

Optimisation for Machine Learning

Ludovic De Matteis
LAAS-CNRS
ldematteis@laas.fr

TABLE OF CONTENTS

1 Motivations in Machine Learning	2
1.a Unconstraint Optimization	2
1.b Regression	2
1.c Classification	3
2 Derivatives and Gradients	4
2.a Reminder on derivatives	4
2.b First order optimality conditions	6
2.c Derivative of classification cost	6
2.d The chain rule	6
3 Gradient Descent	6
3.a The algorithm	6
3.b Convergence Analysis	6
3.c Regularization	6
4 Newton method	6
4.a Comparison the Gradient Descent	6
4.b Advantages	6
4.c Limits	6
5 Stochastic Optimization	6
5.a Problem formulation	6
5.b Batch gradient descent	6
5.c Stochastic gradient descent	6
5.d Backpropagation	6
6 Neural Networks	6
6.a Perceptron	6
6.b Multi-Layer Perceptron (MLP)	6
6.c Additional structures	6
7 Opening	6
7.a Introduction to Large Language Models	6

1 MOTIVATIONS IN MACHINE LEARNING

Optimization describe a mathematical theory focussing on the problem of finding the minimum (or equivalently the maximum) of a function. This problems occurs in almost every domain, from investment portfolios built to maximize the return to minimizing the error in a weather forecast model. This introductory course aims at teaching the basics of optimization and how it is applied to solve machine learning problem.

1.a Unconstraint Optimization

In unconstraint optimization, we consider problems of the form

$$\inf_{x \in \mathbb{R}^n} f(x) \quad (1)$$

we define the set of **global minimizers** of the function f as

$$\arg \min_{x \in \mathbb{R}^n} f(x) \stackrel{\text{def}}{=} \{x_0 \in \mathbb{R}^n \mid \forall x \in \mathbb{R}^n, f(x_0) \leq f(x)\} \quad (2)$$

The minimizer of a function does not necessarily exists and if it does, it can be non unique.

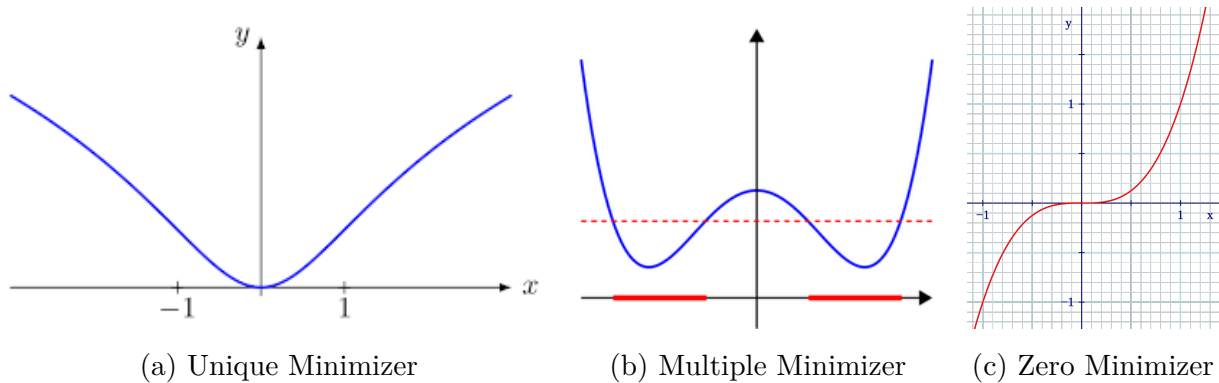


Figure 1: Different functions shapes and corresponding number of minimizers

Figure 1 shows different function shapes corresponding to different number of minimizers. In general, when a minimizer exists - and this will be our focus for the rest of the course - we denote the problem Equation 1 as

$$\min_{x \in \mathbb{R}^n} f(x) \quad (3)$$

1.b Regression

We will start by considering a linear regression problem, in which

$$f(x) = \frac{1}{2} \sum_{i=1}^N (y_i - \langle x, a_i \rangle)^2 = \frac{1}{2} \| Ax - y \|^2 \quad (4)$$

is the least square quadratic risk function.

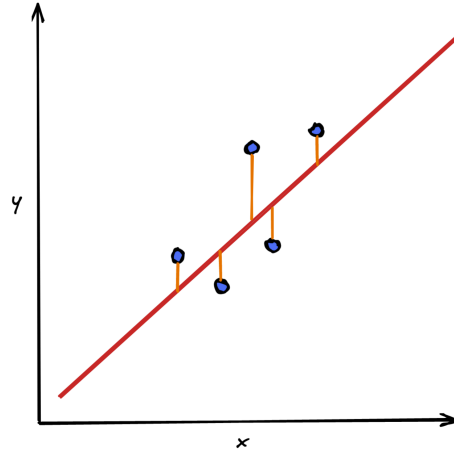


Figure B: Linear regression problem.

The linear function $l(x) = Ax$ in red aims to fit as well as possible the given data points

An illustration of the linear regression problem is shown in Figure B.

The regression problem can be extended to different function, by considering for instance a quadratic function

$$f(x) = \frac{1}{2} \| (x^T Ax + Bx) - y \|^2 \quad (5)$$

or an exponential function

$$f(x) = \frac{1}{2} \| e^{Ax} - y \|^2 \quad (6)$$

1.c Classification

For binary classification, the data points y_i are **labelled** with a value 1 or -1 , defining a class. In this problem, we aim the minimize the function

$$f(x) = \sum_{i=1}^n l(-y_i < x, a_i >) = L(-\text{diag}(y)Ax) \quad (7)$$

where l is the 0-1 loss function $1_{\mathbb{R}^+}$ or a smooth approximation of it, giving a value of 1 if the signs of y_i and $< a_i, x >$ are opposed and zero otherwise.

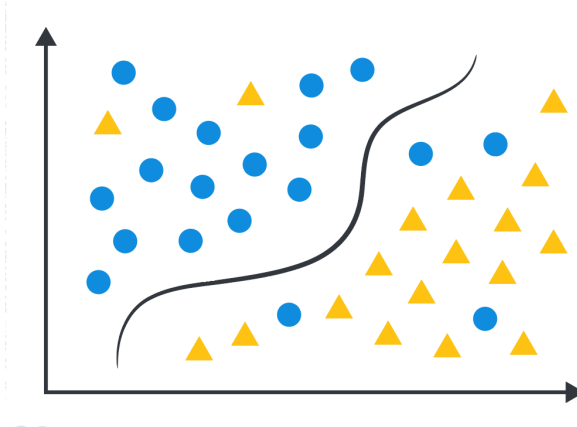


Figure C: Binary Classification problem.

The function in black aims to separate as well as possible the different classes

As for regression, the classification problem can use a different separating function. Figure C presents a binary classification problem with an arbitrary separating function.

2 DERIVATIVES AND GRADIENTS

2.a Reminder on derivatives

A function is said to be differentiable at a point x_0 if the following limit exists:

$$\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} \quad (8)$$

The value of this limit when it exists is called the derivative of f at x_0 and is denoted $f'(x_0)$ or $\frac{df}{dx}(x_0)$.

Examples - Let's illustrate this idea on several classical functions

- $f(x) = 3x^2 - x$

$$\begin{aligned} f(x+h) - f(x) &= 3(x+h)^2 - x - h - 3x^2 + x \\ &= 3(h^2 + 2xh) - h \end{aligned} \quad (9)$$

which gives

$$\frac{f(x+h) - f(x)}{h} = 3h + 6x - 1 \quad (10)$$

the limit when $h \rightarrow 0$ as this expression is defined for all points $x \in \mathbb{R}$ and gives the derivative

$$f'(x) = 6x - 1 \quad (11)$$

- $f(x) = e^{-2x}$

$$\begin{aligned}\frac{f(x+h) - f(x)}{h} &= \frac{e^{-2x-2h} - e^{-2x}}{h} \\ &= \frac{e^{-2x}(e^{-2h} - 1)}{h} \\ &= -2e^{-2x} \frac{(e^{-2h} - 1)}{-2h}\end{aligned}\tag{12}$$

Moreover

2.b First order optimality conditions

2.c Derivative of classification cost

2.d The chain rule

3 GRADIENT DESCENT

3.a The algorithm

3.b Convergence Analysis

3.c Regularization

4 NEWTON METHOD

4.a Comparison the Gradient Descent

4.b Advantages

4.c Limits

5 STOCHASTIC OPTIMIZATION

5.a Problem formulation

5.b Batch gradient descent

5.c Stochastic gradient descent

5.d Backpropagation

6 NEURAL NETWORKS

6.a Perceptron

6.b Multi-Layer Perceptron (MLP)

6.c Additional structures

7 OPENING

7.a Introduction to Large Language Models