

# Optimization for Machine Learning

Pablo Mollá Chárlez

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Optimization Problem</b>	<b>2</b>
2.1	Formal Definition . . . . .	2
2.2	Challenges in Optimization . . . . .	2
<b>3</b>	<b>Sets and Functions</b>	<b>3</b>
3.1	Convex Sets . . . . .	3
3.2	Convex Functions . . . . .	3
3.2.1	Definition 1 . . . . .	3
3.2.2	Definition 2 . . . . .	3
3.2.3	Strong Convexity . . . . .	3
3.2.4	Properties of Convex Functions . . . . .	3
3.3	Gradients . . . . .	3
3.4	Hessians . . . . .	4
3.4.1	Least Square Function: Convexity, Gradient and Hessian . . . . .	4
3.5	Taylor's Expansion . . . . .	5
3.5.1	Taylor's Expansion in $\mathbb{R}$ . . . . .	5
3.5.2	Taylor's Expansion in $\mathbb{R}^n$ . . . . .	5
3.6	Descent Direction . . . . .	5
3.6.1	Descent Direction Lemma . . . . .	6
3.6.2	Steepest Descent Direction . . . . .	6
3.6.3	Gradient Descent Algorithm . . . . .	6
3.6.4	L-Smoothness . . . . .	6
3.6.5	Condition Number . . . . .	7
3.6.6	Level Sets . . . . .	7
<b>4</b>	<b>Convergence Rates</b>	<b>7</b>
4.1	Definition and Importance . . . . .	7
4.2	Types of Convergence . . . . .	7
4.2.1	Linear Convergence: . . . . .	7
4.2.2	Quadratic Convergence: . . . . .	7
4.3	Theorem: Convergence Rate of Gradient Descent for Strongly Convex Functions . . . . .	8
4.4	Theorem: Convergence of Gradient Descent for Smooth and Convex Functions . . . . .	8
<b>5</b>	<b>Continuous Optimization</b>	<b>8</b>
5.1	Unconstrained vs. Constrained Optimization . . . . .	8
5.1.1	Unconstrained Optimization: . . . . .	8
5.1.2	Constrained Optimization: . . . . .	8

# 1 Introduction

**Optimization** is a fundamental aspect of **machine learning** and **data science**, playing a critical role in **model training** and **parameter tuning**. It involves finding the best solution for a given problem by **minimizing** or **maximizing** an objective function **while satisfying specific constraints**. From simple regression models to complex deep learning architectures, **optimization techniques ensure efficient learning and convergence**. This summary outlines the key components of optimization problems and techniques.

## 2 Optimization Problem

### 2.1 Formal Definition

- **Objective Function:** The function  $f(x)$  that we aim to maximize or minimize. For example, in machine learning, this could be a prediction error.
- **Constraints:** Conditions that the solutions must satisfy, including:
  - **Inequality Constraints:**  $g_i(x) \leq 0$
  - **Equality Constraints:**  $h_j(x) = 0$

- **Feasible Set:** The set of all possible solutions that satisfy the constraints:

$$S = \{x \in \text{Dom } f \mid g_i(x) \leq 0 \ \forall i, \ h_j(x) = 0 \ \forall j\}.$$

- **Optimal Solution:** The solution  $x^*$  that maximizes or minimizes the objective function within the feasible set:

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{subject to} \quad \begin{cases} h_j(x) = 0, & \forall j = 1, \dots, \ell, \\ g_i(x) \leq 0, & \forall i = 1, \dots, q. \end{cases}$$

The optimal solution is given by:

$$x^* = \arg \min_{x \in \mathbb{R}^n} f(x).$$

### 2.2 Challenges in Optimization

- **Global vs. Local Optima:** A **global optimum** represents the absolute best solution, whereas a **local optimum** is only the best within a specific neighborhood of the domain of  $f$ . Non-convex problems often cause algorithms to get trapped in local optima.
- **Uniqueness:** While a **unique solution** is ideal, **non-convex problems frequently have multiple optima**, necessitating careful consideration in selecting the most appropriate solution.
- **Multimodality:** Multimodal functions, characterized by **multiple peaks and valleys**, pose significant challenges. Algorithms such as **gradient descent may converge to a local rather than the global optimum**.
- **Convexity:** **Convex functions** simplify optimization, as **any local minimum is also a global minimum**. Non-convex functions, however, lack this property, making them harder to optimize.
- **Differentiability:** **Differentiable functions**, with smooth and continuous slopes, **are well-suited for gradient-based methods**. Non-differentiable points, such as sharp corners or sudden changes, require alternative approaches like subgradient methods.
- **Curse of Dimensionality:** As the **number of variables increases**, the **complexity of the problem grows significantly**, making it more challenging to visualize, analyze, and solve.
- **Non-separability:** **Dependencies between optimization variables prevent them from being optimized independently**, complicating the optimization process and often requiring joint optimization strategies.

- **Ill-posedness:** An ill-posed problem is one where the solution does not depend continuously on the data, meaning small changes in input can lead to large variations in the output, or the solution might not even exist or be unique. Ill-posedness can lead to instability in optimization, especially in numerical methods.
- **Ill-conditioning:** The function may have a poorly scaled gradient, meaning small changes in certain directions lead to disproportionately large changes in others. This imbalance can cause slow or unstable convergence during optimization, as it makes it difficult for gradient-based methods to find a stable path to the optimum.

## 3 Sets and Functions

### 3.1 Convex Sets

A set  $S \subset \mathbb{R}^n$  is **convex** if, for any  $x, y \in S$  and  $\lambda \in [0, 1]$ ,  $\lambda x + (1 - \lambda)y \in S$ . This means that the line segment connecting any two points within the set lies entirely inside the set.

### 3.2 Convex Functions

#### 3.2.1 Definition 1

A function  $f$  is **convex** if, for any  $x, y$  and  $\lambda \in [0, 1]$ ,  $f(\lambda \cdot x + (1 - \lambda) \cdot y) \leq \lambda \cdot f(x) + (1 - \lambda) \cdot f(y)$ . This implies that the line segment connecting any two points on the graph of  $f$  lies above the graph itself. A **twice-differentiable function of a single variable is convex** if and only if its **second derivative is non-negative on its entire domain**.

#### 3.2.2 Definition 2

A differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is said to be **convex** if and only if  $\forall x, \theta \in \mathbb{R}^d$ , we have that:  $f(x) \geq f(\theta) + f'(\theta)^T(x - \theta)$

#### 3.2.3 Strong Convexity

A differentiable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is said to be  **$\mu$ -strongly convex** if, for all  $x, \theta \in \mathbb{R}^d$ ,

$$f(x) \geq f(\theta) + \nabla f(\theta)^T(x - \theta) + \frac{\mu}{2} \|x - \theta\|^2,$$

where  $\mu > 0$  is the strong convexity parameter.

#### 3.2.4 Properties of Convex Functions

The following operations preserve convexity:

- **Non-Negative Weighted Sum:** If  $f_1$  and  $f_2$  are convex functions, then  $\alpha f_1 + \beta f_2$  is also convex for  $\alpha, \beta \geq 0$ .
- **Composition Rules:** If  $f$  is convex and increasing, and  $g$  is convex, then the composition  $f(g(x))$  is convex.
- **Jensen's Inequality:** If  $f$  is a convex function and  $X$  is a random variable, then:  $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$

### 3.3 Gradients

The **gradient of a differentiable function  $f(x)$** , denoted as  $\nabla f(x)$ , is a vector of partial derivatives that represents the rate of change of  $f$  with respect to each variable. Mathematically, for  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,

$$\nabla f(x) = \left[ \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]^T.$$

For a differentiable function  $f(x)$ , the **gradient  $\nabla f(x)$  points in the direction of the steepest ascent**. To **minimize the function, one moves in the direction of  $-\nabla f(x)$** , which is the **direction of steepest descent**.

### 3.4 Hessians

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a twice-differentiable scalar-valued function. The **Hessian** of  $f$  is a **square matrix of second-order partial derivatives**, defined as:

$$H(f) = \nabla^2 f(w) = \begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2} & \frac{\partial^2 f}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_n} \\ \frac{\partial^2 f}{\partial w_2 \partial w_1} & \frac{\partial^2 f}{\partial w_2^2} & \cdots & \frac{\partial^2 f}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_n \partial w_1} & \frac{\partial^2 f}{\partial w_n \partial w_2} & \cdots & \frac{\partial^2 f}{\partial w_n^2} \end{bmatrix}.$$

In compact notation:

$$(H(f))_{ij} = \frac{\partial^2 f}{\partial w_i \partial w_j}, \quad \forall i, j \in \{1, 2, \dots, n\}.$$

The Hessian provides information about the curvature of  $f$ :

- If  $H(f)$  is **positive semi-definite** at  $z$ ,  $f$  is convex at  $z \iff \forall z \in \mathbb{R}^n, z^T \cdot H_f \cdot z \geq 0$ .
- If  $H(f)$  is negative definite at  $z$ ,  $f$  is locally concave at  $z$ .
- If  $H(f)$  is indefinite,  $z$  may be a saddle point.

#### 3.4.1 Least Square Function: Convexity, Gradient and Hessian

Let's prove that the least square function is convex, i.e:

$$f(w) = \frac{1}{2} \|y - Xw\|^2,$$

where  $y \in \mathbb{R}^m$ ,  $X \in \mathbb{R}^{m \times n}$ , and  $w \in \mathbb{R}^n$ . To prove  $f(w)$  is convex, we need to show:

$$f(\lambda w_1 + (1 - \lambda)w_2) \leq \lambda f(w_1) + (1 - \lambda)f(w_2), \quad \forall w_1, w_2 \in \mathbb{R}^n, \lambda \in [0, 1].$$

We have:

$$f(\lambda w_1 + (1 - \lambda)w_2) = \frac{1}{2} \|y - X(\lambda w_1 + (1 - \lambda)w_2)\|^2.$$

Using the linearity of  $X$ , this becomes:

$$f(\lambda w_1 + (1 - \lambda)w_2) = \frac{1}{2} \|\lambda(y - Xw_1) + (1 - \lambda)(y - Xw_2)\|^2.$$

Using the property of the squared norm,  $\|u + v\|^2 = \|u\|^2 + \|v\|^2 + 2u^T v$ , we get:

$$f(\lambda w_1 + (1 - \lambda)w_2) = \frac{1}{2} (\lambda^2 \|y - Xw_1\|^2 + (1 - \lambda)^2 \|y - Xw_2\|^2 + 2\lambda(1 - \lambda)(y - Xw_1)^T (y - Xw_2)).$$

Next, we calculate:

$$\lambda f(w_1) + (1 - \lambda)f(w_2) = \lambda \cdot \frac{1}{2} \|y - Xw_1\|^2 + (1 - \lambda) \cdot \frac{1}{2} \|y - Xw_2\|^2.$$

This simplifies to:

$$\lambda f(w_1) + (1 - \lambda)f(w_2) = \frac{1}{2} (\lambda \|y - Xw_1\|^2 + (1 - \lambda) \|y - Xw_2\|^2).$$

Now, compute the difference:

$$\text{Difference} := f(\lambda w_1 + (1 - \lambda)w_2) - (\lambda f(w_1) + (1 - \lambda)f(w_2)).$$

Substituting the expressions, we get:

$$\text{Difference} = \frac{1}{2} (\lambda^2 \|y - Xw_1\|^2 + (1 - \lambda)^2 \|y - Xw_2\|^2 + 2\lambda(1 - \lambda)(y - Xw_1)^T (y - Xw_2))$$

$$-\frac{1}{2}(\lambda\|y - Xw_1\|^2 + (1-\lambda)\|y - Xw_2\|^2).$$

After some algebra and factoring out  $\lambda(1-\lambda)$ , the difference becomes:

$$\text{Difference} = -\frac{\lambda(1-\lambda)}{2}\|Xw_1 - Xw_2\|^2.$$

Since  $\|Xw_1 - Xw_2\|^2 \geq 0$  and  $\lambda(1-\lambda) \geq 0$  for  $\lambda \in [0, 1]$ , the difference is non-positive. Hence:

$$f(\lambda w_1 + (1-\lambda)w_2) \leq \lambda f(w_1) + (1-\lambda)f(w_2).$$

This proves that  $f(w)$  is convex. The [gradient](#) calculation details of the function can be found on the exercise sheet, although here you can find the results for the gradient and the [hessian](#) of the least square function.

- **Gradient:**  $\nabla f(x) = -X^T \cdot (y - Xw)$
- **Hessian:**  $H_f(w) = X^T X$

## 3.5 Taylor's Expansion

### 3.5.1 Taylor's Expansion in $\mathbb{R}$

Let  $k$  be a natural number,  $x_0 \in \mathbb{R}$ , and  $f$  a function that is  $k$ -times continuously differentiable on an interval containing  $x_0$  and  $x$ . Taylor's theorem states that there exists some  $\xi \in (x_0, x)$  such that:

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \dots + \frac{f^{(k)}(\xi)}{k!}(x - x_0)^k.$$

**Implication:** Taylor's theorem allows us to approximate  $f(x)$  around  $x_0$  using increasingly accurate terms based on the derivatives of  $f$  at  $x_0$ .

### 3.5.2 Taylor's Expansion in $\mathbb{R}^n$

For a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  that is continuously twice differentiable, the Taylor approximation around a point  $x_0 \in \mathbb{R}^n$  is given by:

$$f(x) = f(x_0) + \nabla f(x_0)^T(x - x_0) + \frac{1}{2}(x - x_0)^T \nabla^2 f(x_0)(x - x_0) + R_3(x),$$

where  $R_3(x)$  is the remainder term:

$$R_3(x) = O(\|x - x_0\|^3), \quad \text{which vanishes as } x \rightarrow x_0.$$

## 3.6 Descent Direction

The concept of a descent direction identifies directions  $\mathbf{d}$  in which the function  $f$  decreases locally. Let  $\mathbf{x}$  be a point in the domain of  $f$  such that  $\nabla f(\mathbf{x}) \neq 0$ , meaning  $\mathbf{x}$  is not a critical point of  $f$ . A vector  $\mathbf{d} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  is a [descent direction for  \$f\$](#)  at  $\mathbf{x}$  if there exists  $\bar{\alpha} > 0$  such that:

$$f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x}), \quad \forall \alpha \in (0, \bar{\alpha}).$$

This means that  $f$  strictly decreases along the half-line  $\{\mathbf{x} + \alpha \mathbf{d} : \alpha > 0\}$  for sufficiently small step sizes  $\alpha > 0$ .

### 3.6.1 Descent Direction Lemma

Let  $x$  be a non-critical point of  $f$  ( $\nabla f(x) \neq 0$ ), and  $d \in \mathbb{R}^n - \{0\}$ . If  $\nabla f(x)^T d < 0$ , then  $d$  is a descent direction for  $f$  at  $x$ .

#### Proof

Since  $f$  is differentiable, by the first-order Taylor expansion theorem, we can approximate  $f(x + \alpha d)$  for small  $\alpha > 0$  as:

$$f(x + \alpha d) = f(x) + \alpha \nabla f(x)^T d + o(\alpha),$$

where  $o(\alpha)$  represents higher-order terms that vanish as  $\alpha \rightarrow 0$ .

If  $\nabla f(x)^T d < 0$ , then for small  $\alpha > 0$ , the term  $\alpha \nabla f(x)^T d$  is negative, implying:

$$f(x + \alpha d) < f(x).$$

Therefore,  $d$  is a descent direction for  $f$  at  $x$ .

### 3.6.2 Steepest Descent Direction

The (unnormalized) direction  $\mathbf{d} = -\nabla f(x)$  (anti-gradient) is called the steepest descent direction of  $f$  at  $x$ , as it yields the greatest decrease in  $f$ .

### 3.6.3 Gradient Descent Algorithm

To minimize a differentiable function  $f$ , the Gradient Descent algorithm operates with the following sequence of iterates:

- **Initialization:** Start with an initial point  $x^{(0)}$ .
- **Iteration:** For  $k = 0, 1, 2, \dots$ :

$$x^{(k+1)} = x^{(k)} + \alpha^{(k)} \mathbf{d}^{(k)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}),$$

Where:

- $\mathbf{d}^{(k)} = -\nabla f(x^{(k)})$ : The descent direction (negative gradient).
- $\alpha^{(k)}$ : The step size (learning rate).

The iterations continue until a stopping criterion is reached, such as when  $\|\nabla f(x^{(k)})\|$  is sufficiently small or the change in  $f(x^{(k)})$  becomes negligible.

### 3.6.4 L-Smoothness

A differentiable function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is said to be  $L$ -smooth if and only if,

$$|f(x) - f(\theta) - \nabla f(\theta)^T (x - \theta)| \leq \frac{L}{2} \|x - \theta\|^2 \quad \forall \theta, x \in \mathbb{R}^d$$

where  $L > 0$  is the smoothness parameter.

### 3.6.5 Condition Number

The **condition number**  $\kappa$  measures how "well-conditioned" the optimization problem is. When a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is both  **$L$ -smooth** and  **$\mu$ -strongly convex**, we define its condition number  $\kappa$  as:

$$\kappa = \frac{L}{\mu},$$

where  $L$  is the smoothness constant and  $\mu$  is the strong convexity constant. When  $L = \mu$ , the function is **perfectly conditioned** ( $\sim \kappa = 1$ ). Besides, a small condition number  $\kappa \approx 1$  results in fast convergence and, a large condition number  $\kappa \gg 1$  leads to slow convergence and oscillations (zigzag). In terms of Linear Algebra, the **condition number of a matrix**  $A$  is given by:

$$\kappa(A) = \frac{\lambda_{\max}}{\lambda_{\min}},$$

where  $\lambda_{\max}$  and  $\lambda_{\min}$  are **the largest and smallest eigenvalues of  $A$** .

### 3.6.6 Level Sets

Given a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the **level set of  $f$**  corresponding to a scalar  $c \in \mathbb{R}$  is the set of all points  $x \in \mathbb{R}^n$  such that:

$$L_c = \{x \in \mathbb{R}^n \mid f(x) = c\}.$$

## 4 Convergence Rates

### 4.1 Definition and Importance

The convergence rate describes **how quickly an optimization algorithm approaches the optimal solution**. It is an important metric for comparing algorithms, as faster convergence leads to fewer iterations and faster solutions.

### 4.2 Types of Convergence

#### 4.2.1 Linear Convergence:

**Linear convergence occurs when the error decreases by a constant fraction with each iteration.** Formally, this can be expressed as:

$$\|x^{(k+1)} - x^*\| \leq c \cdot \|x^{(k)} - x^*\|,$$

where  $c$  is a constant that represents the rate of error decrease, satisfying  $0 < c < 1$ . A smaller  $c$  indicates faster convergence.

#### 4.2.2 Quadratic Convergence:

**Quadratic convergence occurs when the error decreases by the square of the previous error at each iteration,** leading to very rapid convergence near the solution. Formally, this is expressed as:

$$\|x^{(k+1)} - x^*\| \leq c \cdot \|x^{(k)} - x^*\|^2,$$

where  $c$  is a positive constant ( $c > 0$ ). **Quadratic convergence is significantly faster than linear convergence,** especially near the optimal solution.

### 4.3 Theorem: Convergence Rate of Gradient Descent for Strongly Convex Functions

Assume  $f$  is  $L$ -smooth and  $\mu$ -strongly convex. For gradient descent with a fixed step size  $\alpha_k = \frac{1}{L}$ , the iterates  $(x_k)_{k \geq 0}$  satisfy:

$$f(x_k) - f(x^*) \leq e^{-\frac{k\mu}{L}} \cdot (f(x_0) - f(x^*)),$$

where:

- $x^*$  is the minimizer of  $f$ ,
- $\frac{\mu}{L}$  determines the rate of convergence and depends on the condition number  $\kappa = \frac{L}{\mu}$ .

Gradient descent therefore achieves exponential (linear in log-scale) convergence rate for strongly convex functions.

### 4.4 Theorem: Convergence of Gradient Descent for Smooth and Convex Functions

For a convex and  $L$ -smooth function  $f$ , gradient descent with a step size  $\alpha = \frac{1}{L}$  satisfies:

$$f(x_k) - f(x^*) = O\left(\frac{1}{k}\right),$$

where  $x^*$  is the minimizer of  $f$ . If  $f$  is only assumed to be smooth and convex, gradient descent with a constant step size  $\alpha = \frac{1}{L}$  still converges, but at a slower rate (sublinear rate).

## 5 Continuous Optimization

### 5.1 Unconstrained vs. Constrained Optimization

#### 5.1.1 Unconstrained Optimization:

In unconstrained optimization, the goal is to minimize a function  $f(x)$  over its domain  $D = \text{dom } f$ , without any explicit constraints on  $x$ :

$$\min_{x \in D} f(x).$$

For these problems, the feasible set is simply  $D$ , the domain of  $f$ .

#### 5.1.2 Constrained Optimization:

If restrictions are imposed on  $x$  (e.g.,  $g_i(x) \leq 0$  for certain constraint functions  $g_i(x)$ ), the problem becomes a constrained optimization problem, where solutions must satisfy these additional conditions.