

Subgradient Method

Indian Institute of Technology, Hyderabad
ai20btech11006@iith.ac.in

April 28, 2022

Abstract

Non-differentiable functions are an important class of functions which often appear in optimization problems, a gradient descent method would fail to optimize such an objective because the gradient might not exist at all points that the method encounters. Methods such as interior point method work great, but it has its own limitations of being computationally inefficient. We will explore the subgradient method which is an iterative first-order method similar to gradient descent.

I. INTRODUCTION

Subgradient method is a simple algorithm used to minimize non-differentiable convex functions. This method is similar to vanilla gradient method which is used to optimize differentiable functions.

i. Algorithm

Lets say we have a nondifferentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. The update rule of subgradient method says

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \alpha_k \underline{g}^{(k)} \quad (1)$$

where $\underline{x}^{(k)}$ is the k^{th} iterate, α_k is the step size at k^{th} iteration and $\underline{g}^{(k)}$ is any subgradient of f at $\underline{x}^{(k)}$. The subgradient $\underline{g}^{(k)}$ is any vector which satisfies

$$f(\underline{y}) \geq f(\underline{x}) + \underline{g}^T(\underline{y} - \underline{x}) \quad (2)$$

At a given point there can be more than one subgradients, we call the set of subgradients as subdifferential.

Theorem I.1. *If the function f is differentiable at $\underline{x}^{(k)}$ then $\underline{g}^{(k)}$ is equal to the gradient of f at $\underline{x}^{(k)}$*

Proof. Substitute $\underline{y} = \underline{x} + \lambda \underline{z}$, $\lambda > 0$ in (2)

$$\frac{f(\underline{x} + \lambda \underline{z}) - f(\underline{x})}{\lambda} \geq \underline{g}^T \underline{z} \quad (3)$$

We can use the limit $\lambda \rightarrow 0$

$$\nabla f(\underline{x})^T \underline{z} \geq \underline{g}^T \underline{z} \quad (4)$$

$$\underline{z}^T (\nabla f(\underline{x}) - \underline{g}) \geq 0 \quad \forall \underline{z} \quad (5)$$

$$\therefore \underline{g} = \nabla f(\underline{x}) \quad (6)$$

□

II. EXAMPLE

Consider the following problem

$$\min \max\{f_1, f_2, \dots, f_n\} \quad (7)$$

This problem is non differentiable but it can be solved using subgradient method. Lets consider a numerical example.

III. CONVERGENCE PROOF

Lets assume x^* is the minimizer of our objective function f . Assume that the norm of subgradients is bounded.

Using Lipschitz condition

$$|f(u) - f(v)| \leq G \|u - v\|_2 \quad (8)$$

for all u, v . Some versions of subgradient method work even when the gradient is not bounded.

IV. SVM USING SUBGRADIENT METHOD

A support vector machine is used for two class classification. The objective is to maximize the slab thickness while still satisfying few constraints. A hard-margin SVM can be formulated as

$$\min_{\underline{w}, b} \underline{w}^T \underline{w} \quad (9)$$

$$\text{s.t: } y_i(\underline{w}^T \underline{x}_i + b) \geq 1 \quad (10)$$

$$(11)$$

We can transform this problem into

$$\min_{\underline{w}, b} \underline{w}^T \underline{w} + \lambda \sum_i \max(0, 1 - y_i(\underline{w}^T \underline{x}_i + b)) \quad (12)$$

Where λ is the trade off factor, the higher is the value of this parameter, the higher is the penalty of violating the given constraints. This problem is essentially a soft-margin SVM. We can now solve this problem by subgradient method.

The first step is to calculate the subgradients,

$$\underline{g}_{\underline{w}} = 2 * \underline{w} + \lambda \sum_i y_i * \underline{x}_i \text{ for } y_i(\underline{w}^T \underline{x} + b) < 1 \quad (13)$$

and

$$g_b = \lambda \sum_i y_i \text{ for } y_i(\underline{w}^T \underline{x} + b) < 1 \quad (14)$$

In practice, we tend to use mini-batch subgradient method because of its several advantages such as computational efficiency, stable convergence and faster learning.

V. RESULTS

i. Maximum of convex functions

I considered a problem of the type maximum of linear functions and compared different step size rules. The results obtained are as follows

1. Constant Step Size

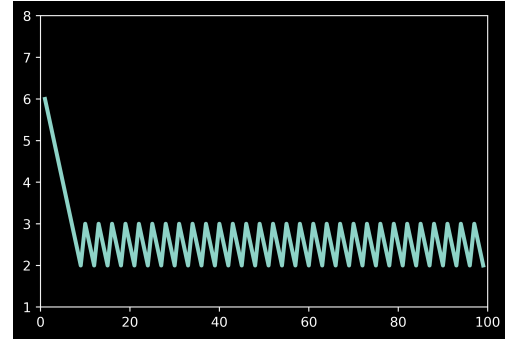


Figure 1: Constant Step Size

2. Constant Step Length

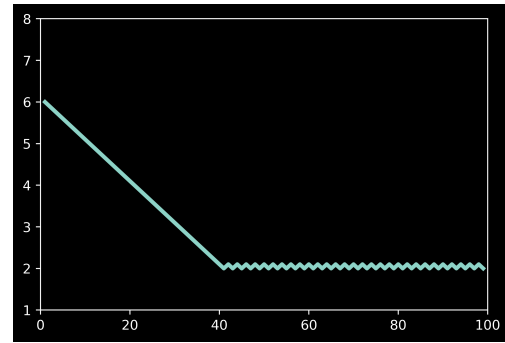


Figure 2: Constant Step Length

3. Square Summable But Not Summable

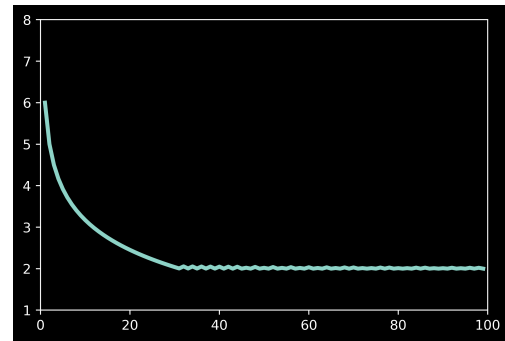


Figure 3: Square Summable But Not Summable

4. Non Summable Diminishing

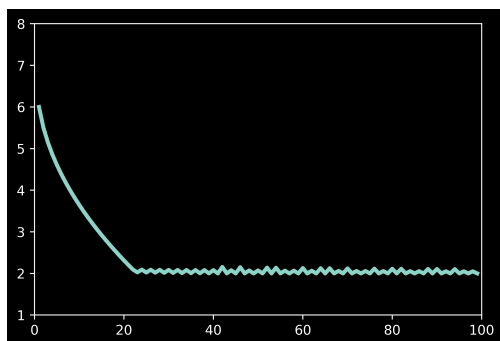


Figure 4: *Non summable diminishing*

VI. PRELIMINARY REFERENCES

1. Boyd, Stephen and Xiao, Lin and Mutapcic, Almir. Subgradient method ¹
2. Boyd, Stephen and Mutapcic, Almir. stochastic Subgradient method ²