

Optimization for Machine Learning

Clément W. Royer

Lecture notes - M2 MIAE ID Apprentissage - 2024/2025

- The last version of these notes can be found at:
<https://www.lamsade.dauphine.fr/~croyer/ensdocs/OID/PolyOID.pdf>.
- Comments, typos, etc, can be sent to clement.royer@lamsade.dauphine.fr.
Thanks to Sébastien Kerleau and Florian Le Bronnec for their feedback on earlier versions of these notes.
- **Version history:**
 - 2024.09.23: First version of the lecture notes with Chapters 1 and 2.
- **Learning goals:**
 - Understand the specifics of optimization problems, and the interest of certain formulations over others.
 - Given an optimization problem, select an algorithm well suited for solving the problem.
 - Analyze the theoretical and practical properties of a particular algorithm.
 - Identify challenges posed by optimization in a data science context, and ways to address these challenges.

Note: This course was previously taught in M2 MIAGE ID Apprentissage under the name "Optimization for Data and Decision Sciences".

Contents

1	Introduction to optimization	4
1.1	About optimization	4
1.1.1	The optimization process	4
1.1.2	Modern optimization	5
1.2	The optimization problem	6
1.2.1	Mathematical background	6
1.2.2	First definitions	10
1.2.3	Convexity	11
1.3	Examples of optimization problems in ML	13
1.3.1	Linear regression	13
1.3.2	Logistic regression	14
1.3.3	Linear SVM	15
1.3.4	Neural networks	15
1.4	Optimization algorithms	16
1.4.1	The algorithmic process	17
1.4.2	Convergence and convergence rates	17
1.4.3	Popular optimization packages	18
2	Unconstrained optimization	19
2.1	Gradient descent	19
2.1.1	Algorithm	19
2.1.2	Choosing the stepsize	21
2.1.3	Theoretical analysis for gradient descent	22
2.2	Acceleration	25
2.2.1	Introduction: the momentum principle	25
2.2.2	Nesterov's accelerated gradient method	26
2.2.3	Other accelerated methods	27
3	Stochastic gradient techniques	29
3.1	Introduction	29
3.2	Stochastic gradient method	30
3.2.1	Algorithm	30
3.2.2	Convergence rate analysis	31
3.3	Variance reduction	34
3.3.1	Batch methods	34

3.3.2	Other variance reduction techniques	35
3.4	Stochastic gradient methods for deep learning	36
3.4.1	Stochastic gradient with momentum	36
3.4.2	AdaGrad	37
3.4.3	RMSProp	37
3.4.4	Adam	38
3.5	Conclusion	38
Appendix A Notations and mathematical tools		41
A.1	Notations	41
A.1.1	Generic notations	41
A.1.2	Scalar and vector notations	41
A.1.3	Matrix notations	42
A.2	Mathematical tools	43
A.2.1	Vector linear algebra	43
A.2.2	Matrix linear algebra	45
A.2.3	Calculus	47

Chapter 1

Introduction to optimization

1.1 About optimization

Optimization is a field of study that is concerned with **making the best possible decision out of a set of alternatives**. *Mathematical* optimization provides a framework to model optimization problems using mathematical language. We provide below a broad definition of optimization problems in the mathematical sense, and connect it to the concept of solving such a problem.

1.1.1 The optimization process

Most optimization problems are readily formulated in plain language, and it can be difficult to extract the core ideas behind the optimization task. On the contrary, a mathematical optimization problem is a well-identified object, that corresponds to the following definition.

Definition 1.1.1 A *mathematical optimization problem* consists of three key components:

- An **objective function**, that quantifies the quality of a decision. If a better decision yields a smaller objective value, the problem is called a minimization problem. On the other hand, if a better decision yields a larger objective value, the problem is called a maximization problem.
- A number of **decision variables**: those are the knobs that can be turned to change the decision and, as a result, the value of the objective function. The goal of optimization is to determine what are the best (optimal) values of the decision variables relatively to the objective function.
- A set of **constraints**, that specify requirements on the decision variables that must be satisfied for the decision to be valid.

An optimization procedure consists in minimizing or maximizing an objective function with respect to the decision variables **subject to** constraints on those variables.

Remark 1.1.1 In practice, many optimization problems are subject to hidden or implicit constraints, that may not appear in the optimization formulation yet can have a significant effect on an algorithm's performance. The absence of these constraints in the formulation may be due to misspecification (not including a positivity constraint for a variable corresponding to a mass, for instance), but this phenomenon often has a more subtle cause. In simulation-based optimization, where the objective function corresponds to running expensive computational codes, an error may be triggered by certain

values of the decision variables that may not be known a priori. Such constraints must then be addressed by an algorithm in real time.

Any concrete problem goes through a **modeling phase** that produces a mathematical optimization problem. In order to solve this problem, we then use **optimization algorithms**, often designed to be implemented and run on a computer. Such methods typically output a candidate solution (or, possibly, that no solution exists), which may or may not form an appropriate answer to the problem. If needed, the optimization process can go over another modeling phase followed by another solve. In practice, it is typical that the optimization process includes a **discussion phase** with experts from the application domain.

Remark 1.1.2 *Numerical optimization typically obeys the following principles:*

- *There is no universal algorithm. Every method can be efficient on a certain class of problems and perform poorly on another class. Studying the structure of a problem of interest helps in choosing the best algorithm for solving this problem.*
- *There may be a big gap between theory and practice. Finite-precision arithmetic can introduce round-off errors that result in worse practical performance than what the theory predicts. Conversely, theory (such as complexity bounds introduced below) can be overly pessimistic, and much better results can be observed in practice.*
- *Theory informs practice, and vice-versa. For most optimization problems, one can use mathematical formulas to assess whether a solution has been found: these expressions are at the heart of most of the methods we will describe in this course. More broadly, theoretical guarantees can help guiding the optimization process beyond pure heuristics and guesses. Meanwhile, some algorithms have demonstrated success without being endowed with theoretical properties: such phenomena are common, and motivate researchers to explain this behavior using mathematical tools.*

1.1.2 Modern optimization

Numerical optimization started during the 1940s. Major theoretical advances were achieved during the 1980s, along with significant algorithmic developments, that leverage the limited computing power of this era. The next two decades saw the rise of computing power, leading to numerous successes of optimization techniques. This trend continues to this day, albeit with a key change of paradigm.

Indeed, as data becomes more prevalent, optimization problems now involve massive amounts of data in the calculation of their objective. In addition to the challenges posed by the use of such datasets, possibly in a distributed fashion, other difficulties arise due to the fact that the solutions of these problems must generalize to yet unseen data.

In this course, we will go on a *tour d'horizon* of optimization problems and algorithms, with a focus on methods that have been successful in classical industrial settings as well as new paradigms such as machine learning. Our presentation will cover theoretical, algorithmic and practical aspects.

1.2 The optimization problem

We now introduce the mathematical foundations behind optimization. As a field of study, optimization is defined as the process of 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

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) \quad \text{subject to} \quad \mathbf{w} \in \mathcal{F}. \quad (1.2.1)$$

In problem (1.2.1), f is the objective function (that we want to minimize), \mathbf{w} is the vector of decision variables defined over a definition set \mathbb{R}^d , and \mathcal{F} is a set encompassing all the constraints on the decision variables. This set is called the feasible set, and is often described using mathematical expressions.

There are multiple ways of formulating an optimization problem given a description of this problem in plain language: constraints can be expressed in the definition set or in the objective, maximization can be preferred over minimization, redundant constraints can be added to the problem, etc. Some formulations will be well suited for theoretical analysis, while others will be efficiently solvable by modern software.

Definition 1.2.1 *Two mathematical optimization problems are called equivalent if the solution of one is readily obtained from the solution of the other, and vice-versa.*

For instance, the minimization problem (1.2.1) has an equivalent reformulation as a maximization problem when $f : \mathbb{R}^d \rightarrow \mathbb{R}$ (i.e. f outputs numerical values):

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

1.2.1 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 normed 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. It also implies that we can measure the distance between two elements of these spaces using norms, the most classical of which is defined below.

Definition 1.2.2 (Euclidean norm on \mathbb{R}^d) *The Euclidean norm (or ℓ_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 1.2.3 (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 $\mathbf{w}^T \mathbf{w} = \|\mathbf{w}\|^2$.

There are natural counterparts to the Euclidean norm and its associated scalar product in a matrix space.

Definition 1.2.4 (Frobenius norm) For any matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, the Frobenius norm is defined by

$$\|\mathbf{A}\|_F := \sqrt{\sum_{i=1}^n \sum_{j=1}^d A_{ij}^2} = \sqrt{\text{trace}(\mathbf{A}^T \mathbf{A})}.$$

The associated scalar product is the mapping

$$(\mathbf{A}, \mathbf{B}) \in (\mathbb{R}^{n \times d})^2 \mapsto \text{trace}(\mathbf{A}^T \mathbf{B}).$$

In this course, we will mainly rely on Euclidean norms to measure distance between vectors, and algorithmic behavior. However, other norms will be used.

Definition 1.2.5 (ℓ_1 norm) The ℓ_1 norm of a vector $\mathbf{w} \in \mathbb{R}^d$ is defined by:

$$\|\mathbf{w}\|_1 := \sum_{i=1}^d |w_i|.$$

Both the ℓ_1 and ℓ_2 norms are special cases of the ℓ_p norm, defined for any $p \in [1, \infty)$ by

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \|\mathbf{w}\|_p := \left[\sum_{i=1}^d |w_i|^p \right]^{1/p}.$$

Definition 1.2.6 (ℓ_∞ norm) The ℓ_∞ norm of a vector $\mathbf{w} \in \mathbb{R}^d$ is defined by:

$$\|\mathbf{w}\|_\infty := \max_{1 \leq i \leq d} |w_i|.$$

Remark 1.2.1 Interestingly, all norms can be defined as optimal values of optimization problems. This is clear from the definition of the ℓ_∞ norm. Below are two optimization problems that yield this optimal value:

$$\begin{aligned} & \underset{i \in \mathbb{R}}{\text{maximize}} |w_i| \quad \text{s.t.} \quad i \in \{1, \dots, d\}. \\ & \underset{v \in \mathbb{R}}{\text{maximize}} |v| \quad \text{s.t.} \quad v \in \{w_1, \dots, w_d\}. \end{aligned}$$

Notice that these two problems do not have the same decision variables or constraints, however they both possess a combinatorial structure. ¹

¹By convention, in these notes, we will always consider optimization problems with real, continuous decision variables. As illustrated by the ℓ_∞ examples, it is indeed possible to model a discrete decision-making problem using continuous formulations. Our focus in these notes is on continuous optimization problems.

Definition 1.2.7 (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 1.2.8 (Positive (semi-)definiteness) A symmetric 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 1.2.9 (Eigenvalues and eigenvectors) Let $\mathbf{A} \in \mathbb{R}^{d \times d}$. A real λ 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 λ .

Theorem 1.2.1 Any symmetric matrix in $\mathbb{R}^{d \times d}$ possesses d real eigenvalues.

Notation 1.2.1 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 positive semi-definite (resp. positive definite) if and only if $\mathbf{A} \succeq 0$ (resp. $\mathbf{A} \succ 0$).

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.

Definition 1.2.10 (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 1.2.11 (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.

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 1.2.12 (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 $f \in \mathcal{C}^1(\mathbb{R}^d)$, its 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 1.2.13 (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 1.2.14 (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 1.2.15 (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 ∇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 1.2.2 (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:

$$f(\mathbf{z}) \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^T(\mathbf{z} - \mathbf{w}) + \frac{L}{2} \|\mathbf{z} - \mathbf{w}\|^2. \quad (1.2.2)$$

1.2.2 First definitions

Having defined an optimization problem as a mathematical object, we are now able to consider solving this problem. For simplicity, we will focus on minimization problems of the form:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) \quad \text{s.t. } \mathbf{w} \in \mathcal{F}, \quad (1.2.3)$$

with $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and $\mathcal{F} \subseteq \mathbb{R}^d$. Similar definitions hold for maximization problems.

Definition 1.2.16 (Feasible point) A point $\mathbf{w} \in \mathbb{R}^d$ is called a *feasible point* of the optimization problem (1.2.3) if $\mathbf{w} \in \mathcal{F}$.

If $\mathbf{w} \notin \mathcal{F}$, we say that \mathbf{w} is infeasible. If \mathcal{F} is empty, the problem (1.2.3) is said to be infeasible. For some optimization problems, the objective function may not be defined at infeasible points.

A solution of problem (1.2.3) must be a feasible point by definition. It should also lead to the best value in terms of cost function, hence the following definition.

Definition 1.2.17 (Global minimum) A point $\mathbf{w}^* \in \mathbb{R}^d$ is called a *solution* or a *global minimum* of problem (1.2.3) if

1. \mathbf{w}^* is feasible, i.e. $\mathbf{w}^* \in \mathcal{F}$;
2. $f(\mathbf{w}^*) \leq f(\mathbf{w})$ for any feasible point $\mathbf{w} \in \mathcal{F}$.

The set of global minima of problem (1.2.3) will be denoted by

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\operatorname{argmin}} \{f(\mathbf{w}) \mid \mathbf{w} \in \mathcal{F}\} \subset \mathbb{R}^d, \quad (1.2.4)$$

and the minimal value of the problem will be defined as

$$\min_{\mathbf{w} \in \mathbb{R}^d} \{f(\mathbf{w}) \mid \mathbf{w} \in \mathcal{F}\} \in \mathbb{R} \cup \{-\infty, +\infty\}. \quad (1.2.5)$$

By convention, if f is unbounded below on $\mathcal{F} \neq \emptyset$, the minimal value is set to $-\infty$, while it is set to $+\infty$ if the problem is infeasible.

Definition 1.2.18 (Local minimum) A point $\mathbf{w}^* \in \mathbb{R}^d$ is called a *local minimum* of problem (1.2.3) if

1. \mathbf{w}^* is feasible, i.e. $\mathbf{w}^* \in \mathcal{F}$;
2. There exists $\epsilon > 0$ such that $f(\mathbf{w}^*) \leq f(\mathbf{w})$ for any feasible point $\mathbf{w} \in \mathcal{F}$ close to \mathbf{w}^* in the sense that $\|\mathbf{w} - \mathbf{w}^*\| \leq \epsilon$.

The notion of local minimum is weaker than that of global minimum. In practice, however, it is often more reasonable to seek local minima, as those can be characterized by mathematical conditions.

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.

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}), \quad (1.2.6)$$

Theorem 1.2.3 (First-order necessary condition) Suppose that the objective function f in problem (1.2.6) belongs to $\mathcal{C}^1(\mathbb{R}^d)$. Then,

$$[\mathbf{w}^* \text{ is a local minimum of } f] \implies \|\nabla f(\mathbf{w}^*)\| = 0. \quad (1.2.7)$$

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 points².

Provided we strengthen our smoothness requirements on f , we can establish stronger optimality conditions for problem (1.2.6).

Theorem 1.2.4 (Second-order necessary condition) Suppose that the objective function f in problem (1.2.6) belongs to $\mathcal{C}^2(\mathbb{R}^d)$. Then,

$$[\mathbf{w}^* \text{ is a local minimum of } f] \implies [\|\nabla f(\mathbf{w}^*)\| = 0 \text{ and } \nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}]. \quad (1.2.8)$$

From Theorem 1.2.3, first-order stationary points that violate the condition $\nabla^2 f(\mathbf{w}^*) \succeq \mathbf{0}$ cannot be local minima. This new condition is thus more precise than the first-order condition, although it remains a necessary. For an arbitrary function, there may exist points with zero gradient and positive semidefinite Hessian (termed **second-order stationary points**) that are not local minima.

On the other hand, and unlike first-order conditions, it is possible to derive a *sufficient* version of the second-order optimality conditions, that can be used to certify local optimality.

Theorem 1.2.5 (Second-order sufficient condition) Suppose that the objective function f in problem (1.2.6) belongs to $\mathcal{C}^2(\mathbb{R}^d)$. Then,

$$[\|\nabla f(\mathbf{w}^*)\| = 0 \text{ and } \nabla^2 f(\mathbf{w}^*) \succ \mathbf{0}] \implies [\mathbf{w}^* \text{ is a local minimum of } f] \quad (1.2.9)$$

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 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.

1.2.3 Convexity

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

²A 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.

Definition 1.2.19 (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 1.2.1 (Examples of convex sets) The following sets are convex:

- \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 $\{\mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|_2^2 = \sum_{i=1}^d [\mathbf{w}]_i^2 \leq 1\}$.

We now provide the basic definition of a convex function.

Definition 1.2.20 (Convex function) A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **convex** if

$$\forall(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \quad f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq t f(\mathbf{u}) + (1 - t) f(\mathbf{v}). \quad (1.2.10)$$

Example 1.2.2 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: $\mathbf{w} \mapsto \|\mathbf{w}\|_2^2 = \mathbf{w}^T \mathbf{w}$.

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

Theorem 1.2.6 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

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}). \quad (1.2.11)$$

Theorem 1.2.7 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

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \nabla^2 f(\mathbf{w}) \succeq \mathbf{0}. \quad (1.2.12)$$

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

Theorem 1.2.8 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 1.2.1 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 1.2.21 (Strongly convex function) A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ in \mathcal{C}^1 is μ -strongly convex (or strongly convex of modulus $\mu > 0$) if for all $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ and $t \in [0, 1]$,

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

Theorem 1.2.9 Any strongly convex function has a unique global minimizer.

As for convex functions, there exist characterizations of strong convexity that involve derivatives.

Theorem 1.2.10 Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be an element of $\mathcal{C}^1(\mathbb{R}^d)$. Then, the function f is μ -strongly convex if and only if

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}) + \frac{\mu}{2}\|\mathbf{v} - \mathbf{u}\|^2. \quad (1.2.13)$$

Theorem 1.2.11 Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be an element of $\mathcal{C}^2(\mathbb{R}^d)$. Then, the function f is μ -strongly convex if and only if

$$\forall \mathbf{w} \in \mathbb{R}^d, \quad \nabla^2 f(\mathbf{w}) \succeq \mu \mathbf{I}. \quad (1.2.14)$$

1.3 Examples of optimization problems in ML

1.3.1 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

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{2} (\mathbf{x}_i^T \mathbf{w} - y_i)^2 + \lambda \|\mathbf{w}\|^2 \right], \quad (1.3.1)$$

where $\lambda > 0$ is a regularization parameter. From an optimizer's point of view, problem (1.3.1) 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 1.1.

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

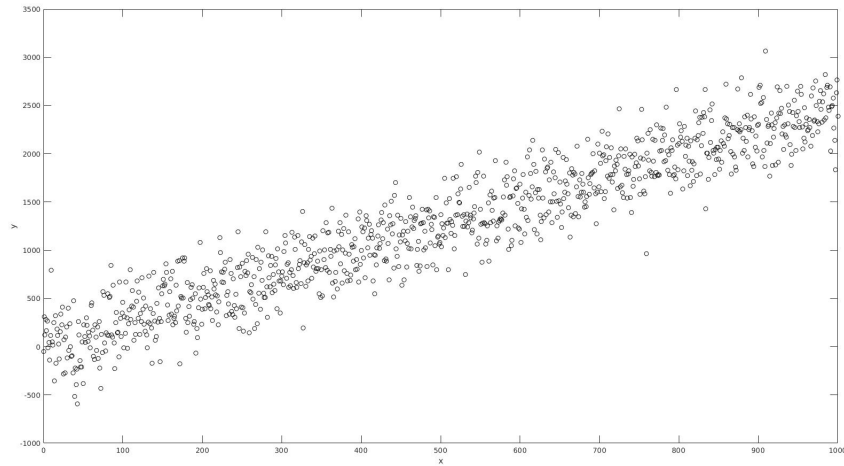


Figure 1.1: Noisy data generated from a linear model with Gaussian noise.

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

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{maximize}} L(y_1, \dots, y_n; \mathbf{w}) := \left[\frac{1}{\sqrt{2\pi}} \right]^m \exp \left(-\frac{1}{2} \sum_{i=1}^m (\mathbf{x}_i^T \mathbf{w} - y_i)^2 - \frac{\lambda}{2} \|\mathbf{w}\|^2 \right). \quad (1.3.2)$$

The solutions of this maximization problem are the same than the solutions of the linear least-squares problem (1.3.1). The resulting solution can be shown to possess very favorable statistical properties: in particular, for λ 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.

1.3.2 Logistic regression

As in Section 1.3.1, 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

$$\forall \mathbf{w} \in \mathbb{R}^d, f(\mathbf{w}) = \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\}. \quad (1.3.3)$$

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

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \left\{ \sum_{y_i=-1} \ln(1 + e^{-\mathbf{x}_i^T \mathbf{w}}) + \sum_{y_i=+1} \ln(1 + e^{\mathbf{x}_i^T \mathbf{w}}) \right\} \quad (1.3.4)$$

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 ??).

1.3.3 Linear SVM

To illustrate the role of optimization in data-related applications, we consider a binary classification problem, illustrated in Figure 1.2.

The red circles and blue squares appear at the same locations on all three figures : they represent data samples identified by their coordinates, while their color or shape represents a certain class they belong to. Our goal is to compute a linear classifier, that is, a separator of the two classes corresponding to a linear function. Each of the three figures shows a separator that achieves the task of classifying the data (the separator is the same for the middle and right plots): in that sense, the task involving the *samples* has been performed. However, if we envision the samples as being part of a (much) larger dataset, represented by the blue and red blobs, it becomes clear that the best classifier is the one on the rightmost figure.

These observations can be modeled and summarized using a mathematical framework. Let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be n vectors of \mathbb{R}^d , and suppose that each vector \mathbf{x}_i is given a label $y(\mathbf{x}_i) \in \{-1, 1\}$ (say, -1 for a red point and 1 for a blue point). Then, one can translate the binary classification problem into the following optimization problem:

$$\underset{\substack{\mathbf{u} \in \mathbb{R}^d \\ v \in \mathbb{R}}}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \max \{1 - y(\mathbf{x}_i)(\mathbf{x}_i^T \mathbf{u} - v), 0\} \quad (1.3.5)$$

Here we seek a linear function defined by $\mathbf{x} \mapsto \mathbf{x}^T \mathbf{u} - v$: it thus suffices to compute its coefficients, given by the vector \mathbf{u} and the scalar v (see Section ?? for a formal definition of these concepts).

1.3.4 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.

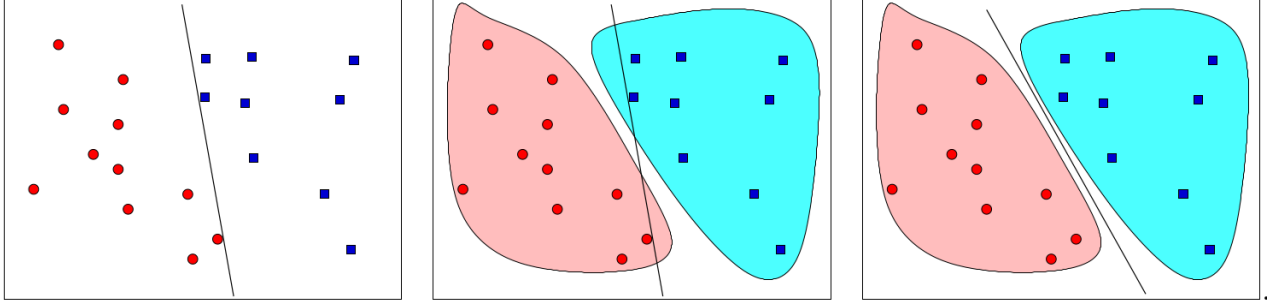


Figure 1.2: A sample for binary classification (red circles/blue squares) and the associated distribution (light red set/light blue set). A same linear classifier is shown on the left and middle plot: although it classifies the samples correctly, its closeness to several data points make it sensitive to the data, and prevents it from correctly classifying the distribution. On the contrary, the linear classifier on the right plot (that has a maximal margin of separation) provides a better classification, and is able to generalize to the distributions. *Source : S. J. Wright and B. Recht, Optimization for Data Analysis [3].*

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:

$$\mathbf{x}_i^{(j)} = \sigma(\mathbf{W}_j \mathbf{x}_i^{(j-1)} + \mathbf{b}_j) \in \mathbb{R}^{d_j}, \quad (1.3.6)$$

where $\mathbf{W}_j \in \mathbb{R}^{d_j \times d_{j-1}}$, $\mathbf{b}_j \in \mathbb{R}^{d_j}$ and $\sigma : \mathbb{R}^{d_j} \rightarrow \mathbb{R}^{d_j}$ is a componentwise nonlinear function, e.g. $\sigma(\mathbf{y}) = \left[\frac{1}{1 + \exp(-y_i)} \right]_i$ (sigmoid function) or $\sigma(\mathbf{y}) = [\max(0, y_i)]_i$. 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 ℓ . It usually results in the following formulation

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i; \mathbf{w}), y_i). \quad (1.3.7)$$

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 (1.3.7) 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.

1.4 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.

1.4.1 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 $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$, resulting in a sequence of iterates $\{\mathbf{w}_k\}_k$. Ideally, these iterates obey one of the scenarios below:

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

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

Although \mathbf{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.

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

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

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

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.

1.4.2 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 1.4.1 (Global convergence rate for the gradient norm) *Given an algorithm applied to $\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ that produces a sequence of iterates $\{\mathbf{w}_k\}$, we say that the method is $\mathcal{O}(1/k)$ for the gradient norm, or $\|\nabla f(\mathbf{w}_k)\| = \mathcal{O}(\frac{1}{k})$ if*

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

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

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

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.

1.4.3 Popular optimization packages

Although a thorough numerical study is out of the scope of this course, we briefly mention popular choices for implementing optimization methods, either for industrial use or as research prototypes.

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/Octave is also widespread throughout the optimization community for prototyping (i.e. rapid and simple validation of an implementation), along with Python and Julia,

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. Some of those solvers are implemented in standard tools such as Microsoft Excel. Open-source codes are also quite popular, again fueled by the massive production of such implementations in the learning community. As far as optimization is concerned, the COIN-OR platform provides a good interface to all of these methods.

Chapter 2

Unconstrained optimization

In this chapter, we study more general **nonlinear optimization problems** of the form

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}). \quad (2.0.1)$$

We make the following assumption on the objective function.

Assumption 2.0.1 *The objective function f in (2.0.1) is $\mathcal{C}_L^{1,1}(\mathbb{R}^d)$ for $L > 0$, and bounded below by $f_{low} \in \mathbb{R}$ (i.e. $f(\mathbf{w}) \geq f_{low} \forall \mathbf{w} \in \mathbb{R}^d$).*

We will design and analyze algorithms that exploit gradient information to move towards better points. Our theoretical results will consist in complexity bounds and convergence rates.

2.1 Gradient descent

The gradient descent algorithm is arguably the most classical technique in unconstrained, smooth optimization. It is based on the following principle, derived from the first-order optimality condition (1.2.7).

For any vector $\mathbf{w} \in \mathbb{R}^d$, two cases can occur:

1. Either $\nabla f(\mathbf{w}) = 0$ and \mathbf{w} is possibly a local minimum (when f is convex, we know that \mathbf{w} is necessarily a global minimum);
2. Or $\nabla f(\mathbf{w}) \neq 0$, and we can show that f must decrease *locally* in the direction of $-\nabla f(\mathbf{w})$.

The second property is formalized below, and is at the core of the gradient descent framework.

2.1.1 Algorithm

The gradient descent algorithm is an iterative process wherein every iteration has the following form:

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla f(\mathbf{w}), \quad (2.1.1)$$

where $\alpha > 0$ is a parameter called **stepsize** or **steplength**. When $\nabla f(\mathbf{w}) = 0$, note that the formula (2.1.1) does not change the value of \mathbf{w} : this is consistent with the fact that the gradient

cannot be used to determine a better point in that case. On the contrary, when $\nabla f(\mathbf{w}) \neq 0$, there will exist values for α leading to a lower function value for the point $\mathbf{w} - \alpha \nabla f(\mathbf{w})$.

Repeated applications of the update (2.1.1) lead to Algorithm 1. This method is called **gradient descent**, or (less frequently) *steepest descent*.¹

Algorithm 1: Gradient descent for minimizing the function f .

```

1 Initialization: Choose  $\mathbf{w}_0 \in \mathbb{R}^d$ .
2 For  $k = 0, 1, \dots$ 
    1. Compute the gradient  $\nabla f(\mathbf{w}_k)$ .
    2. Define a stepsize  $\alpha_k > 0$ .
    3. Set  $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$ .
3 EndFor

```

Algorithm 1 actually describes a framework rather than a specific method. There exist numerous variants upon the gradient descent paradigm: we review the main characteristics of these methods below.

Stopping criterion In practice, an algorithm is often subject to budget requirements (in terms of arithmetic operations, running time, iteration number, etc). These limitations are typically enforced using a stopping criterion. In the context of Algorithm 1, it is typical to terminate the method after k_{\max} iterations, in which case $\mathbf{w}_{k_{\max}}$ is returned as an approximate solution to the problem.

In addition, stopping criteria can be used to check whether the method converged to a solution of the problem (or an approximation thereof). Such criteria are often based on optimality conditions. For gradient descent, the most classical criterion is based on the gradient norm: one stops the algorithm as soon as

$$\|\nabla f(\mathbf{w}_k)\| < \epsilon, \quad (2.1.2)$$

where $\epsilon > 0$ represents a given, desired accuracy level (the condition becomes more expensive to satisfy as ϵ gets smaller).

Finally, it is always possible (and often recommended) to add “safety checks”, that stop the method whenever no visible progress is measured. For instance, if the difference between two successive iterates ($\|\mathbf{w}_{k+1} - \mathbf{w}_k\|$) is of the order of machine precision, the algorithm is likely stalling and no further improvement is expected: in that case, it is often better to stop the method.

Choosing the initial point The performance of a given algorithm can be significantly improved using a suitable starting point. However, finding such a point can be a difficult task without domain expertise, but if such expertise exists, it should be leveraged to obtain a good initial point that the method tries to improve upon. Alternatively, one could use several starting points drawn at random and run a few iterations of gradient descent so as to determine a good starting point.

¹The direction of $-\nabla f(\mathbf{w}_k)$, i.e. the unit vector $\frac{-\nabla f(\mathbf{w}_k)}{\|\nabla f(\mathbf{w}_k)\|}$ (or the zero vector if the gradient is zero) is called the steepest descent direction.

2.1.2 Choosing the stepsize

In this section, we describe the main techniques for selecting the step size sequence in gradient descent. We provide generic principles, and emphasize that any information about a particular problem can be extremely valuable in designing a better step size.

Constant step size One of the most common approaches consists in fixing the step size to a single value for all iterations, i.e. setting $\alpha_k = \alpha > 0$ for all k . Depending on the computational budget, one could run gradient descent with several values of α and select the best value for future use: this practice of *tuning* the step size is commonly adopted in data science problems.

Provided f satisfies Assumption 2.0.1, there exists an interval of values that lead to convergence of gradient descent. In particular, the choice

$$\alpha_k = \alpha = \frac{1}{L}, \quad (2.1.3)$$

where L is the Lipschitz constant for the gradient, is well suited for that problem. Note however that this choice requires the knowledge of the Lipschitz constant: this information is not always available in practice.

Decreasing step size Another classical technique for selecting the stepsize consists in defining a decreasing sequence $\{\alpha_k\}$ such that $\alpha_k \rightarrow 0$ *prior to running the method*. This choice can also lead to converging method, but it risks producing steps that are unnecessarily small in norm. In fact, a good decreasing strategy should drive α_k to 0 quickly enough for convergence, but slowly enough that the norm of the steps do not approach 0 too rapidly.

Adaptive choice using a line search Line-search techniques are widely used in continuous optimization and scientific computing (though less popular in data science, for reasons that we will detail in Chapter 3 of these notes). At a given iteration of index k , we seek a stepsize α_k that leads to a decrease in the objective function along a suitably chosen direction (in the case of Algorithm 1, this would be the direction of $-\nabla f(\mathbf{w}_k)$). An exact line search results in the best possible decrease, but may be costly to perform. In practice, *inexact* approaches are preferred: Algorithm 2 details the most popular of such techniques, called **backtracking**.

Algorithm 2: Backtracking line search in direction d .

- 1 **Inputs:** $\mathbf{w} \in \mathbb{R}^d$, $d \in \mathbb{R}^d$, $\alpha_0 \in \mathbb{R}^d$.
 - 2 **Initialization:** Choose $\alpha = \alpha_0$.
 - 3 **While** $f(\mathbf{w} + \alpha d) > f(\mathbf{w})$
 - 4 $\alpha \rightarrow \frac{\alpha}{2}$.
 - 5 **End**
 - 6 **Output:** α .
-

The line-search procedure of Algorithm 2 can be used at Step 2 of Algorithm 1 using $\mathbf{w} = \mathbf{w}_k$, $d = -\nabla f(\mathbf{w}_k)$ and (for instance) $\alpha_0 = 1$ as inputs. Many variants of this simple idea have been proposed in the literature: those are generally designed to guarantee that a better point (in terms of the objective function) is found. Still, these techniques require at least one additional function evaluation per stepsize value (possibly more), leading to an overall more expensive method.

2.1.3 Theoretical analysis for 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 $f \in C_L^{1,1}(\mathbb{R}^d)$, 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

$$f(\mathbf{w}_k - \frac{1}{L} \nabla f(\mathbf{w}_k)) < f(\mathbf{w}_k) - \frac{1}{2L} \|\nabla f(\mathbf{w}_k)\|^2. \quad (2.1.4)$$

Proof. We use the inequality (1.2.2) with the vectors $(\mathbf{w}_k, \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k))$:

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

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 (2.1.4), 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, and strongly convex.

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 2.0.1. Suppose that Algorithm 1 is applied with $\alpha_k = \frac{1}{L}$. Then, for any $K \geq 1$, we have*

$$\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\| \leq \mathcal{O}\left(\frac{1}{\sqrt{K}}\right). \quad (2.1.5)$$

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 2.0.1) 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 2.0.1, 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 2.0.1. Suppose that Algorithm 1 is applied with $\alpha_k = \frac{1}{L}$. Then, for any $K \geq 1$, the iterate \mathbf{w}_K satisfies

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K}\right). \quad (2.1.6)$$

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 (1.2.11) at \mathbf{w}_k and \mathbf{w}^* gives

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

Combining this property with (2.1.4), we obtain:

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

To proceed onto the next step, one notices that

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

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

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

Hence,

$$f(\mathbf{w}_{k+1}) - f(\mathbf{w}^*) \leq \frac{L}{2} (\|\mathbf{w}_k - \mathbf{w}^*\|^2 - \|\mathbf{w}_{k+1} - \mathbf{w}^*\|^2). \quad (2.1.7)$$

By summing (2.1.7) 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 $C\epsilon^{-1} = \mathcal{O}(\epsilon^{-1})$ 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 μ -strongly convex function satisfying Assumption 2.0.1, with $\mu \in (0, L]$. Suppose that Algorithm 1 is applied with $\alpha_k = \frac{1}{L}$. Then, for any $K \in \mathbb{N}$, we have

$$f(\mathbf{w}_K) - f^* \leq \mathcal{O} \left(\left(1 - \frac{\mu}{L}\right)^K \right). \quad (2.1.8)$$

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 (1.2.13). 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 side². 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

$$\|\nabla f(\mathbf{x})\|^2 \geq 2\mu [f(\mathbf{x}) - f^*], \quad (2.1.9)$$

which is valid for any $\mathbf{x} \in \mathbb{R}^n$. Using (2.1.9) together with (2.1.4) 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 $\mathcal{O}(\frac{L}{\mu} \ln(\frac{1}{\epsilon}))$ 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 $\mathcal{O}((1 - \frac{\mu}{L})^k)$.

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.*

2.2 Acceleration

2.2.1 Introduction: the momentum principle

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

²This is because the gradient of the quadratic function on the right-hand side with respect to \mathbf{y} is $\nabla f(\mathbf{x}) + \mu(\mathbf{y} - \mathbf{x})$. This vector is equal to $\mathbf{0}$ if and only if $\mathbf{y} = \mathbf{x} - \frac{1}{\mu} \nabla f(\mathbf{x})$

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.

1 **Initialization:** $\mathbf{w}_0 \in \mathbb{R}^d$, $\mathbf{w}_{-1} = \mathbf{w}_0$.

2 **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

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k + \beta_k(\mathbf{w}_k - \mathbf{w}_{k-1})) + \beta_k(\mathbf{w}_k - \mathbf{w}_{k-1}). \quad (2.2.1)$$

3 **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 \mathbf{w}_k , but at a combination of this iterate with the previous step $\mathbf{w}_k - \mathbf{w}_{k-1}$: 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 \mathbf{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 α_k and the momentum parameter β_k . The same techniques than those presented in Section 2.1.2 can be considered for the choice of α_k (stepsize parameter). As in the gradient descent case, the choice $\alpha_k = \frac{1}{L}$ is a standard one.

The choice of β_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 μ -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 , β_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 2.0.1, 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*

- i) $f(\mathbf{w}_K) - f^* \leq \mathcal{O}(\frac{1}{K^2})$ for a generic convex function if β_k is set according to the adaptive rule (2.2.4);
- ii) At most $f(\mathbf{w}_K) - f^* \leq \left((1 - \sqrt{\frac{\mu}{L}})^K\right)$ for a μ -strongly convex function, provided β_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 $\mathcal{O}\left(\sqrt{\frac{L}{\mu}} \ln(\epsilon^{-1})\right)$ iterations, which represents an improvement over the $\mathcal{O}\left(\frac{L}{\mu} \ln(\epsilon^{-1})\right)$ complexity over gradient descent. Here the improvement is in terms of problem-dependent constants.

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 μ , because it exploits knowledge from the past iterations. The k -th iteration of conjugate gradient can be written as:

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

In a standard conjugate gradient algorithm, α_k and β_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 μ -strongly convex with μ 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 2.1.3 is a valid one.

If gradient descent is applied, then an ϵ -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 $\mathcal{O}\left(\min\{d, \frac{L}{\mu} \ln\left(\frac{1}{\epsilon}\right)\}\right)$.

Chapter 3

Stochastic gradient techniques

3.1 Introduction

In this chapter, we will develop method that fully exploit the structure of data science problems. To this end, we assume that we are given a dataset of n examples under the form $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^{d_x}$ and $\mathbf{y}_i \in \mathbb{R}^{d_y}$ are obtained from a certain data distribution. As in the linear regression example from the previous chapter, we seek a model function $h : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}$ such that $h(\mathbf{x}_i) \approx \mathbf{y}_i$ for any $i = 1, \dots, n$. We will parameterize this model by a vector $\mathbf{w} \in \mathbb{R}^d$ ($h(\mathbf{x}_i) = h(\mathbf{x}_i; \mathbf{w})$) so that knowing the vector \mathbf{w} suffices to be able to evaluate the model.

In order to quantify the model's ability to represent the data, we define a cost function (or loss function) of the form $\ell : (h, y) \mapsto \ell(h, y)$: this function aims at penalizing values (h, y) such that $h \neq y$. The mean-square or ℓ_2 loss defined by $(h, y) \mapsto \|h - y\|^2$ is a canonical example of a loss. We then define the loss at a given example by $\ell(h(\mathbf{w}; \mathbf{x}_i), \mathbf{y}_i)$.

We wish to compute the best model according to our entire dataset, leading to average the losses over all examples. This finally leads to the following optimization problem.

Definition 3.1.1 (Finite sum problem) Consider a dataset $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^{d_x}$ and $\mathbf{y}_i \in \mathbb{R}^{d_y}$, a model class $\{h(\mathbf{w}; \cdot) : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_y}\}_{\mathbf{w} \in \mathbb{R}^d}$ and a cost function ℓ . The **finite-sum problem** associated with this setup is defined by:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{w}; \mathbf{x}_i), \mathbf{y}_i) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}). \quad (3.1.1)$$

Provided the f_i s are all \mathcal{C}^1 functions, the function f is also \mathcal{C}^1 and we can apply gradient descent to problem (3.1.1). 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 the iteration, it is clear that every gradient descent iteration requires to access **the entire dataset** for a single gradient calculation. In a big data setting, the number of examples n can be extremely large, and gradient descent can be too expensive to be used in practice.

Remark 3.1.1 In a purely stochastic or “online”, examples may only be available in a streaming fashion: in this context, it may be impossible to compute a finite average over all samples, and we

consider instead the stochastic optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimize}} \mathbb{E}_{(\mathbf{x}, \mathbf{y})} [f_{(\mathbf{x}, \mathbf{y})}(\mathbf{w})] . \quad (3.1.2)$$

Even though the gradient of the objective may exist, it can be quite difficult to compute. This formulation is often closer to the real goal of machine learning (i.e. have a good model for all possible examples in the distribution), yet it is often replaced by (3.1.1) to account for the empirical setting.

3.2 Stochastic gradient method

3.2.1 Algorithm

The idea behind stochastic gradient is remarkably simple. It considers the problem (3.1.1) under the assumption that every component function f_i possesses a gradient: every iteration then consists in choosing a random index i_k and taking a step in the direction of the negative gradient of the function f_{i_k} . Algorithm 4 details this process.

Algorithm 4: Stochastic gradient method

```

1 Initialization:  $\mathbf{w}_0 \in \mathbb{R}^d$ .
2 for  $k = 0, 1, \dots$  do
    1. Define a stepsize  $\alpha_k > 0$ .
    2. Draw an index  $i_k \in \{1, \dots, n\}$ .
    3. Compute the new iterate
        
$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k). \quad (3.2.1)$$

3 end
```

The vector $\nabla f_{i_k}(\mathbf{w}_k)$ is called a **stochastic gradient** for \mathbf{w}_k . The key property of iteration (3.2.1) is that it only requires one example from the dataset: as a result, **its cost in terms of accessing the data is n times lower compared to an iteration of gradient descent**. As a result, the comparison between the two methods is often conducted in terms of **epochs**.

Definition 3.2.1 (Epoch) For problem (3.1.1), an epoch represents n calculations of a sample gradient ∇f_i .

As a result, the cost of one iteration of gradient descent is that of one epoch, an iteration of Algorithm 4 only costs a fraction ($\frac{1}{n}$) of an epoch.

To conclude this presentation of stochastic gradient, we make several key remarks on its behavior.

Remark 3.2.1 In general, it is not possible to guarantee convergence of gradient descent. Consider the two-dimensional problem:

$$\underset{\mathbf{w} \in \mathbb{R}}{\text{minimize}} \frac{1}{2}(f_1(w) + f_2(w))$$

with $f_1(w) = 2w^2$ and $f_2 = -w^2$. If $w_0 > 0$ and $i_k = 2$ for any k , then the method will diverge.

It is easy to build problems and draws of indices for which stochastic gradient does not converge. In practice, however, finite-sum problems from data science involve component functions that are quite similar, corresponding to data samples that follow a similar distribution. As a result, improving the loss specific to a given sample may very well improve other losses with respect to similar samples. This observation partly explains the success of stochastic gradient methods in data-related settings.

Remark 3.2.2 Algorithm 4 is often called **Stochastic Gradient Descent**, or SGD. This denomination is inaccurate, in that one cannot guarantee that the stochastic gradient method will behave like a descent method (i.e. that it will decrease the function value at every iteration). For this reason, we will call this method **stochastic gradient** in these notes, thereby following other authors' terminology [1]. However, we point out that many librairies refer to their implementation of Algorithm 4 as SGD.

3.2.2 Convergence rate analysis

In this section, we describe the key steps to deriving convergence rates (and complexity bounds) for stochastic gradient, under a slightly altered version of Assumption 2.0.1.

Assumption 3.2.1 The objective function $f = \frac{1}{n} \sum_{i=1}^n f_i$ of problem (3.1.1) is $C_L^{1,1}(\mathbb{R}^d)$ for some $L > 0$, and bounded below by $f_{low} \in \mathbb{R}$. Moreover, every function f_i is of class C^1 .

The key argument for analyzing gradient descent is the result of Proposition 2.1.1: in particular, the inequality

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

leads to decrease guarantees for a sufficiently small step size.

It is possible to obtain a similar inequality in the case of stochastic gradient (in a probabilistic sense). To this end, assumptions on the random indices (and thus the stochastic gradients used throughout the algorithm) are necessary.

Assumption 3.2.2 (Stochastic gradients' properties) For every iteration k of Algorithm 4, the random index i_k is drawn so that

- i) i_k is independent of i_0, \dots, i_{k-1} ;
- ii) $\mathbb{E}_{i_k} [\nabla f_{i_k}(\mathbf{w}_k)] = \nabla f(\mathbf{w}_k)$;
- iii) $\mathbb{E}_{i_k} [\|\nabla f_{i_k}(\mathbf{w}_k)\|^2] \leq \|\nabla f(\mathbf{w}_k)\|^2 + \sigma^2$ with $\sigma^2 \in (0, \infty)$.

The second property of Assumption 3.2.2 implies that the stochastic gradient $\nabla f_{i_k}(\mathbf{w}_k)$ is an unbiased estimator of the true gradient $\nabla f(\mathbf{w}_k)$. The third property guarantees that the norm of this stochastic estimator cannot deviate too much from the true norm, which is critical in guaranteeing some form of convergence. The most classical strategy to satisfy these assumptions is given below.

Example 3.2.1 (Uniform sampling) Suppose that $\{i_k\}_k$ is a sequence of independent, identically distributed indices uniformly drawn in $\{1, \dots, n\}$. Then, Algorithm 4 satisfies Assumption 3.2.2.

Using Assumption 3.2.2, one can establish a useful inequality for analyzing Algorithm 4.

Proposition 3.2.1 *Let Assumptions 2.0.1 and 3.2.2 hold, and consider the k -th iteration of Algorithm 4. Then,*

$$\mathbb{E}_{i_k} [f(\mathbf{w}_{k+1})] - f(\mathbf{w}_k) \leq \nabla f(\mathbf{w}_k)^\top \mathbb{E}_{i_k} [\mathbf{w}_{k+1} - \mathbf{w}_k] + \frac{L}{2} \mathbb{E}_{i_k} [\|\mathbf{w}_{k+1} - \mathbf{w}_k\|^2]. \quad (3.2.2)$$

As shown by (3.2.2), one can guarantee decrease over an iteration of stochastic gradient **in expectation** (provided the right-hand side of (3.2.2) is negative). This property is sufficient to obtain convergence rates (or complexity bounds) for stochastic gradient applied to strongly convex, convex and nonconvex problems. These results heavily depend on the choice for the stepsize sequence $\{\alpha_k\}_k$: in fact, tuning the stepsize is a major issue in machine learning (in which the stepsize is often referred to as a **learning rate**).

In the rest of this section, we mainly focus on strongly convex problems, and make the following assumption.

Assumption 3.2.3 *The objective function f of problem (3.1.1) is continuous and μ -strongly convex. It possesses a unique global minimum \mathbf{w}^* , and we let $f^* = f(\mathbf{w}^*)$.*

We first consider the case of a constant stepsize, for which we can establish the following result.

Theorem 3.2.1 (Stochastic gradient with constant stepsize) *Let Assumptions 2.0.1, 3.2.2 and 3.2.3 hold. Suppose that we apply Algorithm 4 with a constant stepsize*

$$\alpha_k = \alpha \in (0, \frac{1}{L}] \quad \forall k.$$

Then, for every $K \in \mathbb{N}$,

$$\mathbb{E} [f(\mathbf{w}_K) - f^*] \leq \frac{\alpha L \sigma^2}{2\mu} + (1 - \alpha\mu)^K \left[f(\mathbf{w}_0) - f^* - \frac{\alpha L \sigma^2}{2\mu} \right], \quad (3.2.3)$$

where the expected value is taken on all randomness up to iteration K .

The result of Theorem 3.2.1 shows a linear rate of convergence for stochastic gradient in expectation, which is comparable to that of gradient descent. However, this rate only guarantees that $\{f(\mathbf{w}_K)\}$ converges to a value in $[f^*, f^* + \frac{\alpha L \sigma^2}{2\mu}]$: even on average, stochastic gradient can only be guaranteed to converge towards a neighborhood of the optimal value. Indeed, the result (3.2.3) involves a constant term on its right-hand side, showing a possible gap to the optimum in $\frac{\alpha L \sigma^2}{4\mu}$. It is therefore not possible to guarantee that $\mathbb{E} [f(\mathbf{w}_k) - f^*] \leq \epsilon$ for any $\epsilon > 0$, but only for sufficiently large values of ϵ . However, the use of a constant stepsize guarantees a linear rate of convergence.

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

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

which lead to asymptotic convergence of the method. In order to satisfy these requirements, the sequence α_k must go to 0 as k goes to infinity. Therefore, we now describe a convergence result based on using a decreasing stepsize.

Theorem 3.2.2 (Stochastic gradient with decreasing stepsize) *Under Assumptions 2.0.1, 3.2.2 and 3.2.3, consider Algorithm 4 applied with a decreasing stepsize sequence of the form*

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

with $\beta > \frac{1}{\mu}$ and $\gamma > 0$ chosen so that $\alpha_0 = \frac{\beta}{\gamma} \leq \frac{1}{L}$. Then, for any $K \in \mathbb{N}$, we have

$$\mathbb{E}[f(\mathbf{w}_K) - f^*] \leq \frac{\nu}{\gamma + K}, \quad (3.2.4)$$

where

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

As for gradient descent, using a decreasing stepsize sequence poses the risk of producing tiny steps early in the algorithmic run. However, note that it guarantees convergence of $\{f(\mathbf{w}_k)\}$ to f^* in expectation, which is a strongly result than that of Theorem 3.2.1. Still, the convergence rate in $\mathcal{O}(1/K)$ is worse than the one we established for gradient descent, in $\mathcal{O}((1 - \mu/L)^K)$. Note however that this comparison is related to the number of iterations: if we include the per-iteration cost (in terms of accesses to the data), the stochastic gradient method has a convergence rate of $\mathcal{O}(1/K)$ (1 access per iteration) while the gradient descent method yields a rate in $\mathcal{O}(n(1 - \mu/L)^K)$. The latter can be significantly worse than the former for large n .

Remark 3.2.3 (A hybrid approach) *A common stepsize strategy used in deep learning consists in adopting a learning rate schedule. In this hybrid approach, one runs first stochastic gradient with a constant value α upon a certain point, then replaces α by $\alpha' < \alpha$: the process is repeated until a suitable point has been found, or the budget of epochs is exhausted.*

In the strongly convex setting, the result of Theorem 3.2.1 provides a useful criterion to reduce the stepsize. Indeed, starting with a stepsize α , one can guarantee that a neighborhood of size $\frac{\alpha L \sigma^2}{\mu}$ can be reached in a bounded number of iterations (or epochs). Therefore, once that point is reached, one can reduce α . This approach guarantees to converge to the true optimal value, at a sublinear rate: for any $\epsilon > 0$,

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

In practice, knowing how the neighborhood is reached may not be possible: other criteria, such as detecting stalling in the algorithm (lack of progress in the iterates, in the function values, etc), can help in designing strategies: several learning rate schedules are implemented in modern libraries.

Stepsize choice and nonconvex optimization Stochastic gradient and its variants are the most common methods for training deep neural networks: the associated optimization problem is highly nonconvex, and the above analysis cannot be applied. Still, it is possible to derive convergence rates for the nonconvex case by studying the following quantities:

- $\mathbb{E} \left[\frac{1}{K} \sum_{i=1}^K \|\nabla f(\mathbf{w}_k)\|^2 \right]$ for the variants of Algorithm 4 based on 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 the variants using a decreasing stepsize sequence.

As in the strongly convex case, the iteration results for stochastic gradient will be worse than that of gradient descent. For instance, applying 4 with a constant stepsize guarantees that

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

in at most $\mathcal{O}(\epsilon^{-4})$ iterations. This bound is worse than the $\mathcal{O}(\epsilon^{-2})$ bound for gradient descent. In addition, this result is only valid for sufficiently large ϵ because of the stochastic nature of the method.

Remark 3.2.4 (Advanced stochastic gradient techniques) *The most common stochastic gradient variants for deep learning (SGD WITH MOMENTUM, ADAM, ADAGRAD, RMSPROP) build on the stochastic gradient principle by adding a momentum term to the stochastic gradient and/or a scaling to every component of the step so as to reduce the need for intensive stepsize tuning. Analyzing these methods is significantly more complicated than in the deterministic setting, and theoretical results have yet to be reconciled with practical behavior.*

3.3 Variance reduction

The theory of Section 3.2.2 relies heavily on Assumption 3.2.2, and more precisely on a control over the norm of the stochastic gradient through the quantity σ^2 . A higher value of σ leads to looser neighborhoods of the solution and, as a result, worse theoretical guarantees. Numerically, this translates into high variability of the method's output.

Variance reduction techniques have been designed to reduce the “variance” of gradient estimates used in stochastic gradient. The goal of this section is to summarize key approaches to reducing variance.

3.3.1 Batch methods

As explained in the previous section, the iteration of stochastic gradient is

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

where i_k is a random index drawn within $\{1, \dots, n\}$. This method relies on a **single** example to build a stochastic gradient estimate, and this estimate comes with a certain variance (illustrated by the parameter σ^2 in Assumption 3.2.2). On the contrary, gradient descent relies on an exact gradient computed using **all the samples**, which leads to a deterministic execution and no variance in the (exact) gradient estimate.

To improve the variance of this method, it appears natural to consider stochastic gradient estimates based on **several samples**: this is the underlying principle of **batch stochastic gradient methods**.

At every iteration of a batch method, a (random) set of indices S_k with indices between 1 and n is drawn, and the following update is performed:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k) \quad (3.3.1)$$

When $|S_k| = 1$, we recover the classical stochastic gradient formula. When $|S_k| = n$ and the indices are drawn without replacement, we obtain $S_k = \{1, \dots, n\}$, and iteration (3.3.1) corresponds to that of gradient descent.

More broadly, we can identify two regimes of batch sizes:

- $|S_k| \approx n$: the per-iteration cost of such a variant is close to that of gradient descent, hence these batch methods often exhibit a behavior comparable to that of gradient descent.
- $1 < |S_k| = n_b \ll n$: this regime, called **mini-batching**, is often thought as a good way to reduce variance while keeping the per-iteration cost to an affordable level.

Assuming that the batch size is constant over all iterations, i.e. that $|S_k| = n_b \forall k$, it is possible to show that a mini-batch variant reaches a closer neighborhood than stochastic gradient. Every iteration of the mini-batch method is of course more expensive than an iteration of stochastic gradient (in fact, an epoch corresponds to n/n_b iterations of a batch stochastic gradient method with a fixed batch size of n_b). Note that the cost of batch methods can be mitigated by exploiting parallel computing resources. Indeed, in a distributed data setup, stochastic gradients can often be evaluated in parallel, which reduces the effective cost of batch techniques.

Proposition 3.3.1 *Under Assumptions 2.0.1 and 3.2.2, we have*

$$\mathbb{E}_{S_k} \left[\frac{1}{|S_k|} \sum_{i \in S_k} \nabla f_i(\mathbf{w}_k) \right] = \nabla f(\mathbf{w}_k)$$

and

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

Finally, we mention that (mini)-batch approaches remain more expensive than “vanilla” stochastic gradient, while being more sensitive to redundancies in the data, and introducing an new hyperparameter to be tuned (the batch size). For this reason, the classical stochastic gradient approach can still be observed to be more efficient on some example, and it remains a popular algorithmic choice in practice.

3.3.2 Other variance reduction techniques

Gradient aggregation techniques are variance reduction approaches with well-established convergence rates, that are provably better than that of stochastic gradient. They have attracted significant interest from the academic machine learning and optimization community, though their use in modern data science settings like deep learning has not been successful (these techniques are still efficient on simple models, and are part of libraries like `scikit-learn`). They essentially consist in maintaining a **full gradient estimate** throughout the iterations, which requires at least one full gradient calculation: the gradient estimate used to make a step then combines the full gradient estimator with a new stochastic gradient, leading to a corrected step, and a method with provably smaller variance.

Iterate averaging is another way to reduce the variance in stochastic gradient at no additional cost in terms of accesses to data points. The idea consists in analyzing the properties of a running

average $\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{w}_k$. In certain cases (in particular, when $\alpha = \frac{1}{\mu(k+1)}$ and f is μ -strongly convex), this sequence possesses favorable properties and its behavior is less variable. However, computing this average either requires to store all iterates (which becomes expensive in terms of storage) or necessitates to maintain a running average (subject to numerical errors).

3.4 Stochastic gradient methods for deep learning

In this section, we review the main stochastic gradient techniques that are used to train deep learning models. Our focus remains on finite-sum problems of the form (3.1.1) under Assumption 3.2.1. Our goal is to study several variants on the iterative scheme

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k, \quad (3.4.1)$$

where $\alpha > 0$ is a fixed stepsize (or learning rate), and \mathbf{g}_k is a stochastic gradient estimator, that may correspond to taking a single term from the finite sum (as in stochastic gradient) or to a batch of indices.

To encompass all the variants of interest, we consider a more general iteration of the form

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{m}_k \oslash \mathbf{v}_k, \quad (3.4.2)$$

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

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

To see that (3.4.2) generalizes (3.4.1), set $\mathbf{m}_k = \mathbf{g}_k$ and $\mathbf{v}_k = \mathbf{1}_{\mathbb{R}^d}$. The iteration (3.4.2) is particularly convenient to express the popular methods in deep learning using a single format.

3.4.1 Stochastic gradient with momentum

Most practical implementations of stochastic gradient combine the basic step (3.4.1) together with a momentum term, similarly to the accelerated methods described in Chapter 2. An iteration of gradient descent with momentum reads

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha(1 - \beta)\mathbf{g}_k + \alpha\beta(\mathbf{w}_k - \mathbf{w}_{k-1}), \quad (3.4.3)$$

where $\beta \in (0, 1)$ is a momentum parameter (when $\beta = 0$ we again obtain the standard stochastic gradient iteration). Iteration (3.4.3) can be seen as a version of Polyak's method where the gradient-type step is combined with the previous direction. As for the heavy-ball method, the idea is to incorporate information from the previous step through momentum. In practice, the iteration (3.4.3) often leads to accumulation of good steps (in terms of optimization), whereas bad directions and bad steps tend to cancel out.

The basic method (3.4.3) can be recovered from (3.4.2), by setting $\mathbf{v}_k = \mathbf{1}_{\mathbb{R}^d}$ and defining \mathbf{m}_k in a recursive manner through $\mathbf{m}_{-1} = \mathbf{0}_{\mathbb{R}^d}$ and

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

where $\beta \in (0, 1)$.

Stochastic gradient with momentum is implemented in most deep learning libraries such as PyTorch. It has shown great success in training deep neural networks on computer vision problems, and is partly responsible for the rise of deep learning in the early 2010s.

Remark 3.4.1 *Theoretical guarantees for the iteration (3.4.3) are much more difficult to obtain than for accelerated gradient (in particular, it is not well understood whether such a method can be provably faster than SGD). Nevertheless, stochastic gradient techniques with momentum are widely used in practice, even on nonconvex problems such as neural network training.*

3.4.2 AdaGrad

The adaptive gradient method, or ADAGRAD, was proposed in 2011 to address the issue of setting the learning rate α in stochastic gradient. Rather than using costly procedures such as line search, ADAGRAD scales every coordinate of the stochastic gradient using information from the values of that coordinate in the previous iterations. Mathematically, the method maintains a sequence $\{\mathbf{r}_k\}_k$ defined by

$$\forall i = 1, \dots, d, \quad \begin{cases} [\mathbf{r}_{-1}]_i = 0 \\ [\mathbf{r}_k]_i = [\mathbf{r}_{k-1}]_i + [\mathbf{g}_k]_i^2 \end{cases} \quad \forall k \geq 0, \quad (3.4.4)$$

The ADAGRAD iteration can then be written as

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \mathbf{g}_k \oslash \sqrt{\mathbf{r}_k}, \quad (3.4.5)$$

where the square root is applied componentwise to \mathbf{r}_k . This iteration is a special case of (3.4.2), where $\mathbf{m}_k = \mathbf{g}_k$ and $\mathbf{v}_k = \sqrt{\mathbf{r}_k}$. The novelty in ADAGRAD does not lie in the use of momentum, but in the use of one stepsize per coordinate. The stepsize sequence thus has the form

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

The resulting *diagonal scaling* on the coordinates of \mathbf{g}_k leads to stepsizes that adapt to coordinates that can vary by orders of magnitude (which would require a careful choice of α in basic stochastic gradient). On the other hand, the accumulation process at work in the definition of \mathbf{r}_k results in stepsizes that are monotonically decreasing, and that often converge quickly towards 0.

Remark 3.4.2 *In practice, \mathbf{r}_k is replaced by $\mathbf{r}_k + \eta \mathbf{1}_{\mathbb{R}^d}$ where $\eta > 0$ is a small value that helps with numerical stability.*

The ADAGRAD method is particularly interesting for problems with sparse gradients, in which stochastic gradients tend to have many zero coordinates. In this situation, using \mathbf{r}_k will only modify the stepsizes corresponding to nonzero gradient coordinates. Many problems in recommendation systems have a sparse structure, and ADAGRAD is considered to be an efficient method for this class of problems.

3.4.3 RMSProp

The RMSPROP (*Root Mean Square Propagation*) algorithm is similar in spirit to ADAGRAD, in that it scales the gradient components at every step. This method relies on a sequence $\{\mathbf{r}_k\}_k$ defined by

$$\forall i = 1, \dots, d, \quad \begin{cases} [\mathbf{r}_{-1}]_i = 0 \\ [\mathbf{r}_k]_i = (1 - \lambda)[\mathbf{r}_{k-1}]_i + \lambda[\mathbf{g}_k]_i^2 \end{cases} \quad \forall k \geq 0, \quad (3.4.6)$$

where $\lambda \in (0, 1)$. The value of λ is used to put more weight on the current gradient coordinates than on the coordinates from past iterations (this information being contained in \mathbf{r}_{k-1}). This simple idea slows down the decrease of the stepsizes to 0, compared to the stepsizes of ADAGRAD.

With the definition (3.4.6), the RMSPROP iteration corresponds has the same form as that of ADAGRAD, that is, a special case of (3.4.2) with $\mathbf{m}_k = \mathbf{g}_k$ and $\mathbf{v}_k = \sqrt{\mathbf{r}_k}$.

Remark 3.4.3 As for ADAGRAD, standard practice replaces \mathbf{r}_k by $\mathbf{r}_k + \eta \mathbf{1}_{\mathbb{R}^d}$, where $\eta > 0$ is a small quantity.

The RMSPROP method has been found quite successful for training very deep neural networks.

3.4.4 Adam

The ADAM optimization method was proposed in 2013. This method can be thought as combining the idea of momentum (used in the stochastic gradient method of Section 3.4.1) with the diagonal scaling procedure on which both ADAGRAD and RMSPROP are based. An iteration of ADAM falls into the generic scheme (3.4.2) by setting

$$\mathbf{m}_k = \frac{(1 - \beta_1) \sum_{j=0}^k \beta_1^{k-j} \mathbf{g}_j}{1 - \beta_1^{k+1}} \quad (3.4.7)$$

and

$$\mathbf{v}_k = \sqrt{\frac{(1 - \beta_2) \sum_{j=0}^k \beta_2^{k-j} \mathbf{g}_j \odot \mathbf{g}_j}{1 - \beta_2^{k+1}}}. \quad (3.4.8)$$

Here $\beta_1, \beta_2 \in (0, 1)$, and \odot denotes the Hadamard or componentwise product

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

Remark 3.4.4 In practice, $\mathbf{v}_k + \eta \mathbf{1}_{\mathbb{R}^d}$ (with small $\eta > 0$) is used in lieu of \mathbf{v}_k .

The above formulas describe the two components of ADAM. On one hand, a weighted combination of the previous steps that puts the emphasis on the most recent steps (and the current stochastic gradient) defines the direction of the next step. On the other hand, a diagonal scaling is applied to the coordinates of this direction, again according to a weighted average of the coordinates from the previous iterations. This important feature, that has a statistical motivation, appears to be responsible for the success of ADAM in practice. The impressive performance of ADAM on training neural networks has contributed to its popularity, and it remains the preferred method today in numerous applications. In particular, ADAM and its variant ADAMW (based on regularization principle) are quite efficient on natural language processing models.

3.5 Conclusion

Stochastic gradient methods rely on partial gradient information selected in a random fashion: as such, they are not guaranteed to succeed deterministically, and their convergence guarantees are typically weaker than that of gradient descent. Nevertheless, they have proven very successful in data-driven problems, wherein computing the gradient involves accessing a massive amount of data

points. In this setting, stochastic gradient approaches have an attractive, low per-iteration cost, and typically make fast, significant progress early on in the algorithmic run. The nature of the data is key to the performance of stochastic gradient: on standard datasets from machine and deep learning, it appears beneficial to resort to such techniques.

Multiple variants of stochastic gradient have been developed in academic and industrial contexts. Batch stochastic gradient techniques are among the most popular, as they allow to incorporate more than one sample into the gradient approximation: this leads to variance reduction, a concept that is driving some of the latest advances in stochastic gradient.

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] S. J. Wright and B. Recht. *Optimization for Data Analysis*. Cambridge University Press, 2022.

Appendix A

Notations and mathematical tools

A.1 Notations

A.1.1 Generic notations

- Scalars (i.e. reals) are denoted by lowercase letters: $a, b, c, \alpha, \beta, \gamma$.
- 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 Σ 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 Π operator is used for products. 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 $\prod_{i=1}^m \prod_{j=1}^n$, $\prod_i \prod_j$ and $\prod_{i,j}$ may be used interchangeably.
- The notation $i = 1, \dots, m$ indicates that the variable i takes all integer values between 1 and m .

A.1.2 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. For simplicity, we may write $\mathbf{w} \geq 0$ to indicate that all components of \mathbf{w} are nonnegative.

A.1.3 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} or $[\mathbf{A}]_{ij}$ refers to the coefficient from the i -th row and the j -th column of \mathbf{A} . 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 $[A_{ij}]$ 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.
- The diagonal of a square matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is given by the coefficients A_{ii} . The trace of such a matrix is $\text{trace}(\mathbf{A}) := \sum_{i=1}^d A_{ii}$.
- 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

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

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).

A.2 Mathematical tools

Optimization has mathematical roots in real analysis, mostly through differential calculus. Linear algebra structures also play a major role in optimization (and data science). In this section, we list the basic results that will be used in the course.

For a deeper dive into these notions, the following links are recommended.

- For linear algebra:
 - <https://www.ceremade.dauphine.fr/~carlier/polyalgebre.pdf> (in French);
 - <http://vmls-book.stanford.edu/vmls.pdf> (Chapters 1 to 3, in English).
- For differential calculus:
 - https://www.ceremade.dauphine.fr/~bouin/ens1819/Cours_Bolley.pdf (in French);
 - https://sebastianraschka.com/pdf/books/dlb/appendix_d_calculus.pdf (in English).

A.2.1 Vector linear algebra

We always consider vectors in the normed vector space \mathbb{R}^d , of dimension d . The following operations are defined in this space:

- For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, the sum of \mathbf{x} and \mathbf{y} is denoted by $\mathbf{x} + \mathbf{y} = [x_i + y_i]_{1 \leq i \leq d}$;
- For any $\lambda \in \mathbb{R}$, we define $\lambda \mathbf{x} \stackrel{n}{=} \lambda \cdot \mathbf{x} = [\lambda x_i]_{1 \leq i \leq d}$. In this context, the real value λ is called a *scalar*.

Using these operations, we can build **linear combinations** of vectors in \mathbb{R}^d that produce a vector in \mathbb{R}^d of the form $\sum_{i=1}^p \lambda_i \mathbf{x}_i$, where $\mathbf{x}_i \in \mathbb{R}^d$ and $\lambda_i \in \mathbb{R}$ for any $i = 1, \dots, p$.

The matrix space $\mathbb{R}^{n \times d}$ can also be endowed with a vector space structure of dimension nd :

- For any $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times d}$, the sum of \mathbf{A} and \mathbf{B} is denoted by $\mathbf{A} + \mathbf{B} = [\mathbf{A}_{ij} + \mathbf{B}_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$;
- For any scalar $\lambda \in \mathbb{R}$, we define $\lambda \mathbf{A} \stackrel{n}{=} \lambda \cdot \mathbf{A} = [\lambda \mathbf{A}_{ij}]_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$.

Definition A.2.1 A set $\mathcal{S} \subseteq \mathbb{R}^d$ satisfying the conditions

1. $\mathbf{0}_d \in \mathcal{S}$;
2. $\forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}, \mathbf{x} + \mathbf{y} \in \mathcal{S}$;
3. $\forall \mathbf{x} \in \mathcal{S}, \forall \lambda \in \mathbb{R}, \lambda \mathbf{x} \in \mathcal{S}$.

is called a (linear) subspace of \mathbb{R}^d .

Definition A.2.2 Let $\mathbf{x}_1, \dots, \mathbf{x}_p$ be p vectors in \mathbb{R}^d . The **span** (or linear span) of $\mathbf{x}_1, \dots, \mathbf{x}_p$, denoted by $\text{Span}(\mathbf{x}_1, \dots, \mathbf{x}_p)$, is the subspace of \mathbb{R}^d defined by

$$\text{Span}(\mathbf{x}_1, \dots, \mathbf{x}_p) := \left\{ \mathbf{x} = \sum_{i=1}^p \alpha_i \mathbf{x}_i \mid \alpha_i \in \mathbb{R} \forall i \right\}.$$

We now recall various properties of vector sets.

Definition A.2.3 • The vectors in a set $\{x_i\}_{i=1}^k \subset \mathbb{R}^n$ are called *linearly independent* if for any scalars $\lambda_1, \dots, \lambda_k$ satisfying $\sum_{i=1}^k \lambda_i x_i = 0$, we have $\lambda_1 = \dots = \lambda_k = 0$. In that case, $k \leq n$.

- If the above property does not hold, the vectors are called *linearly dependent*.
- A *spanning set* is a set of vectors $\{x_i\} \subset \mathbb{R}^n$ such that their span is \mathbb{R}^n .
- A set of vectors $\{x_i\}_{i=1}^n \subset \mathbb{R}^n$ is a *basis* if it is both linearly independent and a spanning set. In that case, any vector in \mathbb{R}^n can be written as a uniquely defined linear combination of the x_i s. Any basis in \mathbb{R}^n has exactly n vectors.

Since the size of a basis in \mathbb{R}^n is n , we say that the dimension of the space is n . Consequently, any subspace of \mathbb{R}^n has dimension at most n .

Example A.2.1 Any vector x in \mathbb{R}^n can be written as $x = \sum_{i=1}^n x_i e_i$, where $e_i = [0 \dots 0 \ 1 \ 0 \dots 0]^T$ is the i th vector of the canonical basis (with a 1 in the i th coordinate).

Norm and scalar product Using a Euclidean norm and its associated scalar product allows to compare vectors by measuring the distance between them. This ability is particularly useful to establish that a sequence of vector generated by an optimization method converges toward the solution of a given problem.

Definition A.2.4 The **Euclidean norm** $\|\cdot\|$ on \mathbb{R}^n is defined by

$$\forall x \in \mathbb{R}^n, \quad \|x\| := \sqrt{\sum_{i=1}^n x_i^2}.$$

Remark A.2.1 This is indeed a norm, since it fulfills the four axioms that define what a norm is:

1. $\forall x, y \in \mathbb{R}^n, \|x + y\| \leq \|x\| + \|y\|$;
2. $\|x\| = 0 \Leftrightarrow x = \mathbf{0}_{\mathbb{R}^n}$;
3. $\forall x, \|x\| \geq 0$;
4. $\forall x \in \mathbb{R}^n, \forall \lambda \in \mathbb{R}, \|\lambda x\| = |\lambda| \|x\|$.

A vector $x \in \mathbb{R}^n$ is called a unit vector if $\|x\| = 1$.

Definition A.2.5 For any vectors $x, y \in \mathbb{R}^n$, the **scalar product** derived from the Euclidean norm is a function of x and y , denoted by $x^T y$, defined as follows:

$$x^T y := \sum_{i=1}^n x_i y_i.$$

Two vectors x and y are called *orthogonal* if $x^T y = 0$.

Note that $\mathbf{y}^T \mathbf{x} = \mathbf{x}^T \mathbf{y}$, hence the scalar product defines a “product” between a row vector and a column vector.

Proposition A.2.1 *Let \mathbf{x} and \mathbf{y} be two vectors in \mathbb{R}^n . Then, the following properties hold*

- i) $\|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2\mathbf{x}^T \mathbf{y} + \|\mathbf{y}\|^2$;
- ii) $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 - 2\mathbf{x}^T \mathbf{y} + \|\mathbf{y}\|^2$;
- iii) $\|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 = \frac{1}{4} (\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2)$;
- iv) **Cauchy-Schwarz inequality** :

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad \mathbf{x}^T \mathbf{y} \leq \|\mathbf{x}\| \|\mathbf{y}\|.$$

Remark A.2.2 *The last inequality is a key result in both linear algebra and analysis. In this course, it will play a major role in deriving Taylor-type inequalities.*

A.2.2 Matrix linear algebra

We can define the product of two matrices that have compatible dimensions. More precisely, for any $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$, the product matrix \mathbf{AB} is defined as the matrix $\mathbf{C} \in \mathbb{R}^{m \times p}$ such that

$$\forall i = 1, \dots, m, \forall j = 1, \dots, p, \quad C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}.$$

Using this definition, the product of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with a (column) vector $\mathbf{x} \in \mathbb{R}^n$ is the vector $\mathbf{y} \in \mathbb{R}^m$ given by

$$\forall i = 1, \dots, m, \quad y_i = \sum_{j=1}^n A_{ij} x_j.$$

Remark A.2.3 *Note that the scalar product on \mathbb{R}^n corresponds to the matrix product for matrices of sizes $1 \times n$ and $n \times 1$: the result of this operation is a 1×1 matrix, that is, a scalar.*

When one work with matrices, the following subspaces are of interest.

Definition A.2.6 (Matrix subspaces) *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$.*

- *The null space of \mathbf{A} is the subspace*

$$\text{Null}(\mathbf{A}) := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} = \mathbf{0}_m\}$$

- *The range space of \mathbf{A} is the subspace*

$$\text{Range}(\mathbf{A}) := \{\mathbf{y} \in \mathbb{R}^m \mid \exists \mathbf{x} \in \mathbb{R}^n, \mathbf{y} = \mathbf{Ax}\}$$

*The dimension of this subspace is called the **rank** of \mathbf{A} . We denote it by $\text{rank}(\mathbf{A})$. One always has $\text{rank}(\mathbf{A}) \leq \min\{m, n\}$.*

Theorem A.2.1 (Rank-nullity theorem) Let $A \in \mathbb{R}^{m \times n}$. Then,

$$\dim(\ker(A)) + \text{rang}(A) = n.$$

Definition A.2.7 (Matrix norms) Consider the space $\mathbb{R}^{m \times n}$. The operator norm $\|\cdot\|$ and the Frobenius norm $\|\cdot\|_F$ are defined by

$$\forall A \in \mathbb{R}^{m \times n}, \begin{cases} \|A\| &:= \max_{\substack{x \in \mathbb{R}^n \\ x \neq 0_n}} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\| \\ \|A\|_F &:= \sqrt{\sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} A_{ij}^2}. \end{cases}$$

Definition A.2.8 (Symmetric matrix) A square matrix $A \in \mathbb{R}^{n \times n}$ is called *symmetric* if $A^T = A$. The set of symmetric matrices in $\mathbb{R}^{n \times n}$ is denoted by S^n .

Definition A.2.9 (Invertible matrix) A square matrix $A \in \mathbb{R}^{n \times n}$ is called *invertible* if there exists $B \in \mathbb{R}^{n \times n}$ such that $BA = AB = I_n$ (where we recall that I_n denotes the identity matrix in $\mathbb{R}^{n \times n}$).

When it exists, such a matrix B is unique. It is then called **the inverse of A** and denoted by A^{-1} .

Definition A.2.10 (Positive (semi)definite matrix) A square, symmetric matrix $A \in \mathbb{R}^{n \times n}$ is called *positive semidefinite* if

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

which we write $A \succeq 0$.

Such a matrix is called *positive definite* when $x^T A x > 0$ for any nonzero vector x . We write this as $A \succ 0$.

Definition A.2.11 (Orthogonal matrix) A square matrix $P \in \mathbb{R}^{n \times n}$ is called *orthogonal* if $P^T = P^{-1}$.

More generally, a matrix $Q \in \mathbb{R}^{m \times n}$, where $m \leq n$, is called *orthogonal* if $QQ^T = I_m$ (the columns of Q are orthonormal in \mathbb{R}^m).

When $Q \in \mathbb{R}^{n \times n}$ is orthogonal, then so is its transpose Q^T (this result only applies to square matrices). Orthogonal matrices have the following desirable property.

Lemma A.2.1 Let $A \in \mathbb{R}^{m \times n}$ and $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ be two orthogonal matrices. Then,

$$\|A\| = \|UA\| = \|AV\| \quad \text{and} \quad \|A\|_F = \|UA\|_F = \|AV\|_F,$$

i.e. multiplying by an orthogonal matrix preserves the norm.

As a corollary of the previous lemma, we observe that an orthogonal matrix $Q \in \mathbb{R}^{m \times n}$ with $m \leq n$ must satisfy $\|Q\| = 1$ and $\|Q\|_F = \sqrt{m}$.

Definition A.2.12 (Eigenvalue) Let $A \in \mathbb{R}^{n \times n}$. A scalar $\lambda \in \mathbb{R}$ is called an **eigenvalue of A** if

$$\exists v \in \mathbb{R}^n, v \neq 0_n, \quad Av = \lambda v.$$

The vector v is called an **eigenvector associated with the eigenvalue λ** . The set of eigenvalues of A is the *spectrum of A* .

The span of eigenvectors associated to the same eigenvalue is called the eigenspace. Its dimension corresponds to the multiplicity of the eigenvalue relatively to the matrix.

Proposition A.2.2 For any matrix $A \in \mathbb{R}^{n \times n}$, the following holds:

- A has n complex eigenvalues.
- If A is symmetric positive semidefinite (resp. definite), then its eigenvalues are real nonnegative (resp. real positive).
- The null space of A is spanned by the eigenvectors associated with the 0 eigenvalue.

Theorem A.2.2 (Eigenvalue decomposition theorem) Any symmetric matrix $A \in \mathbb{R}^{n \times n}$ has an **eigenvalue decomposition** of the form

$$A = P\Lambda P^T,$$

where $P \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, and $\Lambda \in \mathbb{R}^{n \times n}$ is a diagonal matrix that contains the n eigenvalues of A $\lambda_1, \dots, \lambda_n$ on its diagonal.

The eigenvalue decomposition is not unique, but the set of eigenvalues that appears in the decomposition is uniquely defined.

Remark A.2.4 There are matrices that possess an eigenvalue decomposition of the form $P\Lambda P^{-1}$, where P is invertible (but not necessarily orthogonal). Those matrices are called diagonalizable.

Link with singular value decomposition Let $A \in \mathbb{R}^{m \times n}$. In general, $m \neq n$ and the notion of eigenvalue that we introduced above does not apply. However, we can always consider the eigenvalues of

$$A^T A \in \mathbb{R}^{n \times n} \quad \text{and} \quad A A^T \in \mathbb{R}^{m \times m}.$$

These matrices are real and symmetric, hence they can be diagonalized. This property is what gives rise to the singular value decomposition (or SVD).

A.2.3 Calculus

Note: This section gives additional background to the concepts and properties introduced in Chapter 1.

Definition A.2.13 (Continuity) A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called **continuous in** $x \in \mathbb{R}^n$ if

$$\forall \epsilon > 0, \exists \delta > 0, \forall y \in \mathbb{R}^n, \quad \|y - x\| < \delta \quad \Rightarrow \quad \|f(y) - f(x)\| < \epsilon.$$

The function f is **continuous** on a set $\mathcal{A} \subseteq \mathbb{R}^n$ if it is continuous at every point of \mathcal{A} . When $\mathcal{A} = \mathbb{R}^n$, we simply say that f is continuous.

Remark A.2.5 In certain textbooks, the notion above is termed uniform continuity. For simplicity of exposure, we will use it as our definition of continuity.

An alternate characterization of continuity based on sequences is given below. Sequences typically appear when considering iterative algorithms, hence the relevance of this notion here.

Definition A.2.14 (Continuity (sequential definition)) A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuous at $\mathbf{x} \in \mathbb{R}^n$ if

$$\forall \{\mathbf{x}_n\} \in (\mathbb{R}^n)^\mathbb{N}, \quad \{\mathbf{x}_n\} \rightarrow \mathbf{x}, \quad \lim_{n \rightarrow \infty} f(\mathbf{x}_n) = f(\mathbf{x}).$$

Example A.2.2 A linear map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, where $f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ for any $\mathbf{x} \in \mathbb{R}^n$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$, is a continuous function on \mathbb{R}^n .

Definition A.2.15 (Differentiability Jacobian matrix) A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called **differentiable** at a point $\mathbf{x} \in \mathbb{R}^n$ if there exists a matrix $\mathbf{J}_f(\mathbf{x}) \in \mathbb{R}^{m \times n}$ such that

$$\lim_{\substack{\mathbf{z} \rightarrow \mathbf{x} \\ \mathbf{z} \neq \mathbf{x}}} \frac{\|f(\mathbf{z}) - f(\mathbf{x}) - \mathbf{J}_f(\mathbf{x})(\mathbf{z} - \mathbf{x})\|}{\|\mathbf{z} - \mathbf{x}\|} = 0.$$

- $\mathbf{J}_f(\mathbf{x})$ is called the **Jacobian** of f at \mathbf{x} , and is uniquely defined.
- If $f(\cdot) = [f_1(\cdot), \dots, f_m(\cdot)]^\top$, then

$$\forall 1 \leq i \leq m, \forall 1 \leq j \leq n, [\mathbf{J}_f(\mathbf{x})]_{ij} = \frac{\partial f_i}{\partial x_j}(\mathbf{x}).$$

The following special cases are instrumental to optimization and basic analysis.

Corollary A.2.1 • When $m = 1$, we define the (column) vector $\nabla f(\mathbf{x}) \equiv \mathbf{J}_f(\mathbf{x})^\top$, called the **gradient** of f at \mathbf{x} . In this case, the gradient is the vector of partial derivatives of f :

$$\forall i = 1, \dots, n, \quad \nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_i}(\mathbf{x}) \right]_{1 \leq i \leq n}.$$

- When $n = m = 1$, both the Jacobian and the gradient are equivalent to a scalar la matrice Jacobienne et le vecteur $f'(\mathbf{x}) \equiv \nabla f(\mathbf{x}) \equiv \mathbf{J}_f(\mathbf{x})^\top$, called the derivative of f at \mathbf{x} .

In these notes, we assume familiarity with the common derivative formulas for functions from \mathbb{R} to \mathbb{R} . More complex formulas are typically obtained thanks to the rule below.

Theorem A.2.3 (Chain rule) If $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ and $g : \mathbb{R}^m \mapsto \mathbb{R}^p$ are both differentiable, respectively on \mathbb{R}^n and \mathbb{R}^m , then $h : \mathbb{R}^n \mapsto \mathbb{R}^p$ is differentiable on \mathbb{R}^n and

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{J}_h(\mathbf{x}) = \mathbf{J}_g(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x}).$$

Remark A.2.6 Special cases of the chain rule:

- $m = p = 1 : \nabla h(\mathbf{x}) = g'(f(\mathbf{x}))\nabla f(\mathbf{x});$
- $n = m = p = 1 : h'(x) = g'(f(x))f'(x).$

Theorem A.2.4 (Mean-value theorem in dimension 1) Let $f : [a, b] \rightarrow \mathbb{R}$. If f is continuous on $[a, b]$ and differentiable on (a, b) , there exists $c \in (a, b)$ such that

$$\frac{f(b) - f(a)}{b - a} = f'(c).$$

Definition A.2.16 (Taylor expansion) Let $f : [a, b] \mapsto \mathbb{R}$ be \mathcal{C}^1 on $[a, b]$, then

$$\begin{aligned} f(b) &= f(a) + f'(c)(b - a) \quad \text{where } c \in [a, b] \\ f(b) &= f(a) + \int_0^1 f'(a + t(b - a))(b - a) dt. \end{aligned}$$

Theorem A.2.5 (Mean-value theorem in dimension d) Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ $f \in \mathcal{C}^1(\mathbb{R}^d)$. For any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, $\mathbf{x} \neq \mathbf{y}$, there exists $t \in (0, 1)$ such that

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T(\mathbf{y} - \mathbf{x}).$$

Definition A.2.17 (Lipschitz continuity) A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is L -**Lipschitz continuous** on $\mathcal{A} \subset \mathbb{R}^n$ if

$$\forall \mathbf{x}, \mathbf{y} \in \mathcal{A}, \quad \|f(\mathbf{x}) - f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Proposition A.2.3 Any Lipschitz continuous function on a set is continuous on this set.

Definition A.2.18 (Function classes) Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$.

- We say that f is $\mathcal{C}^p(\mathbb{R}^d)$ (or simply \mathcal{C}^p) if it is differentiable p times with a continuous p th-order derivative (in which case all derivatives up to order p are continuous). The class of \mathcal{C}^∞ functions is the intersection of all \mathcal{C}^p with $p \in \mathbb{N}$.
- We say that f is $\mathcal{C}_L^{p,p}(\mathbb{R}^d)$ (or simply $\mathcal{C}_L^{p,p}$) if it is differentiable p times and its p th-order derivative is L -Lipschitz continuous.

Theorem A.2.6 (Taylor expansion of order 1) Let $f \in \mathcal{C}^1(\mathbb{R}^d)$. For any vectors \mathbf{x} and \mathbf{y} of \mathbb{R}^d , we have

$$f(\mathbf{y}) = f(\mathbf{x}) + \int_0^1 \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))^T(\mathbf{y} - \mathbf{x}) dt.$$

Moreover, if $f \in \mathcal{C}_L^{1,1}(\mathbb{R}^d)$, then

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|^2. \quad (\text{A.2.1})$$

Theorem A.2.7 (Taylor expansion of order 2) Let $f \in \mathcal{C}^2(\mathbb{R}^d)$. For any vectors \mathbf{x} and \mathbf{y} of \mathbb{R}^d , one has

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2} \int_0^1 (\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}) dt.$$

Moreover, if $f \in \mathcal{C}_L^{2,2}(\mathbb{R}^d)$, then

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|^3. \quad (\text{A.2.2})$$