

Optimization for Machine Learning in Python

Dr. Moritz Wolter

July 14, 2022

High Performance Computing and Analytics Lab, Uni Bonn

Overview

Introduction

The derivative

Optimization in a single dimension

Optimization in many dimensions

Second order optimization

2022-07-14

Optimization for Machine Learning in Python

└ Overview

Overview

Introduction

The derivative

Optimization in a single dimension

Optimization in many dimensions

Second order optimization

TODO

Introduction

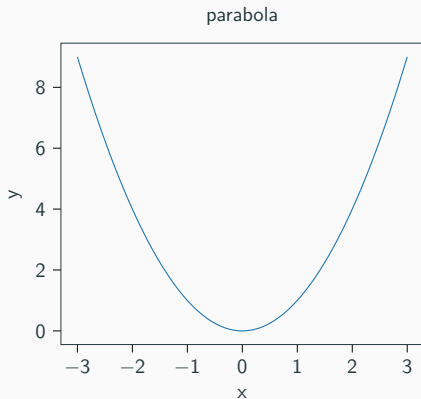
Traditionally, optimization means minimizing using a cost function $f(x)$. Given the cost, we must find the cheapest point x^* on the function, or in other words,

$$x^* = \min_{x \in \mathbb{R}} f(x) \quad (1)$$

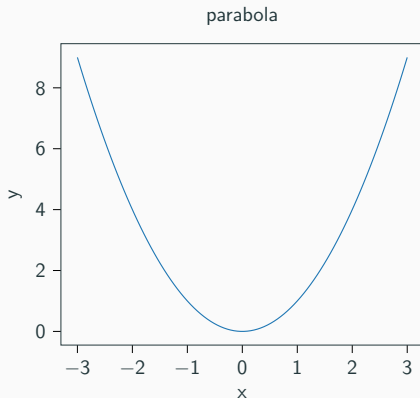
Functions

Functions are mathematical mappings. Consider for example the quadratic function, $f(x) : \mathbb{R} \rightarrow \mathbb{R}$:

$$f(x) = x^2 \quad (2)$$



Where is the minimum?



In this case we immediately see it's at zero. Finding it algorithmically requires derivative information.

Summary

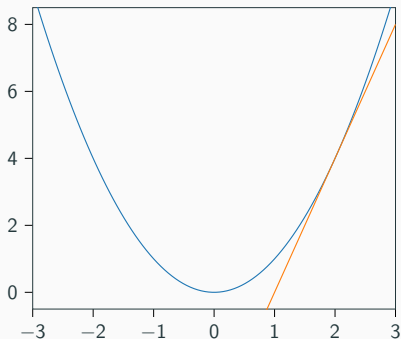
TODO

The derivative

The derivative

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (3)$$

parabola with derivative at two



2022-07-14

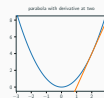
Optimization for Machine Learning in Python

└ The derivative

└ The derivative

The derivative

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (3)$$



TODO

Derivation of the parabola derivative

$$\lim_{h \rightarrow 0} \frac{(x + h)^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h} \quad (4)$$

$$= \lim_{h \rightarrow 0} \frac{2xh + h^2}{h} \quad (5)$$

$$= \lim_{h \rightarrow 0} \frac{h(2x + h)}{h} \quad (6)$$

$$= \lim_{h \rightarrow 0} 2x + h \quad (7)$$

$$= 2x \quad (8)$$

Optimization for Machine Learning in Python

└ The derivative

└ Derivation of the parabola derivative

$$\lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h} \quad (4)$$

$$= \lim_{h \rightarrow 0} \frac{2xh + h^2}{h} \quad (5)$$

$$= \lim_{h \rightarrow 0} \frac{h(2x + h)}{h} \quad (6)$$

$$= \lim_{h \rightarrow 0} 2x + h \quad (7)$$

$$= 2x \quad (8)$$

Derive on the board:

$$\lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h} \quad (9)$$

$$= \lim_{h \rightarrow 0} \frac{2xh + h^2}{h} \quad (10)$$

$$= \lim_{h \rightarrow 0} \frac{h(2x + h)}{h} \quad (11)$$

$$= \lim_{h \rightarrow 0} 2x + h \quad (12)$$

$$= 2x \quad (13)$$

TODO

Optimization in a single dimension

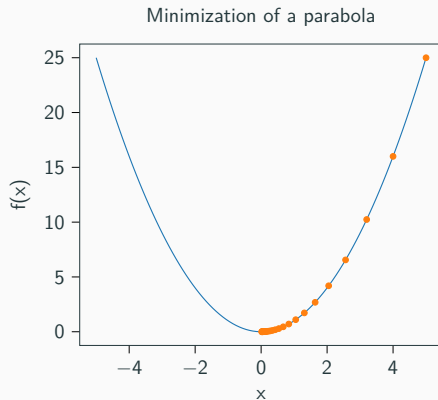
Steepest descent

To find a minimum we descend along the gradient, with n denoting the step number, $\alpha \in \mathbb{R}$ the step size and $\frac{df}{dx}$ the derivative of f along $x \in \mathbb{R}$:

$$x_n = x_{n-1} - \alpha \cdot \frac{df}{dx}. \quad (14)$$

Steepest descent on the parabola

Working with the initial position $x_0 = 5$ and a step size of $\alpha = 0.1$ for 25 steps leads to:



TODO

Optimization in many dimensions

Multidimensional problems

The Rosenbrock test function:

$$f(x, y) = (a - x)^2 + b(y - x^2)^2 \quad (15)$$

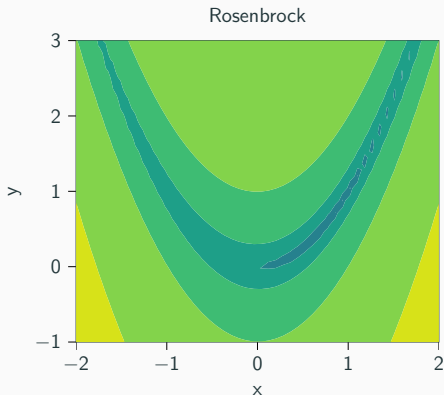


Figure: Rosenbrock function with $a=1$ and $b=100$.

The gradient

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \quad (16)$$

2022-07-14

Optimization for Machine Learning in Python

└ Optimization in many dimensions

└ The gradient

The gradient

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} & \dots & \frac{\partial f}{\partial x_n} \end{pmatrix} \quad (16)$$

TODO

Recall the Rosenbrock function:

$$f(x, y) = (a - x)^2 + b(y - x^2)^2 \quad (17)$$

$$\nabla f(x, y) = \begin{pmatrix} -2a + 2x - 4byx + 4bx^3 \\ 2by - 2bx^2 \end{pmatrix} \quad (18)$$

Optimization for Machine Learning in Python

└ Optimization in many dimensions

└ Rosenbrock gradient

Recall the Rosenbrock function:

$$f(x, y) = (a - x)^2 + b(y - x^2)^2 \quad (17)$$

$$\nabla f(x, y) = \begin{pmatrix} -2a + 2x - 4byx + 4bx^3 \\ 2by - 2bx^2 \end{pmatrix} \quad (18)$$

On the board derive:

$$f(x, y) = (a - x)^2 + b(y - x^2)^2 \quad (19)$$

$$= a^2 - 2ax + x^2 + b(y^2 - 2yx^2 + x^4) \quad (20)$$

$$= a^2 - 2ax + x^2 + by^2 - 2byx^2 + bx^4 \quad (21)$$

$$\Rightarrow \frac{\partial f(x, y)}{\partial x} = -2a + 2x - 4byx + 4bx^3 \quad (22)$$

$$\Rightarrow \frac{\partial f(x, y)}{\partial y} = 2by - 2bx^2 \quad (23)$$

Gradient descent

Initial position: $x_0 = [0.1, 3.]$,

Gradient step size: $\alpha = 0.01$

$$x_n = x_{n-1} - \alpha \cdot \nabla f(\mathbf{x}) \quad (24)$$

n denotes the step number, ∇ the gradient operator, and $f(\mathbf{x})$ a vector valued function.

Gradient descent on the Rosenbrock function

Rosenbrock Optimization

Rosenbrock Optimization

TODO

Second order optimization

Second order optimization

TODO