \in [0, 1]$, if there is a $\tilde{\mathbf{x}}$ such that $\tilde{\mathbf{x}} = \lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2$, then think of y that satisfies $y \geq f(\tilde{\mathbf{x}})$. Recalling the definition of a convex set (0.1), when Y is a set of y , of course, $f(\mathbf{x}_1)$ and $f(\mathbf{x}_2) \in Y$, so $\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \in Y$. That is, $\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2)$ is also included in the value of y , and $\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2)$ holds. \square

Next, I explain the characteristics of the convex function. If the convex function is differentiable, the following necessary and sufficient conditions hold.

*1

Definition 0.4 (Concave function). $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a concave function $\stackrel{\text{def}}{\iff} -f$ is a convex function.

Theorem 0.2 (First-order convexity condition).

Suppose the real-valued function $f : \mathbb{R}^n \rightarrow [-\infty, \infty]$ is differentiable. At this time, the necessary and sufficient conditions for f to be a convex function are as follows:

for $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^n$,

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x})^T \big|_{\mathbf{x}=\mathbf{x}_1} (\mathbf{x}_2 - \mathbf{x}_1)$$

Proof. Since the equal sign holds when $\lambda = 0$ and 1 , prove it for $\lambda \in (0, 1)$. In addition, assuming there are $\forall \mathbf{x}_1$ and $\mathbf{x}_2 \in \mathbb{R}^n$ ($\mathbf{x}_1 \neq \mathbf{x}_2$).

(\Rightarrow)

Since f is a convex function, $\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2)$ holds. Therefore,

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \frac{f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2) - f(\mathbf{x}_1)}{1 - \lambda}$$

Let $\tilde{\lambda} = 1 - \lambda$, the above formula becomes,

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \frac{f(\mathbf{x}_1 + \tilde{\lambda}(\mathbf{x}_2 - \mathbf{x}_1)) - f(\mathbf{x}_1)}{\tilde{\lambda}} (\mathbf{x}_2 - \mathbf{x}_1)$$

Since f is differentiable, perform Taylor expansion around $\Delta \mathbf{x} = \tilde{\lambda}(\mathbf{x}_2 - \mathbf{x}_1)$ on f , the RHS (right-hand side) is,

$$\begin{aligned} \text{RHS} = & f(\mathbf{x}_1) + \nabla f(\mathbf{x})^T \big|_{\mathbf{x}=\mathbf{x}_1} (\mathbf{x}_2 - \mathbf{x}_1) + \frac{1}{2!} (\mathbf{x}_2 - \mathbf{x}_1)^T \nabla^2 f(\mathbf{x}) \big|_{\mathbf{x}=\mathbf{x}_1} (\mathbf{x}_2 - \mathbf{x}_1) + \mathcal{O}(\Delta^3) \end{aligned} \quad (1)$$

Of course, even if the $\Delta \mathbf{x}$ is as close to $\mathbf{0}$ as possible, the above inequality holds, so

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x})^T \big|_{\mathbf{x}=\mathbf{x}_1} (\mathbf{x}_2 - \mathbf{x}_1)$$

(\Leftarrow)

If the inequality (0.2) holds, then the following relational expression can be created.

$$\begin{cases} f(\mathbf{x}_2) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2) + \nabla f(\mathbf{x})^T \big|_{\mathbf{x}=\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2} (\mathbf{x}_2 - (\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2)) \\ f(\mathbf{x}_1) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2) + \nabla f(\mathbf{x})^T \big|_{\mathbf{x}=\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2} (\mathbf{x}_1 - (\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2)) \end{cases}$$

Multiply these relational expressions by $1 - \lambda$ and λ , and take the sum,

$$\lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2) \geq f(\lambda \mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2)$$

Therefore, f is a convex function. □

For example, if there is only one variable, for $\forall x$ and x_0 , the theorem (0.2) becomes,

$$f(x) \geq f(x_0) + \frac{df(x)}{dx} \big|_{x=x_0} (x - x_0)$$

Especially when $\frac{df(x_0)}{dx} = 0$, $f(x)$ is the minimum value at $x = x_0$, as is clear from the formula. So far, I have explained the first-order convexity condition, but the condition of the convex has second, and now I describe the second-order convexity condition.

Theorem 0.3 (Second-order convexity condition).

Let $f : \mathbb{R}^n \rightarrow [-\infty, \infty]$ be a real valued function that can be differentiated twice. At this time, the necessary and sufficient condition for f to be a convex function is to satisfy the following inequality about Hessian matrix H for $\mathbf{x} = \{x_i\}_{i=1,2,\dots,n} \in \mathbb{R}^n$:

for $\forall \mathbf{d} = \{d_i\}_{i=1,2,\dots,n} \in \mathbb{R}^n$,

$$\mathbf{d}^T H(\mathbf{x}) \mathbf{d} \geq 0$$

this inequality is called **positive semi-definite**.

Here, the Hessian matrix is defined as follows:

Definition 0.5 (Hessian matrix).

At $\mathbf{x} = \{x_i\}_{i=1,2,\dots,n}$ and $f(\mathbf{x})$ is differentiable twice, the Hessian matrix defined as

$$H(\mathbf{x})_{ij} := \nabla^2 f(\mathbf{x})_{ij} \quad (i, j = 1, 2, \dots, n)$$

Proof. For $\forall \lambda \in (0, 1)$, utilize the formula used in the proof (0.2).

(\implies)

From equation (1), for $\mathbf{d} := (\mathbf{x}_2 - \mathbf{x}_1)$, perform the Taylor expansion around \mathbf{d} of $f(\mathbf{x})$,

$$f(\mathbf{x} + \mathbf{d}) = f(\mathbf{x}) + \frac{1}{\lambda} \nabla f(\mathbf{x})^T \mathbf{d} + \frac{1}{2!} \frac{1}{\lambda^2} \mathbf{d}^T H(\mathbf{x}) \mathbf{d} + \mathcal{O}(\Delta^3)$$

Also, the following relational expression holds from the theorem (0.2),

$$f(\mathbf{x} + \mathbf{d}) \geq f(\mathbf{x}) + \frac{1}{\lambda} \nabla f(\mathbf{x})^T \mathbf{d}$$

Compare the two formulas above,

$$\frac{1}{2!} \frac{1}{\lambda^2} \mathbf{d}^T H(\mathbf{x}) \mathbf{d} + \mathcal{O}(\Delta^3) \geq 0$$

This formula holds no matter how small $\mathcal{O}(\Delta^3)$ is, so

$$\mathbf{d}^T H(\mathbf{x}) \mathbf{d} \geq 0$$

(\impliedby)

Suppose the Hessian matrix is positive semi-definite. Perform the Taylor expansion around \mathbf{d} of $f(\mathbf{x})$ and move the contents of the formula,

$$f(\mathbf{x} + \mathbf{d}) - \left(f(\mathbf{x}) + \frac{1}{\lambda} \nabla f(\mathbf{x})^T \mathbf{d} \right) = \frac{1}{2!} \frac{1}{\lambda^2} \mathbf{d}^T H(\mathbf{x}) \mathbf{d} + \mathcal{O}(\Delta^3)$$

Since the Hessian matrix is positive semi-definite, that is,

$$f(\mathbf{x} + \mathbf{d}) - \left(f(\mathbf{x}) + \frac{1}{\lambda} \nabla f(\mathbf{x})^T \mathbf{d} \right) \geq \mathcal{O}(\Delta^3)$$

Of course, it holds even if $\mathcal{O}(\Delta^3)$ is small enough,

$$f(\mathbf{x} + \mathbf{d}) - \left(f(\mathbf{x}) + \frac{1}{\lambda} \nabla f(\mathbf{x})^T \mathbf{d} \right) \geq 0$$

Therefore, according to the theorem (0.2), $f(\mathbf{x})$ is a convex function. \square

By utilizing the fact that it is differentiable, it is possible to know whether a function is a convex function without using inequalities. For example, if $f(x) = \exp(x)$, then $f''(x) \geq 0$, so $f(x)$ is a convex function.

0.2 Convex optimization

First, I explain what optimization problem is. When a set S and a function $f : S \rightarrow \mathbb{R}$ are given, optimization problem is a problem described as follows:

$$\begin{cases} \min. & f(\mathbf{x}) \\ \text{s.t.} & \mathbf{x} \in S \end{cases}$$

The function f is called objective function, and \mathbf{x} that maximizes or minimizes the objective function is called the optimal solution. This time, the optimization problem assumes a problem that has a optimal solution the minimizes the objective function. The optimal solution includes global optimal solution and local optimal solution. The global optimal solution is defined as follows:

Definition 0.6 (global optimal solution).

For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and a set $S \subset \mathbb{R}^n$, let $\bar{\mathbf{x}} \in S$ be the solution to the optimization problem (2). When the $\bar{\mathbf{x}}$ satisfies the following condition, $\bar{\mathbf{x}}$ is called a global optimal solution.

$$f(\mathbf{x}) \geq f(\bar{\mathbf{x}}), \quad \forall \mathbf{x} \in S$$

In addition, the local optimal solution is defined as follows:

Definition 0.7 (local optimal solution).

For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and a set $S \subset \mathbb{R}^n$, let $\bar{\mathbf{x}} \in S$ be the solution to the optimization problem (2). When its neighborhood $N(\bar{\mathbf{x}})$ exist and $\bar{\mathbf{x}}$ satisfies the following condition, $\bar{\mathbf{x}}$ is called a local optimal solution.

$$f(\mathbf{x}) \geq f(\bar{\mathbf{x}}), \quad \forall \mathbf{x} \in N(\bar{\mathbf{x}}) \cap S$$

There are many types of optimization problems such as convex optimization problem, linear optimization problem and nonlinear optimization problem, and the solution method and difficulty of finding the optimal solution differ depending on the problem. In particular, the convex optimization problem described this time has the property that the local optimal solution matches the global optimal solution, and it is easy to find the optimal solution.

0.3 Method of optimization

lagrange multiplier

0.4 Log-linear classifier

logistic regression

References

- [1] 奥村 学, 高村 大也 (2010), 言語処理のための機械学習入門 (自然言語処理シリーズ), コロナ社.