

# Optimization for Machine Learning

Clément W. Royer

Lecture notes - M2 IASD Apprentissage - 2023/2024

- The last version of these notes can be found at:  
<https://www.lamsade.dauphine.fr/~croyer/ensdocs/OAA/PolyOAA.pdf>.
- Comments, typos, etc, can be sent to [clement.royer@lamsade.dauphine.fr](mailto:clement.royer@lamsade.dauphine.fr).  
*Thanks to Florentin Goyens for his feedback.*
- **Major updates of the document**
  - 2023.12.03: Augmented version with mock exams.
  - 2023.09.19: First version of the lecture notes.
- **Learning goals:**
  - Understand the nature and structure of optimization problems arising in machine learning.
  - Select an algorithm tailored to solving a particular instance among those seen in class based on theoretical and practical concerns.
  - Know the underlying motivation behind the design of optimization algorithms for machine learning.

# Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>  | <b>4</b>  |
| 1.1      | Motivation . . . . .   | 4         |
| 1.2      | Notations . . . . .  | 4         |
| 1.2.1    | Generic notations . . . . .  | 4         |
| 1.2.2    | Scalar and vector notations . . . . .                              | 5         |
| 1.2.3    | Matrix notations . . . . .   | 6         |
| 1.3      | The optimization problem . . . . .                                 | 6         |
| 1.3.1    | Mathematical background . . . . .                                  | 7         |
| 1.3.2    | Solution and optimality conditions . . . . .                       | 9         |
| 1.3.3    | Convexity . . . . .  | 11        |
| 1.4      | Examples of optimization problems in ML . . . . .                  | 13        |
| 1.4.1    | Linear regression . . . . .  | 13        |
| 1.4.2    | Logistic regression . . . . .                                      | 14        |
| 1.4.3    | Neural networks . . . . .  | 15        |
| 1.5      | Optimization algorithms . . . . .                                  | 16        |
| 1.5.1    | The algorithmic process . . . . .                                  | 16        |
| 1.5.2    | Convergence and convergence rates . . . . .                        | 17        |
| 1.6      | Summary . . . . .  | 18        |
| <b>2</b> | <b>Smooth optimization methods</b>                                 | <b>19</b> |
| 2.1      | Gradient descent . . . . .   | 19        |
| 2.1.1    | Algorithm . . . . .  | 19        |
| 2.1.2    | Choosing the stepsize . . . . .                                    | 21        |
| 2.1.3    | Convergence rate analysis of gradient descent . . . . .            | 22        |
| 2.1.4    | Application: regression with logistic and sigmoid losses . . . . . | 25        |
| 2.2      | Acceleration techniques . . . . .                                  | 27        |
| 2.2.1    | Introduction: the concept of momentum . . . . .                    | 27        |
| 2.2.2    | Nesterov's accelerated gradient method . . . . .                   | 27        |
| 2.2.3    | Other accelerated methods . . . . .                                | 28        |
| 2.3      | Conclusion . . . . .   | 30        |
| <b>3</b> | <b>Regularization</b>  | <b>31</b> |
| 3.1      | Introduction : The perceptron algorithm . . . . .                  | 31        |
| 3.2      | Nonsmooth optimization . . . . .                                   | 32        |
| 3.2.1    | From nonsmooth functions to nonsmooth problems . . . . .           | 32        |

|  |  |           |
|--|--|-----------|
| 3.2.2  | Subgradient methods                                  | 33        |
| 3.3  | Regularization                                       | 34        |
| 3.3.1  | Regularized problems                                 | 34        |
| 3.3.2  | Sparsity-inducing regularizers                       | 34        |
| 3.3.3  | Proximal methods                                     | 35        |
| 3.4  | Conclusion   | 37        |
| <b>4</b>   | <b>Stochastic optimization methods</b>               | <b>38</b> |
| 4.1  | Motivation   | 38        |
| 4.2  | Stochastic gradient algorithm                        | 39        |
| 4.2.1  | Algorithm  | 39        |
| 4.2.2  | Analysis   | 40        |
| 4.3  | Variance reduction                                   | 42        |
| 4.3.1  | Batch variants                                       | 43        |
| 4.3.2  | Other variants                                       | 43        |
| 4.4  | Stochastic gradient methods for deep learning        | 44        |
| 4.4.1  | Stochastic gradient with momentum                    | 44        |
| 4.4.2  | AdaGrad  | 45        |
| 4.4.3  | RMSProp  | 45        |
| 4.4.4  | Adam   | 46        |
| 4.5  | Conclusion   | 47        |
| <b>5</b>   | <b>Large-scale and distributed optimization</b>      | <b>48</b> |
| 5.1  | Coordinate descent methods                           | 48        |
| 5.1.1  | Algorithmic framework                                | 48        |
| 5.1.2  | Theoretical guarantees of coordinate descent methods | 49        |
| 5.1.3  | Applications of coordinate descent methods           | 50        |
| 5.2  | Distributed and constrained optimization             | 51        |
| 5.2.1  | Linear constraints and dual problem                  | 51        |
| 5.3  | Dual algorithms                                      | 52        |
| 5.3.1  | Dual ascent  | 52        |
| 5.3.2  | Augmented Lagrangian                                 | 53        |
| 5.3.3  | ADMM   | 53        |
| 5.4  | Consensus optimization                               | 54        |
| 5.5  | Conclusion   | 55        |
| <b>Appendix A Mock exam: Around the Huber loss</b>   |  | <b>57</b> |
| <b>Appendix B Mock exam: Shallow neural networks</b> |  | <b>65</b> |

# Chapter 1

## Introduction

This course is concerned with optimization problems arising in data-related applications. Such formulations have gained tremendous interest in recent years, due to the increase in computational power that enable significant advances in fields such as image processing. One of the most fundamental tools behind data science is optimization, that combines mathematical formulations and algorithmic procedures. We describe below the motivation behind studying optimization techniques tailored to data-related applications, as well as the characteristics of the associated problems.

### 1.1 Motivation

The words *machine learning* are widely used as a way to characterize any task that involves manipulating data : nevertheless, their precise meaning can be difficult to formalize, as other keywords such as *data mining*, *data analysis*, *artificial intelligence* or *Big Data* also denote fields that involve data and/or a learning process. In these notes, we focus on the link between data-related tasks and optimization; although we will denote our applications of interest as pertaining to machine learning, we point out that a more general, possibly better suited categorization would be that of **data science**. For the purpose of these lectures, we will indeed consider machine learning through two main goals:

- 1) Extract **patterns** from data, possibly in terms of statistical properties;
- 2) Use this information to **infer** or make predictions about yet unseen data.

A number of such machine learning tasks involve an optimization component, as shown Figure 1.1. As a result, for the purpose of these notes, we will view machine learning as a field making use of statistics and optimization, with the latter being our area of interest. Nevertheless, we point out that computer science features such as data management and parallel computing have also been instrumental to the success of machine learning, and thus should eventually be integrated with optimization to form efficient algorithms.

### 1.2 Notations

#### 1.2.1 Generic notations

- Scalars (i.e. reals) are denoted by lowercase letters:  $a, b, c, \alpha, \beta, \gamma$ .

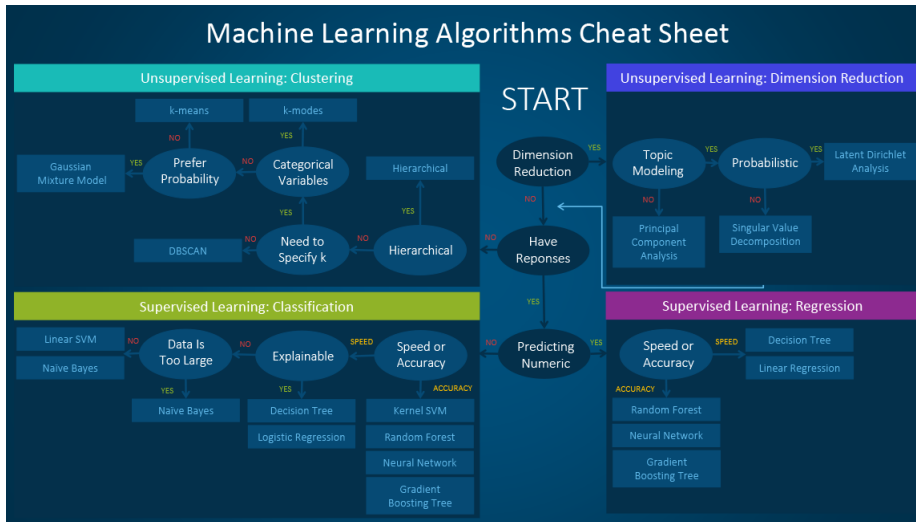


Figure 1.1: A diagram for choosing a machine learning technique appropriate to a given problem; about half of the leaves (Linear SVM, Logistic regression, etc) are directly connected to optimization. Source: <https://blogs.sas.com/content/subconsciousmusings/2017/04/12/machine-learning-algorithm-use/>

- Vectors are denoted by **bold** lowercase letters:  $\mathbf{a}, \mathbf{b}, \mathbf{c}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}$ .
- Matrices are denoted by **bold** uppercase letters:  $\mathbf{A}, \mathbf{B}, \mathbf{C}$ .
- Sets are denoted by **bold** uppercase cursive letters :  $\mathcal{A}, \mathcal{B}, \mathcal{C}$ .
- A new operator or quantity is defined using  $:=$ .
- The following quantifiers are used throughout the notes:  $\forall$  (for every),  $\exists$  (it exists),  $\exists!$  (it exists a unique),  $\in$  (belongs to),  $\subseteq$  (subset of),  $\subset$  (proper subset).
- The  $\Sigma$  operator is used for sums. To lighten the notation, and in the absence of ambiguity, we may omit the first and last indices, or use one sum over multiple indices. As a result, the notations  $\sum_{i=1}^m \sum_{j=1}^n$ ,  $\sum_i \sum_j$  and  $\sum_{i,j}$  may be used interchangeably.
- The notation  $i = 1, \dots, m$  indicates that the variable  $i$  takes all integer values between 1 and  $m$ .

### 1.2.2 Scalar and vector notations

- The set of natural numbers (nonnegative integers) is denoted by  $\mathbb{N}$ ; the set of integers is denoted by  $\mathbb{Z}$ .
- The set of real numbers is denoted by  $\mathbb{R}$ . Our notations for the subset of nonnegative real numbers and the set of positive real numbers are  $\mathbb{R}_+$  and  $\mathbb{R}_{++}$ , respectively. We also define the extended real line  $\overline{\mathbb{R}} := \mathbb{R} \cup \{-\infty, \infty\}$ .
- The notation  $\mathbb{R}^d$  is used for the set of vectors with  $d \in \mathbb{N}$  real components; although we do not explicitly indicate it in the rest of these notes, we always assume that  $d \geq 1$ .

- A vector  $\mathbf{w} \in \mathbb{R}^d$  is thought as a column vector, with  $w_i \in \mathbb{R}$  denoting its  $i$ -th coordinate in the canonical basis of  $\mathbb{R}^d$ . We thus write  $\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$ , or, in a compact form,  $\mathbf{w} = [w_i]_{1 \leq i \leq d}$ .
- Given a column vector  $\mathbf{w} \in \mathbb{R}^d$ , the corresponding row vector is denoted by  $\mathbf{w}^T$ , so that  $\mathbf{w}^T = [w_1 \ \cdots \ w_d]$  and  $[\mathbf{w}^T]^T = \mathbf{w}$ .
- For any integer  $d \geq 1$ , the vectors  $\mathbf{0}_d$  and  $\mathbf{1}_d$  correspond to the vectors of  $\mathbb{R}^d$  for which all elements are 0 or 1, respectively.

### 1.2.3 Matrix notations

- We use  $\mathbb{R}^{m \times n}$  to denote the set of real rectangular matrices with  $m$  rows and  $n$  columns, where  $m$  et  $n$  will always be assumed to be at least 1. If  $m = n$ ,  $\mathbb{R}^{n \times n}$  refers to the set of square matrices of size  $n$ .
- We identify a matrix in  $\mathbb{R}^{m \times 1}$  with its corresponding column vector in  $\mathbb{R}^m$ .
- Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $A_{ij}$  refers to the coefficient from the  $i$ -th row and the  $j$ -th column of  $\mathbf{A}$ : the diagonal of  $\mathbf{A}$  is given by the coefficients  $A_{ii}$ . Provided this notation is not ambiguous, we use the notations  $\mathbf{A}$ ,  $[A_{ij}]_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}}$  and  $[\mathbf{A}_{ij}]$  interchangeably.
- Depending on the context, we may use  $\mathbf{a}_i^T$  to denote the  $i$ -th row of  $\mathbf{A}$  or  $\mathbf{a}_j$  to denote the  $j$ -th column of  $\mathbf{A}$ , leading to  $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_m^T \end{bmatrix}$  or  $\mathbf{A} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_n]$ , respectively.
- Given  $\mathbf{A} = [A_{ij}] \in \mathbb{R}^{m \times n}$ , the *transpose of matrix  $\mathbf{A}$* , denoted by  $\mathbf{A}^T$  (read “ $\mathbf{A}$  transpose”), is defined as the matrix in  $\mathbb{R}^{n \times m}$  (or “ $n$ -by- $m$  matrix”) such that

$$\forall i = 1 \dots m, \forall j = 1 \dots n, \quad \mathbf{A}_{ji}^T = A_{ij}.$$

Note that this generalizes the notation used for row vectors.

- For every  $n \geq 1$ ,  $\mathbf{I}_n$  refers to the identity matrix in  $\mathbb{R}^{n \times n}$  (with 1s on the diagonal and 0s elsewhere).

## 1.3 The optimization problem

We now introduce the mathematical foundations behind optimization.

**Definition 1.3.1** *Optimization is the field of applied mathematics study concerned with making the best decision out of a set of alternatives.*

Mathematically, we write an optimization problem using three components:

- An **objective function**, i.e. a criterion that measures how good a given decision is, that we want to minimize or maximize depending on the context;

- **Decision variables**, that represent the knobs we can turn to change the decision;
- **Constraints**, i.e. conditions that the decision variables must satisfy in order for the decision to be acceptable.

The general form of the optimization problems considered in these notes will be the following

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad f(\mathbf{w}) \quad \text{subject to} \quad \mathbf{w} \in \mathcal{F}. \quad (1.3.1)$$

In problem (1.3.1),  $f$  is the objective function (to be minimized),  $\mathbf{w}$  is the vector of decision variables and  $\mathcal{F}$  is a set encompassing all the constraints on the decision variables. This set is called the feasible set, and is often characterized using mathematical expressions.

### 1.3.1 Mathematical background

Optimization draws from several fields of mathematics, mostly pertaining to linear algebra, topology and differential calculus. We briefly review the key definitions below.

We will always consider  $\mathbb{R}^d$  and  $\mathbb{R}^{n \times d}$  as endowed with their canonical vector space structure; in particular, this means that we will be able to add two vectors (or two matrices), and to multiply a vector (or a matrix) by a scalar value. We will also use the following norm.

**Definition 1.3.2 (Euclidean norm on  $\mathbb{R}^d$ )** *The Euclidean norm (or  $\ell_2$  norm) of a vector  $\mathbf{w} \in \mathbb{R}^d$  is given by:*

$$\|\mathbf{w}\| := \sqrt{\sum_{i=1}^d w_i^2}.$$

**Definition 1.3.3 (Scalar product on  $\mathbb{R}^d$ )** *The scalar product is defined for every  $\mathbf{w}, \mathbf{z} \in \mathbb{R}^d$  by:*

$$\mathbf{w}^T \mathbf{z} := \sum_{i=1}^d w_i z_i.$$

One thus has  $\mathbf{w}^T \mathbf{z} = \mathbf{z}^T \mathbf{w}$  and  $\mathbf{w}^T \mathbf{w} = \|\mathbf{w}\|^2$ .

The notation  $\mathbb{T}$  comes from the concept of transpose in matrix linear algebra.

**Definition 1.3.4 (Transpose matrix)** *Let  $\mathbf{A} = [\mathbf{A}_{ij}] \in \mathbb{R}^{n \times d}$  be a matrix with  $n$  rows and  $d$  columns.*

*The transpose matrix of  $\mathbf{A}$ , denoted by  $\mathbf{A}^T$ , is the matrix with  $d$  rows and  $n$  columns such that*

$$\forall i = 1, \dots, n, \quad \forall j = 1, \dots, d, \quad [\mathbf{A}^T]_{ij} = \mathbf{A}_{ji}.$$

*A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  such that  $\mathbf{A}^T = \mathbf{A}$  is called a symmetric matrix.*

**Definition 1.3.5 (Matrix inversion)** *A matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is invertible if it exists  $\mathbf{B} \in \mathbb{R}^{d \times d}$  such that  $\mathbf{BA} = \mathbf{AB} = \mathbf{I}_d$ , where  $\mathbf{I}_d$  is the identity matrix of  $\mathbb{R}^{d \times d}$ .*

*In this case,  $\mathbf{B}$  is the unique matrix with this property:  $\mathbf{B}$  is called the inverse matrix of  $\mathbf{A}$ , and is denoted by  $\mathbf{A}^{-1}$ .*



**Definition 1.3.6 (Positive (semi-)definiteness)** A matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is positive semidefinite if

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0.$$

It is called positive definite when  $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$  for every nonzero vector  $\mathbf{x}$ .

**Definition 1.3.7 (Eigenvalues and eigenvectors)** Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$ . A real  $\lambda$  is called an eigenvalue of  $\mathbf{A}$  if

$$\exists \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\| \neq 0, \quad \mathbf{A} \mathbf{v} = \lambda \mathbf{v}.$$

The vector  $\mathbf{v}$  is then called an eigenvector of  $\mathbf{A}$  associated to the eigenvalue  $\lambda$ .

**Theorem 1.3.1** Any symmetric matrix in  $\mathbb{R}^{d \times d}$  possesses  $d$  real eigenvalues.

**Notation 1.3.1** Given two symmetric matrices  $(\mathbf{A}, \mathbf{B}) \in \mathbb{R}^{d \times d}$ , we introduce the following notations:

- $\lambda_{\min}(\mathbf{A})/\lambda_{\max}(\mathbf{A})$ : smallest/largest eigenvalue of  $\mathbf{A}$ ;
- $\mathbf{A} \succeq \mathbf{B} \Leftrightarrow \lambda_{\min}(\mathbf{A}) \geq \lambda_{\max}(\mathbf{B})$ ;
- $\mathbf{A} \succ \mathbf{B} \Leftrightarrow \lambda_{\min}(\mathbf{A}) > \lambda_{\max}(\mathbf{B})$ .

Following these notations, a matrix  $\mathbf{A}$  is called **positive semi-definite** (resp. **positive definite**) if and only if  $\mathbf{A} \succeq \mathbf{0}$  (resp.  $\mathbf{A} \succ \mathbf{0}$ ).

**Differential calculus** We will mostly consider minimization problems involving a smooth objective function: the term “smooth” can be loosely defined in the optimization or learning literature, but generally means that the function is as regular as needed for the desired algorithms and analysis to be applicable. In these notes, we will consider that a smooth function is at least continuously differentiable, sometimes twice continuously differentiable. Those concepts are recalled below.

**Definition 1.3.8 (Continuous function)** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$  is continuous at  $\mathbf{w} \in \mathbb{R}^d$  if for every  $\epsilon > 0$ , it exists  $\delta > 0$  such that

$$\forall \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v} - \mathbf{w}\| \leq \delta \implies \|f(\mathbf{v}) - f(\mathbf{w})\| \leq \epsilon.$$

**Definition 1.3.9 (Lipschitz continuous function)** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$  is  $L$ -Lipschitz continuous over  $\mathbb{R}^d$  if

$$\forall (\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \quad \|f(\mathbf{u}) - f(\mathbf{v})\| \leq L \|\mathbf{u} - \mathbf{v}\|,$$

where  $L > 0$  is called a Lipschitz constant.

Lipschitz continuous functions can be sandwiched between two linear functions, which is particularly useful for optimization purposes. Note that every Lipschitz continuous function is continuous.

Derivatives are ubiquitous in continuous optimization, as they allow to characterize the local behavior of a function. We assume that the reader is familiar with the concept of derivative of a function from  $\mathbb{R} \rightarrow \mathbb{R}$ . A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is called *differentiable* at  $\mathbf{w} \in \mathbb{R}^d$  if all its partial derivatives at  $\mathbf{w}$  exist.

**Definition 1.3.10 (Classes of functions)** • A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is continuously differentiable if its first-order derivative exists and is continuous. The set of continuously differentiable functions is denoted by  $\mathcal{C}^1(\mathbb{R}^d)$ .

- A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is twice continuously differentiable if  $f \in \mathcal{C}^1(\mathbb{R}^d)$ , the second-order derivative of  $f$  exists and is continuous. The set of twice continuously differentiable functions is denoted by  $\mathcal{C}^2(\mathbb{R}^d)$ .

**Definition 1.3.11 (First-order derivative)** Let  $f \in \mathcal{C}^1(\mathbb{R}^d)$  be a continuously differentiable function. For any  $\mathbf{w} \in \mathbb{R}^d$ , the **gradient of  $f$  at  $\mathbf{w}$**  is given by

$$\nabla f(\mathbf{w}) := \left[ \frac{\partial f}{\partial w_i}(\mathbf{w}) \right]_{1 \leq i \leq d} \in \mathbb{R}^d.$$

**Definition 1.3.12 (Second-order derivative)** Let  $f \in \mathcal{C}^2(\mathbb{R}^d)$  be a twice continuously differentiable function. For any  $\mathbf{w} \in \mathbb{R}^d$ , the **Hessian of  $f$  at  $\mathbf{w}$**  is given by

$$\nabla^2 f(\mathbf{w}) := \left[ \frac{\partial^2 f}{\partial w_i \partial w_j}(\mathbf{w}) \right]_{1 \leq i, j \leq d} \in \mathbb{R}^{d \times d}.$$

The Hessian matrix is symmetric.

Finally, we define an important class of problems involving a Lipschitz continuity assumption.

**Definition 1.3.13 (Smooth functions with Lipschitz derivatives)** • Given  $L > 0$ , the set  $\mathcal{C}_L^{1,1}(\mathbb{R}^d)$  represents the set of all functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that belong to  $\mathcal{C}^1(\mathbb{R}^d)$  such that  $\nabla f$  is  $L$ -Lipschitz continuous.

- Given  $L > 0$ , the set  $\mathcal{C}_L^{2,2}(\mathbb{R}^d)$  represents the set of all functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  that belong to  $\mathcal{C}^2(\mathbb{R}^d)$  such that  $\nabla^2 f$  is  $L$ -Lipschitz continuous.

An important property of such functions is that one can derive upper approximations on their values, as shown by the following theorem.

**Theorem 1.3.2 (First-order Taylor expansion)** Let  $f \in \mathcal{C}_L^{1,1}(\mathbb{R}^d)$  with  $L > 0$ . For any vectors  $\mathbf{w}, \mathbf{z} \in \mathbb{R}^d$ , one has:

$$f(\mathbf{z}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{z} - \mathbf{w}) + \frac{L}{2} \|\mathbf{z} - \mathbf{w}\|^2. \quad (1.3.2)$$

This expansion is crucial in analyzing the performance of first-order algorithms, as we will do in Chapter 2.

### 1.3.2 Solution and optimality conditions

In the rest of this section, we will focus on unconstrained optimization formulations of the form

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) \quad \text{subject to} \quad \mathbf{w} \in \mathcal{F}, \quad (1.3.3)$$

and characterize properties of solutions of such problems. Since there can be more than one solution, we denote the set of solutions of (1.3.3) by

$$\operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) \mid \mathbf{w} \in \mathcal{F}\} \subseteq \mathbb{R}^d. \quad (1.3.4)$$

The **minimal** value of problem (1.3.6) will be denoted by

$$\min_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) \mid \mathbf{w} \in \mathcal{F}\} \in \mathbb{R} \cup \{-\infty, \infty\}. \quad (1.3.5)$$

If the problem is unbounded (i.e. there always exist a better  $\mathbf{w}$ ), we set the minimum value to be  $-\infty$ , whereas if the feasible set  $\mathcal{F}$  is empty, we set the minimum to be  $+\infty$ .

We now provide two definitions of solutions of (1.3.3), or approximations thereof.

**Definition 1.3.14 (Local minimum)** Given a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $\mathbf{w}^* \in \mathbb{R}^d$  is called a **local minimum** of the problem (1.3.3) if it possesses the lowest value of  $f$  in a neighborhood of feasible points, i.e. if  $\mathbf{w}^* \in \mathcal{F}$  and there exists  $\delta > 0$  such that

$$\forall \mathbf{w} \in \mathcal{B}_\delta(\mathbf{w}^*) \cap \mathcal{F}, \quad f(\mathbf{w}^*) \leq f(\mathbf{w}).$$

Local minima are local approximations of solutions: a stronger notion, much harder to guarantee in practice, is that of global minima.

**Definition 1.3.15 (Global minimum)** Given a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $\mathbf{w}^* \in \mathbb{R}^d$  is called a **global minimum** of  $f$  over  $\mathcal{F}$  if  $\mathbf{w}^* \in \mathcal{F}$

$$\forall \mathbf{w} \in \mathcal{F}, \quad f(\mathbf{w}^*) \leq f(\mathbf{w}).$$

**Optimality conditions** In general, finding global or even local minima is a hard problem. For this reason, researchers in optimization have developed optimality conditions: these are mathematical expressions that can be checked at a given point (unlike the conditions above) and help assessing whether a given point is a local minimum or not.

In this introductory chapter, we will present these conditions in the context of an unconstrained optimization problem

$$\operatorname{minimize}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}). \quad (1.3.6)$$

**Theorem 1.3.3 (First-order necessary condition)** Suppose that the objective function  $f$  in problem (1.3.6) belongs to  $\mathcal{C}^1(\mathbb{R}^d)$ . Then,

$$[\mathbf{w}^* \text{ is a local minimum of } f] \implies \|\nabla f(\mathbf{w}^*)\| = 0. \quad (1.3.7)$$

Note that this condition is only necessary: there may exist points with zero gradient that are not local minima. Indeed, the set of points with zero gradient, called *first-order stationary points*, also includes local maxima and saddle points<sup>1</sup>.

Provided we strengthen our smoothness requirements on  $f$ , we can establish stronger optimality conditions for problem (2.1.1).

<sup>1</sup>A vector is a saddle point of a function if it is a local minimum with respect to certain directions and a local maximum with respect to other directions of the space.

**Theorem 1.3.4 (Second-order necessary condition)** Suppose that the objective function  $f$  in problem (1.3.6) belongs to  $\mathcal{C}^2(\mathbb{R}^d)$ . Then,

$$[\mathbf{w}^* \text{ is a local minimum of } f] \implies [\|\nabla f(\mathbf{w}^*)\| = 0 \text{ and } \nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}]. \quad (1.3.8)$$

From Theorem 1.3.3, first-order stationary points that violate the condition  $\nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}$  cannot be local minima: conversely, a stronger version of this property guarantees that we are in presence of a local minimum.

**Theorem 1.3.5 (Second-order sufficient condition)** Suppose that the objective function  $f$  in problem (1.3.6) belongs to  $\mathcal{C}^2(\mathbb{R}^d)$ . Then,

$$[\|\nabla f(\mathbf{w}^*)\| = 0 \text{ and } \nabla^2 f(\mathbf{w}^*) \succ \mathbf{0}] \implies [\mathbf{w}^* \text{ is a local minimum of } f] \quad (1.3.9)$$

By exploiting the second-order derivative, it is thus possible to certify whether a point is a local minima (note that there could be local or even global minima such that  $\nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}$ ). With further assumptions on the structure of the problem, these optimality conditions can be more informative about minima. This is the case when the objective function is convex: we detail this property in the next section.

### 1.3.3 Convexity

Convexity is at its core a geometric notion: before defining what a convex function is, we describe the corresponding property for a set.

**Definition 1.3.16 (Convex set)** A set  $\mathcal{C} \in \mathbb{R}^d$  is called **convex** if

$$\forall(\mathbf{u}, \mathbf{v}) \in \mathcal{C}^2, \forall t \in [0, 1], \quad t\mathbf{u} + (1 - t)\mathbf{v} \in \mathcal{C}.$$

**Example 1.3.1 (Examples of convex sets)** The following sets are convex:

- The entire space  $\mathbb{R}^d$ ;
- Every line segment of the form  $\{t\mathbf{w} | t \in \mathbb{R}\}$  for some  $\mathbf{w} \in \mathbb{R}^d$ ;
- Every (Euclidean) ball of the form  $\{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|^2 = \sum_{i=1}^d [\mathbf{w}]_i^2 \leq 1\}$ .

We now provide the basic definition of a convex function.

**Definition 1.3.17 (Convex function)** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is **convex** if

$$\forall(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \quad f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq t f(\mathbf{u}) + (1 - t) f(\mathbf{v}).$$

**Example 1.3.2** The following functions are convex :

- Linear functions of the form  $\mathbf{w} \mapsto \mathbf{a}^T \mathbf{w} + b$ , with  $\mathbf{a} \in \mathbb{R}^d$  and  $b \in \mathbb{R}$ ;
- Squared Euclidean norm:  $\mathbf{w} \mapsto \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$ .

If we consider differentiable functions, it is possible to characterize convexity using the derivatives of the function.

**Theorem 1.3.6** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an element of  $\mathcal{C}^1(\mathbb{R}^d)$ . Then, the function  $f$  is convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}). \quad (1.3.10)$$

The inequality (1.3.10) is fundamental in analyzing convex optimization algorithms, as it provides an **underestimator** for the variation of a (convex) objective function.

Convexity can also be characterized using the Hessian matrix (provided the function is sufficiently regular).

**Theorem 1.3.7** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an element of  $\mathcal{C}^2(\mathbb{R}^d)$ . Then, the function  $f$  is convex if and only if

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \nabla^2 f(\mathbf{w}) \succeq \mathbf{0}. \quad (1.3.11)$$

Convex functions are particularly suitable for minimization problems as they satisfy the following property.

**Theorem 1.3.8** If  $f$  is a convex function, then every local minimum of  $f$  is a global minimum.

If the function is differentiable, the optimality conditions as well as the characterization of convexity lead us to the following result.

**Corollary 1.3.1** If  $f$  is continuously differentiable, every point  $\mathbf{w}^*$  such that  $\|\nabla f(\mathbf{w}^*)\| = 0$  is a global minimum of  $f$ .

**Strong convexity** The results above can be further improved by assuming that a convex function is strongly convex, as defined below.

**Definition 1.3.18 (Strongly convex function)** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  in  $\mathcal{C}^1$  is  **$\mu$ -strongly convex** (or strongly convex of modulus  $\mu > 0$ ) if for all  $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$  and  $t \in [0, 1]$ ,

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|^2.$$

Strong convexity leads to an even more desirable property in terms of optimization landscape.

**Theorem 1.3.9** Any strongly convex function has a unique global minimizer.

Similarly to convex functions, it is possible to characterize strong convexity using first- and second-order derivatives.

**Theorem 1.3.10** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an element of  $\mathcal{C}^1(\mathbb{R}^d)$ . Then, the function  $f$  is  $\mu$ -strongly convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}) + \frac{\mu}{2}\|\mathbf{v} - \mathbf{u}\|^2. \quad (1.3.12)$$

**Theorem 1.3.11** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be an element of  $\mathcal{C}^2(\mathbb{R}^d)$ . Then, the function  $f$  is  $\mu$ -strongly convex if and only if

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \nabla^2 f(\mathbf{w}) \succeq \mu \mathbf{I}. \quad (1.3.13)$$

We end this section by giving two examples of strongly convex optimization problems.

**Example 1.3.3 (Convex quadratic problems)** Consider

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}, \quad \mathbf{A} \succeq \mathbf{0}.$$

The function  $f$  belongs to  $\mathcal{C}^2(\mathbb{R}^d)$ , with  $\nabla^2 f(\mathbf{w}) = \mathbf{A}$  for every  $\mathbf{w} \in \mathbb{R}^d$ . As a result, this function is convex. Moreover, if we assume that  $\mathbf{A} \succ \mathbf{0}$ , then the function is  $\lambda_{\min}(\mathbf{A})$ -strongly convex.

**Example 1.3.4 (Projection onto a closed, convex set)** Let  $\mathcal{X} \subseteq \mathbb{R}^d$  be a convex, closed<sup>2</sup> set, and  $\mathbf{a} \in \mathbb{R}^d$ . The problem of computing the projection of  $\mathbf{a}$  onto  $\mathcal{X}$  is formulated as

$$\underset{\mathbf{w} \in \mathcal{X}}{\text{minimize}} \frac{1}{2} \|\mathbf{w} - \mathbf{a}\|^2.$$

The objective function of this problem is 1-strongly convex, which implies that the problem has a unique solution (i. e. the projection is unique).

## 1.4 Examples of optimization problems in ML

### 1.4.1 Linear regression

Linear least squares is arguably the most classical problem in data analysis. We consider a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . Our goal is to compute a linear model that best fits (or explains) the data. We define this model as a function  $h : \mathbb{R}^d \rightarrow \mathbb{R}$ , and we parameterize it through a vector  $\mathbf{w} \in \mathbb{R}^d$ , so that for any  $\mathbf{x} \in \mathbb{R}^d$ , we have  $h(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$ . For every example  $(\mathbf{x}_i, y_i)$  in the dataset, we evaluate how we fit the data based on the squared error  $(\mathbf{x}_i^T \mathbf{w} - y_i)^2$ . We then compute a model by solving the following optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2n} \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} [(\mathbf{x}_i^T \mathbf{w} - y_i)^2 + \lambda \|\mathbf{w}\|^2], \quad (1.4.1)$$

where  $\lambda > 0$  is a regularization parameter. From an optimizer's point of view, problem (1.4.1) is well understood: this is a strongly convex, quadratic problem, and its solution can be computed in close form.

In a typical **linear regression** setting, one assumes that there exists an underlying truth but that the measurements are noisy, i.e.

$$\mathbf{y} = \mathbf{X} \mathbf{w}^* + \boldsymbol{\epsilon},$$

where  $\boldsymbol{\epsilon} \in \mathcal{N}(\mathbf{0}, \mathbf{I})$  is a vector with i.i.d. entries following a standard normal distribution: this is illustrated in Figure 1.2.

In this setting, we wish to compute the most likely value for  $\mathbf{w}^*$ , while being robust to variance in the data. To this end, we suppose that  $\mathbf{y}$  follows a Gaussian distribution of mean  $\mathbf{X} \mathbf{w}$  and of covariance matrix  $\mathbf{I}$ . We also assume a prior Gaussian distribution on the entries of  $\mathbf{w}$ , in order to

<sup>2</sup>A set  $\mathcal{X} \subseteq \mathbb{R}^d$  is closed if for every converging subsequence of  $\{\mathbf{x}_n\}_n$ , the limit of this sequence belongs to  $\mathcal{X}$ .

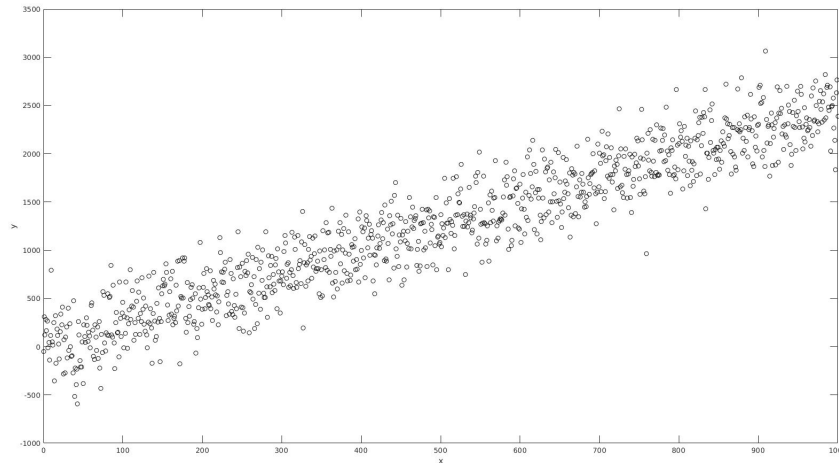


Figure 1.2: Noisy data generated from a linear model with Gaussian noise.

reduce the variance with respect to the data. As a result, an estimate of  $\mathbf{w}^*$ , called the maximum a posteriori estimator, can be computed by solving

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{maximize}} L(y_1, \dots, y_n; \mathbf{w}) := \left[ \frac{1}{\sqrt{2\pi}} \right]^m \exp \left( -\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \mathbf{w} - y_i)^2 - \frac{\lambda}{2} \|\mathbf{w}\|^2 \right). \quad (1.4.2)$$

The solutions of this maximization problem are the same than the solutions of the linear least-squares problem (1.4.1). The resulting solution can be shown to possess very favorable statistical properties: in particular, for  $\lambda$  close to 0, its expected value is close to  $\mathbf{w}^*$ .

Linear regression (with or without regularization) has been extensively studied in optimization and statistics; however, when the number of samples is extremely large, it still poses a number of challenges in practice, as the solution of the problem cannot be computed exactly.

### 1.4.2 Logistic regression

As in Section 1.4.1, we consider a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where  $\mathbf{x}_i \in \mathbb{R}^d$  are feature vectors, and the  $y_i$ s represent binary labels. We wish to build a linear classifier  $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$  to perform this classification, i. e. identify the correct label from the feature. We first suppose that  $y_i \in \{-1, +1\}$ . To model these discrete-valued labels, we introduce an *odds-like* function

$$p(\mathbf{x}; \mathbf{w}) = (1 + e^{\mathbf{x}^T \mathbf{w}})^{-1} \in (0, 1).$$

Given this function, our goal is to choose the model  $\mathbf{w}$  such that

$$\begin{cases} p(\mathbf{x}_i; \mathbf{w}) \approx 1 & \text{if } y_i = +1; \\ p(\mathbf{x}_i; \mathbf{w}) \approx 0 & \text{if } y_i = -1. \end{cases}$$

Given this goal, we want to build an objective function that measures the error between our model and the labels according to the property above. Therefore, we penalize situations in which  $y_i = +1$

and  $p(\mathbf{x}_i; \mathbf{w})$  is close to 0, or  $y_i = -1$  and  $p(\mathbf{x}_i; \mathbf{w})$  is close to 1. This results in the so-called logistic loss, which is a function from  $\mathbb{R}^d$  to  $\mathbb{R}$  defined by

$$\forall \mathbf{w} \in \mathbb{R}^d, f(\mathbf{w}) = \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\}. \quad (1.4.3)$$

The motivation behind introducing the logarithm of the function  $p$  is twofold. On the one hand, it provides a statistical interpretation of the loss as a joint distribution; on the other hand, the derivatives of this function have a more favorable structure.

Given this objective function, the **logistic regression** problem is given by

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\} \quad (1.4.4)$$

This is a convex, smooth problem (though not a strongly convex one), that can be made strongly convex by adding a regularizing term (see Chapter 3).

### 1.4.3 Neural networks

Neural networks have enabled the most impressive, recent advances in perceptual tasks such as image recognition and classification. Thanks to the increase in computational capabilities over the past decade, it is now possible to train extremely deep and wide neural networks, so that they can learn efficient representations of the data.

Given an input vector  $\mathbf{x}_i \in \mathbb{R}^{d_0}$ , a neural network represents a prediction function  $h : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_J}$ , which applies a series of transformations in layers  $\mathbf{x}_i = \mathbf{x}_i^{(0)} \mapsto \mathbf{x}_i^{(1)} \mapsto \dots \mapsto \mathbf{x}_i^{(J-1)} \mapsto \mathbf{x}_i^{(J)}$ . The  $j$ -th layer typically performs the following transformation:

$$\mathbf{x}_i^{(j)} = \boldsymbol{\sigma}(\mathbf{W}_j \mathbf{x}_i^{(j-1)} + \mathbf{b}_j) \in \mathbb{R}^{d_j}, \quad (1.4.5)$$

where  $\mathbf{W}_j \in \mathbb{R}^{d_j \times d_{j-1}}$ ,  $\mathbf{b}_j \in \mathbb{R}^{d_j}$  and  $\boldsymbol{\sigma} : \mathbb{R}^{d_j} \rightarrow \mathbb{R}^{d_j}$  is a componentwise nonlinear function, e.g.  $\boldsymbol{\sigma}(\mathbf{y}) = \left[ \frac{1}{1 + \exp(-y_i)} \right]_i$  (sigmoid function) or  $\boldsymbol{\sigma}(\mathbf{y}) = [\max(0, y_i)]_i$ . As a result, we have  $\mathbf{x}_i^{(J)} = h(\mathbf{x}_i; \mathbf{w})$ , where  $\mathbf{w} \in \mathbb{R}^d$  gathers all the parameters  $\{(\mathbf{W}_1, \mathbf{b}_1), \dots, (\mathbf{W}_J, \mathbf{b}_J)\}$  of the layers.

The optimization problem corresponding to training this neural network architecture involves a training set  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and the choice of a loss function  $\ell$ . It usually results in the following formulation

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i; \mathbf{w}), y_i). \quad (1.4.6)$$

This optimization problem is highly nonlinear and nonconvex in nature, which makes it particularly difficult to solve using algorithms such as gradient descent. Moreover, it typically involves costly algebraic operations, as the number of layers and/or parameters is tremendously large in modern deep neural network architectures. Therefore, problem (1.4.6) also possesses characteristics that are not accounted for in its formulation. The optimization algorithms that efficiently tackle this problem are those that can both guarantee convergence and perform well in practice.



## 1.5 Optimization algorithms

The field of optimization can be broadly divided into three categories:

- **Mathematical** optimization is concerned with the theoretical study of complex optimization formulations, and the proof of well-posedness of such problems (for instance, prove that their exist solutions);
- **Computational** optimization deals with the development of software that can solve a family of optimization problems, through careful implementation of efficient methods;
- **Algorithmic** optimization lies in-between the previous two categories, and aims at proposing new algorithms that address a particular issue, with theoretical guarantees and/or validation of their practical interest.

These notes cover material from the third category of optimization activities. The design of optimization algorithms (also called methods, or schemes) is a particularly subtle process, as an algorithm must exploit the theoretical properties of the problem while being amenable to implementation on a computer.

### 1.5.1 The algorithmic process

Most numerical optimization algorithms do not attempt to find a solution of a problem in a direct way, and rather proceed in an *iterative* fashion. Given a current point, that represents the current approximation to the solution, an optimization procedure attempts to move towards a (potentially) better point: to this end, the method generally requires a certain amount of calculation.

Suppose we apply such a process to the problem minimize  $w \in \mathbb{R}^d$   $f(w)$ , resulting in a sequence of iterates  $\{w_k\}_k$ . Ideally, these iterates obey one of the scenarios below:

1. The iterates get increasingly close to a solution, i. e.

$$\|w_k - w^*\| \rightarrow 0 \quad \text{when } k \rightarrow \infty.$$

Although  $w^*$  is generally not known in practice, such results can be guaranteed by the theory, for instance on strongly convex problems.

2. The function values associated with the iterates get increasingly close to the optimum, i. e.

$$f(w_k) \rightarrow f^* \quad \text{when } k \rightarrow \infty,$$

As for the case above,  $f^*$  may not be known, but it can still be possible to prove convergence for certain algorithms and function classes (typically strongly convex, smooth functions).

3. The first-order optimality condition gets close to being satisfied, that is,  $f \in \mathcal{C}^1(\mathbb{R}^d)$  and

$$\|\nabla f(w_k)\| \rightarrow 0 \quad \text{when } k \rightarrow \infty.$$

Out of the three conditions, the last one is the easiest to track as the algorithm unfolds: it is, however, only a necessary condition, and does not guarantee convergence to a local minimum for generic, nonconvex functions. On the other hand, the first two conditions can only be measured approximately (by looking at the behavior of the iterates and enforcing decrease in the function values), but lead to stronger guarantees.

### 1.5.2 Convergence and convergence rates

The typical theoretical results that optimizers aim at proving for algorithms are asymptotic, as shown above: they only provide a guarantee in the limit. In practice, one may want to obtain more precise guarantees, that relate to a certain accuracy target that the practitioner would like to achieve. This led to the development of **global convergence rates**.

**Example 1.5.1 (Global convergence rate for the gradient norm)** *Given an algorithm applied to minimize  $w \in \mathbb{R}^d$   $f(w)$  that produces a sequence of iterates  $\{w_k\}$ , we say that the method is  $\mathcal{O}(1/k)$  for the gradient norm, or  $\|\nabla f(w_k)\| = \mathcal{O}\left(\frac{1}{k}\right)$  if*

$$\exists C > 0, \quad \|\nabla f(w_k)\| \leq \frac{C}{k} \quad \forall k.$$

Such rates allow to quantify how much effort (in terms of iterations) is needed to reach a certain target accuracy  $\epsilon > 0$ . This leads to the companion notion of **worst-case complexity bound**.

**Example 1.5.2 (Worst-case complexity for the gradient norm)** *Given an algorithm applied to minimize  $w \in \mathbb{R}^d$   $f(w)$  that produces a sequence of iterates  $\{w_k\}$ , we say that the method has a worst-case complexity of  $\mathcal{O}(\epsilon^{-1})$  for the gradient norm if*

$$\exists C > 0, \quad \|\nabla f(w_k)\| \leq \epsilon \quad \text{when } k \geq \frac{C}{\epsilon}.$$

Such results are quite common in theoretical computer science or statistics, which partly explain their popularity in machine learning. In optimization, they have been developed for a number of years in the context of convex optimization but have only gained momentum in general optimization over the last decade.

**Remark 1.5.1 (The computational side of optimization)** *The most popular programming languages for optimization are C/C++/Fortran for high performance implementations, with Python and Julia raising increasing interest. The use of MATLAB is also widespread throughout the optimization community.*

*In addition to programming languages, optimizers have developed **modeling** languages that help bringing the code and the mathematical formulation of a problem closer. The broad-spectrum languages GAMS/AMPL/CVX are reknown examples; other languages, that are more domain-oriented, include MATPOWER and PyTorch.*

*Finally, there are many commercial solvers available (with CPLEX and Gurobi being arguably some of the most efficient for certain classes of problems), along with open-source codes (the COIN-OR platform provides a good interface to all of these methods).*

## 1.6 Summary

Optimization is a key component of modern science, with many tasks in machine learning and related fields involving an optimization problem of some form. The specifics of dealing with massive amounts of data, yet possibly not enough to perfectly model the task at hand, poses a challenge to optimizers. Still, optimization algorithms can prove quite useful to help practitioners in data science (and beyond) in making better decisions.

Optimization begins by a modeling phase, in which a given problem must be stated in terms of objective, variable and constraints. This allows to characterize the properties of the problem, and most importantly its solutions. Properties such as differentiability or convexity lead to specific conditions that one can exploit to identify solutions of this problem.

In general, it is not possible to directly compute a solution of an optimization problem from its formulation; one must thus design a method that will try to compute an approximate solution of the problem. By analyzing this method, it is often possible to identify how fast a method can be at getting close to a solution.

## Chapter 2

# Smooth optimization methods

In this chapter, we review the main methods for solving smooth unconstrained optimization problems. Our starting point will be the Gradient Descent (GD) algorithm, which we study from a theoretical and computational viewpoint in Section 2.1. We will then focus on convex problems and investigate **accelerated** techniques in Section 2.2.

### 2.1 Gradient descent

In this section, we investigate more general, nonlinear unconstrained problems of the form

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}). \quad (2.1.1)$$

We will assume that  $f \in \mathcal{C}^1(\mathbb{R}^d)$ , therefore the gradient mapping for  $f$  exists, is continuous: we will also assume that it can be used in an algorithm. We will develop an algorithm that primarily relies on the use of gradient information, termed **gradient descent**. For such a method, we will derive theoretical guarantees with and without the assumption of convexity: in the latter case, we will see that better results are obtained compared to the general, nonconvex setting.

#### 2.1.1 Algorithm

Because we consider a problem with a continuously differentiable function, we know from the optimality conditions that for any local minimum  $\mathbf{w}^*$ , we necessarily have  $\nabla f(\mathbf{w}^*) = 0$ . As a result, given any point  $\mathbf{w} \in \mathbb{R}^d$ , only one of the two properties below holds:

1. Either  $\nabla f(\mathbf{w}) = 0$ , and  $\mathbf{w}$  can be a local minimum;
2. Or  $\nabla f(\mathbf{w}) \neq 0$  and the function  $f$  decreases *locally* from  $\mathbf{w}$  in the direction of  $-\nabla f(\mathbf{w})$ .

We will formally establish the second property in the next section, thanks to the Taylor expansions we derived in Section 1.3.1. Using this result, we can design the update rule

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla f(\mathbf{w}), \quad (2.1.2)$$

where  $\alpha > 0$  is a stepsize parameter. If  $\nabla f(\mathbf{w}) = 0$ , the vector  $\mathbf{w}$  does not change: this is consistent with the notion of first-order stationarity (we cannot get more information by using the gradient).

On the contrary, when  $\nabla f(\mathbf{w}) \neq 0$ , we expect that there exists a range of values for  $\alpha > 0$  for which such an update leads to a point with a lower objective value.

Using the updating rule (2.1.2), we can design an algorithm for the minimization of the function  $f$ : this method is called **gradient descent**<sup>1</sup> and described in Algorithm 1.

---

**Algorithm 1:** Gradient descent algorithm.

---

**Initialization:**  $\mathbf{w}_0 \in \mathbb{R}^d$ .  
**for**  $k = 0, 1, \dots$  **do**  
    1. Compute the gradient  $\nabla f(\mathbf{w}_k)$ .  
    2. Compute a steplength  $\alpha_k > 0$ .  
    3. Set  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ .  
**end**

---

As written, Algorithm 1 does not have any stopping criterion, and a number of variants can be derived depending on the choice of this stopping criterion, that of the initial point and that of the sequence  $\{\alpha_k\}_k$ . We comment on these aspects below.

**Stopping criterion** In general numerical algorithms operate under a certain budget (of floating-point operations, time, number of iterations), thus any reasonable numerical algorithm will have an embedded stopping criterion, that forces the method to terminate if this budget is reached. In Algorithm 1, for instance, we could have stopped the method after  $k_{\max}$  iterations.

In addition to these practical concerns, algorithms are run in the hope of reaching a prescribed level of accuracy, corresponding to the metrics we described in Section 1.3. For instance, a typical stopping criterion (also called **convergence criterion**) for gradient descent is

$$\|\nabla f(\mathbf{w}_k)\| < \epsilon, \quad (2.1.3)$$

where  $\epsilon > 0$  is a prescribed tolerance, convergence being supposedly harder to achieve as  $\epsilon$  gets smaller.

Finally, additional safety checks can be added to the algorithm. For instance, if the difference between two successive points falls below machine precision, it may not be worth running the method for more iterations.

**Choosing the initial point** Good initialization can lead to significant gains in performance, that must however be put in perspective with the cost of this initialization. For general problems, there could be no incentive to choose one point over another: in this case, random multistart (i.e. running multiple versions of the method with randomly generated starting points) can be used with a small budget to determine a suitable initial point. However, in many applications, the practitioner might already have a reference point, or take an educated guess at what values the decision variables could take: using this as a starting point can be quite valuable, as it will represent a reference value the method is trying to improve upon.

---

<sup>1</sup>Although “gradient descent” is the most common terminology in data science, the historical name used in optimization is “steepest descent”, because the gradient is the direction of steepest change at a given point.

### 2.1.2 Choosing the stepsize

There are numerous techniques used to select the stepsize<sup>2</sup>. We review the most general below, but point out that those are generally combined with knowledge about the problem in practice.

**Constant stepsize** One possible strategy is to maintain a constant step size throughout the entire algorithmic run, i. e. set  $\alpha_k = \alpha > 0$ . If the budget allows for it, several values of  $\alpha$  can be tested for comparison. Under regularity assumptions on  $f$ , one can guarantee that there exists a value below which a constant stepsize will lead to complexity guarantees (see Section 2.1.3). For instance, when  $f \in \mathcal{C}_L^{1,1}(\mathbb{R}^d)$ , the choice

$$\alpha_k = \alpha = \frac{1}{L} \quad (2.1.4)$$

leads to such guarantees. Because of its dependence in  $L$ , this choice is tailored to the problem at hand. Note that the rule (2.1.4) requires knowledge of the Lipschitz constant, but this information may not be available in practice.

**Decreasing stepsize** Another popular choice consist in choosing the entire sequence  $\{\alpha_k\}$  in advance so as to guarantee that  $\alpha_k \rightarrow 0$  as  $k \rightarrow \infty$ . This also enables the derivation of theoretical results, under some conditions that can help designing the formula for the  $\alpha_k$ s. However, this process forces the steps to get increasingly smaller, which may prevent fast progress towards the end of the algorithm.

**Adaptive choice with line search** Line-search techniques have been widely used in continuous optimization: at every iteration, they aim at computing the value of  $\alpha_k$  that leads to the largest decrease in the function value in the direction  $-\nabla f(\mathbf{w}_k)$ . In general, such exact line searches are not practical, and thus an inexact process is preferred. The most popular method is backtracking, that proceeds by testing a set of decreasing values: a simple version of a backtracking line search is described in Algorithm 2.

---

**Algorithm 2:** Basic backtracking line search in direction  $d$ .

---

**Inputs:**  $\mathbf{w} \in \mathbb{R}^d$ ,  $\mathbf{d} \in \mathbb{R}^d$ ,  $\alpha_0 \in \mathbb{R}^d$ .  
**Initialization:** Set  $\alpha = \alpha_0$  and  $j = 0$ .  
**while**  $f(\mathbf{w} + \alpha_j \mathbf{d}) > f(\mathbf{w})$  **do**  
  | Set  $\alpha_j = \frac{\alpha_j}{2}$  and  $j = j + 1$ .  
**end**  
**Output:**  $\alpha_j$ .

---

We can thus incorporate this line-search technique in step 2 of Algorithm 1 by calling the method with  $\mathbf{w} = \mathbf{w}_k$ ,  $\mathbf{d} = -\nabla f(\mathbf{w}_k)$  and (for instance)  $\alpha_0 = 1$ . Many variants can be build upon this simple framework. One drawback of line-search methods is that they require to evaluate the objective function, which can be deemed too expensive in certain applications.

---

<sup>2</sup>Or *learning rate* in machine learning.

### 2.1.3 Convergence rate analysis of gradient descent

In this section, we present several convergence rates for gradient descent, in the case of a smooth objective function. We will see that the nonconvex, convex and strongly convex cases exhibit different behavior.

**Proposition 2.1.1** Consider the  $k$ -th iteration of Algorithm 1 applied to  $f \in \mathcal{C}_L^{1,1}(\mathbb{R}^d)$ , and suppose that  $\nabla f(\mathbf{w}_k) \neq 0$ . Then, if  $0 < \alpha_k < \frac{2}{L}$ , we have

$$f(\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k).$$

In particular, choosing  $\alpha_k = \frac{1}{L}$  leads to

$$f(\mathbf{w}_k - \frac{1}{L} \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2. \quad (2.1.5)$$

**Proof.** We use the inequality (1.3.2) with the vectors  $(\mathbf{w}_k, \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k))$  :

$$\begin{aligned} f(\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)) &\leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top [-\alpha_k \nabla f(\mathbf{w}_k)] + \frac{L}{2} \|\alpha_k \nabla f(\mathbf{w}_k)\|^2 \\ &= f(\mathbf{w}_k) - \alpha_k \nabla f(\mathbf{w}_k)^\top \nabla f(\mathbf{w}_k) + \frac{L}{2} \alpha_k^2 \|\nabla f(\mathbf{w}_k)\|^2 \\ &= f(\mathbf{w}_k) + \left(-\alpha_k + \frac{L}{2} \alpha_k^2\right) \|\nabla f(\mathbf{w}_k)\|^2. \end{aligned}$$

If  $-\alpha_k + \frac{L}{2} \alpha_k^2 < 0$ , the second term on the right-hand side will be negative, thus we will have  $f(\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k)$ . Since  $-\alpha_k + \frac{L}{2} \alpha_k^2 < 0 \Leftrightarrow \alpha_k < \frac{2}{L}$  and  $\alpha_k > 0$  by definition, this proves the first part of the result.

To obtain (2.1.5), one simply needs to use  $\alpha_k = \frac{1}{L}$  in the series of equations above.  $\square$

The result of Proposition 2.1.1 will be instrumental to obtain complexity guarantees on Algorithm 1 in three different settings (nonconvex, convex, strongly convex): this analysis will be performed under the following assumption.

**Assumption 2.1.1** The objective function  $f$  belongs to  $\mathcal{C}_L^{1,1}(\mathbb{R}^d)$  for  $L > 0$  and there exists  $f_{low} \in \mathbb{R}$  such that for every  $\mathbf{w} \in \mathbb{R}^d$ ,  $f(\mathbf{w}) \geq f_{low}$  (i. e.  $f$  is bounded below on  $\mathbb{R}^d$ ).

**Nonconvex case** In the nonconvex case, we aim at bounding the number of iterations required to drive the gradient norm below some threshold  $\epsilon > 0$ : this means that we should be able to show that the gradient norm actually goes below this threshold, which is a guarantee of convergence.

**Theorem 2.1.1 (Complexity of gradient descent for nonconvex functions)** Let  $f$  be a nonconvex function satisfying Assumption 2.1.1. Suppose that Algorithm 1 is applied with  $\alpha_k = \frac{1}{L}$ . Then, for any  $K \geq 1$ , we have

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right). \quad (2.1.6)$$

**Proof.** Let  $K$  be an iteration index such that for every  $k = 0, \dots, K-1$ , we have  $\|\nabla f(\mathbf{w}_k)\| > \epsilon$ . From Proposition 2.1.1, we have that

$$\forall k = 0, \dots, K-1, \quad f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 \leq f(\mathbf{w}_k) - \frac{1}{2L} \left(\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\|\right)^2.$$

By summing across all such iterations, we obtain :

$$\sum_{k=0}^{K-1} f(\mathbf{w}_{k+1}) \leq \sum_{k=0}^{K-1} f(\mathbf{w}_k) - \frac{K}{2L} \left( \min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \right)^2.$$

Removing identical terms on both sides yields

$$f(\mathbf{w}_K) \leq f(\mathbf{w}_0) - \frac{K}{2L} \left( \min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \right)^2.$$

Using  $f(\mathbf{w}_K) \geq f_{low}$  (which holds by Assumption 2.1.1) and re-arranging the terms leads to

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \left[ \frac{2L(f(\mathbf{w}_0) - f_{low})}{K} \right]^{1/2} = \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

□

Equivalently, we say that the worst-case complexity of gradient descent is  $\mathcal{O}(\epsilon^{-2})$ , because for any  $\epsilon > 0$ , a reasoning similar to the proof of Theorem 2.1.1 guarantees that  $\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \epsilon$  after at most

$$\lceil 2L(f(\mathbf{w}_0) - f_{low})\epsilon^{-2} \rceil = \mathcal{O}(\epsilon^{-2})$$

iterations.

**Convex/Strongly convex case** In addition to Assumption 2.1.1, if we further assume that the objective is convex or strongly convex, we can show that stronger guarantees than that of the nonconvex case can be obtained at a lower cost. This improvement illustrates the interest of convex functions in optimization.

In this paragraph, we let  $f^* = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$  denote the minimal value of  $f$  (note that  $f^* \geq f_{low}$ ) and we assume that there exists  $\mathbf{w}^* \in \mathbb{R}^d$  such that  $f(\mathbf{w}^*) = f^*$  (i.e. the set of minima is not empty). Given an accuracy threshold  $\epsilon > 0$ , we are interested in bounding the number of iterations necessary to reach an iterate  $\mathbf{w}_k$  such that  $f(\mathbf{w}_k) - f^* \leq \epsilon$ .

**Theorem 2.1.2** *Convergence of gradient descent for convex functions* Let  $f$  be a convex function satisfying Assumption 2.1.1. Suppose that Algorithm 1 is applied with  $\alpha_k = \frac{1}{L}$ . Then, for any  $K \geq 1$ , the iterate  $\mathbf{w}_K$  satisfies

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K}\right). \quad (2.1.7)$$

method runs for at most  $\mathcal{O}(\epsilon^{-1})$  iterations before computing  $\mathbf{w}_k$  such that  $f(\mathbf{w}_k) - f^* \leq \epsilon$ .

**Proof.** Let  $K$  be an index such that for every  $k = 0, \dots, K-1$ ,  $f(\mathbf{w}_k) - f^* > \epsilon$ .

For any  $k = 0, \dots, K-1$ , the characterization of convexity (1.3.10) at  $\mathbf{w}_k$  and  $\mathbf{w}^*$  gives

$$f(\mathbf{w}^*) \geq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}^* - \mathbf{w}_k).$$

Combining this property with (2.1.5), we obtain:

$$\begin{aligned} f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 \\ &\leq f(\mathbf{w}^*) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_k - \mathbf{w}^*) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2. \end{aligned}$$



To proceed onto the next step, one notices that

$$\nabla f(\mathbf{w}_k)^\top (\mathbf{w}_k - \mathbf{w}^*) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 = \frac{L}{2} \left( \|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_k - \mathbf{w}^* - \frac{1}{L} \nabla f(\mathbf{w}_k)\|^2 \right).$$

Thus, recalling that  $\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{1}{L} \nabla f(\mathbf{w}_k)$ , we arrive at

$$\begin{aligned} f(\mathbf{w}_{k+1}) &\leq f(\mathbf{w}^*) + \frac{L}{2} \left( \|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_k - \mathbf{w}^* - \frac{1}{L} \nabla f(\mathbf{w}_k)\|^2 \right) \\ &= f(\mathbf{w}^*) + \frac{L}{2} (\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2). \end{aligned}$$

Hence,

$$f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \leq \frac{L}{2} (\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2). \quad (2.1.8)$$

By summing (2.1.8) on all indices  $k$  between 0 and  $K - 1$ , we obtain

$$\sum_{k=0}^{K-1} f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \leq \frac{L}{2} (\|\mathbf{w}_0 - \mathbf{w}^*\|^2 - \|\mathbf{w}_K - \mathbf{w}^*\|^2) \leq \frac{L}{2} \|\mathbf{w}_0 - \mathbf{w}^*\|^2.$$

Finally, using  $f(\mathbf{w}_0) \geq f(\mathbf{w}_1) \geq \dots \geq f(\mathbf{w}_K)$  (a consequence of Proposition 2.1.1, we obtain that

$$\sum_{k=0}^{K-1} f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \geq K (f(\mathbf{w}_K) - f^*).$$

Injecting this formula into the previous equation finally yields the desired outcome:

$$f(\mathbf{w}_k) - f(\mathbf{w}^*) \leq \frac{L \|\mathbf{w}_0 - \mathbf{w}^*\|^2}{2} \frac{1}{K}.$$

□

Equivalently, we say that the worst-case complexity of gradient descent is  $\mathcal{O}(\epsilon^{-1})$ , which means here that there exist a positive constant  $C$  (that depends on  $\|\mathbf{w}_0 - \mathbf{w}^*\|$  and  $L$ ) such that

$$f(\mathbf{w}_K) - f_{low} \leq \epsilon.$$

after at most  $C\epsilon^{-1}$  iterations.

We now turn to the strongly convex case.

**Theorem 2.1.3** *Convergence of gradient descent for strongly convex functions* Let  $f$  be a  $\mu$ -strongly convex function satisfying Assumption 2.1.1, with  $\mu \in (0, L]$ . Suppose that Algorithm 1 is applied with  $\alpha_k = \frac{1}{L}$  and let  $\epsilon > 0$ . Then, for any  $K \in \mathbb{N}$ , we have

$$f(\mathbf{w}_k) - f^* \leq \mathcal{O} \left( \left(1 - \frac{\mu}{L}\right)^k \right) \quad (2.1.9)$$

for at most  $\mathcal{O} \left( \frac{L}{\mu} \ln \left( \frac{1}{\epsilon} \right) \right)$  iterations before computing  $\mathbf{w}_k$  such that  $f(\mathbf{w}_k) - f^* \leq \epsilon$ .

Equivalently, we say that the convergence rate of gradient descent is  $\mathcal{O} \left( \left(1 - \frac{\mu}{L}\right)^k \right)$ .

**Proof.** We exploit the strong convexity property (1.3.12). For any  $(\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^n)^2$ , we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Minimizing both sides with respect to  $\mathbf{y}$  lead to  $\mathbf{y} = \mathbf{w}^*$  on the left-hand side, and  $\mathbf{y} = \mathbf{x} - \frac{1}{\mu} \nabla f(\mathbf{x})$  on the right-hand side (see Example ??). As a result, we obtain

$$\begin{aligned} f^* &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top \left[ -\frac{1}{\mu} \nabla f(\mathbf{x}) \right] + \frac{\mu}{2} \left\| -\frac{1}{\mu} \nabla f(\mathbf{x}) \right\|^2 \\ f^* &\geq f(\mathbf{x}) - \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|^2. \end{aligned}$$

By re-arranging the terms, we arrive at

$$\|\nabla f(\mathbf{x})\|^2 \geq 2\mu [f(\mathbf{x}) - f^*], \tag{2.1.10}$$

which is valid for any  $\mathbf{x} \in \mathbb{R}^n$ . Using (2.1.10) together with (2.1.5) thus gives

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2 \leq f(\mathbf{w}_k) - \frac{\mu}{L} (f(\mathbf{w}_k) - f^*).$$

This leads to

$$f(\mathbf{w}_{k+1}) - f^* \leq \left(1 - \frac{\mu}{L}\right) (f(\mathbf{w}_k) - f^*),$$

which we can iterate in order to obtain

$$f(\mathbf{w}_K) - f^* \leq \left(1 - \frac{\mu}{L}\right)^K (f(\mathbf{w}_0) - f^*).$$

It then suffices to note that the bound is also valid for  $K = 0$ . □

Equivalently, we can show a worst-case complexity result: the method computes  $\mathbf{w}_k$  such that  $f(\mathbf{w}_k) - f^* \leq \epsilon$  in at most  $\mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right)$  iterations.

Similar results can be shown for the criterion  $\|\mathbf{w}_k - \mathbf{w}^*\|$ : in other words, the distance between the current iterate and the (unique) global optimum decreases at a rate  $\mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^k\right)$ .

**Remark 2.1.1** *Proofs of convergence rates are typically more technical for convex and strongly convex problems: in order to obtain better bounds than in the nonconvex setting, one must make careful use of the (strong) convexity inequalities. In this course, we do not focus on these aspects, but rather draw insights from the final complexity bounds or convergence rates.*

### 2.1.4 Application: regression with logistic and sigmoid losses

As in Section ??, we consider a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where  $\mathbf{x}_i \in \mathbb{R}^d$  are feature vectors, and the  $y_i$  represent binary labels. We wish to build a linear classifier  $\mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x}$  to perform this classification, i. e. identify the correct label from the feature.

**Logistic loss** We first suppose that  $y_i \in \{-1, +1\}$ ; to model these discrete-valued labels, we introduce an *odds-like* function

$$p(\mathbf{x}; \mathbf{w}) = (1 + e^{\mathbf{x}^T \mathbf{w}})^{-1} \in (0, 1).$$

Given this function, our goal is to choose the model  $\mathbf{w}$  such that

$$\begin{cases} p(\mathbf{x}_i; \mathbf{w}) \approx 1 & \text{if } y_i = +1; \\ p(\mathbf{x}_i; \mathbf{w}) \approx 0 & \text{if } y_i = -1. \end{cases}$$

Given this goal, we want to build an objective function that measures the error between our model and the labels according to the property above. Therefore, we penalize situations in which  $y_i = +1$  and  $p(\mathbf{x}_i; \mathbf{w})$  is close to 0, or  $y_i = -1$  and  $p(\mathbf{x}_i; \mathbf{w})$  is close to 1. This results in the so-called logistic loss, which is a function from  $\mathbb{R}^d$  to  $\mathbb{R}$  defined by

$$\forall \mathbf{w} \in \mathbb{R}^d, f(\mathbf{w}) = \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\}. \quad (2.1.11)$$

The motivation behind introducing the logarithm of the function  $p$  is twofold. On the one hand, it provides a statistical interpretation of the loss as a joint distribution; on the other hand, the derivatives of this function have a more favorable structure.

Given this objective function, the **logistic regression** problem is given by

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\} \quad (2.1.12)$$

This is a convex, smooth problem (though not a strongly convex one), that can be made strongly convex by adding a regularizing term, which will be done in a subsequent chapter. In both cases, we can apply gradient descent with guaranteed convergence rates.

**Sigmoid loss** We now assume that  $y_i \in \{0, 1\}$  for every  $i$ . In this case, and for similar reasons than in the case of the logistic loss, we can measure agreement between the model and the label for example  $i$  by looking at the sigmoid function

$$\phi(\mathbf{x}_i; \mathbf{w}) = \left(1 + e^{-\mathbf{x}_i^T \mathbf{w}}\right)^{-1};$$

Drawing inspiration from Section ??, we may want to penalize the average of the squared errors  $(y_i - \phi(\mathbf{x}_i; \mathbf{w}))^2$ . This is the philosophy behind the nonlinear regression problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \left( y_i - \frac{1}{1 + e^{-\mathbf{x}_i^T \mathbf{w}}} \right)^2. \quad (2.1.13)$$

This problem is a *nonlinear* least-squares problems: it is twice continuously differentiable, but non-convex. Therefore, we can apply gradient descent to this problem, but we will only be guaranteed to reach a first-order stationary point.

## 2.2 Acceleration techniques

### 2.2.1 Introduction: the concept of momentum

In Section 2.1.3, we derive complexity bounds for the gradient descent algorithm, and we saw in particular that assuming that the function was convex (respectively, strongly convex) improved the complexity. These results are called *upper* complexity bounds, in the sense that they reflect the worst possible convergence rate that this algorithm could exhibit on a given problem. The issue of *lower* bounds, that show a rate that cannot be improved upon, has been the subject to a lot of attention, particularly in the convex optimization community.

For nonconvex optimization, it is known that there exists a function for which gradient descent converges exactly at the  $\mathcal{O}(\frac{1}{\sqrt{K}})$  rate: in this case, the lower bound matches the upper bound. On the contrary, for convex functions, the lower bound is actually  $\mathcal{O}(\frac{1}{K^2})$ , which is a sensible improvement over the bound in  $\mathcal{O}(\frac{1}{K})$  of Theorem 2.1.2. There are methods that can achieve this bound, thanks to an algorithmic technique called **acceleration**.

The underlying idea of acceleration is that, at a given iteration and given the available information from previous iterations (in particular, the latest displacement), one can move along a better step than that given by the current gradient.

### 2.2.2 Nesterov's accelerated gradient method

Among the existing methods based on acceleration, the accelerated gradient algorithm proposed by Yurii Nesterov in 1983 is the most famous, to the point that it has been termed “Nesterov’s algorithm”.

---

**Algorithm 3:** Accelerated gradient method.

---

**Initialization:**  $w_0 \in \mathbb{R}^d$ ,  $w_{-1} = w_0$ .

**for**  $k = 0, 1, \dots$  **do**

1. Compute a steplength  $\alpha_k > 0$  and a parameter  $\beta_k > 0$ .
2. Compute the new iterate as

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k + \beta_k(w_k - w_{k-1})) + \beta_k(w_k - w_{k-1}). \quad (2.2.1)$$

**end**

---

Algorithm 3 provides a description of the method. Like the gradient descent method of Section 2.1, it requires a single gradient calculation per iteration; however, unlike in gradient descent, the gradient is not evaluated at the current iterate  $w_k$ , but at a combination of this iterate with the previous step  $w_k - w_{k-1}$ : this term is called the **momentum term**, and is key to the performance of accelerated gradient techniques.

Another view of the accelerated gradient descent is that of a two-loop recursion: given  $w_0$  and

$\mathbf{z}_0 = \mathbf{w}_0$ , the update (2.2.1) can be rewritten as

$$\begin{cases} \mathbf{w}_{k+1} &= \mathbf{z}_k - \alpha_k \nabla f(\mathbf{z}_k) \\ \mathbf{z}_{k+1} &= \mathbf{w}_{k+1} + \beta_{k+1}(\mathbf{w}_{k+1} - \mathbf{w}_k). \end{cases} \quad (2.2.2)$$

This formulation decouples the two steps behind the accelerated gradient update: a gradient step on  $\mathbf{z}_k$ , combined with a momentum step on  $\mathbf{w}_{k+1}$ .

**Choosing the parameters** We now comment on the choice of the stepsize  $\alpha_k$  and the momentum parameter  $\beta_k$ . The same techniques than those presented in Section 2.1.2 can be considered for the choice of  $\alpha_k$  (stepsize parameter). As in the gradient descent case, the choice  $\alpha_k = \frac{1}{L}$  is a standard one.

The choice of  $\beta_k$  is most crucial to obtaining the improved complexity bound. The standard values proposed by Nesterov depend on the nature of the objective function:

- If  $f$  is a  $\mu$ -strongly convex, we set

$$\beta_k = \beta = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} \quad (2.2.3)$$

for every  $k$ . Note that this requires the knowledge of both the Lipschitz constant of the gradient and the strong convexity constant.

- For a general convex function  $f$ ,  $\beta_k$  is computed in an adaptive way using two sequences, as follows:

$$t_{k+1} = \frac{1}{2}(1 + \sqrt{1 + 4t_k^2}), t_0 = 0, \quad \beta_k = \frac{t_k - 1}{t_{k+1}}. \quad (2.2.4)$$

The following informal theorem summarizes the complexity results that can be proven for Algorithm 3.

**Theorem 2.2.1** Consider Algorithm 3 applied to a convex function  $f$  satisfying Assumption 2.1.1, with  $\alpha_k = \frac{1}{L}$ , and let  $\epsilon > 0$ . Then, for any  $K \geq 1$ , the iterate  $\mathbf{w}_K$  computed by Algorithm 3 satisfies

- $f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K^2}\right)$  for a generic convex function if  $\beta_k$  is set according to the adaptive rule (2.2.4);
- At most  $f(\mathbf{w}_K) - f^* \leq \left(1 - \sqrt{\frac{\mu}{L}}\right)^K$  for a  $\mu$ -strongly convex function, provided  $\beta_k$  is set to the constant value given by (2.2.3).

Note that we can also derive worst-case complexity bounds for the accelerated gradient method, that show the same improvement. For instance, for strongly convex functions, we can establish that  $f(\mathbf{w}_k) - f^* \leq \epsilon$  after at most  $\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \ln(\epsilon^{-1})\right) \mathcal{O}\left(\frac{L}{\mu} \ln(\epsilon^{-1})\right)$ .

### 2.2.3 Other accelerated methods

**Heavy ball method** The heavy ball method is a precursor of the accelerated gradient algorithm, that was proposed by Boris T. Polyak in 1964. Its  $k$ -th iteration can be written as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \nabla f(\mathbf{w}_k) + \beta(\mathbf{w}_k - \mathbf{w}_{k+1}),$$

where the stepsize and momentum parameters are chosen to be constant values. The key difference between this iteration and Nesterov's lies in the gradient evaluation, which the heavy ball method performs at the current point: in that sense, the heavy ball method performs first the gradient update, then the momentum step, while Nesterov's method adopts the inverse approach. This method achieves the optimal rate of convergence on strongly convex quadratic functions, but can fail on general strongly convex functions.

**Conjugate gradient** The (linear) conjugate gradient method, proposed by Hestenes and Stiefel in 1952, has remained to this day one of the preferred methods to solve linear systems of equations and strongly convex quadratic minimization problems. Unlike Polyak's method, the conjugate gradient algorithm does not require knowledge of the Lipschitz constant  $L$  nor the parameter  $\mu$ , because it exploits knowledge from the past iterations. The  $k$ -th iteration of conjugate gradient can be written as:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha_k p_k, \quad p_k = -\nabla f(x_k) + \beta_k p_{k-1}.$$

In a standard conjugate gradient algorithm,  $\alpha_k$  and  $\beta_k$  are computed using formulas tailored to the problem: this contributes to their convergence rate analysis, which leads to a rate similar to that of accelerated gradient. However, unlike accelerated gradient, the conjugate gradient is guaranteed to terminate after  $d$  iterations on a  $d$ -dimensional problem. When  $d$  is very large, the bound for conjugate gradient matches that of the other methods, and in that sense does not depend on the problem dimension.

**Example 2.2.1 (Strongly convex quadratic minimization)** *A strongly convex quadratic minimization problem is an optimization problem of the form*

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad q(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} - \mathbf{b}^T \mathbf{w}$$

where  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is a symmetric positive definite matrix and  $\mathbf{b} \in \mathbb{R}^d$ . This problem is smooth (because the objective is polynomial in all of the decision variables) and  $\nabla^2 f(\mathbf{w}) \succ \mathbf{0}$  for every  $\mathbf{w}$ , meaning that the problem is  $\mu$ -strongly convex with  $\mu$  denoting the minimum eigenvalue of  $\mathbf{A}$ . As a result, there exist a unique global minimum given by the solution of  $\nabla q(\mathbf{w}) = \mathbf{A} \mathbf{w} - \mathbf{b} = \mathbf{0}$ . This equation is a linear system but the cost of inverting this system and computing a solution can be prohibitive. For this reason, one can replace the exact solve by an iterative, gradient-based approach, and apply Algorithm 1 or Algorithm 3. Note that  $q \in \mathcal{C}_{\|\mathbf{A}\|}^{1,1}(\mathbb{R}^d)$ , hence the choice of steplength 2.1.4 is a valid one.

If gradient descent is applied, then an  $\epsilon$ -accuracy in the objective value can be reached in at most  $\mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right)$  iterations, while if one applies the accelerated gradient or the heavy ball method with appropriately chosen parameters, this bound improves to  $\mathcal{O}\left(\frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right)$ . Finally, if we aim at using conjugate gradient, the result bound will be in  $\mathcal{O}\left(\min\left\{d, \frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\right\}\right)$ .

## 2.3 Conclusion

The most classical optimization problems involve linear algebra: this is the case for linear least squares as well as eigenvalue and singular value calculations, that can be viewed as solutions of optimization problems. For these problems, it is possible to compute the solution explicitly (or in closed form). Linear least squares is a particular case of such instances.

For general unconstrained optimization problems, it is not possible to obtain a closed-form expression of the solution(s). As a result, one must construct algorithms that proceed iteratively to move from a starting point towards a solution. The gradient descent method is the canonical example of such a framework: many variants have been built on this paradigm, especially regarding the choice of the stepsize (or learning rate in machine learning applications). To analyze the behavior of gradient descent, one can establish global convergence rates (or, equivalently, global complexity bounds) that can be refined depending on the nature of the objective function. Indeed, gradient descent can be shown to converge faster on convex problems than on nonconvex ones, and even faster on strongly convex problems.

A natural question arising from these convergence rates results is whether those are optimal. For gradient descent applied to nonconvex, differentiable functions, it is not possible to improve over the rate established in Section 2.1.3. However, one can design accelerated methods for strongly convex and convex functions that possess better rates, a fact that reflects on the practical performance. These methods all rely on the concept of momentum, which is also exploited in state-of-the-art algorithms used to learn complex models in machine learning (e. g. Adagrad).

## Chapter 3

# Regularization

In this chapter, we investigate several challenges that can be posed while trying to apply stochastic gradient techniques to machine learning problems. To motivate these issues further, we will begin with an introductory example and method.

### 3.1 Introduction : The perceptron algorithm

Recall that in section 1.1, we introduced a linear SVM problem of the following form :

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \max\{1 - y_i \mathbf{x}_i^T \mathbf{w}, 0\} + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (3.1.1)$$

where  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  represents the dataset, and  $\lambda > 0$ .

One of the earliest methods that was proposed to solve this algorithm is the **perceptron algorithm**, given in Algorithm 4.

---

**Algorithm 4:** Perceptron algorithm for problem 3.1.2.

---

**Initialization:**  $\mathbf{w}_0 \in \mathbb{R}^d$ ,  $\alpha > 0$ .

**for**  $k = 0, 1, \dots$  **do**

1. Draw an index  $i_k \in \{1, \dots, n\}$  at random.

2. Compute the new iterate as

$$\mathbf{w}_{k+1} = \left(1 - \frac{\alpha\lambda}{n}\right) \mathbf{w}_k + \begin{cases} \alpha y_{i_k} \mathbf{x}_{i_k} & \text{if } 1 - y_{i_k} \mathbf{x}_{i_k}^T \mathbf{w}_k > 0 \\ 0 & \text{otherwise,} \end{cases} \quad (3.1.2)$$

**end**

---

In its basic form, the perceptron algorithm is quite similar to stochastic gradient with a constant step size, in that it selects a single sample at every iteration and performs an update based on this



value. In fact, this would be exactly the stochastic gradient method if the problem were

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (1 - y_i \mathbf{x}_i^T \mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

However, this choice of loss function would not satisfy our desired requirements (see section 1.1). The *hinge loss* is a more meaningful quantity, however it is **nonsmooth**, i.e. the gradient does not exist at every point. In this situation, and with structured functions such as the hinge loss, it is possible to define quantities that act as a proxy for the gradient, and can thus drive the optimization process : we detail these aspects in Section 3.2.

Another interesting property of the problem (3.1.1) is that the objective function involves two terms: the hinge loss term, which depends on the data and a regularizing term, which does not depend on the data and serves to enforce structural properties on the solution. We will address this topic and the associated algorithms in Section 3.3.

## 3.2 Nonsmooth optimization

### 3.2.1 From nonsmooth functions to nonsmooth problems

Problems such as (3.1.1), that involve a function possibly not differentiable, are termed *nonsmooth problems*. They involve functions that we will call nonsmooth (by opposition with smooth) : for the purpose of these notes, we will define nonsmooth functions as follows.

**Definition 3.2.1 (Nonsmooth functions)** A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is called **nonsmooth** if it is not differentiable everywhere.

**Remark 3.2.1** A nonsmooth function can be continuous (this is the case for the hinge loss above).

**Example 3.2.1** Examples of nonsmooth functions

- $w \mapsto |w|$  from  $\mathbb{R}$  to  $\mathbb{R}$ ;
- $w \mapsto \|w\|_1$  from  $\mathbb{R}^d$  to  $\mathbb{R}$ ;
- *ReLU*:  $w \mapsto \max\{w, 0\}$  from  $\mathbb{R}^d$  to  $\mathbb{R}$ .

Since nonsmooth functions are not differentiable everywhere, optimization problems that involve nonsmooth functions may be impossible to solve via gradient-based methods. Still, several approaches can be used to tackle these problems.

One useful technique consists in reformulating a nonsmooth problem as a smooth one when possible. For instance, the problem  $\min_{w \in \mathbb{R}} |w|$  is equivalent to

$$\min_{w, t^+, t^- \in \mathbb{R}} t^+ + t^- \quad \text{s. t.} \quad w = t^+ - t^-, t^+ \geq 0, t^- \geq 0.$$

This reformulation is a smooth problem involving only linear objective and constraints, which is easily solvable by smooth solvers.

Another technique, frequently employed in practice, consists in working with functions that are nonsmooth but Lipschitz continuous (denoted by  $\mathcal{C}_L^{0,0}$ , by analogy with  $\mathcal{C}_L^{1,1}$ ) and using a gradient-based scheme. This approach is motivated by the following property.

**Theorem 3.2.1** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a Lipschitz continuous function. Then it is differentiable at almost every point in  $\mathbb{R}^d$ .

For instance, the ReLU function is Lipschitz continuous (not differentiable at 0) thus most constructions involving ReLU (such as neural networks) would not be differentiable everywhere. However, most algorithms will operate under the assumption that the function is indeed differentiable. This is the case for most points (in fact, almost every point), but nonsmooth functions are likely to be non-differentiable at their minima, should they possess one.

### 3.2.2 Subgradient methods

In the case of convex functions, one can define a proxy for the gradient called the subgradient.

**Definition 3.2.2 (Subgradient and subdifferential)** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a convex function. A vector  $\mathbf{g} \in \mathbb{R}^d$  is called a *subgradient* of  $f$  at  $\mathbf{w} \in \mathbb{R}^d$  if

$$\forall \mathbf{z} \in \mathbb{R}^n, \quad f(\mathbf{z}) \geq f(\mathbf{w}) + \mathbf{g}^T(\mathbf{z} - \mathbf{w}).$$

The set of all subgradients of  $f$  at  $\mathbf{w}$  is called the *subdifferential* of  $f$  at  $\mathbf{w}$ , and denoted by  $\partial f(\mathbf{w})$ .

Note that when the function  $f$  is differentiable at  $\mathbf{w}$ , we have  $\partial f(\mathbf{w}) = \{\nabla f(\mathbf{w})\}$ , thus the notion of subdifferential matches that of the gradient for differentiable functions.

The interest of subgradients is further illustrated by the following result.

**Theorem 3.2.2** Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a convex function, and  $\mathbf{w} \in \mathbb{R}^d$ .

$$\mathbf{0} \in \partial f(\mathbf{w}) \Leftrightarrow \mathbf{w} \text{ minimum of } f.$$

**Example 3.2.2** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $f(w) = |w|$ .

$$\partial f(w) = \begin{cases} -1 & \text{if } w < 0 \\ 1 & \text{if } w > 0 \\ [-1, 1] & \text{if } w = 0. \end{cases}$$

The set  $[-1, 1]$  contains 0, which confirms that  $w^* = 0$  is the minimum of  $f$ .

**Remark 3.2.2** Subgradients can also be defined for nonconvex functions, however in that case the subdifferential may be empty (typically at local maxima of the function).

By analogy with gradient descent, we can design a subgradient method, as shown by Algorithm 5.

Such a method offers a flexibility in choosing the subgradient, which can be an issue. Moreover, choosing the stepsize is more difficult than for gradient descent, due to the nonsmooth nature of the problem. In fact, a subgradient can lead to increase in the function value for any stepsize, hence the choice of subgradient is critical to the success of this method.

**Variants of subgradient method** Based on the existing variants on the gradient descent paradigm, one can build algorithms that incorporate momentum and/or stochastic aspects; however, their analysis is also more intricate.

**Algorithm 5:** Subgradient descent method.

---

**Initialization:**  $\mathbf{w}_0 \in \mathbb{R}^d$ .  
**for**  $k = 0, 1, \dots$  **do**  
    1. Compute a subgradient  $\mathbf{g}_k \in \partial f(\mathbf{w}_k)$ .  
    2. Compute a steplength  $\alpha_k > 0$ .  
    3. Set  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{g}_k$ .  
**end**

---

### 3.3 Regularization

#### 3.3.1 Regularized problems

As we mentioned in introduction, a common practice in machine learning problems consists in enforcing a specific structure of the machine learning model **through the objective function**. Such regularized problems have the following form :

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \underbrace{f(\mathbf{w})}_{\text{loss function}} + \underbrace{\lambda \Omega(\mathbf{w})}_{\text{regularization term}} .$$

where  $\lambda > 0$  is called a regularization parameter.

**Example 3.3.1 (Ridge regularization)** A problem with ridge regularization has the following form:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 .$$

The ridge regularizer  $\mathbf{w} \mapsto \frac{1}{2} \|\mathbf{w}\|^2$  has several interpretations. It effectively penalizes  $\mathbf{w}$ s with large components, and can be shown to be equivalent to a constraint on the squared norm  $\|\mathbf{w}\|^2$ . In addition, a ridge regularizer has the effect to reduce the variance of the problem solution with respect to the data. Finally, when the regularizer  $\lambda > 0$  is big enough, this often turns the objective function into a strongly convex one, with the positive implications in terms of convergence speed and uniqueness of the (global) minimum.

#### 3.3.2 Sparsity-inducing regularizers

While computing a model to explain some data, we might want to compute a model that explains the data using as few features as possible<sup>1</sup>. Mathematically speaking, if our model is parameterized by a vector  $\mathbf{w} \in \mathbb{R}^d$ , our goal is to compute a vector that explains the data with as few nonzero coordinates as possible.

There exists a regularizer that penalized vectors with nonzero components (not just large as opposed to the ridge regularizer), called the  $\ell_0$  norm<sup>2</sup>. An  $\ell_0$ -regularized problem has the form

$$\underset{\mathbf{w}}{\text{minimize}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_0, \quad \|\mathbf{v}\|_0 = |\{i | [\mathbf{v}]_i \neq 0\}| .$$

<sup>1</sup>The goal of this process is *feature selection*.

<sup>2</sup>Though technically this function defines a semi-norm.

However, this function is nonsmooth and discontinuous; its combinatorial nature also introduces more complexity to the original problem. As a result, researchers have turned to an intermediate regularization term, the  $\ell_1$  norm defined by

$$\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|. \quad (3.3.1)$$

This function is continuous and convex; moreover, it is a norm function, which endows it with many desirable properties.

An illustration of this method is given below.

**Example 3.3.2** *LASSO (Least Absolute Shrinkage and Selection Operator)* Consider the setting of linear regression with data  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and  $\mathbf{y} \in \mathbb{R}^n$ . With an  $\ell_1$  regularizer, the problem becomes:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1.$$

The solution of this problem is known to possess fewer nonzero elements than the un-regularized, least-squares solution.

### 3.3.3 Proximal methods

Following our introduction of regularized problems in the previous section, we now describe optimization algorithms tailored to such formulations.

We begin by describing our problem class of interest.

**Definition 3.3.1 (Composite optimization)** A composite optimization problem is of the form:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \lambda \Omega(\mathbf{w}),$$

where  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is a smooth,  $C^{1,1}$  function,  $\lambda > 0$  and  $\Omega : \mathbb{R}^d \rightarrow \mathbb{R}$  is a convex, nonsmooth regularizer.

The **proximal approach** follows a classical optimization paradigm, in which a given problem is replaced by a sequence of easier problems called subproblems (note that all methods that we covered in these notes implicitly rely on these techniques). In the case of proximal methods, one aims at exploiting the smoothness of  $f$  to obtain easier problems, while using the structure of  $\Omega$  directly into the subproblems.

Algorithm 6 gives a sketch of a proximal gradient method. The cost of an iteration of this algorithm is clearly more than that of other methods we have seen so far, given that it includes a gradient calculation as well as solving an auxiliary optimization problem (3.3.2), called the **proximal subproblem**.

**Remark 3.3.1** If  $\Omega \equiv 0$  (i. e.  $\Omega$  is the zero function and the problem is un-regularized), one can show that the solution of (3.3.2) is given by

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k).$$

We thus recognize the gradient iteration of Algorithm 1.

**Algorithm 6:** Proximal gradient method.**Initialization:**  $\mathbf{w}_0 \in \mathbb{R}^d$ .**for**  $k = 0, 1, \dots$  **do**

1. Compute the gradient of the smooth part  $\nabla f(\mathbf{w}_k)$ .
2. Compute a steplength  $\alpha_k > 0$ .
3. Compute  $\mathbf{w}_{k+1}$  such that

$$\mathbf{w}_{k+1} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|_2^2 + \lambda \Omega(\mathbf{w}) \right\}. \quad (3.3.2)$$

**end**

Proximal gradient methods can be designed using most of the tools that can be applied to gradient descent : this includes stepsize choices, acceleration as well as stochastic aspects. Moreover, complexity results exist for nonconvex and convex  $f$ , though the latter has attracted more attention in the literature.

**Example of proximal method: ISTA** We end this section on proximal methods by a instance of Algorithm 6 that has proven successful in signal and image processing. This method is dedicated to solving problems with an  $\ell_1$  regularization term, of the form:

$$\operatorname{minimize}_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1.$$

Unlike for general regularizers, one can obtain a closed-form solution of the subproblem (3.3.2). Indeed, the proximal subproblem, given by

$$\operatorname{minimize}_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|_2^2 + \lambda \|\mathbf{w}\|_1 \right\},$$

has a unique solution. To obtain it, one computes the usual gradient step  $\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ , then one applies the **soft-thresholding function**  $s_{\alpha_k \lambda}(\bullet)$  to each component, where this function is given by

$$\forall \mu > 0, \forall t \in \mathbb{R}, \quad s_\mu(t) = \begin{cases} t + \mu & \text{if } t < -\mu \\ t - \mu & \text{if } t > \mu \\ 0 & \text{otherwise.} \end{cases}$$

As a result, the solution of the proximal subproblem is defined component-wise according to the components of the gradient step. The resulting update is at the heart of the corresponding proximal algorithm, called ISTA (Iterative Soft-Thresholding Algorithm): a description of ISTA is given in Algorithm 7.

It can be shown that the use of the soft-thresholding function does promote zero components in the new iterates, which results in sparser solutions at the end of the algorithmic run.

**Algorithm 7:** ISTA: Iterative Soft-Thresholding Algorithm.**Initialization:**  $\mathbf{w}_0 \in \mathbb{R}^d$ .**for**  $k = 0, 1, \dots$  **do**

1. Compute the gradient of the smooth part  $\nabla f(\mathbf{w}_k)$ .
2. Compute a steplength  $\alpha_k > 0$ .
3. Compute  $\mathbf{w}_{k+1}$  component-wise through the following rule

$$[\mathbf{w}_{k+1}]_i = \begin{cases} [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i + \alpha_k \lambda & \text{if } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i < -\alpha_k \lambda \\ [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i - \alpha_k \lambda & \text{if } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i > \alpha_k \lambda \\ 0 & \text{if } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i \in [-\alpha_k \lambda, \alpha_k \lambda]. \end{cases} \quad (3.3.3)$$

**end**

**Remark 3.3.2** A notable improvement on ISTA was the inclusion of momentum, which resulted in a new algorithm called FISTA (Fast ISTA): this method is now the most widely used instance of ISTA.

### 3.4 Conclusion

Nonsmoothness is a very common property in optimization, that can lead to mild or major challenges in implementing algorithms to minimize nonsmooth functions. In certain cases, the structure and the impact of nonsmoothness are well understood; in other cases, generalized notions of derivative such as subgradients may have to come into play.

Nonsmoothness frequently arises in regularized problem, where the goal is to enforce properties for a model, that do not depend on the data. The optimization schemes of choice for these problems are proximal gradient methods, that proceed by solving subproblems involving the regularizer. For instance, the  $\ell_1$  regularizer, that promotes sparsity of the solution, can be tackled using the ISTA method. Note that a regularizer need not be nonsmooth, in which case a classical gradient method could be applied. This is for instance the case with the  $\ell_2$  regularizer, that aims at reducing variance with respect to the data, and leads to a smooth, possibly strongly convex problem.

## Chapter 4

# Stochastic optimization methods

### 4.1 Motivation

In this chapter, we will leverage the structure inherent to data science problems. More formally, we suppose that we have access to data samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , that are drawn from an unknown distribution. As in the regression examples studied above, we seek a predictor function or a model  $h$  such that  $h(\mathbf{x}_i) \approx y_i$  for every  $i = 1, \dots, n$ . Rather than optimizing over a space of models, we assume that a given model is defined by means of a vector  $\mathbf{w} \in \mathbb{R}^d$  (i.e.  $h(\mathbf{x}_i) = h(\mathbf{w}; \mathbf{x}_i)$ ). Therefore, we only need to determine the vector  $\mathbf{w}$  in order to obtain the model.

To assess the accuracy of our model in predicting the data, we define a loss function, i.e. a mapping  $\ell : (h, y) \mapsto \ell(h, y)$ , that penalize pairs  $(h, y)$  such that  $h \neq y$ . We have already seen several examples of such losses (least-squares loss, sigmoid loss, etc). The loss at a given sample of the dataset thus is  $\ell(h(\mathbf{w}; \mathbf{x}_i), y_i)$ : in order to account for all samples, we consider the average of all losses as our objective to be minimized. This gives rise to the following optimization problem.

**Definition 4.1.1 (Finite-sum optimization problem)** *Given a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , a class of predictor functions  $\{h(\mathbf{w}; \cdot)\}_{\mathbf{w} \in \mathbb{R}^d}$  and a loss function  $\ell$ , we define the corresponding optimization problem:*

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{w}; \mathbf{x}_i), y_i) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}). \quad (4.1.1)$$

Suppose that we apply gradient descent (Algorithm 1) to that problem, assuming all  $f_i$  are differentiable. The  $k$ -th iteration of this method is

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k) = \mathbf{w}_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_k).$$

From this update, we see that one iteration of gradient descent requires to look over the **entire dataset** in order to compute the gradient vector. In a big data setting where the number of samples  $n$  is very large, this cost can be prohibitive.

**Remark 4.1.1** *In stochastic optimization, the data samples might be generated directly from the distribution, and be available in a streaming fashion. Instead of involving a discrete average on the*

sample, the resulting optimization problem would involve a mathematical expectation of the form

$$\min_{\mathbf{w} \in \mathbb{R}^d} \mathbb{E}_{(\mathbf{x}, y)} [f_{(\mathbf{x}, y)}(\mathbf{w})].$$

In such a context, the full gradient cannot be computed exactly. However, most of the reasoning of stochastic gradient will still be applicable.

## 4.2 Stochastic gradient algorithm

### 4.2.1 Algorithm

At its core, the idea of the stochastic gradient method is remarkably simple. Starting from the problem  $\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w})$ , and assuming each component function  $f_i$  is differentiable, the method picks an index  $i$  at random and takes a step in the direction of the negative gradient of the component function  $f_i$ .

---

**Algorithm 8:** Stochastic gradient method.

---

**Initialization:**  $\mathbf{w}_0 \in \mathbb{R}^d$ .

**for**  $k = 0, 1, \dots$  **do**

1. Compute a steplength  $\alpha_k > 0$ .
2. Draw a random index  $i_k \in \{1, \dots, n\}$ .
3. Compute the new iterate as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k). \quad (4.2.1)$$

**end**

---

The key motivation for this process is that using a single data point at a time results in updates that are  $n$  times cheaper than a full gradient step.

**Remark 4.2.1** In general, considering independent updates may not be desirable. Consider for instance the problem  $\min_{w \in \mathbb{R}} \frac{1}{2}(f_1(w) + f_2(w))$  with  $f_1(w) = 2w^2$  and  $f_2 = -w^2$ . Starting from  $w_k > 0$ , drawing  $i_k = 2$  will necessarily lead to an increase in the function value.

In finite-sum problems arising from machine learning, the data samples are correlated enough that an update according to one sample might lead to improvement with respect to other samples as well: this is a key reason for the success of stochastic gradient methods in this setting.

**Remark 4.2.2** Algorithm 8 is often referred to as Stochastic Gradient Descent, or SGD, by analogy with Gradient Descent. However, for the reason mentioned in the previous remark, the stochastic gradient algorithm is not a descent method in general (as we will see in the next section, it can however produce descent in expectation). In these notes, we will adopt the terminology *stochastic gradient*.



### 4.2.2 Analysis

We now describe the main arguments in deriving convergence rates for stochastic gradient, under a slightly modified version of Assumption 2.1.1.

**Assumption 4.2.1** *The objective function  $f = \frac{1}{n} \sum_{i=1}^n f_i$  belongs to  $\mathcal{C}_L^{1,1}(\mathbb{R}^d)$  for  $L > 0$  and there exists  $f_{low} \in \mathbb{R}$  such that for every  $\mathbf{w} \in \mathbb{R}^d$ ,  $f(\mathbf{w}) \geq f_{low}$ . Moreover, every function  $f_i$  belongs to  $\mathcal{C}^1(\mathbb{R}^d)$ .*

Recall that, for gradient descent, the key result was Proposition 2.1.1, which gave

$$f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^\top (\mathbf{w}_{k+1} - \mathbf{w}_k) + \frac{L}{2} \|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2.$$

A similar result can be shown for stochastic gradient under certain assumptions on how the random components are drawn. Those are summarized below.

**Assumption 4.2.2 (Assumptions on stochastic gradient)** *At any iteration of Algorithm 8 of index  $k$ , the index  $i_k$  is drawn independently from the previous indices  $i_0, \dots, i_{k-1}$  so that the following properties are satisfied:*

1.  $\mathbb{E}_{i_k} [\nabla f_{i_k}(\mathbf{w}_k)] = \nabla f(\mathbf{w}_k)$ ;
2.  $\mathbb{E}_{i_k} [\|\nabla f_{i_k}(\mathbf{w}_k)\|^2] \leq \sigma^2 + \|\nabla f(\mathbf{w}_k)\|^2$  with  $\sigma^2 > 0$ .

The first property of Assumption 4.2.2 forces the stochastic gradient  $\nabla f_{i_k}(\mathbf{w}_k)$  to be an unbiased estimate of the true gradient  $\nabla f(\mathbf{w}_k)$ . The second property controls the variance of the norm of this stochastic gradient, so as to control the variations in its magnitude due to noise. Several strategies can be designed to draw an index  $i_k$  that satisfies these properties, the most classical of which is given below.

**Example 4.2.1 (Uniform sampling)** *Suppose that the  $k$ -th iteration of stochastic gradient draws the index  $i_k$  uniformly at random in  $\{1, \dots, n\}$ . Then Algorithm 8 satisfies Assumption 4.2.2.*

**Proposition 4.2.1** *Under Assumptions 2.1.1 and 4.2.2, consider the  $k$ -th iteration of Algorithm 8. Then,*

$$\mathbb{E}_{i_k} [f(\mathbf{w}_{k+1})] - f(\mathbf{w}_k) \leq \nabla f(\mathbf{w}_k)^\top \mathbb{E}_{i_k} [\mathbf{w}_{k+1} - \mathbf{w}_k] + \frac{L}{2} \mathbb{E}_{i_k} [\|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2].$$

A stochastic gradient update will thus lead to decrease **in expectation**. Such a property suffices to derive convergence rates (or complexity results) for stochastic gradient applied to strongly convex, convex or nonconvex problems. Those results heavily depend upon the formula for the step sizes  $\{\alpha_k\}_k$ . In fact, one of the major problems in machine learning consists in **tuning the learning rate**, which corresponds to choosing the step size in stochastic gradient. We will illustrate the various challenges posed by this choice in the context of strongly convex functions.

**Assumption 4.2.3** *The objective function is  $\mu$ -strongly convex and possesses a unique global minimizer  $\mathbf{w}^*$ . We let  $f^* = f(\mathbf{w}^*)$ .*

We first provide a global rate result in the case of a constant step size.

**Theorem 4.2.1 (SG with constant stepsize)** *Let Assumptions 2.1.1, 4.2.2 and 4.2.3, and consider Algorithm 8 applied with a constant stepsize*

$$\alpha_k = \alpha \in (0, \frac{1}{2\mu}) \forall k.$$

Then,

$$\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \frac{\alpha L \sigma^2}{2\mu} + (1 - 2\alpha\mu)^k \left[ f(\mathbf{w}_0) - f^* - \frac{\alpha L \sigma^2}{2\mu} \right]. \quad (4.2.2)$$

We note that this convergence rate corresponds to guaranteeing  $\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \epsilon$  after at most  $\mathcal{O}(\ln(1/\epsilon))$  iterations. However, unlike in the gradient descent case, the tolerance  $\epsilon$  cannot be arbitrarily close to zero. In fact, the use of stochastic gradients introduces an additional (bias) term  $\frac{\alpha L \sigma^2}{2\mu}$ . As a result, SG with constant stepsize can only be guaranteed to converge towards a **neighborhood** of the optimal function value  $f^*$ . On the other hand, such a method is capable of taking long steps, as opposed to the next technique based on decreasing step sizes.

In the original stochastic gradient method (proposed by Robbins and Monro in 1951), the stepsize sequence was required to satisfy

$$\sum_{k=0}^{\infty} \alpha_k = \infty \quad \text{and} \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty,$$

which implies that  $\alpha_k \rightarrow 0$ . In our next result, we thus consider the case of diminishing stepsizes.

**Theorem 4.2.2 (SG with diminishing stepsize)** *Let Assumptions 2.1.1, 4.2.2 and 4.2.3, and consider Algorithm 8 applied with a decreasing stepsize sequence  $\{\alpha_k\}_k$  satisfying*

$$\alpha_k = \frac{\beta}{k + \gamma},$$

where  $\beta > \frac{1}{\mu}$  and  $\gamma > 0$  is chosen such that  $\alpha_0 = \frac{\beta}{\gamma} \leq \frac{1}{L}$ . Then,

$$\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \frac{\nu}{\gamma + k}, \quad (4.2.3)$$

where

$$\nu = \max \left\{ \gamma(f(\mathbf{w}_0) - f^*), \frac{\beta^2 L \sigma^2}{2(\beta\mu - 1)} \right\}.$$

The decreasing stepsize choice possesses the same drawbacks than for gradient descent, namely that it results in increasingly small steps. It also provides a global convergence rate that is sublinear, as opposed to linear with a constant stepsize. Note, however, that SG with a decreasing stepsize is guaranteed to reach any neighborhood of a solution, unlike its variant with a constant stepsize.

**Remark 4.2.3 (A practical constant stepsize approach)** *A common practical strategy in machine learning consists in running the algorithm with a value  $\alpha$  until the method stalls (which can indicate that the smallest neighborhood attainable with this stepsize choice has been reached). When that occurs, the stepsize can be reduced, and the algorithmic run can continue until it stalls again, then the stepsize will be further reduced, etc (say  $\alpha, \alpha/2, \alpha/4$ , etc). This process can lead to*

convergence guarantees, however the convergence is slower than that produced by constant stepsize SG:

$$\mathbb{E}[f(\mathbf{w}_k) - f^*] \leq \epsilon \quad \text{after } \mathcal{O}(1/\epsilon) \text{ iterations.}$$

This choice of stepsize is adaptive, in that it is designed to reach closer and closer neighborhoods as the algorithm proceeds. However, it requires the method to be able to detect stalling, and act upon it.

**Stepsize choice in the nonconvex setting** Stochastic gradient (or some variant thereof) is the method of choice for training neural networks, which is usually a nonconvex problem. It is thus natural to ask whether global rates can be obtained for stochastic gradient in the nonconvex setting. The situation is significantly more complicated, as we get guarantees on

- $\mathbb{E}\left[\frac{1}{K} \sum_{i=1}^K \|\nabla f(\mathbf{w}_k)\|^2\right]$  for constant stepsizes;
- $\mathbb{E}\left[\frac{1}{\sum_{i=1}^K \alpha_k} \sum_{i=1}^K \alpha_k \|\nabla f(\mathbf{w}_k)\|^2\right]$  for decreasing stepsizes.

Similarly to the strongly convex case, the complexity bounds are affected by a residual term which in turns lead to worse rates than in the deterministic setting.

**Example 4.2.2** A typical stochastic gradient method with constant stepsize will satisfy

$$\mathbb{E}\left[\frac{1}{K} \sum_{i=1}^K \|\nabla f(\mathbf{w}_k)\|^2\right] \leq \epsilon$$

in at most  $\mathcal{O}(\epsilon^{-4})$  iterations, where  $\epsilon$  is a sufficient large threshold of accuracy.

**Remark 4.2.4 (What about momentum?)** The most successful implementations of stochastic gradient, such as ADAM, rely on some form of momentum incorporated in the stochastic gradient update. Intuitively, the hope is that incorporating momentum will allow the method to promote moves along good directions of decrease, while steps in bad directions will eventually cancel out. Some theory has been developed in the recent years to accelerate stochastic gradient yet, unlike in the deterministic setting, theory is still decorrelated from practice.

### 4.3 Variance reduction

As we saw in the previous section, the theory for stochastic gradient is based on Assumption 4.2.2, and in particular on the fact that the variance of stochastic gradient estimates is bounded (by  $\sigma^2$ ). It can clearly be seen from bounds such as (4.2.2) that the bigger  $\sigma$  is, the looser the bound becomes. More practically, this means that gradient estimates with high variance are unlikely to yield fast convergence.

Variance reduction techniques have precisely been developed in the aim of diminishing the variance of traditional stochastic gradient estimates. They can be categorized in two families, that either exploit more sampled gradients at every iteration, or use past history of the method. In these notes, we will focus on the former category.

### 4.3.1 Batch variants

We recall that the main part of Algorithm 8 consists in the update

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k),$$

where the index  $i_k$  is drawn at random. The use of a **single** sample is partially responsible for the importance of the variance term  $\sigma^2$  in Assumption 4.2.2. One can thus consider stochastic gradient estimates that are built using *several* samples at once : this is the idea behind **batch stochastic gradient**.

Formally, the update of a batch stochastic gradient method is given by

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k) \quad (4.3.1)$$

where  $S_k \subset \{1, \dots, n\}$  is drawn at random. When  $S_k$  consists in a single index, we recover the usual stochastic gradient algorithm; conceptually, one could also consider a set  $S_k$  of cardinality  $n$ , in which case we would recover the usual gradient method.

Overall, two batch regimes can be distinguished:

- $|S_k| \approx n$ , which has a cost essentially equivalent to that of a full gradient update;
- $|S_k| = n_b \ll n$ , also called **mini-batching**, which may be advantageous in theory and variance reduction while still being affordable in practice. The resulting method is called **mini-batch SG**.

In fact, if we assume that  $|S_k| = n_b \forall k$ , it is possible to show that with the same stepsize, mini-batch SG requires  $n_b$  less iterations than SG. Moreover, mini-batch SGD can exploit parallel computing, by computing the  $n_b$  stochastic gradients on distributed processors. Moreover, we have the following property.

**Proposition 4.3.1** *Under Assumptions 2.1.1 and 4.2.2, the variance of a mini-batch stochastic gradient estimate is given by*

$$\text{Var}_{S_k} \left[ \left\| \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k) \right\|_2 \right] \leq \frac{\sigma^2}{n_b}.$$

As a final note, we mention that batch techniques are still more expensive than stochastic gradient, while being more sensitive to redundancies in the data. Tuning the best batch size is not necessarily an easy task. These concerns partly explain why stochastic gradient (or other schemes based on his sampling paradigm) remains the preferred approach.

### 4.3.2 Other variants

**Gradient aggregation** methods have attracted a lot of attention in the learning and optimization theory, because of the nice theory and algorithms that have been proposed and guarantee linear convergence rates. Their main principle consists in computing a **full gradient step** once in a while during the algorithmic run, in order to correct high-variance components. Despite their strong guarantees, they have not been widely exploited in practice, due to the cost of full gradient evaluations, that is still too prohibitive in certain applications.

**Iterate averaging** is another popular technique, that can be easier to implement. The underlying idea consists in analyzing (and possibly returning as output) the average iterate of a run of stochastic gradient, given by  $\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{w}_k$ . In certain contexts (e.g.  $\alpha = \frac{1}{\mu(k+1)}$  and  $f$   $\mu$ -strongly convex), this average has good properties with respect to the optimization, and is also a more robust solution than the last iterate obtained. However, returning this average either requires to store the history of iterates, or to maintain an average which can be prone to cancellation or numerical errors.

## 4.4 Stochastic gradient methods for deep learning

In this section, we focus on stochastic gradient algorithms that have proven useful in training deep learning models (though the methods we will present are not tailored to a particular architecture).

We again consider a finite-sum problem of the form (4.1.1) under Assumption 4.2.1. Our objective is to analyze several variants on the basic scheme

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k, \quad (4.4.1)$$

where  $\alpha > 0$  is a stepsize (also known as learning rate in the machine learning community) and  $\mathbf{g}_k$  is a stochastic gradient estimator, that either corresponds to a single gradient component (as in vanilla stochastic gradient) or a batch of indices.

We will present all our variants within a unified framework that highlights the key features of these methods: this framework is given by the iteration

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{m}_k \oslash \mathbf{v}_k, \quad (4.4.2)$$

where  $\alpha > 0$ ,  $\mathbf{m}_k, \mathbf{v}_k \in \mathbb{R}^d$  and  $\oslash$  denotes the componentwise division, i. e.

$$\mathbf{m}_k \oslash \mathbf{v}_k := \left[ \frac{[\mathbf{m}_k]_i}{[\mathbf{v}_k]_i} \right]_{i=1, \dots, d}.$$

Note that by letting  $\mathbf{m}_k = \mathbf{g}_k$  and  $\mathbf{v}_k = \mathbf{1}_{\mathbb{R}^d}$ , we recover the classical stochastic gradient iteration (4.4.1).

### 4.4.1 Stochastic gradient with momentum

Inspired by the accelerated methods that we investigated in Chapter 2, we first consider adding momentum to the basic iteration (4.4.1). The most common approach, called **stochastic gradient with momentum**, corresponds to the iteration:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha(1 - \beta_1) \mathbf{g}_k + \alpha \beta_1 (\mathbf{w}_k - \mathbf{w}_{k-1}), \quad (4.4.3)$$

where  $\beta_1 \in (0, 1)$  is a constant parameter ( $\beta_1 = 0$  would correspond to the classical stochastic gradient method). This method is a (stochastic) variant on Polyak's heavy-ball method, for which the gradient step is combined with the previous displacement. As in momentum-based methods, the idea consists in accumulating information from the previous iteration. In practice, the iteration (4.4.3) tends to accumulate good directions (in the optimization sense) while "bad" directions tend to cancel out.

The method (4.4.3) is a special case of (4.4.2), corresponding to  $\mathbf{v}_k = \mathbf{1}_{\mathbb{R}^d}$  and  $\mathbf{m}_k$  defined recursively by  $\mathbf{m}_{-1} = \mathbf{0}_{\mathbb{R}^d}$  and

$$\mathbf{m}_k = (1 - \beta_1) \mathbf{g}_k - \beta_1 \mathbf{m}_{k-1} \quad \forall k \in \mathbb{N}.$$

where  $\beta_1$  is the constant defined in (4.4.3).

Stochastic gradient with momentum is implemented in standard deep learning libraries such as PyTorch. It is particularly useful in training deep neural networks on computer vision tasks, and, as such, played a role in the outbreak of deep learning circa 2012.

**Remark 4.4.1** *It is less straightforward to derive theoretical guarantees for the method (4.4.3) than for accelerated gradient descent, even in a strongly convex setting. Nevertheless, adding momentum to the stochastic gradient iteration is a popular practice in solving nonconvex problems such as those arising from training neural networks.*

#### 4.4.2 AdaGrad

The adaptive gradient method, or ADAGRAD, was proposed in 2011 to address the issue of selecting the learning rate  $\alpha$  in stochastic gradient without relying on adaptive approaches like line searches. In ADAGRAD, every component of the stochastic gradient is scaled according to a running average of the values taken by that component over all iterations. The method maintains a sequence  $\{\mathbf{r}_k\}_k$  given by

$$\forall i = 1, \dots, d, \quad \begin{cases} [\mathbf{r}_{-1}]_i = 0 \\ [\mathbf{r}_k]_i = [\mathbf{r}_{k-1}]_i + [\mathbf{g}_k]_i^2 \quad \forall k \geq 0, \end{cases} \quad (4.4.4)$$

The ADAGRAD iteration is thus

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k \odot \sqrt{\mathbf{r}_k}, \quad (4.4.5)$$

where the square root is applied to every component of  $\mathbf{r}_k$ . This iteration matches (4.4.2) with  $\mathbf{m}_k = \mathbf{g}_k$  and  $\mathbf{v}_k = \sqrt{\mathbf{r}_k}$ . The contribution of ADAGRAD thus consists in using a different stepsize for each coordinate, leading the sequence :

$$\left\{ \left[ \frac{\alpha}{\sqrt{[\mathbf{r}_k]_i}} \right]_{i=1}^d \right\}_k.$$

The method performs a *diagonal scaling* of the components of the stochastic gradient  $\mathbf{g}_k$ , which is particularly well suited for ill-conditioned problems where the components have a high variance. However, such stepsizes typically decrease very quickly towards 0.

**Remark 4.4.2** *In practice, we replace  $\mathbf{r}_k$  by  $\mathbf{r}_k + \eta \mathbf{1}_{\mathbb{R}^d}$  where  $\eta > 0$  is a small quantity, so that the algorithm is numerically stable.*

ADAGRAD is particularly suited for problems with *sparse gradients*, for which stochastic gradients also tend to have many zero components. In this situation, computing  $\mathbf{r}_k$  will only change the stepsize for the nonzero coordinates. Problems from recommender systems typically come with sparse gradients, which explains the popularity of ADAGRAD in this setting.

#### 4.4.3 RMSProp

The Root Mean Square Propagation algorithm, or RMSPROP, is similar to ADAGRAD in that it scales the stochastic gradient components. To this end, the method computes a vector sequence  $\{\mathbf{r}_k\}_k$  as follows:

$$\forall i = 1, \dots, d, \quad \begin{cases} [\mathbf{r}_{-1}]_i = 0 \\ [\mathbf{r}_k]_i = (1 - \lambda)[\mathbf{r}_{k-1}]_i + \lambda[\mathbf{g}_k]_i^2 \quad \forall k \geq 0, \end{cases} \quad (4.4.6)$$

where  $\lambda \in (0, 1)$ . The value of  $\lambda$  controls how much weight is given to the past stochastic gradient components over the current stochastic gradient components. This idea leads to a slower decrease in the stepsizes compared to the values of ADAGRAD.

As for ADAGRAD, the iteration of RMSPROP corresponds to a special case of (4.4.2) using  $\mathbf{m}_k = \mathbf{g}_k$  and  $\mathbf{v}_k = \sqrt{\mathbf{r}_k}$ .

**Remark 4.4.3** *In practice, and as in ADAGRAD, the vector  $\mathbf{r}_k$  is replaced by  $\mathbf{r}_k + \eta \mathbf{1}_{\mathbb{R}^d}$  for a small value  $\eta > 0$ .*

The RMSPROP algorithm has been successfully applied to training very deep neural networks.

#### 4.4.4 Adam

The ADAM algorithm<sup>1</sup> was proposed in 2013, and has been one of the most popular stochastic gradient technique in practice. This method can be viewed as combining the idea of momentum together with scaling: scaling will be performed according to the past gradients, and the search direction will also include pas gradient information. The ADAM iteration corresponds to applying (4.4.2) with

$$\mathbf{m}_k = \frac{(1 - \beta_1) \sum_{j=0}^k \beta_1^{k-j} \mathbf{g}_j}{1 - \beta_1^{k+1}} \quad (4.4.7)$$

for  $\beta_1 \in (0, 1)$ . This is indeed a momentum-type iteration, since we can obtain  $\mathbf{m}_k$  from  $\mathbf{m}_{k-1}$  and  $\mathbf{g}_k$  through the formula

$$\mathbf{m}_k = \beta_1 \frac{1 - \beta_1^k}{1 - \beta_1^{k+1}} \mathbf{m}_{k-1} + \frac{1 - \beta_1}{1 - \beta_1^{k+1}} \mathbf{g}_k.$$

The other component of the ADAM update is given by

$$\mathbf{v}_k = \sqrt{\frac{(1 - \beta_2) \sum_{j=0}^k \beta_2^{k-j} \mathbf{g}_j \odot \mathbf{g}_j}{1 - \beta_2^{k+1}}}. \quad (4.4.8)$$

where  $\beta_2 \in (0, 1)$  and  $\odot$  denoting the componentwise or Hadamard product given by

$$\mathbf{g}_k \odot \mathbf{g}_k = \left[ [\mathbf{g}_k]_i^2 \right]_{i=1}^d.$$

**Remark 4.4.4** *In practice, a vector of the form  $\mathbf{v}_k + \eta \mathbf{1}_{\mathbb{R}^d}$  will be used in lieu of  $\mathbf{v}_k$ , with  $\eta$  being a small positive number.*

The above formulae amount to combining previously employed directions with the latest stochastic gradient vector, and normalizing the components of the obtained vector according to the history of these components. In both cases, more importance is given to the latest values that have been computed. This is a key feature of the method, that has statistical motivations, and may explain the impressive performance of ADAM. In practice, ADAM (and its variant ADAMW based on regularization) are among the most efficient methods for training architectures on Natural Language Processing tasks.

<sup>1</sup>The name Adam is derived from ADaptive Momentum estimation.

## 4.5 Conclusion

From a pure optimization perspective, stochastic gradient methods may not seem so attractive, as they only rely on partial information from the gradient and possess worse convergence guarantees than gradient descent. However, they have encountered tremendous success in data-related applications, where computing gradients involves looking at the entire data and is thus too prohibitive. On the contrary, using stochastic gradient estimates represents a significantly cheaper cost per iteration; in a data science setting, where there can be redundancies (or even underlying randomness) in the data, such updates do not necessarily hinder the progress of the algorithm, but rather lead to faster convergence in practice.

Still, the stochastic gradient approach suffers from high-variance estimates. For this reason, practical variants typically incorporate enhancements to reduce the variance. The most prominent technique for finite-sum and stochastic problems consist in using a batch of samples, which provably reduces the variance and can improve the performance. Meanwhile, the most efficient stochastic gradient techniques, such as those used in deep learning, employ both momentum terms and diagonal scaling to improve the quality of the steps. These techniques may not be endowed with better (if any) theoretical guarantees, especially when applied to nonconvex training problems. However, methods such as Stochastic Gradient with Momentum or `ADAM` have been widely adopted by the learning community because of their practical efficiency.



## Chapter 5

# Large-scale and distributed optimization

In this last chapter, we dive into a increasingly important area of focus in optimization methods for data science. As we witness a growth in both the model complexity (i.e. the number of parameters) and the amount of data available (i.e. the size of the dataset), standard optimization techniques may suffer from the curse of dimensionality and their performance may deteriorate as dimensions grow. The goal of this chapter is to present some algorithmic ideas that can reduce the impact of large dimensions, either in terms of parameters or data points.

### 5.1 Coordinate descent methods

In this section, we address the treatment of large-scale optimization problems, where the number of parameters to be optimized over is extremely large. In general, due to the curse of dimensionality, the difficulty of the problem increases with the dimension, simply because there are more variables to consider. However, on structured problems such as those arising in data science, the problem may possess a low-dimensional or separable structure that allows for optimization steps to be taken over a subset of variables. This is the underlying idea of **coordinate descent methods**, that have regained interest in the early 2000s due to their applicability in certain data science settings.

#### 5.1.1 Algorithmic framework

Consider the unconstrained optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}), \quad (5.1.1)$$

where  $f \in \mathcal{C}^1(\mathbb{R}^d)$ . The idea of coordinate descent methods consist in taking a gradient step with respect to a single decision variable at every iteration. To this end, we observe that for every  $\mathbf{w} \in \mathbb{R}^d$ , the gradient of  $f$  at  $\mathbf{w}$  can be decomposed as

$$\nabla f(\mathbf{w}) = \sum_{j=1}^d \nabla_j f(\mathbf{w}) \mathbf{e}_j,$$

where  $\nabla_j$  denotes the partial derivative with respect to the  $j$ -th variable of the function  $f$  (that is, the  $j$ th coordinate of  $f$ ) and  $\mathbf{e}_j \in \mathbb{R}^d$  is the  $j$ th coordinate vector of the canonical basis in  $\mathbb{R}^d$ . Not

unlike the stochastic gradient paradigm<sup>1</sup>, the coordinate descent approach replaces the full gradient by a step along a coordinate gradient, as formalized in Algorithm 9.

---

**Algorithm 9:** Coordinate descent method.

---

**Initialization:**  $w_0 \in \mathbb{R}^d$ .

**for**  $k = 0, 1, \dots$  **do**

1. Select a coordinate index  $j_k \in \{1, \dots, d\}$ .

2. Compute a steplength  $\alpha_k > 0$ .

3. Set

$$w_{k+1} = w_k - \alpha_k \nabla_{j_k} f(w_k) e_{j_k}. \quad (5.1.2)$$

**end**

---

The variants of coordinate descent are mainly identified by the way they select the coordinate sequence  $\{j_k\}$ . There exist numerous rules for choosing the coordinate index, among which:

- **Cyclic:** Select the indices by cycling over  $\{1, \dots, d\}$  in that order. After  $d$  iterations, all indices have been selected.
- **Randomized cyclic:** Cycle through a random ordering of  $\{1, \dots, d\}$ , that changes every  $d$  steps.
- **Randomized:** Draw  $j_k$  at random in  $\{1, \dots, d\}$  at every iteration.

The last two strategies are those for which the strongest results can be obtained.

**Block coordinate descent** Rather than using a single index, it is possible to select a subset of the coordinates (called “block” in the literature). The  $k$ th iteration of such a *block coordinate descent* algorithm thus is

$$w_{k+1} = w_k - \alpha_k \sum_{j \in \mathcal{B}_k} \nabla_j f(w_k) e_j, \quad (5.1.3)$$

where  $\mathcal{B}_k \subset \{1, \dots, d\}$ .

### 5.1.2 Theoretical guarantees of coordinate descent methods

A famous 3-dimensional example designed by M. J. D. Powell in 1973 shows that coordinate descent methods do not necessarily converge. Nevertheless, it is possible to provide guarantees on coordinate descent methods under appropriate assumptions. In particular, a linear rate of convergence can be obtained for coordinate descent methods on strongly convex problems: we provide below the necessary assumptions to arrive at such a result.

---

<sup>1</sup>But not to be confused with it!

**Assumption 5.1.1** The objective function  $f$  in (5.1.1) is  $\mathcal{C}^1$  and  $\mu$ -strongly convex, with  $f^* = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ . Moreover, for every  $j = 1, \dots, d$ , the partial derivative  $\nabla_j f$  is  $L_j$ -Lipschitz continuous, i.e.

$$\forall \mathbf{w} \in \mathbb{R}^d, \forall h \in \mathbb{R}, \quad |\nabla_j f(\mathbf{w} + h\mathbf{e}_j) - \nabla_j f(\mathbf{w})| \leq L_j |h|. \quad (5.1.4)$$

We let  $L_{\max} = \max_{1 \leq j \leq d} L_j$ .

**Theorem 5.1.1** Suppose that Assumption 5.1.1 holds, and that Algorithm 9 is applied to problem (5.1.1) with  $\alpha_k = \frac{1}{L_{\max}}$  for all  $k$  and  $j_k$  being drawn uniformly at random in  $\{1, \dots, d\}$ . Then, for any  $K \in \mathbb{N}$ , we have

$$\mathbb{E}[f(\mathbf{w}_K) - f^*] \leq \left(1 - \frac{\mu}{dL_{\max}}\right)^K (f(\mathbf{w}_0) - f^*). \quad (5.1.5)$$

Other results have been established in the convex and nonconvex settings, under additional assumptions. In all cases, properties on the partial derivatives are required.

**Remark 5.1.1** As described in the lab session, it is possible to combine randomized coordinate descent with Nesterov's acceleration technique to yield improved theoretical guarantees. However, this raises implementation issues that may alleviate the practical interest of coordinate descent approaches.

### 5.1.3 Applications of coordinate descent methods

Coordinate descent techniques are particularly useful for large-scale sparse optimization. Consider a regularized problem of the form

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \tilde{f}_i(\mathbf{x}_i^T \mathbf{w}) + \sum_{j=1}^d \Omega(w_j), \quad (5.1.6)$$

where  $\tilde{f}_i : \mathbb{R} \rightarrow \mathbb{R}$  is (possibly) data-dependent,  $\mathbf{x}_i \in \mathbb{R}^d$  is a **sparse** data vector, and  $\Omega : \mathbb{R} \rightarrow \mathbb{R}$  is a regularization function applied componentwise to the vector  $\mathbf{w}$ .

**Example 5.1.1 (Regularized least squares with sparse data)** Given  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with sparse rows and  $\mathbf{y} \in \mathbb{R}^n$ , consider the problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad f(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n \|\mathbf{X} \mathbf{w} - \mathbf{y}\|^2 + \lambda \sum_{j=1}^d w_j^2.$$

Apply Algorithm 9 to this problem. For any iteration  $k$ , if we move along the  $j_k$ th coordinate, the partial derivative under consideration is

$$\nabla_{j_k} f(\mathbf{w}_k) = \mathbf{x}_{j_k}^T (\mathbf{X} \mathbf{w}_k - \mathbf{y}) + 2\lambda [\mathbf{w}_k]_{j_k}.$$

By storing the vector  $\{\mathbf{X} \mathbf{w}_k\}$  across all iterations, the calculation of  $\nabla_{j_k} f(\mathbf{w}_k)$  can be greatly reduced when  $\mathbf{x}_{j_k}$  is sparse, to the point that the cost of a coordinate descent iteration will be of the order of the number of nonzero elements in  $\mathbf{x}_{j_k}$ .

Coordinate descent techniques are quite prominent in parallel optimization algorithms. In this setting, several cores are cooperating to solve problem (5.1.1): each core can then run *its own coordinate descent method* and all cores update the same shared iterate vector. The most efficient parallel coordinate descent techniques perform these iterations in an asynchronous fashion, which does not prevent from guaranteeing convergence of this framework!

**Link with stochastic gradient** Consider a finite-sum problem of the form

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}), \quad f_i(\mathbf{w}) := \ell_i(\mathbf{x}_i^T \mathbf{w}), \quad (5.1.7)$$

where  $\mathbf{x}_i^T \mathbf{w}$  is a linear model of the data vector  $\mathbf{x}_i$ , and  $\ell_i : \mathbb{R} \rightarrow \mathbb{R}$  is a convex loss function specific to the  $i$ th data point (such as  $\ell_i(h) = \frac{1}{2}(h - y_i)^2$  for linear least squares). In Chapter 4, we saw how to apply stochastic gradient to this problem. Another approach consists in considering an equivalent formulation of (5.1.7) through duality, given by

$$\underset{\mathbf{v} \in \mathbb{R}^n}{\text{maximize}} g(\mathbf{v}) := -\frac{1}{n} \sum_{i=1}^n f_i^*(v_i) \quad (5.1.8)$$

where for any convex function  $\phi : \mathbb{R}^m \rightarrow \mathbb{R}$ , the conjugate function  $\phi^*$  is defined by

$$\phi^*(\mathbf{a}) = \sup_{\mathbf{b} \in \mathbb{R}^m} \{\mathbf{a}^T \mathbf{b} - \phi(\mathbf{b})\}.$$

The so-called dual problem (5.1.8) has a finite-sum, separable form. It can thus be tackled using (dual) *coordinate ascent*, the counterpart of coordinate descent for minimization: the iteration of this method is given by

$$\mathbf{v}_{k+1} = \mathbf{v}_k + \alpha_k \nabla_i g(\mathbf{v}_k), \quad (5.1.9)$$

leading to updating the iterate one coordinate at a time. Under the appropriate assumptions on the problem, the iteration (5.1.9) is equivalent to the original stochastic gradient iteration: this is why stochastic gradient is sometimes viewed as applying coordinate ascent to the dual problem. We will come back to this notion of duality in the next section.

## 5.2 Distributed and constrained optimization

In this section, we describe the theoretical insights behind distributed optimization formulations, in which several agents collaborate to solve an optimization problem. This paradigm can be modeled using a constrained optimization formulation, leading to a dual view of certain algorithms.

### 5.2.1 Linear constraints and dual problem

Consider the following optimization problem with linear equality constraints:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) \quad \text{subject to} \quad \mathbf{A}\mathbf{w} = \mathbf{b}, \quad (5.2.1)$$

where  $\mathbf{A} \in \mathbb{R}^{m \times d}$  and  $\mathbf{b} \in \mathbb{R}^m$ . For simplicity, we will assume that the feasible set  $\{\mathbf{w} \in \mathbb{R}^d \mid \mathbf{A}\mathbf{w} = \mathbf{b}\}$  is not empty.

Duality theory consists in handling constraints formulations by reformulating the problem into an unconstrained optimization problem. We present the theoretical arguments for the special case of problem (5.2.1), which yields a much simpler analysis.

**Definition 5.2.1** The Lagrangian function of problem (5.2.1) is given by

$$\mathcal{L}(\mathbf{w}, \mathbf{z}) := f(\mathbf{w}) + \mathbf{z}^T (\mathbf{A}\mathbf{w} - \mathbf{b}). \quad (5.2.2)$$

The Lagrangian function combines the objective function and the constraints, and allows to restate the original problem as an unconstrained one, called the primal problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \max_{\mathbf{z} \in \mathbb{R}^m} \mathcal{L}(\mathbf{w}, \mathbf{z}). \quad (5.2.3)$$

The solutions of the primal problem are identical to that of problem (5.2.3) in our case. The difficulty of solving problem (5.2.3) lies in the definition of its objective function as the optimal value of a maximization problem.

**Definition 5.2.2** The **dual problem** of (5.2.1) is the maximization problem

$$\underset{\mathbf{z} \in \mathbb{R}^m}{\text{maximize}} \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{z}), \quad (5.2.4)$$

where the function  $\mathbf{z} \mapsto \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{z})$  is called the dual function of the problem.

Unlike the primal problem, the dual problem is always concave (i.e. the opposite of the dual function is convex), which facilitates its resolution by standard optimization techniques. The goal is then to solve the dual problem in order to get the solution of the primal problem, thanks to properties such as the one below.

**Assumption 5.2.1** We suppose that **strong duality** holds between problem (5.2.1) and its dual, that is,

$$\min_{\mathbf{w} \in \mathbb{R}^d} \max_{\mathbf{z} \in \mathbb{R}^m} \mathcal{L}(\mathbf{w}, \mathbf{z}) = \max_{\mathbf{z} \in \mathbb{R}^m} \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{z}).$$

A sufficient condition for Assumption 5.2.1 is that  $f$  be convex, but this is not necessary.

## 5.3 Dual algorithms

We are now concerned with solving the dual problem (5.2.4), and we will present three methods for this purpose.

### 5.3.1 Dual ascent

The **dual ascent** method is implicitly a subgradient method applied to the dual problem (which we recall is a maximization problem). At every iteration, it starts from a primal-dual pair  $(\mathbf{w}_k, \mathbf{z}_k)$  and performs the following iteration:

$$\begin{cases} \mathbf{w}_{k+1} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{z}_k) \\ \mathbf{z}_{k+1} = \mathbf{z}_k + \alpha_k (\mathbf{A}\mathbf{w}_{k+1} - \mathbf{b}), \end{cases} \quad (5.3.1)$$

where  $\alpha_k > 0$  is a stepsize for the dual ascent step, and  $\mathbf{A}\mathbf{w}_{k+1} - \mathbf{b}$  is a subgradient for the dual function  $\mathbf{z} \mapsto \min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{z})$  at  $\mathbf{z}_k$ .

### 5.3.2 Augmented Lagrangian

The dual ascent method generally has weak convergence guarantees. For this reason, the optimization literature has introduced other frameworks based on a regularized version of the Lagrangian function.

**Definition 5.3.1** The **augmented Lagrangian** of problem (5.2.1) is the function on  $\mathbb{R}^d \times \mathbb{R}^m \times \mathbb{R}_{++}$  by

$$\mathcal{L}^a(\mathbf{w}, \mathbf{z}; \lambda) := f(\mathbf{w}) + \mathbf{z}^T (\mathbf{A}\mathbf{w} - \mathbf{b}) + \frac{\lambda}{2} \|\mathbf{A}\mathbf{w} - \mathbf{b}\|^2. \quad (5.3.2)$$

Augmented Lagrangians thus are a family of functions parameterized by  $\lambda > 0$ , that put more emphasis on the constraint violation as  $\lambda$  grows.

The **augmented Lagrangian** algorithm, also called method of multipliers, performs the following iteration:

$$\begin{cases} \mathbf{w}_{k+1} \in \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}^a(\mathbf{w}, \mathbf{z}_k; \lambda) \\ \mathbf{z}_{k+1} = \mathbf{z}_k + \lambda(\mathbf{A}\mathbf{w}_{k+1} - \mathbf{b}). \end{cases} \quad (5.3.3)$$

In this algorithm,  $\lambda$  is constant and used as a constant stepsize: many more sophisticated choices of both the augmented Lagrangian function and the stepsizes have been proposed. In general, the advantages of augmented Lagrangian techniques are that the subproblems defining  $\mathbf{w}_{k+1}$  become easier to solve (thanks to regularization) and that the overall guarantees on the primal-dual pair are stronger.

### 5.3.3 ADMM

The **Alternated Direction Method of Multipliers**, or **ADMM**, is an increasingly popular variation on the augmented Lagrangian paradigm that bears some connection with coordinate descent approaches, in that it splits the problem in two sets of variables.

Suppose that we consider a linearly constrained problem with a separable form:

$$\begin{cases} \operatorname{minimize}_{\mathbf{u} \in \mathbb{R}^{d_1}, \mathbf{v} \in \mathbb{R}^{d_2}} f(\mathbf{u}) + g(\mathbf{v}) \\ \text{subject to} \quad \mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} = \mathbf{c}, \end{cases} \quad (5.3.4)$$

where  $\mathbf{A} \in \mathbb{R}^{d_1 \times m}$ ,  $\mathbf{B} \in \mathbb{R}^{d_2 \times m}$  and  $\mathbf{c} \in \mathbb{R}^m$ . In that case, for any  $\lambda > 0$ , the augmented Lagrangian of problem (5.3.4) has the form

$$\mathcal{L}^a(\mathbf{u}, \mathbf{v}, \mathbf{z}; \lambda) = f(\mathbf{u}) + g(\mathbf{v}) + \mathbf{z}^T (\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} - \mathbf{c}) + \frac{\lambda}{2} \|\mathbf{A}\mathbf{u} + \mathbf{B}\mathbf{v} - \mathbf{c}\|^2.$$

The ADMM iteration exploits the separable nature of the problem by computing the values  $\mathbf{u}$  and  $\mathbf{v}$  independently. Starting from  $(\mathbf{u}_k, \mathbf{v}_k, \mathbf{z}_k)$ , the ADMM counterpart to iteration (5.3.3) is

$$\begin{cases} \mathbf{u}_{k+1} \in \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^{d_1}} \mathcal{L}^a(\mathbf{u}, \mathbf{v}_k, \mathbf{z}_k; \lambda) \\ \mathbf{v}_{k+1} \in \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^{d_2}} \mathcal{L}^a(\mathbf{u}_{k+1}, \mathbf{v}, \mathbf{z}_k; \lambda) \\ \mathbf{z}_{k+1} = \mathbf{z}_k + \lambda(\mathbf{A}\mathbf{u}_{k+1} + \mathbf{B}\mathbf{v}_{k+1} - \mathbf{c}). \end{cases} \quad (5.3.5)$$

The two-subproblem process of (5.3.5) corresponds to two iterations of block coordinate descent, which is often beneficial to the optimization process compared to a joint iteration in  $\mathbf{u}$  and  $\mathbf{v}$ .

**Remark 5.3.1** The idea of splitting the objective and the constraints across two groups of variables can be declined into as many groups of variables as possible, depending on the structure of the problem.

To end this section, we briefly mention that there exist convergence results for ADMM-type frameworks, typically under convexity assumptions on the problem [?]. A typical result consist in showing that

$$\begin{cases} \| \mathbf{A}\mathbf{u}_k + \mathbf{B}\mathbf{v}_k - \mathbf{c} \| & \rightarrow 0 \\ f(\mathbf{u}_k) + g(\mathbf{v}_k) & \rightarrow \min_{\mathbf{u}, \mathbf{v}} f(\mathbf{u}) + g(\mathbf{v}) \\ \mathbf{z}_k & \rightarrow \mathbf{z}^*, \end{cases}$$

where  $\mathbf{z}^*$  is a solution of the dual problem.

## 5.4 Consensus optimization

We end this chapter by describing an increasingly common setup in optimization over large datasets, often termed **consensus optimization** or **decentralized optimization**. In this setup, we consider a dataset that is split across  $m$  entities called *agents*. Every agent uses its own data to train a certain learning model parameterized by a vector in  $\mathbb{R}^d$ . To this end, each agent not only has its own function  $f^{(i)}$ , but also its own copy of the model parameters  $\mathbf{w}^{(i)}$ . The optimization problem at hand considers a master iterate  $\mathbf{w}$ , and attempts to reach consensus between all the agents. This leads to the following formulation:

$$\begin{aligned} & \text{minimize}_{\mathbf{w}, \mathbf{w}^{(1)}, \dots, \mathbf{w}^{(m)} \in \mathbb{R}^d} \sum_{i=1}^m f^{(i)}(\mathbf{w}^{(i)}) \\ & \text{subject to} \quad \mathbf{w} = \mathbf{w}^{(i)} \quad \forall i = 1, \dots, m. \end{aligned} \quad (5.4.1)$$

This problem is a proxy for  $\text{minimize}_{\mathbf{w} \in \mathbb{R}^d} \sum_{i=1}^m f^{(i)}(\mathbf{w})$ , but the latter problem cannot be solved by a single agent since every agent has exclusive access to its data by design. The formulation (5.4.1) models the fact that all agents are involved in computing  $\mathbf{w}$  by acting on  $\mathbf{w}_i$ . It is possible to apply ADMM to problem (5.4.1) by setting

$$\mathbf{u} = \begin{bmatrix} \mathbf{w}^{(1)} \\ \vdots \\ \mathbf{w}^{(m)} \end{bmatrix} \in \mathbb{R}^{md}, \quad \mathbf{v} = \mathbf{w} \in \mathbb{R}^d.$$

**Generalization** The idea behind the formulation (5.4.1) can be extended to the case of data spread over a network, represented by a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ : every vertex  $s \in \mathcal{V}$  of the graph represents an agent, while every edge  $(s, s') \in \mathcal{E}$  represents a channel of communication between two agents in the graph. Letting  $\mathbf{w}^{(s)} \in \mathbb{R}^d$  and  $f^{(s)} : \mathbb{R}^d \rightarrow \mathbb{R}$  represent the parameter copy and objective function for agent  $s \in \mathcal{V}$ , respectively, the consensus optimization problem can be written as:

$$\begin{aligned} & \text{minimize}_{\{\mathbf{w}^{(s)}\}_{s \in \mathcal{V}} \in (\mathbb{R}^d)^{|\mathcal{V}|}} \sum_{s \in \mathcal{V}} f^{(s)}(\mathbf{w}^{(s)}) \\ & \text{subject to} \quad \mathbf{w}^{(s)} = \mathbf{w}^{(s')} \quad \forall (s, s') \in \mathcal{E}. \end{aligned} \quad (5.4.2)$$

When the graph is fully connected, i.e. all agents communicate, this problem reduces to an unconstrained problem. However, in general, the solutions of this problem are much difficult to identify, and one must work through minimizing the objective and satisfying the so-called consensus constraints.

## 5.5 Conclusion

Large-scale problems have always pushed optimization algorithms to their limits, and have led to reconsidering certain algorithms in light of their applicability to large-scale settings. Coordinate descent methods are the perfect example of classical techniques that regained popularity because of their efficiency in data science settings. On some instances, randomized coordinate descent techniques bear a close connection with stochastic gradient methods, but are more amenable to large-dimensional problems, particularly those exhibiting sparsity in the problem data.

In modern data science tasks, the amount of data available requires distributed storage, and possibly agents cooperating in order to solve the optimization problem at hand. Linearly constrained formulations can capture this behavior, and ad hoc algorithms such as ADMM are perfectly suited for distributing the optimization effort among all agents.



# Bibliography

- [1] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization Methods for Large-Scale Machine Learning. *SIAM Rev.*, 60:223–311, 2018.
- [2] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, United Kingdom, 2004.
- [3] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer-Verlag, New York, second edition, 2006.
- [4] S. J. Wright and B. Recht. *Optimization for Data Analysis*. Cambridge University Press, 2022.

## Appendix A

# Mock exam: Around the Huber loss

This appendix contains the material of the 2022-2023 exam of the course. It investigated various optimization problems arising from the use of the Huber loss (named after the Swiss statistician Peter J. Huber), a loss function commonly used in robust statistics. Throughout, we define the *Huber loss* as the function  $\ell$  from  $\mathbb{R}$  to  $\mathbb{R}$  such that

$$\ell(t) = \begin{cases} \frac{1}{2}t^2 & \text{if } |t| < 1 \\ |t| - \frac{1}{2} & \text{otherwise.} \end{cases} \quad (\text{A.0.1})$$

This function behaves like  $t \mapsto \frac{t^2}{2}$  for  $|t| < 1$ , and looks like  $t \mapsto |t|$  when  $|t|$  becomes large. Unlike what its expression might suggest, it is continuously differentiable (i.e.  $\ell \in \mathcal{C}^1$ ).

### Part 1: Huber loss and linear models

We consider a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  where  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ . Our goal is to find a linear model that predicts every  $y_i$  given the corresponding  $\mathbf{x}_i$  as best as possible. We thus define a family of model functions parameterized by a vector  $\mathbf{w}$  as follows:

$$\begin{aligned} h_{\mathbf{w}} : \mathbb{R}^d &\rightarrow \mathbb{R} \\ \mathbf{x} &\mapsto \mathbf{x}^T \mathbf{w} = \sum_{i=1}^d [\mathbf{x}]_i [\mathbf{w}]_i. \end{aligned}$$

Given  $h_{\mathbf{w}}$ , we will consider that the model function correctly predicts  $y_i$  from  $\mathbf{x}_i$  if

$$\ell(h_{\mathbf{w}}(\mathbf{x}_i) - y_i) = \ell(\mathbf{x}_i^T \mathbf{w} - y_i) = 0.$$

The value  $\ell(\mathbf{x}_i^T \mathbf{w} - y_i)$  represents the error of the model at  $(\mathbf{x}_i, y_i)$ . Therefore, we are interested in selecting a model (i.e. a vector  $\mathbf{w} \in \mathbb{R}^d$ ) such that the sum of the errors is minimized. As a result, we arrive at the following optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i^T \mathbf{w} - y_i). \quad (\text{A.0.2})$$

a) Justify that 0 is a lower bound for problem (A.0.2). Is it necessarily the minimum value of (A.0.2)?

b) The gradient of  $f$  at  $\mathbf{w} \in \mathbb{R}^d$  is given by

$$\nabla f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell'(\mathbf{x}_i^T \mathbf{w} - y_i) \mathbf{x}_i, \quad (\text{A.0.3})$$

with

$$\ell'(t) = \begin{cases} 1 & \text{if } t > 1 \\ t & \text{if } |t| \leq 1 \\ -1 & \text{if } t < -1. \end{cases}$$

- i) Write down the gradient descent iteration with a constant stepsize  $\alpha$  and using the formula (A.0.3).
- ii) What happens to this iteration if the current point is a local minimum?
- c) A Lipschitz constant for the gradient is given by  $L = \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i\|^2$ .
- i) How can this constant be used to choose the stepsize?
- ii) Provide two additional ways of choosing the stepsize when the value of  $L$  is unknown.
- d) The function  $f$  can be decomposed as  $f = \frac{1}{n} \sum_{i=1}^n f_i$ , with  $f_i(\mathbf{w}) = \ell(\mathbf{x}_i^T \mathbf{w} - y_i)$ . The gradient of  $f_i$  at  $\mathbf{w}$  is

$$\nabla f_i(\mathbf{w}) = \ell'(\mathbf{x}_i^T \mathbf{w} - y_i) \mathbf{x}_i.$$

Write down the stochastic gradient iteration for this problem with a generic stepsize choice.

- e) We suppose that our unit of cost is an access to one  $\mathbf{x}_i$ . What is the cost of a gradient iteration and that of a stochastic gradient iteration?
- f) We now consider a batch stochastic gradient method in which we select a subset of  $n_b$  components.
- i) Write the resulting iteration using a constant stepsize.
- ii) If  $m$  is the number of processors available for computation, what can be the interest of choosing  $n_b = m$ ?
- iii) Practical situation: suppose that we use several batch sizes and we observe that going from  $n_b = 1$  to  $n_b = n/4$  constantly gives better results in terms of convergence speed. Suppose that we also see a degradation in the convergence speed for batch values greater than  $n/4$ . How can this behavior be explained?
- g) Using batches of gradients is one way of reducing the variance in the stochastic gradient estimates. Name one other variance reduction technique among those seen in class.
- h) Consider an instance of problem (A.0.2) for which the components of the stochastic gradient estimates differ by orders of magnitude. Propose (with justification) an advanced stochastic gradient method among those seen in class that could prove efficient given that property.

## Solutions

a) The function  $\ell$  is nonnegative on  $\mathbb{R}$ . Thus, for any  $\mathbf{w} \in \mathbb{R}^d$ ,

$$f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i^T \mathbf{w} - y_i) \geq \frac{1}{n} \sum_{i=1}^n 0 = 0.$$

This shows that the value 0 is a lower bound for problem (??). This value is only attained if there exists a  $\mathbf{w}$  such that  $\mathbf{x}_i^T \mathbf{w} - y_i = 0$  for every  $i$ : this might not always be the case (for instance with  $n = 2, d = 1, \mathbf{x}_1 = 1, \mathbf{x}_2 = -1, y_1 = y_2 = 1$ ), therefore 0 is not necessarily the minimum value of the problem.

b)

i) At a point  $\mathbf{w}_k \in \mathbb{R}^d$ , the gradient descent iteration with a constant stepsize  $\alpha$  is written as follows:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\alpha}{n} \sum_{i=1}^n \ell'(\mathbf{x}_i^T \mathbf{w}_k - y_i) \mathbf{x}_i.$$

ii) If  $\mathbf{w}_k$  is a local minimum, then  $\nabla f(\mathbf{w}_k) = 0$ , and the gradient descent iteration reduces to  $\mathbf{w}_{k+1} = \mathbf{w}_k$ .

c)

i) If a Lipschitz constant  $L$  for the gradient is known, an appropriate choice for a constant stepsize is  $\alpha = \frac{1}{L}$ .

ii) When such a value is unknown, one can use a decreasing stepsize (e.g.  $\alpha_k = \frac{1}{k+1}$ ) or perform a line search at every iteration in order to find an appropriate stepsize value.

d) At  $\mathbf{w}_k \in \mathbb{R}^d$ , the stochastic gradient iteration with generic stepsize  $\alpha_k$  consists in two steps. First, an index  $i_k$  is drawn at random in  $\{1, \dots, n\}$ ; secondly, the new iterate is computed using the formula:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k) = \mathbf{w}_k - \alpha_k \ell'(\mathbf{x}_{i_k}^T \mathbf{w}_k - y_{i_k}) \mathbf{x}_{i_k}.$$

e) Every gradient descent iteration must access all the data in order to compute the gradient: if our unit of cost is an access to one  $\mathbf{x}_i$ , the cost of a gradient descent iteration is  $n$ . As for the stochastic gradient iteration, it only accesses one data point ( $\mathbf{x}_{i_k}$  where  $i_k$  is drawn at random): its cost is thus 1.

f)

i) The batch stochastic gradient iteration at  $\mathbf{w}_k \in \mathbb{R}^d$  consists in two steps. First, a random index set  $S_k \subset \{1, \dots, n\}$  of size  $|S_k| = n_b$  is drawn; then, the following iteration is performed:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\alpha}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k),$$

where  $\alpha > 0$  is the constant stepsize.

ii) If  $m$  processors are available for computation, and the gradients  $\nabla f_i$  can be computed in parallel, the cost of a batch stochastic gradient can be distributed over these  $m$  processors.

- iii) If improvement is observed while using a small batch size, this means that the data is sufficiently correlated that considering a subset of it at every iteration is enough to converge. Using more than one data point at every iteration also leads to steps with a lower variance, and this explains why  $n_b = n/4$  can give better performance than  $n_b = 1$  (classical stochastic gradient iteration). When the batch size gets closer to  $n$ , its cost also gets closer to that of a full gradient iteration, and the method is also at risk of suffering from redundancies in the data. This explains why the behavior of the method worsens when  $n_b > n/4$ .
- g) Gradient aggregation is another technique that reduces the variance. *Another valid answer was iterate averaging.*
- h) The use of a variant with diagonal scaling is appropriate in this setting, since it uses different stepsizes for every coordinate. As such, it will be less sensitive to differences of magnitude between gradient components.

## Part 2: Pseudo-Huber loss

The goal of this part is to replace the Huber loss by a (smoothed) *pseudo-Huber loss function*, namely:

$$\begin{aligned} p: \mathbb{R} &\rightarrow \mathbb{R} \\ t &\mapsto p(t) := \sqrt{1+t^2} - 1. \end{aligned} \quad (\text{A.0.4})$$

It can be shown that the function  $p$  is twice continuously differentiable (whereas  $\ell$  is only once continuously differentiable). Its derivatives are given for every  $t \in \mathbb{R}$  by

$$\nabla p(t) = \frac{t}{\sqrt{1+t^2}} \quad \text{and} \quad \nabla^2 p(t) = \frac{1}{1+t^2}.$$

- a) Justify that the function  $p$  is convex. What can be said of its local minima ?
- b) Show that  $\operatorname{argmin}_{t \in \mathbb{R}} p(t) = \{0\}$ .
- c) Using the same data as that of Exercise 1, consider the problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad g(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n p(\mathbf{x}_i^\top \mathbf{w} - y_i). \quad (\text{A.0.5})$$

The function  $g$  is convex but not strongly convex in general.

- i) What is the convergence rate of gradient descent applied to problem (A.0.5)? What quantity does this rate apply to?
- ii) Give the convergence rate of accelerated gradient on problem (A.0.5). Is it better or worse than that of gradient descent?
- d) Suppose that we apply a stochastic gradient method to problem (A.0.5) by exploiting its finite-sum structure, and that the method we apply has a convergence rate (in expectation) in  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$  for the same quantity than that considered in question c), with  $K$  being the iteration count.
  - i) In terms of iterations, justify that this rate is worse than that of gradient descent.

- ii) Is there another metric in which this rate can be better than that of gradient descent? If so, which metric?
- iii) Does the result from the previous question apply when we compare the rates for stochastic gradient and accelerated gradient?

### Solutions

a) From the formula for the second derivative of  $p$ , we observe that

$$\nabla^2 p(t) = \frac{1}{1+t^2} > 0 \quad \forall t \in \mathbb{R}.$$

Therefore, the function  $p$  has a positive definite Hessian, which means that it is convex. This implies that any local minimum of  $u$  is a global minimum.

b) We have  $p(t) \geq 0$  for every  $t \in \mathbb{R}$ , and  $p(0) = 0$ : this allows us to conclude that 0 is a global minimum of the function. In addition, for any value  $t > 0$ , we have  $p(t) > 0$ , showing that 0 is the unique global minimum.

c)

i) After  $K \geq 1$  iterations of gradient descent, the convergence rate guarantee is

$$g(\mathbf{w}_k) - \min_{\mathbf{w} \in \mathbb{R}^d} g(\mathbf{w}) \leq \mathcal{O}\left(\frac{1}{K}\right).$$

ii) The convergence rate of accelerated gradient on this problem is  $\mathcal{O}\left(\frac{1}{K^2}\right)$ , which is better than that for gradient descent.

d)

i) The sequence  $\left\{\frac{1}{\sqrt{K}}\right\}$  converges to 0 more slowly than  $\left\{\frac{1}{K}\right\}$ , which is the rate of gradient descent. As a result, the rate for stochastic gradient is worse in terms of iterations.

ii) By comparing the convergence rates in terms of epochs rather than iterations, we obtain (for a fixed epoch number  $N_E \geq 1$ ) a rate of  $\mathcal{O}\left(\frac{1}{\sqrt{nN_E}}\right)$  for stochastic gradient and a rate of  $\mathcal{O}\left(\frac{1}{N_E}\right)$  for gradient descent. The former is better when  $n \gg N_E$ .

iii) The same observation applies to the rates for stochastic gradient and accelerated gradient, which become  $\mathcal{O}\left(\frac{1}{\sqrt{nN_E}}\right)$  and  $\mathcal{O}\left(\frac{1}{N_E^2}\right)$ , respectively. For stochastic gradient to yield a better convergence rate, one needs  $n \gg N_E^2$ , which is a stronger requirement.

### Part 3: Reversed Huber loss

In this part, we consider the reverse philosophy of the Huber loss, that is, we propose to use a loss function that looks like the absolute value on  $[-1, 1]$  and like a quadratic everywhere else.

The *reversed Huber loss* is thus defined as:

$$\begin{aligned} r : \mathbb{R} &\rightarrow \mathbb{R} \\ t &\mapsto r(t) := \begin{cases} |t| & \text{if } |t| < 1 \\ \frac{t^2+1}{2} & \text{otherwise.} \end{cases} \end{aligned} \quad (\text{A.0.6})$$

This function is convex but nonsmooth, since it is not differentiable at 0.

As in the first part, we consider linear models  $\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x}$  and a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ .

a) We first consider the convex, nonsmooth problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i^T \mathbf{w} - y_i). \quad (\text{A.0.7})$$

- i) What mathematical tool can we use to design algorithms applicable to problem (A.0.7)?
- ii) Using this tool, how can the solutions of (A.0.7) be characterized?

b) We now study the family of problems:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \lambda \sum_{i=1}^d r([\mathbf{w}]_i), \quad (\text{A.0.8})$$

where  $f$  is the objective function of (A.0.2), and  $\lambda > 0$ .

- i) How is this type of problem called? What is the purpose of the second term?
- ii) Write the generic proximal gradient iteration for this problem.
- iii) When is this algorithm worthy of consideration in practice?

## Solutions

a)

- i) Since  $\mathbf{w} \mapsto \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i^T \mathbf{w} - y_i)$  is convex, it is possible to define the subdifferential of  $v$  at any point: the elements of the subdifferential, called the subgradients, can be used in lieu of the gradient to construct optimization methods for solving problem (A.0.7).
- ii) Let  $\phi_r : \mathbf{w} \mapsto \frac{1}{n} \sum_{i=1}^n r(\mathbf{x}_i^T \mathbf{w} - y_i)$ . A point  $\bar{\mathbf{w}} \in \mathbb{R}^d$  is a global minimum of  $\phi_r$  if and only if

$$\mathbf{0} \in \partial \phi_r(\bar{\mathbf{w}}),$$

where  $\partial \phi_r(\cdot)$  denotes the subdifferential of  $\phi_r$ .

b)

- i) Problem (A.0.8) is a regularized (or composite) optimization problem. The goal of the second term, that does not depend on data, is to enforce desired properties for the solution.
- ii) At a point  $\mathbf{w}_k$ , the generic proximal gradient iteration (with a generic stepsize  $\alpha_k$ ) for this problem is:

$$\mathbf{w}_{k+1} \in \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|_2^2 + \lambda \sum_{i=1}^d r([\mathbf{w}]_i) \right\}.$$

- iii) The proximal gradient algorithm is only interesting when the cost of solving the subproblem is cheaper than that of solving the original problem.

## Part 4: Large-scale reversed Huber loss

In this part, we consider an instance of the problem family (A.0.8). More precisely, we focus on

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) + \gamma \Omega(\mathbf{w}), \tag{A.0.9}$$

for a fixed value  $\gamma \geq 0$ , where  $f = \frac{1}{n} \sum_{i=1}^n f_i$  is defined as in Exercise 1 and  $\Omega(\mathbf{w}) := \sum_{i=1}^d r([\mathbf{w}]_i)$  with  $r$  being the function (A.0.6) defined in Exercise 3.

- a) Suppose first that the number of parameters  $d$  is quite large.
  - i) Assuming  $\gamma = 0$ , write down a block coordinate descent iteration for problem (A.0.9).
  - ii) How can you combine this iteration with other algorithms seen in class to develop a method based on coordinate updates that can be applied to problem (A.0.9) with  $\gamma > 0$ ?
- b) We now introduce an auxiliary variable to the problem, which we then rewrite as

$$\begin{aligned} &\underset{\substack{\mathbf{u} \in \mathbb{R}^d \\ \mathbf{v} \in \mathbb{R}^d}}{\text{minimize}} && f(\mathbf{u}) + \gamma \Omega(\mathbf{v}) \\ &\text{subject to} && \mathbf{u} - \mathbf{v} = \mathbf{0}_{\mathbb{R}^d}. \end{aligned} \tag{A.0.10}$$

- i) Write down the augmented Lagrangian formula for problem (A.0.10).
  - ii) Which method is based on the augmented Lagrangian and can exploit the structure of problem (A.0.10)? What is the main idea behind exploiting such a structure?
- c) Suppose now that the number of data points used in defining  $f$  is so large that all  $f_i$  are spread across several agents, each of which has only access to its own  $f_i$  and maintains its own copy  $\mathbf{u}^{(i)}$  of  $\mathbf{u}$ . All agents share knowledge of  $\mathbf{v}, \gamma$  and  $\Omega$ .
  - i) Rewrite problem (A.0.10) to model the distributed aspect of the problem as described above.
  - ii) How can the method from question b)ii) be adapted to this new setting?

### Solutions

- a)
  - i) Starting from  $\mathbf{w}_k \in \mathbb{R}^d$ , a block coordinate descent iteration with a generic stepsize  $\alpha_k > 0$  can be written as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \sum_{j \in \mathcal{B}_k} [\nabla f(\mathbf{w}_k)]_j \mathbf{e}_j,$$

where  $\mathcal{B}_k \subset \{1, \dots, d\}$  is a block of coordinate indices

- ii) Since  $\Omega$  is separable, we can combine proximal gradient and coordinate descent to perform updates only concerned with the coordinates in the block, leading to the iteration

$$\mathbf{w}_{k+1} \in \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \left\{ f(\mathbf{w}_k) + \sum_{j \in \mathcal{B}_k} [\nabla f(\mathbf{w}_k)]_j \mathbf{e}_j^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|^2 + \sum_{j \in \mathcal{B}_k} r([\mathbf{w}]_j) \right\}$$

*An answer without a formula could have sufficed.*



b)

i) The augmented Lagrangian for problem (A.0.10) is

$$\mathcal{L}^a(\mathbf{u}, \mathbf{v}, \mathbf{z}; \lambda) = f(\mathbf{u}) + \gamma\Omega(\mathbf{v}) + \mathbf{z}^T(\mathbf{u} - \mathbf{v}) + \frac{\lambda}{2}\|\mathbf{u} - \mathbf{v}\|^2.$$

ii) The Alternating Direction Method of Multipliers, or ADMM. Its goal consists in exploiting separable structure in the problem to perform possibly cheaper updates consecutively rather than jointly.

c)

i) The problem can be rewritten by adding the copies as variables, so that we obtain

$$\begin{aligned} & \underset{\substack{\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n)} \in \mathbb{R}^d \\ \mathbf{v} \in \mathbb{R}^d}}{\text{minimize}} && \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{u}^{(i)}) + \gamma\Omega(\mathbf{v}) \\ & \text{subject to} && \mathbf{u}^{(i)} - \mathbf{v} = \mathbf{0}_{\mathbb{R}^d} \quad \forall i = 1, \dots, n. \end{aligned}$$

ii) To adapt ADMM, one would consider concurrent updates on the variables  $\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n)}$  so as to exploit the separability of the objective further.

## Appendix B

# Mock exam: Shallow neural networks

In this appendix, we review the material from the lectures through a study of (basic) shallow neural network architectures.

### Part 1: Two-layer linear neural networks

We consider a dataset  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  where  $\mathbf{x}_i \in \mathbb{R}^{d_x}$  and  $\mathbf{y}_i \in \mathbb{R}^{d_y}$ . We wish to learn a mapping from  $\mathbb{R}^{d_x}$  to  $\mathbb{R}^{d_y}$  that correctly outputs  $\mathbf{y}_i$  when given  $\mathbf{x}_i$  as an input. Our model will be that of a two-layer linear neural network :

$$\begin{aligned} \mathbf{h}(\cdot; \mathbf{w}) : \mathbb{R}^{d_x} &\longrightarrow \mathbb{R}^{d_y} \\ \mathbf{x} &\longmapsto \mathbf{W}_2(\mathbf{W}_1\mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2, \end{aligned} \quad (\text{B.0.1})$$

where  $\mathbf{W}_1 \in \mathbb{R}^{d_x \times m}$ ,  $\mathbf{b}_1 \in \mathbb{R}^m$ ,  $\mathbf{W}_2 \in \mathbb{R}^{m \times d_y}$  and  $\mathbf{b}_2 \in \mathbb{R}^{d_y}$ . We will consider  $\mathbf{h}$  as being parameterized by  $\mathbf{w} \in \mathbb{R}^d$ , with  $d = d_x m + m + m d_y + d_y$  and  $\mathbf{w}$  concatenating all coefficients from  $\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2$ . Our goal is to determine a value of  $\mathbf{w}$  so that  $\mathbf{h}(\mathbf{x}_i; \mathbf{w}) \approx \mathbf{y}_i$ , which we formalize using the squared loss  $(\mathbf{h}, \mathbf{y}) \mapsto \frac{1}{2} \|\mathbf{h} - \mathbf{y}\|^2$ .

Overall, we obtain the following problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n \|\mathbf{h}(\mathbf{x}_i; \mathbf{w}) - \mathbf{y}_i\|^2. \quad (\text{B.0.2})$$

- a) Give a lower bound on the objective function of problem (B.0.2).
- b) In general, problem (B.0.2) is nonconvex. What does this imply about its local minima?
- c) The function  $f$  is continuously differentiable, or  $\mathcal{C}^1$ .
  - i) Suppose that  $\mathbf{w}^*$  is a solution of (B.0.2). What can be said about the derivative of  $f$  at  $\mathbf{w}^*$ ?
  - ii) Write down the gradient descent iteration for problem (B.0.2) with an arbitrary stepsize.
  - iii) Given that the problem is nonconvex, what is the theoretical convergence rate of gradient descent applied to (B.0.2)?

d) We now exploit the fact that  $f$  has a finite-sum structure :

$$f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) \quad \text{with} \quad f_i(\mathbf{w}) = \frac{1}{2} \|\mathbf{h}(\mathbf{x}_i; \mathbf{w}) - \mathbf{y}_i\|^2 \quad \forall i = 1, \dots, n.$$

Every  $f_i$  is  $\mathcal{C}^1$ .

- i) Write down the iteration of stochastic gradient for this problem with a constant stepsize.
- ii) Can we guarantee that a run of stochastic gradient will converge on this problem? Justify your answer.
- iii) What is the main computational advantage of stochastic gradient over gradient descent?
- iv) Recall the definition of an epoch. How many iterations of stochastic gradient does an epoch correspond to?
- v) Write down the iteration of a batch stochastic gradient method with fixed batch size  $n_b \in \{1, \dots, n\}$  and an arbitrary stepsize.
- vi) Compare the cost of one iteration of the batch stochastic gradient from the previous question with that of one iteration of stochastic gradient.
- vii) In a parallel environment, how can the cost of a batch approach be reduced?

### Solutions

- a) Any value less than or equal to 0 is a lower bound on the objective function of (B.0.2). *Note that 0 is the only attainable value, but that is not necessarily attained by the function.*
- b) The local minima of problem (B.0.2) may not be global minima.
- c) (Using the  $\mathcal{C}^1$  nature of the objective.)
  - i) The derivative, or gradient of  $f$  at  $\mathbf{w}^*$  is zero:  $\nabla f(\mathbf{w}^*) = \mathbf{0}$ .
  - ii) The  $k$ th iteration of gradient descent for problem (B.0.2) is

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k),$$

where  $\alpha_k > 0$ .

- iii) For nonconvex problems, the convergence rate of gradient descent is  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$ : that is, after  $K \geq 1$  iterations, we can guarantee that

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right).$$

- d) (Exploiting the finite-sum structure.)
  - i) The  $k$ th iteration of stochastic gradient is

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \nabla f_{i_k}(\mathbf{w}_k),$$

where  $\alpha > 0$  is the chosen constant stepsize, and  $i_k$  is an index drawn at random in  $\{1, \dots, n\}$ .

- ii) Because of the intrinsic randomness within stochastic gradient, we cannot guarantee that it will converge to a solution of the problem.
- iii) Stochastic gradient only accesses one data point per iteration, whereas an iteration of gradient descent must access all data points in order to compute a gradient.
- iv) An epoch is a unit of cost corresponding to  $n$  access to one element in the dataset. Since every iteration of stochastic gradient accesses a single element, one epoch corresponds (in terms of cost) to  $n$  iterations of stochastic gradient.
- v) The  $k$ th iteration of the proposed batch stochastic gradient method is

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \frac{\alpha}{n_b} \sum_{i \in \mathcal{S}_k} \nabla f_i(\mathbf{w}_k),$$

where  $\mathcal{S}_k$  is a set of  $n_b$  random indices drawn in  $\{1, \dots, n\}$  with or without replacement, and  $\alpha > 0$  is the chosen constant stepsize.

- vi) In terms of accesses to data points, one iteration of the proposed batch stochastic gradient has a cost of  $n_b$ , which is  $n_b$  times more expensive than that of one iteration of stochastic gradient, where only a single index is accessed.
- vii) A parallel environment can enable parallel accesses to data points, thereby reducing the cost of a batch stochastic gradient approach. *Note: The cost will only be reduced according to the number of tasks that can be performed in parallel.*

## Part 2: Two-layer ReLU neural network

Building on Part 1, we consider a dataset  $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^{d_x}$  and  $\mathbf{y}_i \in \mathbb{R}^{d_y}$ . Our model will now consist in a two-layer neural network with nonsmooth rectified linear unit (ReLU) activation:

$$\begin{aligned} \mathbf{h}^{ReLU}(\cdot; \mathbf{w}) : \mathbb{R}^{d_x} &\longrightarrow \mathbb{R}^{d_y} \\ \mathbf{x} &\longmapsto \mathbf{W}_2 \boldsymbol{\sigma}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2, \end{aligned} \tag{B.0.3}$$

where  $\mathbf{W}_1 \in \mathbb{R}^{d_x \times m}$ ,  $\mathbf{b}_1 \in \mathbb{R}^m$ ,  $\mathbf{W}_2 \in \mathbb{R}^{m \times d_y}$ ,  $\mathbf{b}_2 \in \mathbb{R}^{d_y}$ , and  $\boldsymbol{\sigma} : \mathbb{R}^m \rightarrow \mathbb{R}^m$  is defined componentwise by

$$\forall \mathbf{v} \in \mathbb{R}^m, \quad \boldsymbol{\sigma}(\mathbf{v}) := [\max\{v_i, 0\}]_{i=1}^m.$$

As before,  $\mathbf{w} \in \mathbb{R}^d$  with  $d = d_x m + m + m d_y + d_y$  concatenates all the coefficients of the  $\mathbf{W}_i$  matrices and the  $\mathbf{b}_i$  vectors.

Using a square loss  $\frac{1}{2} \|\cdot\|^2$ , our training problem thus becomes

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f^{ReLU}(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n \|\mathbf{h}^{ReLU}(\mathbf{x}_i; \mathbf{w}) - \mathbf{y}_i\|^2. \tag{B.0.4}$$

- a) Justify that gradient-based methods (such as gradient descent) cannot be directly applied to that problem.
- b) What other approaches can be used to tackle problem (B.0.4)?
- c) We now look at the activation function  $r : t \mapsto \max\{t, 0\}$ , which is convex on  $\mathbb{R}$ .

i) Using the expression of  $r(t)$ , justify that

$$\operatorname{argmin}_{t \in \mathbb{R}} \{r(t)\} = \{t \in \mathbb{R}, t \leq 0\},$$

that is, every  $t \leq 0$  is a global minimum of  $r$ .

ii) For any  $t \in \mathbb{R}$ , we define the set  $\mathcal{S}(t)$  as

$$\mathcal{S}(t) := \begin{cases} 1 & \text{if } t > 0, \\ 0 & \text{if } t < 0, \\ [0, 1] & \text{if } t = 0. \end{cases}$$

How does this set confirm the result of question i)?

### Solutions

- a) The objective function of problem (B.0.4) is not differentiable at every point. Therefore, the gradient does not exist at certain points, and this prevents from applying gradient-based methods like gradient descent.
- b) Subgradient methods can be applied to problem (B.0.4)?
- c) (Activation function  $r : t \mapsto \max\{t, 0\}$ .)

i) For every  $t \leq 0$ , we have

$$r(t) = 0 \leq r(s) \quad \forall s \in \mathbb{R}.$$

By definition, this means that  $t$  is a global minimum of the function  $r$ , leading to

$$\operatorname{argmin}_{t \in \mathbb{R}} \{r(t)\} = \{t \in \mathbb{R}, t \leq 0\},$$

- ii) The set  $\mathcal{S}(t)$  represents the subdifferential of the function  $r$  at  $t$ . From its expression, we see that  $0 \in \mathcal{S}(t)$  if and only if  $t \leq 0$ , therefore any  $t \leq 0$  is a global minimum of the convex function  $r$ .

### Part 3: One-layer linear neural network

In this exercise, we consider the special case of a dataset with scalar labels/outputs, i.e. of the form  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $\mathbf{x}_i \in \mathbb{R}^{d_x}$  and  $y_i \in \mathbb{R}$  for every  $i = 1, \dots, n$ . We build a simple neural network with no activation function and one homogeneous linear layer to predict the value  $y_i$  from the vector  $\mathbf{x}_i$ , resulting in the model

$$\begin{aligned} h^{lin}(\cdot; \mathbf{w}) : \mathbb{R}^{d_x} &\longrightarrow \mathbb{R} \\ \mathbf{x} &\longmapsto \mathbf{W}_1 \mathbf{x}, \end{aligned} \tag{B.0.5}$$

with  $\mathbf{W}_1 \in \mathbb{R}^{1 \times d_x}$ . Letting  $d = d_x$  and  $\mathbf{w} = \mathbf{W}_1^T \in \mathbb{R}^d$ , finding the best model amounts to solving

$$\operatorname{minimize}_{\mathbf{w} \in \mathbb{R}^d} f^{lin}(\mathbf{w}) := \frac{1}{2n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2. \tag{B.0.6}$$

- a) What class of problems does problem (B.0.6) belong to?
- b) The objective function  $f^{lin}$  is  $\mathcal{C}_L^{1,1}$ , i.e. its gradient is  $L$ -Lipschitz continuous. If  $L$  is known, how can its value be used in an algorithm such as gradient descent?
- c) Problem (B.0.6) is convex with a  $\mathcal{C}^1$  objective function.
  - i) What can then be said about a point  $\bar{\mathbf{w}}$  such that  $\nabla f^{lin}(\bar{\mathbf{w}}) = \mathbf{0}_{\mathbb{R}^d}$ ?
  - ii) What is the convergence rate of gradient descent on this problem?
  - iii) What is the convergence rate of accelerated descent on a convex problem? Is it better or worse than that of the previous question ?
- d) Suppose that the data is such that the objective  $f^{lin}$  is  $\mu$ -strongly convex, in addition to the properties already mentioned above.
  - i) Let  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$  be two points such that  $\nabla f^{lin}(\mathbf{w}) = \nabla f^{lin}(\mathbf{v}) = \mathbf{0}_{\mathbb{R}^d}$ . What can we say about  $\mathbf{v}$  and  $\mathbf{w}$ ?
  - ii) What is the convergence rate of accelerated gradient on this problem?
- e) If  $f^{lin}$  is convex but not strongly convex, how can problem (B.0.6) be modified into a strongly convex one without changing its problem class?

### Solutions

- a) Problem (B.0.6) is a linear least-squares problems (it also belongs to the class of quadratic optimization problems).
- b) If  $f^{lin}$  is  $\mathcal{C}_L^{1,1}$  with  $L$  is known, any positive value less than  $\frac{2}{L}$  can be used as a constant stepsize: the value  $\frac{1}{L}$  gives precise decrease guarantees at every iteration. *Note: The answer  $\frac{1}{L}$  would suffice here.*
- c) (Convexity of problem (B.0.6).)

- i) If  $\nabla f^{lin}(\bar{\mathbf{w}}) = \mathbf{0}_{\mathbb{R}^d}$ , we know that this point is a global minimum because the function is convex.
- ii) Since the objective function is convex, the convergence rate of gradient descent is  $\mathcal{O}(\frac{1}{K})$ : that is, for every  $K \geq 1$ , if  $\mathbf{w}_K$  is the  $K$ th iterate of gradient descent, we have

$$f^{lin}(\mathbf{w}_K) - \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \leq \mathcal{O}\left(\frac{1}{K}\right).$$

- iii) The convergence rate of accelerated descent on a convex problem is  $\mathcal{O}(\frac{1}{K^2})$ : this is better than the rate for gradient descent, in that it converges more rapidly towards 0 as  $K$  increases.

- d) ( $f^{lin}$   $\mu$ -strongly convex.)

- i) If  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$  are such that  $\nabla f^{lin}(\mathbf{w}) = \nabla f^{lin}(\mathbf{v}) = \mathbf{0}_{\mathbb{R}^d}$ , then they are both global minima (recall that  $f^{lin}$  is convex). But since  $f^{lin}$  is strongly convex, it has a unique global minimum, from which we conclude that  $\mathbf{v} = \mathbf{w}$ .
- ii) Accelerated gradient has a convergence rate in  $\mathcal{O}\left((1-t)^K\right)$  on this problem with  $t \in (0, 1)$ .  
With standard assumptions, it is possible to establish this rate with  $t = \sqrt{\frac{\mu}{L}}$ .

- e) One way to turn problem (B.0.6) into a strongly convex one without changing its problem class consists in adding an  $\ell_2$  regularization term, so that (B.0.6) becomes

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f^{\text{lin}}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

for some  $\lambda > 0$ . This problem is still a quadratic optimization problem (and can even be expressed as a linear least-squares problems), but the objective function is now  $\lambda$ -strongly convex.

## Part 4: Revised one-layer linear neural network

Building on Part 3, we finally consider the linear neural network model  $h^{\text{lin}}(\mathbf{x}; \mathbf{w})$  defined in (B.0.5). In order to train this model on a dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , we consider the optimization problem

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f_{\ell_1}(\mathbf{w}) := f^{\text{lin}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1, \quad (\text{B.0.7})$$

where  $f^{\text{lin}}$  is the objective function of problem (B.0.6),  $\lambda > 0$  and  $\|\mathbf{w}\|_1 = \sum_{i=1}^d |[\mathbf{w}]_i|$ .

- What is the purpose of adding the  $\lambda \|\mathbf{w}\|_1$  term to the objective? How are problems of this form called?
- We recall that  $f^{\text{lin}} \in \mathcal{C}^1$ . Write down the proximal gradient iteration for problem (B.0.7) in its generic form, using an arbitrary stepsize.
- In the specific case of problem (B.0.7), the proximal gradient iteration corresponds to the ISTA iteration. What is the interest of using the ISTA formula compared to that of the previous question?
- A possible reformulation of (B.0.7) as a constrained optimization problem is

$$\begin{aligned} & \underset{\mathbf{u}, \mathbf{v} \in \mathbb{R}^d}{\text{minimize}} f^{\text{lin}}(\mathbf{u}) + \lambda \|\mathbf{u}\|_1 \\ & \text{s. t.} \quad \mathbf{u} - \mathbf{v} = \mathbf{0}. \end{aligned} \quad (\text{B.0.8})$$

- Justify that the set of solutions of problems (B.0.7) and (B.0.8) are identical.
- For any  $\rho > 0$ , form the augmented Lagrangian associated to problem (B.0.8) with parameter  $\rho$ .
- Write down the ADMM iteration for problem (B.0.8) using  $\rho_k > 0$  as a stepsize. What is the interest of such an approach?

## Solutions

- The term  $\lambda \|\mathbf{w}\|_1$  is added to the objective function to promote sparse solutions. As a result, the problem becomes a regularized optimization problem, in a composite form.
- The  $k$ th iteration of proximal gradient for problem (B.0.6) is given by

$$\mathbf{w}_{k+1} \in \underset{\mathbf{w} \in \mathbb{R}^d}{\text{argmin}} \left\{ f^{\text{lin}}(\mathbf{w}_k) + \nabla f^{\text{lin}}(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|^2 + \lambda \|\mathbf{w}\|_1 \right\},$$

where  $\alpha_k > 0$  is the stepsize for iteration  $k$ .

c) The ISTA iteration gives an explicit formula for  $\mathbf{w}_{k+1}$  given  $\mathbf{w}_k$  and  $\alpha_k$ , which means that the next iterate is computed explicitly without the need to solve a subproblem.

d) (Constrained reformulation.)

i) For any feasible point of problem (B.0.8), we have  $\mathbf{u} = \mathbf{v}$ , and thus the objective value is equal to  $f^{lin}(\mathbf{u}) + \lambda\|\mathbf{u}\|_1$ . With this observation in mind, we have

$$\operatorname{argmin}_{(\mathbf{u}, \mathbf{v})} \left\{ f^{lin}(\mathbf{u}) + \lambda\|\mathbf{v}\|_1 \mid \mathbf{u} - \mathbf{v} = \mathbf{0} \right\} = \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^d} \left\{ f^{lin}(\mathbf{u}) + \lambda\|\mathbf{u}\|_1 \right\}.$$

The set of the right-hand side is the set of solutions of problem (B.0.7) by definition, which concludes the argument.

ii) For any  $\rho > 0$ , the augmented Lagrangian of problem (B.0.8) with parameter  $\rho$  is given by

$$\mathcal{L}^a(\mathbf{u}, \mathbf{v}, \mathbf{z}; \rho) = f^{lin}(\mathbf{u}) + \lambda\|\mathbf{v}\|_1 + \mathbf{z}^T(\mathbf{u} - \mathbf{v}) + \frac{\rho}{2}\|\mathbf{u} - \mathbf{v}\|^2.$$

iii) The  $k$ th iteration of ADMM applied to problem (B.0.8) is

$$\begin{cases} \mathbf{u}_{k+1} \in \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^d} \mathcal{L}^a(\mathbf{u}, \mathbf{v}_k, \mathbf{z}_k; \rho_k) \\ \mathbf{v}_{k+1} \in \operatorname{argmin}_{\mathbf{v} \in \mathbb{R}^d} \mathcal{L}^a(\mathbf{u}_{k+1}, \mathbf{v}, \mathbf{z}_k; \rho_k) \\ \mathbf{z}_{k+1} = \mathbf{z}_k + \rho_k(\mathbf{A}\mathbf{u}_{k+1} - \mathbf{v}_{k+1}), \end{cases}$$

with  $\rho_k > 0$  being the augmented Lagrangian parameter. This approach exploits the separable nature of the objective by solving two subproblems at every iteration, akin to two iterations of block coordinate descent.