# Basics of Applied Mathematics - Part III: Optimization

# Preliminary Draft Version, from January 20, 2025

Léo Simpson and Moritz Diehl Department of Microsystems Engineering and Department of Mathematics, University of Freiburg, Germany moritz.diehl@imtek.uni-freiburg.de

# Preface

This script is designed for the course Basics of Applied Mathematics (BAM) - Optimization, for students pursuing a master's of mathematics that focuses on data and technologies.

The script is quite self-contained. Nevertheless, to follow this course, one must already have good mathematics fundamentals.

An important note is that this script is largely based on the book [1].

# Contents

	Pref	ace				
1	Optimization problems: Definitions and Applications 4					
	$1.1^{-1}$	Optimization in the real world				
	1.2	General definitions and notations				
		1.2.1 Formulation of optimization problems				
		1.2.2 Minimizers				
		1.2.3 Local optimality				
	1.3	Examples of optimization problems in data analysis				
		1.3.1 Regression problems				
		1.3.2 Ridge regression				
		1.3.3 LASSO regression				
		1.3.4 Cross-entropy for classification tasks				
	1.4	Basic properties of optimization problems 11				
		1.4.1 Existence of solutions 11				
		1.4.2 First order optimality conditions for smooth functions				
		1.4.3 Second order optimality conditions for smooth functions				
	1.5	Application of the properties to the regression examples 16				
		1.5.1 Quadratic programs with $Q \succ 0$				
		1.5.2 Quadratic Programs with $Q \succeq 0$				
		1.5.3 LASSO regression: existence of minimizers 19				
<b>2</b>	Convexity 20					
	2.1					
	2.2	Convex functions				
		2.2.1 General case				
		2.2.2 Convexity of smooth functions 23				
		2.2.3 Strictly convex Functions				
		2.2.4 Strong convexity				
	2.3	Convex optimization problems				
	2.4	Examples of convex optimization problems in data analysis				
3	$\mathbf{Des}$	cent Methods for Solving Optimization Problems 38				
	3.1	Generalities about Descent methods				
	3.2	The gradient descent method				
	3.3	Convergence properties				
		3.3.1 The important assumption				

		3.3.2	Convergence for a general smooth function	44		
		3.3.3	Convergence for a strongly convex function	46		
		3.3.4	Convergence for a weakly convex function	48		
3.4 Glob		Globa	lization techniques	49		
3.5 Other		Other	examples of descent methods	53		
		3.5.1	Quasi-Newton methods	53		
		3.5.2	Stochastic gradient methods	55		
		3.5.3	Coordinate descent methods	55		
4	Des	cent N	fethods with Momentum	56		
	4.1	Deriva	tion of descent methods with momentum	56		
		4.1.1	Motivation from differential equations	56		
		4.1.2	Derivation of the heavy-ball method	57		
		4.1.3	Nesterov's accelerated gradient method	58		
	4.2	Conve	rgence analysis of Nesterov's accelerated gradient method	59		
	4.3	The co	onjugate gradient method	63		
		4.3.1	Motivation	64		
		4.3.2	Construction of the CG method	65		
		4.3.3	Some properties of the CG method	66		
		4.3.4	Termination of the CG method	68		
		4.3.5	More efficient formulation	71		
		4.3.6	Computing $Q^{-1}$ with the CG method $\ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$	72		
$\mathbf{A}$	Very basics of mathematics 74					
	A.1	Basics	of linear algebra	74		
	A.2	Basics	of differential calculus	77		
	A.3		of topology	78		
Bi	Bibliography					

# Chapter 1

# Optimization problems: Definitions and Applications

# 1.1 Optimization in the real world

There are two contexts where optimization problems arise: Decision making, and model learning:

- In *decision making*, we are faced with a set of possible decisions, and we want to find the decision that minimizes some cost. This decision can be made based on the solution to some optimization problem.
- In *model learning*, a set of data is available, and we want to find a model that fits the data. This model can be found by solving an optimization problem.

In this course, we will mainly focus on the applications arising in model learning. In the rest of the section, we will introduce the basic idea of formulating a model learning problem as an optimization problem.

The typical optimization problem in data analysis is to find a model that agrees with some collected data but adheres to some structural constraints that reflect our beliefs about what a good model should be. The data set is typically a collection of inputs and outputs corresponding to different samples:

$$\mathcal{D} \coloneqq \{(a_1, y_1), \dots, (a_m, y_m)\},\tag{1.1}$$

where  $a_i$  is the input (also sometimes called *features*) and  $y_i$  is the output (also sometimes called *measurements*).

The goal is to find a model that predicts the output y given the input a. This typically takes the form of a function  $\varphi$  that maps inputs to outputs:

$$y \approx \varphi(a).$$
 (1.2)

To define the problem mathematically, we need to parameterize the possible model functions  $\varphi$  with some unknown parameters  $x \in \mathbb{R}^n$ . The fitting problem can then be formulated as finding some  $x \in \mathbb{R}^n$  such that for any input a, the input can be predicted by the model:

$$y \approx \varphi(a; x). \tag{1.3}$$

To formulate this as an optimization problem, we define a loss function  $\mathcal{L}_{\mathcal{D}}(x)$  that measures how well the model fits the data:

$$\mathcal{L}_{\mathcal{D}}(x) \coloneqq \frac{1}{m} \sum_{i=1}^{m} l\left(y_i, \varphi(a_i; x)\right), \qquad (1.4)$$

where  $l(y, \bar{y})$  represents some distance between the true output y and the predicted output  $\bar{y}$ .

The goal is to find the parameter x that best fits the data while meeting some prior assumptions on the model. This is typically done by solving the optimization problem:

$$\min_{x \in \mathbb{R}^n} \mathcal{L}_{\mathcal{D}}(x) + \lambda \ \operatorname{pen}(x) \tag{1.5}$$

where pen(x) is a penalty function that measures how well the model meets the constraints. The parameter  $\lambda \geq 0$  is the *regularization parameter*, a hyperparameter that controls the trade-off between fitting the data and meeting the constraints.

Depending on the nature of the labels  $y_j$ , the model-fitting task takes different names:

- When  $y_j$  are real numbers, or vectors of real numbers, the task is called a regression problem.
- When  $y_j$  are *labels*, i.e. integers in a set  $\{1, \ldots, q\}$ , the task is called a *classification problem*.

Later, we will look at several examples of regression and classification problems.

# **1.2** General definitions and notations

### **1.2.1** Formulation of optimization problems

### **Definition 1.1: Optimization Problems**

An optimization problem is mathematically formulated as follows:

$$\min_{x \in \mathcal{X}} f(x) \tag{1.6}$$

In (1.6), three "ingredients" are present:

- The decision variable  $x \in \mathbb{R}^n$  that can be chosen, and that may contain several components.
- The feasible set  $\mathcal{X}$  in which the decision variable x is imposed to be. Often, we will choose  $\mathcal{X} = \mathbb{R}^n$ . In this case, the optimization is qualified as an unconstrained problem.
- An objective function,  $f(x) : \mathcal{X} \to \mathbb{R}$ , that shall be minimized. Note that when a function  $\tilde{f}(x)$  shall be maximized, one can minimize the function  $f(x) \equiv -\tilde{f}(x)$ .

*Remark.* Sometimes, an alternative formulation might be found, that accounts more explicitly for the constraint  $x \in \mathcal{X}$ :

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\min i x \in \mathbb{R}^n} & f(x) \\ \text{subject to} & x \in \mathcal{X}. \end{array}$$
(1.7)

### 1.2.2 Minimizers

### **Definition 1.2: Global optimality**

The point  $x^* \in \mathbb{R}^n$  is called a "global minimizer" (often also called a "global minimum") when  $x^* \in \mathcal{X}$  and  $\forall x \in \mathcal{X} : f(x) \ge f(x^*)$ .

### **Definition 1.3: Strict optimality**

The point  $x^* \in \mathbb{R}^n$  is called a "strict global minimizer" when  $x^* \in \mathcal{X}$  and  $\forall x \in \mathcal{X} \setminus \{x^*\}: f(x) > f(x^*)$ .

### 1.2.3 Local optimality

### **Definition 1.4: Local optimality**

The point  $x^* \in \mathbb{R}^n$  is called a "local minimizer" when  $x^* \in \mathcal{X}$  and there exists a neighborhood of  $\mathcal{N}$  of  $x^*$  such that  $\forall x \in \mathcal{X} \cap \mathcal{N} : f(x) \geq f(x^*)$ .

Note that this neighborhood can be chosen to be in the form of an open ball:  $\mathcal{N} \coloneqq \{x \mid ||x - x^*|| < \varepsilon\}$  for some  $\varepsilon > 0$ .

### **Definition 1.5: Strict local optimality**

The point  $x^* \in \mathbb{R}^n$  is called a "strict local minimizer" when  $x^* \in \mathcal{X}$  and there exists a neighborhood  $\mathcal{N}$  of  $x^*$  so that  $\forall x \in (\mathcal{X} \cap \mathcal{N}) \setminus \{x^*\} : f(x) > f(x^*)$ .

*Remark.* Note that a global minimizer is also a local minimizer, but the converse is not necessarily true.

### Example 1.1: (Unconstrained) Quadratic Program

A Quadratic Program (QP) is an optimization problem where the objective function is quadratic. In the unconstrained case (i.e.  $\mathcal{X} = \mathbb{R}^n$ ), the problem is formulated as follows:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} x^T Q x - c^T x + r, \tag{1.8}$$

where  $Q \in \mathbb{R}^{n \times n}$  is a symmetric matrix,  $c \in \mathbb{R}^n$  is a vector, and  $r \in \mathbb{R}$  is a scalar.

*Remark.* The Hessian of the objective function of (1.8) is constant, and equal to Q:

$$\forall x \in \mathbb{R}^n \quad \nabla^2 f(x) = Q. \tag{1.9}$$

# **1.3** Examples of optimization problems in data analysis

### 1.3.1 Regression problems

The most common optimization problem in data analysis is the least squares problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \|y_j - \varphi(a_j; x)\|^2, \qquad (1.10)$$

where the outputs  $y_i \in \mathbb{R}^p$  are vectors, and the inputs  $a_i$  can be some matrices or vectors.

For a general function  $\varphi$ , the optimization problem (1.10) is often called "non-linear least squares". Since it is quite general, it is quite difficult to analyze: it might have multiple local minima, it might have no global minimum at all, etc.

There is, however, a special case where the analysis is way easier: the *linear least squares problem*. This corresponds to the case where the model  $\varphi$  is affine in the parameters x.

### Example 1.2: The linear least squares problem

When the model takes the form:  $\varphi(a_j; x) = A(a_j)x + b(a_j)$ , then the optimization problem (1.10) takes the following form:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \|\tilde{y}_j - A_j x\|^2 \eqqcolon f(x), \tag{1.11}$$

where  $\tilde{y}_j \coloneqq y_j - b(a_j)$  and  $A_j \coloneqq A(a_j)$ . As a small abuse of notation, we might write  $y_j$  instead of  $\tilde{y}_j$  in the next parts.

### Proposition 1.1: Linear least squares is a QP

The linear least squares problem (1.11) is a quadratic optimization problem in the form (1.8), with:

$$Q \coloneqq \frac{1}{m} \sum_{j=1}^{m} A_j^{\top} A_j,$$
$$c \coloneqq \frac{1}{m} \sum_{j=1}^{m} A_j^{\top} \tilde{y}_j,$$
$$r \coloneqq \frac{1}{2m} \sum_{j=1}^{m} \|\tilde{y}_j\|^2.$$

In the next sections, we will create new mathematical tools for the analysis of optimization problems. In particular, these tools will help us to study the solutions of the optimization problem (1.11).

### 1.3.2 Ridge regression

As we saw earlier, some regularization is often used in practice. There can be different reasons for that; one of them is called *overfitting*. The idea is that if there is a lot of degrees of freedom in the model, and yet not enough data points, the procedure might fits "too well" the data, and might not generalize well to new data points. An extreme case is when there are more parameters than data points, i.e., n > m.

To prevent this effect, one can add a regularization term to the optimization problem. This will force the model to stay "simple" while still fitting the data well.

When the regularization term is a penalty on the squared norm of the parameters, the optimization problem is called *ridge regression*. The optimization problem (1.11) becomes the following:

### Example 1.3: Ridge regression

The ridge regression problem is the following optimization problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \|\tilde{y}_j - A_j x\|^2 + \frac{\lambda}{2} \|x\|^2 \eqqcolon f(x), \tag{1.12}$$

where  $\lambda > 0$  is a hyperparameter.

Later, we will see that the ridge regression problem has the nice property that there always exists a unique local minimizer, which is also the unique global minimizer.

### Proposition 1.2: Ridge regression is a QP

The ridge regression problem (1.12) is a quadratic optimization problem in the form (1.8), with:

$$Q \coloneqq \lambda I_n + \frac{1}{m} \sum_{j=1}^m A_j^\top A_j,$$
  

$$c \coloneqq \frac{1}{m} \sum_{j=1}^m A_j^\top \tilde{y}_j,$$
  

$$r \coloneqq \frac{1}{2m} \sum_{j=1}^m \|\tilde{y}_j\|^2.$$
  
(1.13)

where  $I_n$  is the identity matrix of size n.

*Remark.* When  $\lambda > 0$  we have:

$$Q \succeq \lambda I_n \succ 0 \tag{1.14}$$

which implies that Q is non-singular.

## 1.3.3 LASSO regression

In the case one looks for a sparse solution to the linear least squares problem, one would need to penalize the number of non-zero entries in x. Since this would imply a non-continuous (hence

non-convex) penalization, the resulting problem would be extremely difficult to solve. Instead, one can use the LASSO regression, which penalizes the  $l_1$  norm of the parameters:

### Example 1.4: LASSO regression

The following optimization problem is called the LASSO regression:

$$\underset{x \in \mathbb{R}^{n}}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^{m} \|y_{j} - A_{j}x\|^{2} + \lambda \|x\|_{1}, \qquad (1.15)$$

*Remark.* In (1.15), we used the  $l_1$  norm of a vector, which is defined as the sum of the absolute values of its components:

$$\|x\|_{1} \coloneqq \sum_{i=1}^{n} |x_{i}|.$$
 (1.16)

Note that the optimization problem (1.15) is not a QP, but it still has a rather simple structure. With the tools from the next chapter (convexity), we will be able to characterize the solution-set of this problem.

### 1.3.4 Cross-entropy for classification tasks

Now that we have seen several examples of regression tasks, let us mention another common type of optimization problem in data analysis: classification tasks.

Here, the outputs  $y_i$  are discrete labels: integers in a set  $\{1, \ldots, q\}$ . Typically, the labels represent different classes to which the input  $a_i$  can belong. In this case, it is useful to define  $p^{y_j} \in \{0, 1\}^q$  as follows:

for 
$$j = 1, ..., m$$
 and  $l = 1, ..., q$ ,  $(p^{y_j})_l = \begin{cases} 1 & \text{if } y_j = l, \\ 0 & \text{otherwise.} \end{cases}$  (1.17)

Here, the vector  $p^{y_j} \in \{0,1\}^q$  represents the membership of the label  $y_j$  to each of the q classes. We will learn a model  $\varphi(a_j; x) \in [0,1]^q$  that approximates the vectors  $p^{y_j}$ . One could interpret  $\varphi_l(a_j; x)$  as the probability that the label  $y_j$  belongs to class l. In the previous section, we have seen that the squared norm error is a common loss function for regression tasks. For classification problems, instead of the least squares loss, one often uses the cross-entropy loss between the probability distribution associated with  $p^{y_j}$  and the one associated with  $\varphi(a_j; x)$ :

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{m} \sum_{j=1}^m L\left(p^{y_j}, \varphi(a_j; x)\right).$$
(1.18)

where the function  $L(p, \bar{p})$  is the cross-entropy between two distributions, which is defined as follows:

$$L(p,\bar{p}) \coloneqq -\sum_{l=1}^{q} p_l \log(\bar{p}_l).$$

$$(1.19)$$

While the simplest structure for  $\varphi$  was a linear one in the regression case, in the present case, we need  $\varphi(a_j; x) \in [0, 1]^q$ . Therefore, the standard choice uses the softmax function, as in the following example:

### Example 1.5: Logistic regression

The logistic regression is the optimization problem (1.18) with the following model for  $\varphi(a; x)$ :

$$\varphi(a;x) \coloneqq \operatorname{softmax}(A(a)x) \tag{1.20}$$

for some matrix  $A(a) \in \mathbb{R}^{q \times n}$  of  $\mathbb{R}^n$  that depend on the input a, and the function softmax is defined as follows:

$$\operatorname{softmax}(z)_{l} \coloneqq \frac{e^{z_{l}}}{\sum_{k=1}^{q} e^{z_{k}}}.$$
(1.21)

# 1.4 Basic properties of optimization problems

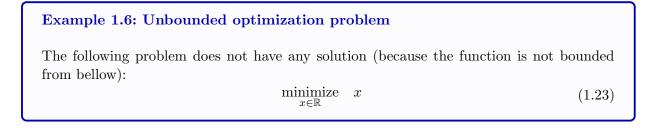
In this section, we will discuss some properties of optimization problems of the form:

$$\min_{x \in \mathcal{X}} \quad f(x) \tag{1.22}$$

### 1.4.1 Existence of solutions

Since we study the solutions of optimization problems, a natural question is whether such solutions exist.

First, let us keep in mind the following two examples where no solution exists.



### Example 1.7: Bounded optimization problem but with no solution

The following opitmization problem is buonded from bellow, but still, it does not have any solution:

 $\underset{x \in \mathbb{R}}{\text{minimize}} \quad e^{-x}$ 

(1.24)

In the following theorems, we will see some conditions that are easy to check that guarantee the existence of a solution.

Theorem 1.1: Existence of a minimizer for a compact feasible set

If the feasible set  $\mathcal{X} \subset \mathbb{R}^n$  is non-empty and compact (i.e., bounded and closed) and  $f : \mathcal{X} \to \mathbb{R}$  is continuous, then there exists a global minimizer to the optimization problem (1.22).

*Proof.* Let us write  $f^* := \inf_{x \in \mathcal{X}} f(x) \in \mathbb{R}^n \cup \{-\infty\}$ . Almost by definition of the infimum, there exists a sequence  $(x_k)_{k \in \mathbb{N}}$  in  $\mathcal{X}$  such that:

$$f(x_k) \xrightarrow[k \to +\infty]{} f^* \tag{1.25}$$

Since  $\mathcal{X}$  is compact, the sequence  $x_k$  has at least one accumulation point  $\bar{x} \in \mathcal{X}$ , i.e. for some increasing sequence of integers  $(k_j)_{j \in \mathbb{N}}$ , we have  $x_{k_j} \xrightarrow[j \to +\infty]{} \bar{x}$ .

By continuity of f, we have:

$$f(x_{k_j}) \xrightarrow[k_j \to +\infty]{} f(\bar{x}) \tag{1.26}$$

By combining (1.25) and (1.26), we obtain  $f(\bar{x}) = f^* = \inf_{x \in \mathcal{X}} f(x)$ . This proves that  $\bar{x}$  is a minimizer of f over  $\mathcal{X}$ .

### **Definition 1.6: Coercive functions**

A function  $f : \mathcal{X} \to \mathbb{R}$  is called *coercive* if there exists a function  $\kappa : \mathbb{R} \to \mathbb{R}$  such that  $\kappa(t) \xrightarrow[t \to +\infty]{} +\infty$  and such that  $\forall x \in \mathcal{X}, f(x) \ge \kappa(||x||)$ .

### Theorem 1.2: Existence of a minimizer for a coercive function

Let f be a continuous and coercive function and  $\mathcal{X}$  a closed and non-empty set. Then, there exists a global minimizer of the optimization problem (1.22). *Proof.* Since  $\mathcal{X}$  is non-empty, there exists some  $x_0 \in \mathcal{X}$ .

Using the fact that  $\kappa(t) \xrightarrow[t \to +\infty]{t \to +\infty} +\infty$ , there exists some  $c \in \mathbb{R}$  such that if t > c, then  $\kappa(t) > f(x_0) + 1$ . The set  $\tilde{\mathcal{X}}_c := \{x \in \mathcal{X} | \|x\| \le c\}$  is closed and bounded (hence it is compact). Using the previous theorem, there exists a vector  $x^*$  that minimizes f on  $\tilde{\mathcal{X}}_c$ . The following also holds:

$$\forall x \in \mathcal{X} \setminus \mathcal{X}_c \qquad f(x) > f(x_0) + 1 \ge f(x^*) + 1 > f(x^*)$$

Hence,  $x^*$  also minimizes f on  $\mathcal{X} \setminus \tilde{\mathcal{X}}_c$ . We can conclude that  $x^*$  is a minimizer of f over the whole set  $\mathcal{X}$ .

### 1.4.2 First order optimality conditions for smooth functions

In this part, we assume the function  $f : \mathcal{X} \to \mathbb{R}$  to be continuously differentiable.

### **Definition 1.7: Stationary points**

Let  $\bar{x}$  be a point in the interior of  $\mathcal{X}$ . We say that  $\bar{x}$  is a *stationary point* of the optimization problem (1.22) if it satisfies:

$$\nabla f\left(\bar{x}\right) = 0 \tag{1.27}$$

*Remark.* A more general definition exists for points that are not in the interior of the feasible set. In this course, we do not focus on the cases concerning the points at the border of the feasible set.

#### Theorem 1.3: First-order condition

Let  $\bar{x}$  be in the interior of the feasible set  $\mathcal{X}$ . If  $\bar{x}$  is a local minimizer of the optimization problem (1.22), then it is a stationary point:

$$\nabla f\left(\bar{x}\right) = 0 \tag{1.28}$$

*Proof.* Let p be a vector of  $\mathbb{R}^n$ .

Since  $\bar{x}$  is in the interior of  $\mathcal{X}$ ,  $\bar{x} + \alpha p \in \mathcal{X}$  for  $\alpha$  small enough. Furthermore, since  $\bar{x}$  is a local minimizer,  $f(\bar{x} + \alpha p) \geq f(\bar{x})$  for for  $\alpha$  small enough. This implies that for  $\alpha$  small enough:

$$0 \le \frac{f(\bar{x} + \alpha p) - f(\bar{x})}{\alpha} \xrightarrow[\alpha \to 0]{} \nabla f(\bar{x})^{\top} p \qquad (1.29)$$

Hence, we have  $\nabla f(\bar{x})^{\top} p \ge 0$  for all  $p \in \mathbb{R}^n$ . By choosing  $p = -\nabla f(\bar{x})$ , we obtain  $- \|\nabla f(\bar{x})\|^2 \ge 0$ , which implies  $\nabla f(\bar{x}) = 0$ .

### 1.4.3 Second order optimality conditions for smooth functions

In this part, we now assume the function  $f : \mathcal{X} \to \mathbb{R}$  to be *twice* continuously differentiable.

Theorem 1.4: Second-order necessary conditions for unconstrained optimization

Let  $\bar{x}$  be in the interior of the feasible set  $\mathcal{X}$ . If  $\bar{x}$  is a local minimizer of the optimization problem (1.22), then not only it satisfies  $\nabla f(\bar{x}) = 0$  but also:

$$\nabla^2 f(\bar{x}) \succeq 0 \tag{1.30}$$

*Proof.* Let  $\bar{x}$  be a local minimizer.

First,  $\bar{x}$  is a stationary point from the previous theorem. Let p be an element of  $\mathbb{R}^n$ . Using the same argument as in the proof above,  $f(\bar{x} + \alpha p) \ge f(\bar{x})$  for  $\alpha$  small enough. Let us now use the first-order Taylor expansion of f at  $\bar{x}$ :

$$f(\bar{x} + \alpha p) = f(\bar{x}) + \nabla f(\bar{x})^{\top} (\alpha p) + \int_0^1 s(\alpha p)^{\top} \nabla^2 f(\bar{x} + s\alpha p)(\alpha p) ds$$
  
=  $f(\bar{x}) + \alpha^2 \int_0^1 s p^{\top} \nabla^2 f(\bar{x} + s\alpha p) p ds$  (1.31)

After a rearrangement, we obtain:

$$0 \leq \frac{f(\bar{x} + \alpha p) - f(\bar{x})}{\alpha^2} = \int_0^1 s p^\top \nabla^2 f(\bar{x} + s\alpha p) \, p \mathrm{d}s \xrightarrow[\alpha \to 0]{} \int_0^1 s p^\top \nabla^2 f(\bar{x}) \, p \mathrm{d}s = \frac{1}{2} p^$$

This proves  $p^{\top} \nabla^2 f(\bar{x}) p \ge 0$ . Since this is true for any vector  $p \in \mathbb{R}^n$ , we have  $\nabla^2 f(\bar{x}) \ge 0$ .  $\Box$ 

### Theorem 1.5: Second-order sufficient conditions for unconstrained optimization

Let  $\bar{x}$  be in the interior of the feasible set  $\mathcal{X}$ . If  $\bar{x}$  is a stationary point of the optimization problem (1.22), and in addition,  $\bar{x}$  satisfies the following condition:

$$\nabla^2 f\left(\bar{x}\right) \succ 0,\tag{1.33}$$

then x is a strict local minimizer of the optimization problem (1.22).

*Proof.* Let  $\bar{x}$  be a stationary point that satisfies the condition (1.33). Using (1.33), there exists  $\lambda > 0$  such that  $\nabla^2 f(\bar{x}) \succeq \lambda I_n$ .

The function  $\nabla^2 f(x)$  is continuous. This implies that for any  $\delta > 0$ , there exists a neighborhood  $\mathcal{N} \subset \mathcal{X}$  of  $\bar{x}$  such that:

$$\forall x \in \mathcal{N}, \quad -\delta I_n \preccurlyeq \nabla^2 f(x) - \nabla^2 f(\bar{x}) \preccurlyeq \delta I_n \tag{1.34}$$

In particular, choosing  $\delta < \lambda$ , we have:

$$\forall x \in \mathcal{N}, \quad \nabla^2 f(x) \succ 0 \tag{1.35}$$

Let x' be an element of  $\mathcal{N} \setminus \{\bar{x}\}$ . Let us define  $p := x - \bar{x} \neq 0$ . Now let us use the first-order Taylor expansion of f at  $\bar{x}$ :

$$f(x') = f(\bar{x} + p) = f(\bar{x}) + \nabla f(\bar{x})^{\top} p + \int_{0}^{1} sp^{\top} \nabla^{2} f(\bar{x} + sp) p ds$$
  
=  $f(\bar{x}) + \int_{0}^{1} sp^{\top} \nabla^{2} f(\bar{x} + sp) p ds$  (1.36)

Furthermore, for all  $s \in [0, 1]$ ,  $\bar{x} + sp \in \mathcal{N}$  (assuming that  $\mathcal{N}$  is convex, which is the case if  $\mathcal{N}$  is a ball for example). From (1.35), this implies  $\nabla^2 f(\bar{x} + sp) \succ 0$ , and therefore:  $p^\top \nabla^2 f(\bar{x} + sp) p > 0$ . Injecting this into (1.36), we obtain:

$$f(x') > f(\bar{x}) \tag{1.37}$$

This is true for any  $x' \in \mathcal{N} \setminus \{\bar{x}\}$ , where  $\mathcal{N}$  is a neighborhood of  $\bar{x}$  in  $\mathcal{X}$ . Hence,  $\bar{x}$  is a strict local minimum of the optimization problem (1.22).

15

### Example 1.8: Illustrative example

For p = 1, 2, 3, 4, consider the following optimization problem:

$$\underset{x \in \mathbb{R}}{\text{minimize}} \quad x^p \tag{1.38}$$

Let us detail the status of the point x = 0 for each value of p:

- If p = 1, the condition (1.28) is not satisfied, hence x can not be a (local) minimizer.
- If p = 2, the conditions (1.28) and (1.33) are satisfied, hence x is a local minimizer. In this specific case, it is also the unique global minimizer.
- If p = 3, the conditions (1.28) and (1.30) are satisfied, but not the conditions (1.33), hence, the second order conditions do not allow us to conclude. In this specific case, x is not a (local) minimizer.
- If p = 4, the conditions (1.28) and (1.30) are satisfied, but not the conditions (1.33), hence, the second order conditions do not allow us to conclude. In this specific case, x is the unique local, and even global minimizer.

# **1.5** Application of the properties to the regression examples

In this section, we will apply the properties that were proven in Section 1.4 to the regression problems that were introduced in Section 1.3.

# **1.5.1** Quadratic programs with $Q \succ 0$

In this part, we consider the loss of the quadratic program (1.8):

$$f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x + r.$$
 (1.39)

In particular, we focus on the case where Q is positive definite:  $Q \succ 0$ .

*Remark.* As noted before, the loss of the Ridge regression problem (1.12) with  $\lambda > 0$  falls into this category.

*Remark.* Some of the linear least square problems (1.11) also fall into this category, under the assumption that  $\frac{1}{m} \sum_{j=1}^{m} A_j^{\top} A_j \succ 0$ .

Proposition 1.3: Existence of a global minimizer

The function f(x) is coercive (see Definition 1.6). Hence, using Theorem 1.2, it has at least one global minimizer in  $\mathbb{R}^n$ .

*Proof.* Define  $\lambda_{\min}$  as the lowest eigenvalue of Q. Then we have:  $Q - \lambda_{\min} I_n \succeq 0$  (where  $I_n$  is the identity matrix of size n). This implies that for all  $x \in \mathbb{R}^n$ :  $x^\top Q x \ge \lambda_{\min} ||x||^2$ . This implies that:

$$f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x + r$$
$$\geq \frac{\lambda_{\min}}{2} \|x\|^2 - \|c\| \|x\| + r$$
$$=: \kappa(\|x\|)$$

where  $\kappa(t) := \frac{\lambda_{\min}}{2} t^2 - \|c\| t + r \xrightarrow[t \to +\infty]{} +\infty.$ 

By definition, this implies that f(x) is coercive.

### Proposition 1.4: Unique stationary point

The function f(x) has a unique stationary point, given by the following formula:

$$x^{\star} = Q^{-1}c \tag{1.40}$$

*Proof.* Since the problem is unconstrained, i.e.  $\mathcal{X} = \mathbb{R}^n$ , the stationary points all verify the formula

$$\nabla f(x) = Qx - c = 0 \tag{1.41}$$

Since Q is positive definite, it is also invertible. Hence, the equation (1.41) has a unique solution, given by (1.40).  $\Box$ 

### Proposition 1.5: Unique minimizer

The stationary point  $x^*$  is the unique global minimizer of the function f(x). It is also the unique local minimizer.

*Proof.* All global (resp. local) minimizers are stationary points (cf. Theorem 1.3). Using Proposition 1.4, there exists not more than one stationary point, given by  $x^*$ . This implies that there exists no more than one global (resp.) minimizer.

On the other hand, using Proposition 1.3, there exists at least one global minimizer (which is also a local minimizer).

Combining these two facts,  $x^*$  is the unique global minimizer, and the unique local minimizer.  $\Box$ 

### **1.5.2** Quadratic Programs with $Q \succeq 0$

In the case where Q is positive semi-definite, the function f(x) is no longer coercive. It could even be that no global minimizer exists. For example, this is the case if Q = 0 and  $c \neq 0$ .

However, we still have some interesting results. Unfortunately, we cannot prove them using only the properties from Section 1.4. We still derive them for completeness.

### Proposition 1.6: Stationary points are global minimizers for QP

Assume that there exists some stationary point  $x^*$  of  $f(x) = \frac{1}{2}x^\top Qx - c^\top x + r$ (with  $Q \succ 0$ ) over  $\mathbb{R}^n$ . Then  $x^*$  is a global minimizer of the function f(x).

*Proof.* Since  $\mathbb{R}^n$  is an open-set, the stationary point  $x^*$  verifies  $\nabla f(x^*) = Qx^* - c = 0$ . Let  $x \in \mathbb{R}^n$ .

We want to prove that  $f(x) \ge f(x^*)$ . Let us define  $z := x - x^*$ . Then the following holds:

$$f(x) = f(z + x^{\star})$$
  
=  $\frac{1}{2}(x^{\star} + z)^{\top}Q(x^{\star} + z) - c^{\top}(x^{\star} + z) + r$   
=  $\frac{1}{2}z^{\top}Qz + \frac{1}{2}(z^{\top}Qx^{\star} + x^{\star}Qz) - c^{\top}z + (\frac{1}{2}x^{\star}Qx^{\star} - c^{\top}x^{\star} + r)$   
=  $\frac{1}{2}z^{\top}Qz + x^{\star}Qz - c^{\top}z + f(x^{\star})$   
=  $\frac{1}{2}z^{\top}Qz + (Qx^{\star} - c)^{\top}z + f(x^{\star})$   
=  $\frac{1}{2}z^{\top}Qz + f(x^{\star})$   
 $\geq f(x^{\star})$ 

*Remark.* In the next chapter, we will generalize the result from Proposition 1.6 to a broader class of optimization problems: *the convex optimization problems.* 

Proposition 1.7: Lower-bounded quadratic programs

Assume that  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x + r$  is lower-bounded over  $\mathbb{R}^n$ , i.e.  $\min_{x \in \mathbb{R}^n} f(x) > -\infty$ . Then the function f(x) has at least one stationary point over  $\mathbb{R}^n$ .

*Proof.* Let  $z \in \text{Ker}(Q)$ , i.e. Qz = 0. Then we have  $f(tz) = r - t(c^{\top}z)$ . Since f is lower bounded over  $\mathbb{R}^n$ , the function  $t \mapsto f(tz)$  is also lower bounded. This implies that  $c^{\top}z = 0$ . Since this is true for any  $z \in \text{Ker}(Q)$ ,  $\text{Ker}(Q) \subset \text{Ker}(c^{\top})$ .

Now, let us use some linear algebra results (which can be found in Appendix A.1):

$$c \in \operatorname{Im}(c) = \operatorname{Ker}(c^{\top})^{\perp} \subset \operatorname{Ker}(Q)^{\perp} = \operatorname{Im}(Q^{\top}) = \operatorname{Im}(Q).$$

This implies that  $c \in \text{Im}(Q)$ , i.e. there exists  $x^* \in \mathbb{R}^n$  such that  $Qx^* = c$ . Hence,  $\nabla f(x^*) = Qx^* - c = 0$ , i.e.  $x^*$  is a stationary point.

### Corollary 1.1

Lower-bounded QPs have at least one global minimizer. Furthermore, these are all the points that satisfy the equation  $Qx^* = c$ .

*Proof.* This is a direct application of the two previous theorems.

*Remark.* The linear least squares problem (1.11) falls into that category.

18

### 1.5.3 LASSO regression: existence of minimizers

We recall that the LASSO regression problem is associated with the following loss function:

$$f(x) \coloneqq \frac{1}{2m} \sum_{j=1}^{m} \|\tilde{y}_j - A_j x\|^2 + \lambda \|x\|_1.$$

Proposition 1.8: Existence of minimizers for LASSO regression

Assume that  $\lambda > 0$ . Then, the loss function of the LASSO regression problem is coercive (see Definition 1.6). Hence, using Theorem 1.2, the optimization problem (1.15) has at least one solution

*Proof.* We have:

$$f(x) \ge \lambda \|x\|_1 \ge \underbrace{\lambda \|x\|_2}_{=:\kappa(\|x\|_2)} \xrightarrow[\|x\|_2 \to +\infty]{} +\infty$$

*Remark.*  $||x||_1 \ge ||x||_2$  because:  $||x||_1^2 = \left(\sum_{i=1}^n |x_i|\right)^2 = \sum_{\substack{i=1\\ = ||x||_2^2}}^n |x_i|^2 + \sum_{\substack{i\neq j\\ \ge 0}} |x_i| |x_j| \ge ||x||_2^2$ 

# Chapter 2

# Convexity

# 2.1 Convex sets

In this section, we define what convexity means for sets, and discuss the properties of convex sets.

**Definition 2.1: Convex Set** 

A set  $\mathcal{X} \subset \mathbb{R}^n$  is convex if

$$\forall x, y \in \mathcal{X}, \alpha \in [0, 1]: (1 - \alpha)x + \alpha y \in \mathcal{X}.$$
(2.1)

*Remark.* Intuitively, one could translate this definition into "all connecting lines lie inside the set."

### Proposition 2.1: Intersection of convex sets

If S is a set of convex sets, then their intersection  $\bigcap_{\mathcal{X}\in S} \mathcal{X}$  is also convex.

*Proof.* Let  $x, y \in \bigcap_{\mathcal{X} \in S} \mathcal{X}$  and  $\alpha \in [0, 1]$ . For all  $\mathcal{X} \in S$ , we have  $x, y \in \mathcal{X}$ . Since  $\mathcal{X}$  is convex,  $(1 - \alpha)x + \alpha y \in \mathcal{X}$ . This holds for any  $\mathcal{X} \in S$ , hence:

$$(1-\alpha)x + \alpha y \in \bigcap_{\mathcal{X} \in S} \mathcal{X}.$$
(2.2)

This holds for all  $x, y \in \bigcap_{\mathcal{X} \in S} \mathcal{X}$  and  $\alpha \in [0, 1]$ . Therefore,  $\bigcap_{\mathcal{X} \in S} \mathcal{X}$  is convex.  $\Box$ 

*Remark.* In particular, if  $\mathcal{X}_1$  and  $\mathcal{X}_2$  are convex sets, then their intersection  $\mathcal{X}_1 \cap \mathcal{X}_2$  is also convex. *Remark.* The union of convex sets is not necessarily convex.

### Proposition 2.2: Cartesian product of convex sets

Let  $\mathcal{X}_1$  and  $\mathcal{X}_2$  be convex sets. Then their Cartesian product  $\mathcal{X}_1 \times \mathcal{X}_2 := \{(x_1, x_2) \text{ such that } x_1 \in \mathcal{X}_1 \text{ and } x_2 \in \mathcal{X}_2\}$  is also convex.

*Proof.* Let  $x, y \in \mathcal{X}_1 \times \mathcal{X}_2$ , and  $\alpha \in [0, 1]$ . Then  $x = (x_1, x_2)$  and  $y = (y_1, y_2)$ , where  $x_1, y_1 \in \mathcal{X}_1$  and  $x_2, y_2 \in \mathcal{X}_2$ . This implies that:

This implies that:

$$(1-\alpha)x + \alpha y = \left(\underbrace{(1-\alpha)x_1 + \alpha y_1}_{\in \mathcal{X}_1}, \underbrace{(1-\alpha)x_2 + \alpha y_2}_{\in \mathcal{X}_2}\right) \in \mathcal{X}_1 \times \mathcal{X}_2.$$
(2.3)

### Proposition 2.3: Affine transformation on convex sets

Let  $\mathcal{X} \in \mathbb{R}^n$  be a convex set. Let  $A \in \mathbb{R}^{m \times n}$  be a matrix, and  $b \in \mathbb{R}^m$  be a vector. Then the set  $A\mathcal{X} + b = \{Ax + b \text{ for } x \in \mathcal{X}\}$  is convex.

*Proof.* Left as an exercise.

# 2.2 Convex functions

### 2.2.1 General case

### **Definition 2.2: Convex Function**

A function  $f : \mathcal{X} \to \mathbb{R}$  is convex, if  $\mathcal{X}$  is convex and if

$$\forall x, y \in \mathcal{X}, \alpha \in [0, 1]: f\left((1 - \alpha)x + \alpha y\right) \le (1 - \alpha)f(x) + \alpha f(y) \tag{2.4}$$

*Remark.* In words, a function is convex when all secants are above the graph.

### Proposition 2.4: Convex over affine is convex

Assume that  $f : \mathbb{R}^n \to \mathbb{R}$  is convex. Then for any  $A \in \mathbb{R}^{n \times m}$  and  $b \in \mathbb{R}^m$ , the function  $g : \mathbb{R}^m \to \mathbb{R}$  defined by g(x) = f(Ax + b) is also convex.

*Proof.* Left as an exercise.

22

### Proposition 2.5: Increasing affine function over convex is convex

Assume that  $f : \mathbb{R}^n \to \mathbb{R}$  is convex. Let a > 0 and c be real numbers. Then, the function  $g : \mathbb{R} \to \mathbb{R}$  defined by g(x) = af(x) + c is also convex.

*Proof.* Left as an exercise.

### Proposition 2.6: Sum of convex functions is convex

Assume that  $f : \mathbb{R}^n \to \mathbb{R}$  and  $g : \mathbb{R}^n \to \mathbb{R}$  are convex. Then the function h(x) = f(x) + g(x) is also convex.

*Proof.* Left as an exercise.

### Proposition 2.7: Characterization of convex functions with epigraph

A function  $f : \mathcal{X} \to \mathbb{R}$  is convex if and only if its epigraph, i.e., the set  $\{(x,s) \in \mathcal{X} \times \mathbb{R} | x \in \mathcal{X}, s \ge f(x)\}$ , is a convex set.

*Proof.* Left as an exercise.

### **Definition 2.3: Sublevel sets**

The set  $\{x \in \mathbb{R}^n | f(x) \le c\}$  is the "sublevel set" of f for the value c.

Proposition 2.8: Convexity of sublevel Sets

The sublevel sets of a convex function  $f : \mathcal{X} \to \mathbb{R}$  are convex.

*Proof.* If  $f(x) \leq c$  and  $f(y) \leq c$  then for any  $\alpha \in [0, 1]$  it holds also

$$f((1-\alpha)x + \alpha y) \le (1-\alpha)f(x) + \alpha f(y) \le (1-\alpha)c + \alpha c = c$$

### 2.2.2 Convexity of smooth functions

# **Proposition 2.9: Convexity for** $C^1$ **Functions**

Assume that  $f : \mathcal{X} \to \mathbb{R}$  is continuously differentiable and  $\mathcal{X}$  is convex. Then the following properties are equivalent

f is convex (2.5a)

$$\forall x, y \in \mathcal{X}, \quad f(x) + \nabla f(x)^{\top} (y - x) \le f(y)$$
(2.5b)

$$\forall x, y \in \mathcal{X}, \quad (\nabla f(y) - \nabla f(x))^{\top} (y - x) \ge 0$$
(2.5c)

*Remark.* The property (2.5b) means that the graph of f is above its tangents. *Remark.* The property (2.5c) is a multi-dimensional equivalent of saying " $\nabla f(x)$  is a non-decreasing function".

*Proof.* We recall that by definition, (2.5a) is equivalent to:

$$\forall x, y \in \mathcal{X}, \alpha \in [0, 1]: \quad f\left((1 - \alpha)x + \alpha y\right) \le (1 - \alpha)f(x) + \alpha f(y) \tag{2.6}$$

In order to prove the equivalences  $(2.5a) \iff (2.5b) \iff (2.5c)$ , we will show the following chain of implications:

$$((2.5a) \iff) (2.6) \implies (2.5b) \implies (2.5c) \implies (2.6) (\iff (2.5a))$$

•  $(2.6) \implies (2.5b):$ 

A rearrangement of (2.6) gives:

$$f(y) - f(x) \ge \frac{f(x + \alpha(y - x)) - f(x)}{\alpha} \xrightarrow[\alpha \to 0]{} \nabla f(x)^{\top}(y - x)$$
(2.7)

which proves (2.5b).

•  $(2.5b) \implies (2.5c)$ :

Let x, y be in  $\mathcal{X}$ . Consider both equation (2.5b), and the one where we swap x and y:

$$f(x) + \nabla f(x)^{\top} (y - x) \le f(y)$$
  
$$f(y) + \nabla f(y)^{\top} (x - y) \le f(x)$$

Let us add these two inequalities:

$$f(y) + f(x) + (\nabla f(x) - \nabla f(y))^{\top} (y - x) \le f(y) + f(x)$$

Now, subsract f(x) + f(y) from both sides:

$$\left(\nabla f(x) - \nabla f(y)\right)^{\top} (y - x) \le 0$$

which proves (2.5c) after multiplication by -1.

•  $(2.5c) \implies (2.6):$ 

Let x, y be in  $\mathcal{X}$  and  $\alpha \in [0, 1]$ . The property (2.4) being trivial for  $\alpha = 0, 1$ , we will assume  $\alpha \in (0, 1)$ . Let us define the function  $g(t) \coloneqq f(x + t(y - x))$ . We have g(0) = f(x) and g(1) = f(y).

Furthermore, from (2.5c), we have that if  $t_1 > t_2$ , then  $g'(t_1) \ge g'(t_2)$ :

$$g'(t_1) - g'(t_2) = \nabla f(x + t_1(y - x))^\top (y - x) - \nabla f(x + t_2(y - x))^\top (y - x)$$
  
=  $\frac{1}{t_1 - t_2} (\nabla f(x_1) - \nabla f(x_2))^\top (x_1 - x_2)$   
\ge 0

with  $x_1 = x + t_1(y - x)$  and  $x_2 = x + t_2(y - x)$ .

Now let us prove (2.4):

$$g(\alpha) = g(0) + \int_0^{\alpha} g'(t) dt$$
  
=  $g(0) + \alpha \int_0^1 g'(s\alpha) ds$   
=  $g(0) + \alpha \left(g(1) - g(0) - \int_0^1 g'(s) ds\right) + \alpha \int_0^1 g'(s\alpha) ds$   
=  $(1 - \alpha)g(0) + \alpha g(1) - \alpha \int_0^1 \underbrace{\left(g'(\alpha) - g'(s\alpha)\right)}_{\geq 0} ds$   
 $\leq (1 - \alpha)g(0) + \alpha g(1)$ 

Finally, remarking that  $g(\alpha) = f((1 - \alpha)x + \alpha y)$ , g(0) = f(x) and g(1) = f(y), we have proved (2.6). This concludes the proof that f is convex.

# **Theorem 2.1: Convexity for** $C^2$ Functions

Assume that  $f : \mathcal{X} \to \mathbb{R}$  is twice continuously differentiable and  $\mathcal{X}$  convex. Then f is convex if and only if the following holds:

$$\forall x, y \in \mathcal{X}: \quad (x-y)^{\top} \nabla^2 f(x)(x-y) \ge 0.$$
(2.8)

*Proof.* Let us use the following formula:

$$\nabla f(y) - \nabla f(x) = \int_0^1 \nabla^2 f(x + s(y - x))(y - x) \,\mathrm{d}s$$
(2.9)

If we combine (2.9) with the equivalence (2.5a)  $\iff$  (2.5c) from Proposition 2.9, we obtain that f is convex if and only if:

$$\forall x, y \in \mathcal{X}, \quad \int_0^1 (y-x)^\top \nabla^2 f\left(x + s(y-x)\right) \left(y-x\right) \mathrm{d}s \ge 0 \tag{2.10}$$

Now let us show that (2.10) is equivalent to (2.8).

• (2.10)  $\implies$  (2.8): Assume that (2.10) holds. Let  $x, y \in \mathcal{X}$ . Apply (2.10) to x' = x and  $y' = x + \alpha(y - x)$ :

$$0 \leq \int_0^1 (y' - x')^\top \nabla^2 f(x' + s(y' - x'))(y' - x') ds$$
$$= \int_0^1 (\alpha(y - x))^\top \nabla^2 f(x + s\alpha(y - x))(\alpha(y - x)) ds$$

Then divide by  $\alpha^2$  and take the limit  $\alpha \to 0$ :

$$0 \le \int_0^1 \left(\alpha(y-x)\right)^\top \nabla^2 f\left(x + s\alpha(y-x)\right) \left(\alpha(y-x)\right) \,\mathrm{d}s \xrightarrow[\alpha \to 0]{} (y-x)^\top \nabla^2 f(x)(y-x)$$

This proves (2.8).

• (2.8)  $\implies$  (2.10): Assume that (2.8) holds. Let  $x, y \in \mathcal{X}$ . Apply (2.8) to x' = x + s(y-x) and y' = y + s(x-y):

$$0 \le (y' - x')^{\top} \nabla^2 f(x') (y' - x') = (1 - s)^2 (y - x)^{\top} \nabla^2 f(x + s(y - x)) (y - x)$$

By dividing by  $(1-s)^2$  and integrating with respect to s from 0 to 1, we find implies (2.10).

*Remark.* For the points x in the interior of  $\mathcal{X}$ , equation (2.8) is equivalent to:

$$\nabla^2 f(x) \succcurlyeq 0 \tag{2.11}$$

### **Example 2.1: Quadratic Function**

The function  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x + r$  is convex if and only if  $Q \succeq 0$ , because  $\forall x \in \mathbb{R}^n : \nabla^2 f(x) = Q$ .

### 2.2.3 Strictly convex Functions

### **Definition 2.4: Strict Convexity**

A function  $f : \mathcal{X} \to \mathbb{R}$  is said to be *strictly convex* if:

 $\forall x, y \in \mathcal{X} \text{ such that } x \neq y, \forall \alpha \in (0, 1): \quad f\left((1 - \alpha)x + \alpha y\right) < (1 - \alpha)f(x) + \alpha f(y).$ (2.12)

**Proposition 2.10: Strict convexity for**  $C^1$  **Functions** 

Assume that  $f : \mathcal{X} \to \mathbb{R}$  is continuously differentiable and  $\mathcal{X}$  is convex. Then the following properties are equivalent

f is strictly convex (2.13a)

 $\forall x, y \in \mathcal{X}, \text{ such that } x \neq y, \quad f(x) + \nabla f(x)^{\top} (y - x) > f(y)$  (2.13b)

 $\forall x, y \in \mathcal{X}, \text{ such that } x \neq y, \quad (\nabla f(y) - \nabla f(x))^{\top} (y - x) > 0$  (2.13c)

*Proof.* Simply take the proof of Proposition 2.9 and replace the inequalities by strict inequalities.

### Theorem 2.2: Strict convexity of smooth functions

Let f be a twice continuously differentiable function on a convex set  $\mathcal{X}$ . Assume that the following holds:

$$\forall x \in \mathcal{X} : \quad \nabla^2 f(x) \succ 0. \tag{2.14}$$

Then f is strictly convex.

*Proof.* Assume that (2.14) holds. Let  $x, y \in \mathcal{X}$ , such that  $x \neq y$ . Using (2.9) again, we have:

$$\left(\nabla f(y) - \nabla f(x)\right)^{\top} (y - x) = \int_0^1 (y - x) \nabla^2 f(x + s(y - x))(y - x) \,\mathrm{d}s > 0 \tag{2.15}$$

This allows us to conclude that f is strictly convex using the equivalence (2.13a)  $\iff$  (2.13c) from Proposition 2.10.

*Remark.* The converse is not necessarily true. For example, the function  $f(x) = x^4$  is strictly convex, but  $\nabla^2 f(x) = 13x^2$  is zero for x = 0.

### Example 2.2: Strongly Convex Quadratic

The quadratic function  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x + r$  is strictly convex if and only if  $Q \succ 0$ .

## 2.2.4 Strong convexity

### **Definition 2.5: Strongly convex function**

Let  $\mu > 0$  be a positive scalar. We say that f is  $\mu$ -strongly convex when the function  $f(x) - \frac{\mu}{2} ||x||^2$  is convex.

### **Proposition 2.11: Characterization of** $\mu$ -strongly convex functions

A function f being  $\mu$ -strongly convex is equivalent to each of the following properties (when f is sufficiently differentiable):

$$\forall x, y \in \mathbb{R}^{n}, \alpha \in [0, 1], \quad f((1 - \alpha)x + \alpha y) \le (1 - \alpha)g(x) + \alpha g(y) - \frac{\mu}{2}\alpha(1 - \alpha) \|y - x\|^{2}$$
(2.16a)

$$\forall x, y \in \mathcal{X}, \quad f(y) \ge f(x) + \nabla f(x)^{\top} (y - x) + \frac{\mu}{2} ||y - x||^2$$
 (2.16b)

$$\forall x, y \in \mathcal{X}, \quad (\nabla f(x) - \nabla f(y))^{\top} (x - y) \ge \mu ||x - y||^2$$
(2.16c)

$$\forall x, y \in \mathcal{X}, \quad (x-y)^{\top} \nabla^2 f(x) (x-y) \ge \mu \left\| x - y \right\|^2$$
(2.16d)

*Proof.* Using the characterizations (2.4), (2.5b), (2.5c), and (2.8) to the function  $h(x) \coloneqq f(x) - \frac{\mu}{2} ||x||^2$ , we find that f being  $\mu$ -strongly convex is equivalent to each of the following properties (when f is sufficiently differentiable):

$$\forall x, y \in \mathbb{R}^n, \alpha \in [0, 1], \quad h\left((1 - \alpha)x + \alpha y\right) \le (1 - \alpha)h(x) + \alpha h(y) \tag{2.17a}$$

$$\forall x, y \in \mathcal{X}, \quad h(y) \ge h(x) + \nabla h(x)^{\top} (y - x)$$
(2.17b)

$$\forall x, y \in \mathcal{X}, \quad (\nabla h(x) - \nabla h(y))^{\top} (x - y) \ge 0$$
(2.17c)

$$\forall x, y \in \mathcal{X}, \quad \forall x, y \in \mathcal{X}, \quad (x - y)^{\top} \nabla^2 h(x)(x - y) \ge 0$$
(2.17d)

Replacing h(x) with  $f(x) - \frac{\mu}{2} ||x||^2$ , we find that (2.17a), (2.17b), (2.17c) and (2.17d) are respectively equivalent to the following:

$$\forall x, y \in \mathbb{R}^{n}, \alpha \in [0, 1],$$

$$f((1 - \alpha)x + \alpha y)^{2} \leq (1 - \alpha)f(x) + \alpha f(y) - \frac{\mu}{2} \left((1 - \alpha) ||x||^{2} + \alpha ||y||^{2}\right)$$
(2.18a)

$$f\left((1-\alpha)x + \alpha y\right) - \frac{1}{2} \left\| (1-\alpha)x + \alpha y \right\| \le (1-\alpha)f(x) + \alpha f(y) - \frac{1}{2} \left( (1-\alpha) \|x\| + \alpha \|y\| \right)$$
  
$$\forall x, y \in \mathcal{X}, \quad f(y) - \frac{\mu}{2} \|y\|^2 \ge f(x) + \nabla f(x)^\top (y-x) - \frac{\mu}{2} \left( \|x\|^2 + \left(\nabla \|x\|^2\right)^\top (y-x) \right)$$
(2.18b)

$$\forall x, y \in \mathcal{X}, \quad (\nabla f(x) - \nabla f(y))^{\top} (x - y) - \frac{\mu}{2} \left( \nabla \|x\|^2 - \nabla \|y\|^2 \right)^{\top} (x - y) \ge 0$$
(2.18c)

$$\forall x, y \in \mathcal{X}, \quad \forall x, y \in \mathcal{X}, \quad (x-y)^{\top} \nabla^2 f(x) (x-y) - \frac{\mu}{2} (x-y)^{\top} \left( \nabla^2 \left\| x \right\|^2 \right) (x-y) \ge 0 \quad (2.18d)$$

After some rearrangements, these properties can be rewritten as:

$$\forall x, y \in \mathbb{R}^n, \alpha \in [0, 1], \tag{2.19a}$$

$$f((1-\alpha)x + \alpha y) \le (1-\alpha)g(x) + \alpha g(y) - \frac{\mu}{2} \left( \underbrace{\|(1-\alpha)x + \alpha y\|^2 - (1-\alpha)\|x\|^2 - \alpha\|y\|^2}_{=\alpha(1-\alpha)\|y-x\|^2} \right)$$

$$\forall x, y \in \mathcal{X}, \quad f(y) \ge f(x) + \nabla f(x)^{\top} (y - x) + \frac{\mu}{2} \left( \underbrace{\|y\|^2 - \left(\|x\|^2 + \left(\nabla \|x\|^2\right)^{\top} (y - x)\right)}_{=\|x - y\|^2} \right) \quad (2.19b)$$

$$\forall x, y \in \mathcal{X}, \quad (\nabla f(x) - \nabla f(y))^{\top} (x - y) \ge \frac{\mu}{2} \underbrace{\left( \nabla \|x\|^2 - \nabla \|y\|^2 \right)^{\top} (x - y)}_{=2\|x - y\|^2}$$
(2.19c)

$$\forall x, y \in \mathcal{X}, \quad \forall x, y \in \mathcal{X}, \quad (x-y)^{\top} \nabla^2 f(x)(x-y) \ge \frac{\mu}{2} \underbrace{(x-y)^{\top} \left(\nabla^2 \|x\|^2\right)(x-y)}_{=2\|x-y\|^2} \tag{2.19d}$$

The properties (2.19a), (2.19b), (2.19c), and (2.19d) are respectively equivalent to (2.16a), (2.16b), (2.16c), and (2.16d).

This concludes the proof that f being  $\mu$ -strongly convex is equivalent to each of the properties (2.16a), (2.16b), (2.16c), and (2.16d).

*Remark.* For the points x in the interior of  $\mathcal{X}$ , equation (2.16d) is equivalent to:

$$\nabla^2 f(x) \succcurlyeq \mu I_n \tag{2.20}$$

### Proposition 2.12: Strong convexity implies strict convexity

If f is  $\mu$ -strongly convex, then it is also strictly convex.

*Proof.* Direct consequence from the characterization (2.16a) from Proposition 2.11.

### Example 2.3: Strongly Convex Quadratic

The quadratic function  $f(x) = \frac{1}{2}x^{\top}Qx - c^{\top}x + r$  is  $\mu$ -strongly convex if and only if  $Q \succeq \mu I_n$ .

# 2.3 Convex optimization problems

An important class of optimization problems is the convex optimization problem.

"The great watershed in optimization is not between linearity and nonlinearity, but convexity and nonconvexity" R. Tyrrell Rockafellar

### **Definition 2.6: Convex Optimization Problem**

If  $\mathcal{X}$  is a convex set and  $f : \mathcal{X} \to \mathbb{R}$  is a convex function, then the optimization problem (1.6) is called a "convex optimization problem".

Theorem 2.3: Local Implies Global Optimality for Convex Problems

For a convex optimization problem, every local minimum is also a global one.

*Proof.* Regard a local minimum  $x^*$  of the convex optimization problem (1.6). This means that there exists  $\varepsilon > 0$  such that for all  $x \in \mathcal{N}_{\varepsilon} := \{x \mid \|x - x^*\| \le \varepsilon\}$  we have  $f(x) \ge f(x^*)$ . Now let  $y \in \mathcal{X} \setminus \{x^*\}$ . Let us define  $\alpha = \frac{\varepsilon}{\|x^* - y\|}$ . Then  $(1 - \alpha)x^* + \alpha y \in \mathcal{N}_{\varepsilon}$ . Hence, using convexity of f, we have that

$$f(x^{\star}) \le f\left((1-\alpha)x^{\star} + \alpha y\right) \le (1-\alpha)f(x^{\star}) + \alpha f(y)$$

This implies that  $f(y) \ge f(x^*)$ . Since this is true for any  $y \in \mathcal{X}$ , we can conclude that  $x^*$  is a global minimum.

### Theorem 2.4: Solution set

For a convex optimization problem, the set of minimizers is convex.

*Proof.* The set of minimizers is the set  $\{x \in \mathcal{X} \mid f(x) = \min_{x} f(x)\} = \{x \in \mathcal{X} \mid f(x) \le \min_{x} f(x)\}$ . As it is a sublevel set of a convex function, it is convex. *Remark.* The function f is constant on the set of global minimizers.

### Theorem 2.5: First order optimality condition for convex problems

Let f be a convex and continuously differentiable function on a convex set  $\mathcal{X}$ . Let  $\bar{x}$  be in the interior of  $\mathcal{X}$ . Then  $\bar{x}$  is a global optimizer of (1.6) if and only if it is a stationary point.

Proof.

- $\Rightarrow$ : Global optimality implies stationarity from Theorem 1.3.
- $\Leftarrow$ : Let  $\bar{x}$  be a stationary point, i.e.  $\nabla f(\bar{x}) = 0$  Let us use the characterization (2.5b) of convexity of smooth functions, from Proposition 2.9 with  $x = \bar{x}$ :

$$\forall y \in \mathcal{X}, \quad f(y) \ge f(\bar{x}) + \nabla f(\bar{x})^{\top} (y - \bar{x}) = f(y)$$

This implies that  $\bar{x}$  is a global optimizer.

#### Theorem 2.6: Unicity of minimizer for strictly convex functions

Let f be a strictly convex function on a convex set  $\mathcal{X}$ . Then f has at most one minimizer in  $\mathcal{X}$  (which is the unique stationary point in case of existence).

*Proof.* Let us prove this by contradiction. Assume that there exists two distinct minimizers x and y. Then apply the strict convexity property (2.12) with  $\alpha = \frac{1}{2}$ :

$$f\left(\frac{x+y}{2}\right) < \frac{f(x)+f(y)}{2}f(x) + \frac{1}{2}f(y) = \min_{x \in \mathcal{X}} f(x),$$

meaning that f evaluated at  $\frac{x+y}{2}$  has a value lower than its minimum, and yet  $\frac{x+y}{2}$  is in  $\mathcal{X}$ , since  $\mathcal{X}$  is convex. This is a contradiction.

*Remark.* In the case where  $\mathcal{X}$  is open and f is also continuously differentiable, this result, combined with Theorem 2.5, implies that if the minimizer exists, it is also the unique stationary point. Furthermore, the minimizer exists if and only if a stationary point exists.

### Lemma 2.1: Strong convexity implies coercive

Let f be a  $\mu$ -strongly convex function. Then it is coercive.

*Proof.* Let us define the set  $B = \{x \in \mathcal{X} \mid ||x|| \le 1\}$ . Since B is compact,  $\overline{f} \coloneqq \inf_{x \in B} f(x)$  is finite.

Let y be such that  $||y|| \ge 1$ . Let us use the property (2.16a) for x = 0:

$$f(\alpha y) \le (1-\alpha) f(0) + \alpha f(y) - \frac{\mu}{2} \alpha (1-\alpha) ||y||^2$$
 (2.21)

Now, apply (2.21) to  $\alpha = \frac{1}{\|y\|}$ :

$$\bar{f} \le f\left(\frac{y}{\|y\|}\right) \le \left(1 - \frac{1}{\|y\|}\right) f(0) + \frac{f(y)}{\|y\|} - \frac{\mu}{2} \left(1 - \frac{1}{\|y\|}\right) \|y\|$$
(2.22)

After rearranging the terms, we get:

$$f(y) \ge \underbrace{\frac{\mu}{2} \|y\|^2 + \left(\bar{f} - \frac{\mu}{2}\right) \|y\| - f(0) + \frac{f(0)}{\|y\|}}_{=:\kappa(\|y\|)}$$
(2.23)

Clearly,  $\kappa(\|y\|) \xrightarrow{\|y\| \to +\infty} +\infty$ . This implies that f is coercive.

# 

## Theorem 2.7: Existence and unicity theorem for strongly convex functions

Let f be a  $\mu$ -strongly convex function on a convex set  $\mathcal{X}$ . Also, assume that  $\mathcal{X}$  is closed and non-empty. Then f has a unique global minimizer in  $\mathcal{X}$  (which is also the unique stationary point).

*Proof.* Using Lemma 2.1 and Proposition 2.12, f is coercive and strictly convex. Using Theorem 1.2, f has at least one global minimizer. Using Theorem 2.6, it has at most one minimizer. This implies that f has a unique global minimizer.

*Remark.* In the case where  $\mathcal{X}$  is open, and f is also continuously differentiable, this result combined with Theorem 2.5 implies that the unique global minimizer is also the unique stationary point.

# 2.4 Examples of convex optimization problems in data analysis

### Example 2.4: Quadratic Programs

Consider the QP (1.8):

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2}x^TQx - c^Tx + r_s$$

Then the following holds:

- The problem is convex if and only if  $Q \succcurlyeq 0$ .
- The problem is strictly convex if and only if  $Q \succ 0$ .
- The problem is  $\mu$ -strongly convex if and only  $Q \succcurlyeq \mu I_n$

### Example 2.5: Linear least squares problems

Consider the linear least squares optimization problem (1.11):

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \|y_j - A_j x\|^2$$

Then the following holds:

- The problem is convex.
- The problem is strictly convex if and only if  $\frac{1}{m} \sum_{j=1}^{m} A_j^{\top} A_j \succ 0$ .
- The problem is  $\mu$ -strongly convex if and only if  $\frac{1}{m} \sum_{j=1}^{m} A_j^{\top} A_j \succeq \mu I_n$ .

### Example 2.6: Ridge Regression

Consider the ridge regression optimization problem (1.12):

minimize 
$$\frac{1}{2m} \sum_{j=1}^{m} ||y_j - A_j x||^2 + \frac{\lambda}{2} ||x||^2$$

The problem is always  $\lambda$ -strongly convex; hence, it has a unique solution, as we saw before.

### Example 2.7: LASSO Regression

Consider the LASSO regression optimization problem (1.15):

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \|y_j - A_j x\|^2 + \lambda \|x\|_1$$

The problem is convex, and even strongly convex if the matrix  $\frac{1}{m} \sum_{j=1}^{m} A_j^{\top} A_j$  is invertible.

### Example 2.8: Robust regression problems

Another variant of linear least squares is the robust regression problem, where the goal is to be robust against potential outliers in the data set. More precisely, there might be a couple of data points that are highly corrupted with noise, and we do not want these data points to affect the solution too much. The robust regression problem can be formulated as follows:

$$\underset{x \in \mathbb{R}^{n}, O \in \mathbb{R}^{m \times p}}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^{m} \|y_{j} - (A_{j}x + o_{j})\|^{2} + 2\rho \|o_{j}\|_{1}$$
(2.24)

where  $o_j$  are some additional errors, typically high only for a couple of samples, corresponding to the outliers. This is translated by the  $l_1$  penalization on  $o_j$  in the objective function.

The optimization problem (2.24) is non-differentiable, but it can be transformed into another form, by explicitly optimizing over  $o_j$ :

### Proposition 2.13: Equivalent form of the robust regression problem

The optimization problem (2.24) is equivalent to the following optimization problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{m} \sum_{j=1}^m h_{\rho}^{\text{vec}} \left( y_j - A_j x \right)$$
(2.25)

where  $h_{\rho}^{\text{vec}}(e) \coloneqq \sum_{k=1}^{p} h_{\rho}(e_k)$ , and where  $h_{\rho}$  is the Huber function defined as:

$$h_{\rho}(e) \coloneqq \begin{cases} \frac{1}{2}e^2 & \text{if } |e| \le \rho, \\ \rho |e| - \frac{1}{2}\rho^2 & \text{otherwise.} \end{cases}$$
(2.26)

Furthermore, the objective function of the optimization problem (2.25) is convex and continuously differentiable.

*Proof.* Let us define  $e_{ij}(x) \coloneqq (y_j - A_j x)_i$ . Then the following holds:

$$\frac{1}{2m}\sum_{j=1}^{m}\|y_j - (A_jx + o_j)\|^2 + 2\rho \|o_j\|_1 = \frac{1}{2m}\sum_{j=1}^{m}\sum_{i=1}^{p}(e_{ij}(x) - o_{ij})^2 + 2\rho \sum_{i=1}^{p}|o_{ij}|$$
(2.27)

$$= \frac{1}{m} \sum_{j=1}^{n} \sum_{i=1}^{n} f_{if}(o_{ij}; x), \qquad (2.28)$$

where we defined  $f_{ij}(o;x) \coloneqq \frac{1}{2}(e_{ij}(x) - o)^2 + \rho |o|.$ 

By explicitly minimizing over o, we can rewrite the optimization problem (2.24) as follows:

$$\underset{x \in \mathbb{R}^{n}}{\text{minimize}} \left( \min_{O \in \mathbb{R}^{m \times p}} \frac{1}{2m} \sum_{j=1}^{m} \|y_{j} - (A_{j}x + o_{j})\|^{2} + 2\rho \|o_{j}\|_{1} \right)$$
(2.29)

$$= \left( \min_{O \in \mathbb{R}^{m \times p}} \frac{1}{m} \sum_{j=1}^{m} \sum_{i=1}^{p} f_{if}(o_{ij}; x) \right)$$
(2.30)

Now, let us use the following identity:

$$\min_{O \in \mathbb{R}^{m \times p}} \frac{1}{m} \sum_{j=1}^{m} \sum_{i=1}^{p} f_{ij}(o_{ij}; x) = \frac{1}{m} \sum_{j=1}^{m} \sum_{i=1}^{p} \min_{o \in \mathbb{R}} f_{ij}(o; x)$$
(2.31)

This implies that the optimization problem (2.32) is equivalent to the following optimization problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^p \min_{o \in \mathbb{R}} f_{if}(o; x)$$
(2.32)

Now a simple analyse of the 1D function  $f_{ij}(o; x)$  shows that the minimum is reached for:

$$o_{ij}^{\star} = \begin{cases} e_{ij}(x) - \rho & \text{if } e_{ij}(x) > \rho, \\ e_{ij}(x) + \rho & \text{if } e_{ij}(x) < -\rho, \\ 0 & \text{otherwise.} \end{cases}$$
(2.33)

This yields the following value for the minimum of  $f_{ij}(o; x)$ :

$$\min_{o \in \mathbb{R}} f_{ij}(o; x) = \begin{cases} \rho e_{ij}(x) - \frac{1}{2}\rho^2 & \text{if } |e_{ij}(x)| > \rho, \\ \frac{1}{2}e_{ij}(x)^2 & \text{otherwise.} \end{cases}$$
(2.34)

Note that this matches with the definition of the Huber function  $h_{\rho}$  in (2.26). Hence, we have shown that  $\min_{o \in \mathbb{R}} f_{ij}(o; x) = h_{\rho}(e_{ij}(x))$ . Plugging this equality into (2.32), we find that the optimization problem (2.24) is equivalent to the following:

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\text{minimize}} & \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^p h_\rho(e_{ij}(x)) \\ \iff \underset{x \in \mathbb{R}^n}{\text{minimize}} & \frac{1}{m} \sum_{j=1}^m \sum_{i=1}^p h_\rho\left((y_j - A_j x)_i\right) \\ \iff \underset{x \in \mathbb{R}^n}{\text{minimize}} & \frac{1}{m} \sum_{j=1}^m h_\rho^{\text{vec}}\left(y_j - A_j x\right) \end{array}$$

as desired.

Example 2.9: Nonlinear least squares problems

In the case of the nonlinear least squares problem (1.10):

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \|y_j - \varphi(a_j; x)\|^2,$$

the problem is, in general, non-convex.

Because of the non-convexity, the analysis of the solutions of (1.10) is quite difficult. However, there are some algorithms that can be used to find at least stationary points of the problem. The next chapter will be dedicated to optimization algorithms. To give an idea of how one could approach the problem, one very natural idea is to linearize the model around some initial guess  $\bar{x}$ , and solve the resulting linear least squares problem. Such a problem would take the following form:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \left\| y_j - \left( \varphi(a_j; \bar{x}) + \nabla \varphi(a_j; \bar{x})^\top (x - \bar{x}) \right) \right\|^2.$$
(2.35)

Iterating over the described procedure results in a method called the *Gauss-Newton method*, which is a popular method to solve nonlinear least squares problems. In this course, we will not study this method specifically but rather similar methods.

#### Example 2.10: Logistic regression

In the previous chapter, we introduced the logistic regression problem, in Example 1.5.

This optimization problem can be written more explicitly as follows:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{m} \sum_{j=1}^m \log\left(\sum_{l=1}^q e^{a_j^\top x_l}\right) - \frac{1}{m} \sum_{j=1}^m \sum_{l=1}^q p_l^{y_j}(a_j^\top x_l)$$
(2.36)

This optimization problem is convex as we will see in the next proposition.

To probe the convexity of the logistic loss, we first need to prove the following lemma.

Lemma 2.2: Convexity of log-sum-exp functions  
The function 
$$h(z) \coloneqq \log\left(\sum_{l=1}^{q} e^{z_l}\right)$$
 is convex.

*Proof.* First, let us explicitly write the derivatives of h:

$$(\nabla h(z))_i = \frac{e^{z_i}}{\sum\limits_{l=1}^q e^{z_l}}$$

Now, regarding the second derivatives:

$$\left(\nabla^2 h(z)\right)_{ij} = \frac{1}{\left(\sum_{l=1}^q e^{z_l}\right)^2} \left(-e^{z_i}e^{z_j} + \begin{cases} e^{z_i}\left(\sum_{l=1}^q e^{z_l}\right) & \text{if } i=j\\ 0 & \text{else} \end{cases}\right)$$

Let z and d be vectors of  $\in \mathbb{R}^q$ . The following holds:

$$d^{\top} \nabla^2 h(z) d = \sum_{i=1}^q \sum_{j=1}^q d_i d_j \left( \nabla^2 h(z) \right)_{ij}$$
  
=  $\frac{1}{\left( \sum_{l=1}^q e^{z_l} \right)^2} \left( \left( \sum_{i=1}^q d_i^2 e^{z_i} \right) \left( \sum_{i=1}^q e^{z_l} \right) - \sum_{i=1}^q \sum_{j=1}^q d_i d_j e^{z_i} e^{z_j} \right)$   
=  $\frac{1}{\left( \sum_{l=1}^q e^{z_l} \right)^2} \left( \left( \sum_{i=1}^q d_i^2 e^{z_i} \right) \left( \sum_{i=1}^q e^{z_l} \right) - \left( \sum_{i=1}^q d_i e^{z_i} \right)^2 \right)$ 

Furthermore, the following inequality holds, using the Cauchy-Schwarz inequality:

$$\left(\sum_{i=1}^{q} d_i e^{z_i}\right)^2 = \left(\sum_{i=1}^{q} \left(d_i \sqrt{e^{z_i}}\right) \left(\sqrt{e^{z_i}}\right)\right)^2$$
$$\leq \left(\sum_{i=1}^{q} \left(d_i \sqrt{e^{z_i}}\right)^2\right) \left(\sum_{i=1}^{q} \left(\sqrt{e^{z_i}}\right)^2\right)$$
$$= \left(\sum_{i=1}^{q} d_i^2 e^{z_i}\right) \left(\sum_{i=1}^{q} e^{z_i}\right)$$

This allows us to conclude that  $d^{\top}\nabla^2 h(z)d \ge 0$ . Since this holds for any point z and any direction d, we can conclude that h is convex.

Proposition 2.14: Convexity of the logistic regression loss

The logistic regression optimization problem (2.36) is convex.

*Proof.* The objective function in (2.36) is the sum of functions that are either linear, either in the form of  $h(a_j^{\top}x)$ .

Using Lemma 2.2, we know that h is convex. Since a convex-over-linear function is convex, the term  $h(a_j^{\top}x)$  is also convex.

Moreover, linear functions are convex.

Therefore, the objective function is a sum of convex functions. Hence, it is itself a convex function.

## Chapter 3

# Descent Methods for Solving Optimization Problems

In this section, we will assume that  $\mathcal{X} = \mathbb{R}^n$  and that f is a continuously differentiable function. This is usually referred to as *smooth and unconstrained optimization*. The optimization problem that we will treat in this section takes the following form:

$$\min_{x \in \mathbb{R}^n} f(x) \tag{3.1}$$

## 3.1 Generalities about Descent methods

If a closed-form of the solution of the optimization problem (3.1) is not available, one needs to approximate the solution with some algorithm. Such an algorithm typically takes the following iterative form:

Initialize  $x_0$  to some initial guess For k = 0, ..., T or until some convergence criterion is satisfied (3.2) compute  $x_{k+1}$  according to rule

Furthermore, the iterative rule for  $x_k$  will take the form:

$$x_{k+1} = x_k + \alpha_k d_k \tag{3.3}$$

where  $\alpha_k \in (0, 1]$  is called the *step-length* and  $d_k \in \mathbb{R}^n$  represents the direction in which we update the solution point  $x_k$ .

In most of the algorithms that are going to be seen in this chapter, the direction  $d_k$  is (only) a function of the current point  $x_k$ :

$$d_k = \phi(x_k) \tag{3.4}$$

where  $\phi : \mathbb{R}^n \to \mathbb{R}^n$  is some function. In the next section, we will see that a very natural choice is  $\phi(x) = -\nabla f(x)$ .

#### **Definition 3.1: Descent directions**

A vector d is called a *descent direction* for f at the point x if:

 $\exists \varepsilon > 0 \text{ such that } \forall \alpha \in (0, \varepsilon], \ f(x + \alpha d) < f(x)$ (3.5)

Proposition 3.1: Conditions for descent directions

Let f be an L-smooth function. Then, for d to be a descent direction at point x:

- it is necessary that  $\nabla f(x)^{\top} d \leq 0$ ,
- it is sufficient that  $\nabla f(x)^{\top} d < 0$

Remark. One could maybe relax the L-smooth assumption here, but it makes the proof easier.

*Proof.* Let us define the function  $g(\alpha) \coloneqq \frac{f(x+\alpha d)-f(x)}{\alpha}$ . Note that  $\lim_{\alpha \to 0} g(\alpha) = \nabla f(x)^{\top} d$ .

Now we prove each of the points.

- $\nabla f(x)^{\top} d \leq 0$  is necessary: If d is a descent direction, then  $g(\alpha) < 0$  for  $\alpha$  small enough. This implies that  $\nabla f(x)^{\top} d = \lim_{\alpha \to 0} g(\alpha) \leq 0$
- $\nabla f(x)^{\top} d < 0$  is sufficient: If  $\nabla f(x)^{\top} d < 0$ , then  $\lim_{\alpha \to 0} g(\alpha) < 0$ . This implies that  $g(\alpha) < 0$  for  $\alpha$  small enough, hence d is a descent direction.

Regularity of the function f

#### **Definition 3.2:** *L*-smooth function

Let L > 0 be a positive scalar. We say the function f is L-smooth with its gradient is L-Lipschitz-continuous:

$$\forall x, y \in \mathbb{R}^n, \quad \|\nabla f(x) - \nabla f(y)\| \le L \|x - y\|$$
(3.6)

Proposition 3.2: Inequality for L-smooth functions

Assume that f is L-smooth. Then the following holds:

$$\forall x, y \in \mathbb{R}^n, \quad f(y) = f(x) + \nabla f(x)^\top (y - x) + r(x, y)$$
(3.7)

with r(x, y) satisfying:

$$|r(x,y)| \le \frac{L}{2} ||x-y||^2$$
(3.8)

*Proof.* Using the integral form of the first-order Taylor expansion, we find:

$$r(x,y) = \int_0^1 \left( \nabla f(x + t(y - x)) - \nabla f(x) \right)^\top (y - x) \, dt \tag{3.9}$$

Using the L-smoothness of f, we find:

$$\begin{aligned} |r(x,y)| &= \left| \int_0^1 \left( \nabla f \left( x + t(y-x) \right) - \nabla f(x) \right)^\top (y-x) \, dt \right| \\ &\leq \int_0^1 \left| \left( \nabla f \left( x + t(y-x) \right) - \nabla f(x) \right)^\top (y-x) \right| \, dt \\ &\leq \int_0^1 \| \nabla f \left( x + t(y-x) \right) - \nabla f(x) \| \, \| x - y \| \, dt \quad \text{using the Cauchy-Schwarz inequality} \\ &\leq \int_0^1 L \| x + t(y-x) - x \| \, \| x - y \| \, dt \quad \text{using (3.6)} \\ &= L \int_0^1 t \, dt \, \| x - y \|^2 \\ &= \frac{L}{2} \, \| x - y \|^2 \end{aligned}$$

#### Proposition 3.3: Characterization of L-smooth function

Let f be a function twice-differentiable. Then f is L-smooth if and only if

$$\forall x \in \mathbb{R}^n, \quad -LI_n \preccurlyeq \nabla^2 f(x) \preccurlyeq LI_n \tag{3.10}$$

*Proof.* Now let us prove that f is L-smooth  $\iff$  (3.10) holds:

 $\Rightarrow$ : Let  $x \in \mathbb{R}^n$  and  $d \in \mathbb{R}^n$ . Using (3.11) and (3.6), we find:

$$\begin{aligned} \left| d^{\top} \nabla^{2} f(x) d \right| &= \left| \lim_{\alpha \to 0} \frac{1}{\alpha} \left( \nabla f(x + \alpha d) - \nabla f(x) \right)^{\top} d \right| \\ &= \lim_{\alpha \to 0} \frac{1}{\alpha} \left| \left( \nabla f(x + \alpha d) - \nabla f(x) \right)^{\top} d \right| \\ &\leq \lim_{\alpha \to 0} \frac{1}{\alpha} \left\| \left( \nabla f(x + \alpha d) - \nabla f(x) \right)^{\top} \right\| \|d\| \\ &\leq \lim_{\alpha \to 0} \frac{1}{\alpha} L \|\alpha d\| \|d\| \\ &\leq \lim_{\alpha \to 0} L \|d\|^{2} \end{aligned}$$

which proves that  $-L \|d\|^2 \leq d^\top \nabla^2 f(x) d \leq L \|d\|^2$ . Since this holds for all  $d \in \mathbb{R}^n$ , we have that  $-LI_n \preccurlyeq \nabla^2 f(x) \preccurlyeq LI_n$ .

 $\Leftarrow$ : Using the fundamental theorem of calculus applied to  $\nabla f(\cdot)$ , we have that:

$$\nabla f(y) - \nabla f(x) = \int_0^1 \nabla^2 f(x + t(y - x))(y - x) \, dt \tag{3.11}$$

Since implies:

$$\begin{aligned} \|\nabla f(y) - \nabla f(x)\| &= \left\| \int_0^1 \nabla^2 f(x + t(y - x))(y - x) \, dt \right\| \\ &= \int_0^1 \left\| \nabla^2 f(x + t(y - x))(y - x) \right\| \, dt \\ &= \int_0^1 \sqrt{(y - x)^\top (\nabla^2 f(x + t(y - x)))^2 (y - x)} \, dt \\ &\leq \int_0^1 \sqrt{L^2 \|x - y\|^2} \, dt \end{aligned} = L \|x - y\|$$

which proves that f is L-smooth.

Note that we used the fact that (3.10) implies that  $(\nabla^2 f(x))^2 \preccurlyeq LI_n$ .

#### A useful inequality for descent methods

#### **Proposition 3.4**

Let f be an L-smooth function. Let  $x_0, \ldots, x_T$  be updated according to the general rule (3.3). Then the following holds:

$$f(x_{k+1}) \le f(x_k) + \alpha_k \nabla f(x_k)^\top d_k + \alpha_k^2 \frac{L}{2} \|d_k\|^2$$
(3.12)

*Proof.* This comes directly from Proposition 3.2.

*Remark.* From the inequality (3.12), we see that to minimize the right-hand side of (3.12), the choice  $d_k = -\nabla f(x_k)$  and  $\alpha_k = \frac{1}{L}$  is optimal. This choice results in a method called *the gradient descent method*. This will be studied more in-depth in the next section.

### 3.2 The gradient descent method

In this section, one very classical, widely used, algorithm: the gradient descent algorithm.

Definition 3.3: The gradient descent algorithm The gradient descent algorithm is the following iterative method: Initialize  $x_0$  to some initial guess For k = 0, ..., T or until some convergence criterion is satisfied (3.13)  $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$ 

*Remark.* As mentionned above, this corresponds to the choice  $d_k = \phi(x_k) = -\nabla f(x_k)$ *Remark.* We have not yet defined how the step length is chosen. There are several ways to choose it, either via fixing a value  $\alpha_k = \alpha$ , or via a dedicated procedure, as we will see later.

42

#### **Proposition 3.5**

Assume that f is L-smooth. Let  $\alpha_{\min}$  and  $\alpha_{\max}$  be such that  $0 < \alpha_{\min} \le \alpha_{\max} < \frac{2}{L}$ . Let  $x_0, \ldots, x_T$  be updated according to the gradient descent algorithm (3.13) with  $\alpha_k \in [\alpha_{\min}, \alpha_{\max}]$ . Then the following holds:

$$f(x_{k+1}) \le f(x_k) - C \|\nabla f(x_k)\|^2$$
(3.14)

with  $C = \alpha_{\min}(1 - \alpha_{\max}\frac{L}{2}) > 0.$ 

*Proof.* Using equation (3.12) with  $d_k = -\nabla f(x_k)$ :

$$f(x_{k+1}) \leq f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \alpha_k^2 \frac{L}{2} \|\nabla f(x_k)\|^2$$
  
=  $f(x_k) - \alpha_k \left(1 - \alpha_k \frac{L}{2}\right) \|\nabla f(x_k)\|^2$   
$$\leq f(x_k) - \underbrace{\alpha_{\min} \left(1 - \alpha_{\max} \frac{L}{2}\right)}_{=C} \|\nabla f(x_k)\|^2$$

*Remark.* We see that to get the best bound on the decrease of the function, we should choose:

$$\alpha_{\min} = \alpha_{\max} = \frac{1}{L}.$$
(3.15)

This choice is called the *steepest descent method*, and it results in the following constant C:

$$C = \frac{1}{2L} \tag{3.16}$$

## 3.3 Convergence properties

#### 3.3.1 The important assumption

In this section, we go back to the general case of descent methods. However, we make the assumption that a property similar to the property (3.14) is verified:

Assumption 3.1

Assume that there exists a constant C such that the following holds:

$$\forall k \in \mathbb{N} \quad f(x_{k+1}) \le f(x_k) - C \left\| \nabla f(x_k) \right\|^2 \tag{3.17}$$

*Remark.* As we saw above, this property is verified for the gradient method under the assumption that the step lengths  $\alpha_k$  are in a certain interval.

#### 3.3.2 Convergence for a general smooth function

#### **Proposition 3.6: Convergence of** $\nabla f(x_k)$ to 0

Assume that f is L-smooth and bounded from below. Also, assume that  $x_0, \ldots, x_t$  is a sequence of points such that the Assumption 3.2 is fulfilled. Then we have

$$\nabla f(x_t) \xrightarrow[t \to +\infty]{} 0 \tag{3.18}$$

*Proof.* From the inequality (3.17), we have, for all k:

$$\|\nabla f(x_k)\|^2 \le \frac{1}{C} \left( f(x_k) - f(x_{k+1}) \right)$$
(3.19)

Summing this inequality from k = 0 to t - 1, we find:

$$\underbrace{\sum_{k=0}^{t-1} \|\nabla f(x_k)\|^2}_{=:u_t} \le \frac{1}{C} \left( f(x_0) - f(x_t) \right) \le \frac{1}{C} \left( f(x_0) - \inf_x f(x) \right)$$
(3.20)

The sequence  $u_t$  is non-decreasing and bounded; hence it converges. Therefore, the following holds:

$$\|\nabla f(x_t)\| = u_{t+1} - u_t \xrightarrow[t \to +\infty]{} \lim_{t \to +\infty} u_{t+1} - \lim_{t \to +\infty} u_t = 0,$$

which proves (3.18).

*Remark.* This proposition is important, but it does not necessarily mean that  $x_k$  will approach a stationary point of f. For example, if  $f(x) = e^{-x}$ , we still have f bounded from below, yet it does not admit any stationary point. In fact, such a point does not necessarily exist.

#### Theorem 3.1: Convergence theorem for a general smooth function

Assume that the set  $K := \{x \mid f(x) \leq f(x_0)\}$  is bounded. Also, assume that  $x_0, \ldots, x_t$  is a sequence of points such that the Assumption 3.2 is fulfilled. Then:

 $f(x_t) \xrightarrow[t \to +\infty]{} f(\bar{x}) \tag{3.21}$ 

where  $\bar{x}$  is a stationary point of f.

*Remark.* The set K is always bounded if f is coercive.

*Proof.* Because of Assumption 3.2,  $f(x_k)$  is decreasing, which implies that  $x_k \in K$  for all k. Using Bolzano-Weierstrass theorem, this implies that it has at least one accumulation point  $\bar{x}$ , i.e. there exists a sequence  $k_j$  that goes to  $+\infty$  such that:

$$x_{k_j} \xrightarrow[j \to +\infty]{} \bar{x} \tag{3.22}$$

Moreover, the sequence  $f_k \coloneqq f(x_k)$  is non-increasing, hence  $f_k \xrightarrow[k \to +\infty]{} \bar{f}$  for some  $\bar{f} \in \mathbb{R} \cup \{+\infty\}$ . Let us write  $\bar{f}$  its limit. The following holds:

$$\bar{f} = \lim_{k \to +\infty} f(x_k) = \lim_{j \to +\infty} f(x_{k_j}) = f(\bar{x})$$
(3.23)

where the last equality comes from the continuity of  $f(\cdot)$ . This proves  $f(x_k) \xrightarrow[k \to +\infty]{} f(\bar{x})$ .

Now, we only have to show that  $\bar{x}$  is a stationary point of f, i.e.  $\nabla f(\bar{x}) = 0$ . Let us use (3.17) from Assumption 3.2:

$$\|\nabla f(x_k)\|^2 \le \frac{1}{C} \left( f(x_k) - f(x_{k+1}) \right) \xrightarrow[k \to +\infty]{} \frac{1}{C} \left( \bar{f} - \bar{f} \right) = 0$$
(3.24)

This implies that  $\nabla f(x_k) \xrightarrow[k \to +\infty]{} 0.$ 

On the other hand,  $\nabla f(\cdot)$  is continuous, which implies that  $\nabla f(x_{k_j}) \xrightarrow{j \to +\infty} \nabla f(\bar{x})$ . Combining these two limits, we find  $\nabla f(\bar{x}) = 0$ .

*Remark.* While Theorem 3.1 provides a good guarantee, it does not provide any information on the speed of convergence. In fact, the convergence might be very slow, as we see in the following example.

#### Example 3.1: A function where CG converges slowly

Let us define the function  $f : \mathbb{R} \to \mathbb{R}$  as:

$$f(x) = \begin{cases} 1 - e^{-x} & \text{if } x \ge 0\\ \frac{1}{2}(x+1)^2 - 1 & \text{if } x < 0 \end{cases}$$
(3.25)

This function is L-smooth with L = 1. It has a unique stationary point and minimum at x = -1. The steepest gradient method for f is:

$$x_{k+1} = x_k - \nabla f(x_k).$$
 (3.26)

While the sequence  $x_k$  converges to -1, the convergence is very slow. Indeed, if  $x_0 > 0$ , it will take more than  $e^{x_0} - 1$  iterations to reach a negative point  $x_t \le 0$ . *Proof.* Let t be the first iteration such that  $x_t < 0$ . Then, for all  $k \le t - 1$ , we have:

$$x_{k+1} = x_k + d_k$$

with  $d_k = -\nabla e^{-x_k}$ .

Using the convexity of  $e^x$ , we have  $e^{d_k} \ge 1 + d_k$ . This implies:

$$e^{x_{k+1}} = e^{x_k} e^{d_k} \ge e^{x_k} (1+d_k) = e^{x_k} - 1$$

By repeating the inequality found for k = 0 to t - 1, we find:

$$e^{x_t} \ge e^{x_0} - t.$$

Finally, since  $x_t \leq 0$ , we have  $e^{x_t} \leq 1$ , which implies that  $t \geq e^{x_0} - 1$ .

#### 3.3.3 Convergence for a strongly convex function

#### Theorem 3.2: Convergence theorem for a strongly convex function

Assume that f is  $\mu$ -strongly convex. Let  $x^*$  be the solution of the optimization problem (3.1). Also, assume that  $x_0, \ldots, x_t$  is a sequence of points such that the Assumption 3.2 is fulfilled. Then the following holds for all  $k \in \mathbb{N}$ :

$$f(x_k) - f(x^*) \le (1 - 2\mu C)^k \left( f(x_0) - f(x^*) \right)$$
(3.27)

*Proof.* Using (2.16b) from Proposition 2.11, we have:

$$\forall x, y \in \mathbb{R}^n, \quad f(y) \ge f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|^2$$

By minimizing both sides over y, we find:

$$f(x^{\star}) = \min_{y} f(y) \ge \min_{y} f(x) + \nabla f(x)^{\top} (y - x) + \frac{\mu}{2} \|y - x\|^{2} = f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^{2}$$

Now, apply this inequality for  $x = x_k$ :

$$\|\nabla f(x_k)\|^2 \ge 2\mu (f(x_k) - f(x^*))$$

Now, plugging this inequality into (3.17), we find:

$$\forall k \in \mathbb{N} \quad f(x_{k+1}) \le f(x_k) - C \|\nabla f(x_k)\|^2 \le f(x_k) - 2\mu C(f(x_k) - f(x^*))$$

Hence, the sequence  $u_k := f(x_k) - f(x^*)$  verifies  $u_{k+1} \le (1 - 2\mu C)u_k$ . Repeating this inequality iteratively leads to  $u_k \le (1 - 2\mu C)^k u_0$ , which is the inequality (3.27).

*Remark.* The speed of convergence here is very satisfying. For example, it ensures that the number of steps that are required is linear in the number of digits of precision that is desired. For example, if a solution with  $\left|f(x) - \min_{x} f(x)\right| \leq 10^{-M}$  is required, then one needs less that  $k = \frac{M + \log_{10}(f(x_0) - f(x^*))}{-\log_{10}(1 - 2\mu C)}$  number of steps.

Corollary 3.1: Convergence rate for the steepest descent method

For the steepest descent method  $x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k)$ , the following holds:

$$f(x_k) - f(x^*) \le \left(1 - \frac{\mu}{L}\right)^k (f(x_0) - f(x^*))$$
 (3.28)

*Proof.* From (3.16), we saw that  $C = \frac{1}{2L}$  is achieved for the steepest descent method. Plugging this into (3.27), we find the desired result.

Corollary 3.2: Bound on  $x_k - x^*$ 

If the assumptions of Theorem 3.2 are fulfilled, then the following also holds:

$$\|x_k - x^\star\|^2 \le \frac{2}{\mu} (1 - 2\mu C)^k \left( f(x_0) - f(x^\star) \right)$$
(3.29)

*Proof.* Directly follow from (3.27) and from:

$$f(x_k) \ge f(x^*) + \frac{\mu}{2} \|x_k - x^*\|^2, \qquad (3.30)$$

which is derived from the strong convex inequality (2.16b) for  $y = x_k$  and  $x = x^*$ .

#### **Corollary 3.3**

If in addition to the assumption in Theorem 3.2, f is L-smooth, then the following holds

$$\|x_k - x^\star\|^2 \le \frac{L}{\mu} (1 - 2\mu C)^k \|x_0 - x^\star\|^2$$
(3.31)

*Proof.* Directly follow from (3.29) and from:

$$f(x_k) \le f(x^*) + \frac{L}{2} ||x_0 - x^*||^2$$
(3.32)

which is derived from the L-smooth inequality (3.7) for  $y = x_k$  and  $x = x^*$ .

#### 3.3.4 Convergence for a weakly convex function

Regarding the convergence of gradient descent methods for convex functions, Theorem 3.2 is very encouraging. Indeed, it provides a strong convergence guarantee and ensures a fast convergence when the function is strongly convex.

Furthermore, a function that is convex is not so far from being strongly convex. For examples, if f is convex, then  $f + \frac{\epsilon}{2} ||x||^2$  is strongly convex. That means that by perturbing the function f a little, the gradient method would converge rather quickly to the solution.

The following theorem provides convergence guarantee for weakly convex functions under reasonable assumptions. The convergence speed is however slower than for strongly convex functions.

Before we state the theorem, let us prove the following lemma:

#### Lemma 3.1: A small mathematical point

Let  $u_0, \ldots, u_t$  be a sequence of non-negative real numbers such that for some  $\lambda > 0$ :

$$\forall k \in \mathbb{N}: \quad u_{k+1} \le u_k - \lambda u_k^2 \tag{3.33}$$

Then  $u_k \leq \frac{1}{\lambda k}$  for all  $k \in \mathbb{N}$ .

*Proof.* First, note that  $u_{k+1} \leq u_k$  for all  $k \in \mathbb{N}$ . This permits to rearrange (3.33) into:

$$\frac{1}{u_{k+1}} - \frac{1}{u_k} = \frac{u_k - u_{k+1}}{u_k u_{k+1}} \ge \frac{u_k - u_{k+1}}{u_k^2} \ge \lambda \tag{3.34}$$

where the last inequality comes from (3.33) directly.

Now, summing the inequality (3.34) from k = 0 to t - 1, we find:

$$\frac{1}{u_t} - \frac{1}{u_0} \ge \lambda t \tag{3.35}$$

which implies  $u_t \leq \frac{1}{\lambda t}$ .

#### Theorem 3.3: Convergence theorem for a convex function

Assume that f is L-smooth and convex. Let  $x^*$  be a solution of the optimization problem (3.1). Also, assume that  $x_0, \ldots, x_t$  is a sequence of points such that the Assumption 3.2 is fulfilled. Moreover, like in Theorem 3.1, assume that the set  $K := \{x \mid f(x) \leq f(x_0)\}$  is bounded. Then define:

$$R_{0} \coloneqq \max\left\{ \|x - x^{\star}\| \mid x \in K \right\} = \max\left\{ \|x - x^{\star}\| \mid f(x) \le f(x_{0}) \right\} < \infty \quad (3.36)$$

Then the following holds for all k:

$$f(x_k) - f(x^*) \le \frac{R_0^2}{C} \frac{1}{k}$$
(3.37)

*Proof.* First,  $f(x_k)$  is decreasing because of Assumption 3.2. This implies  $x_k \in K$  for all k, which implies:

$$\forall k \in \mathbb{N}, \quad \|x_k - x^\star\| \le R_0. \tag{3.38}$$

Furthermore, since f is convex, we have:

$$f(x_k) \leq f(x^*) + \nabla f(x^*)^\top (x_k - x^*) \quad \text{from (2.5b)}$$
  
$$\leq f(x^*) + \|\nabla f(x^*)\| \|x_k - x^*\| \quad \text{from the Cauchy-Schwarz inequality}$$
  
$$\leq f(x^*) + \|\nabla f(x^*)\| R_0 \quad \text{from (3.38)}$$

which can be rearranged as:

$$\|\nabla f(x^{\star})\| \ge \frac{f(x_k) - f(x^{\star})}{R_0}$$
 (3.39)

Now, using equation (3.17) combined with (3.39), we have:

$$f(x_{k+1}) \le f(x_k) - \frac{C}{R_0^2} (f(x_k) - f(x^*))^2$$
(3.40)

Now, define  $\lambda \coloneqq \frac{C}{R_0^2}$  and  $u_k \coloneqq f(x_{k+1}) - f(x^*)$ . Note that  $u_k \leq 0$ . Substracting  $f(x^*)$  in both sides of equation (3.40):

$$u_{k+1} \le u_k - \lambda u_k^2 \tag{3.41}$$

Using Lemma 3.1, this implies that  $u_k \leq \frac{1}{\lambda k}$  for all  $k \in \mathbb{N}$ , which proves (3.37).

## 3.4 Globalization techniques

As we saw for the gradient descent method, to ensure convergence, we need to ensure some constraints on the step-lengths  $\alpha_k$ . The critical part of the constraint that we saw was  $\alpha_{\max} < \frac{2}{L}$ .

However, the constant L depends on global properties of the function f, while we can only access local properties of f.

We make the distinction between three types of method to choose the step-length  $\alpha_k$ :

- Fixed step-length:  $\alpha_k = \alpha \approx \frac{1}{L}$  for all k. This is the simplest method, but it requires a good guess on the value of L.
- Exact line search:  $\alpha_k$  can be chosen according to the optimization problem:

$$\alpha_k = \arg\min_{\alpha>0} f\left(x_k + \alpha d_k\right)$$

This method require to solve an additional 1D optimization problem at each iteration, so it is in general computationally expensive.

• Backtracking line search: one finds a value on the form  $\alpha_k = \beta^i \bar{\alpha}$  for some  $\bar{\alpha}, \beta \in (0, 1)$  by iteratively decreasing *i* until the Armijo criterion (cf. definition below) is satisfied. This is summarized in the following algorithm:

$$\alpha_k \leftarrow \bar{\alpha}$$
  
while the Armijo Criterion is not satisfied for  $\alpha_k$ : (3.42)  
 $\alpha_k \leftarrow \beta \alpha_k$ 

#### **Definition 3.4: Armijo criterion**

The Armijo criterion is a test on  $\alpha$  regarding the veracity of the following inequality:

$$f(x_k + \alpha d_k) \le f(x_k) + c_A \alpha \nabla f(x_k)^{\top} d_k$$
(3.43)

for some  $c_A \in (0, 1)$ .

In the following theorem, we will prove a strong convergence result for the backtracking line search method. Before stating that theorem, however, we need to introduce the following assumption:

#### Assumption 3.2

Assume that the direction is chosen as  $d_k = \phi(x_k)$ , and assume that  $\phi(x)$  is such that, for some  $\varepsilon > 0$  and somme  $\gamma > 0$ , the two following conditions hold:

$$\forall x \in \mathbb{R}^{n}, \quad -\nabla f(x)^{\top} \phi(x) \ge \varepsilon \left\| \nabla f(x) \right\|^{2}$$
(3.44a)

$$\forall x \in \mathbb{R}^n, \qquad \|\phi(x)\| \le \gamma \|\nabla f(x)\| \qquad (3.44b)$$

*Remark.* The Assumption 3.2 is verified for the gradient descent method with  $\gamma = \varepsilon = 1$ .

#### Proposition 3.7: Termination of backtracking line-search

Assume that f is an L-smooth function. Under Assumption 3.2, the backtracking line-search method (3.42) with Armijo criterion (3.43) terminates in a finite number of iterations.

*Proof.* Using the L-smooth property of f and equation (3.7), we have for all  $\alpha$ :

$$f(x_k + \alpha d_k) \le f(x_k) + \alpha \nabla f(x_k)^\top d_k + \alpha^2 \frac{L}{2} ||d_k||^2$$
 (3.45)

This implies:

$$f(x_{k} + \alpha d_{k}) \leq f(x_{k}) + \alpha \nabla f(x_{k})^{\top} d_{k} + \alpha^{2} \frac{L}{2} ||d_{k}||^{2}$$
  
=  $f(x_{k}) + c_{A} \alpha \nabla f(x_{k})^{\top} d_{k} + (1 - c_{A}) \alpha \nabla f(x_{k})^{\top} d_{k} + \frac{L}{2} \alpha^{2} ||d_{k}||^{2}$   
 $\leq f(x_{k}) + c_{A} \alpha \nabla f(x_{k})^{\top} d_{k} - (1 - c_{A}) \alpha \varepsilon ||\nabla f(x_{k})||^{2} + \frac{L}{2} \alpha^{2} ||d_{k}||^{2}$   
 $\leq f(x_{k}) + c_{A} \alpha \nabla f(x_{k})^{\top} d_{k} + \alpha \frac{L}{2} ||d_{k}||^{2} \left(\alpha - \frac{2(1 - c_{A})\varepsilon ||\nabla f(x_{k})||^{2}}{L ||d_{k}||^{2}}\right)$ 

On the other hand, since  $\beta^i \bar{\alpha} \xrightarrow[i \to +\infty]{i \to +\infty} 0$ , there exists an iteration i' such that  $\alpha' \coloneqq \beta^{i'} \bar{\alpha} \leq \frac{2(1-c_A)\varepsilon \|\nabla f(x_k)\|^2}{L\|d_k\|^2}$ . This implies, using the inequality derived above:

$$f(x_k + \alpha' d_k) \le f(x_k) + c_A \alpha' \nabla f(x_k)^\top d_k$$

This proves that the Armijo criterion at iteration i', therefore, the backtracking line-search method terminates in a finite number of iterations (maximum i').

#### Theorem 3.4: Convergence result for backtracking line search

Let f be an L-smooth function. Assume that the chosen direction  $d_k$  follows Assumption 3.2. Furthermore, consider that the backtracking line-search method is used with the Armijo criterion (3.43). Then, Assumption 3.2 holds for  $C = \max\left(c_A \varepsilon \bar{\alpha}, c_A (1-c_A) \frac{2\beta \varepsilon^2}{L\gamma^2}\right)$ , i.e.:

 $f(x_{k+1}) \le f(x_k) - C \|\nabla f(x_k)\|^2$ (3.46)

*Proof.* We have to distinguish between two cases: the case where backtracking is uncessary, and the case where at least one backtracking was performed.

• For the iterations where no backtracking is needed, the Armijo criterion is satisfied for  $\alpha_k = \bar{\alpha}$ . Using (3.43) from the Armijo condition and (3.44a) from Assumption 3.2:

$$f(x_{k+1}) \leq f(x_k) + c_A \bar{\alpha} \nabla f(x_k)^\top d_k,$$
  
$$\leq f(x_k) - c_A \bar{\alpha} \varepsilon \|\nabla f(x_k)\|^2,$$
  
$$\leq f(x_k) - C \|\nabla f(x_k)\|^2$$

• For the iterations where backtracking is needed, the condition is satisfied for  $\alpha = \alpha_k$ , but not for  $\alpha = \beta^{-1} \alpha_k$ . This is summarized as follows:

$$f\left(x_{k} + \beta^{-1}\alpha_{k}d_{k}\right) > f\left(x_{k}\right) + c_{A}\beta^{-1}\alpha_{k}\nabla f\left(x_{k}\right)^{\top}d_{k}$$
(3.47a)

$$f(x_{k+1}) = f(x_k + \alpha_k d_k) \le f(x_k) + c_A \alpha_k \nabla f(x_k)^{\top} d_k$$
(3.47b)

On the other hand, the inequality (3.7) for L-smooths functions gives:

$$f(x_{k} + \beta^{-1}\alpha_{k}d_{k}) \leq f(x_{k}) + \beta^{-1}\alpha_{k}\nabla f(x_{k})^{\top}d_{k} + \frac{L}{2}\beta^{-2}\alpha_{k}^{2} ||d_{k}||^{2}$$
(3.48)

Combining (3.47a) and (3.48), we have:

$$f(x_k) + c_A \beta^{-1} \alpha_k \nabla f(x_k)^\top d_k \le f(x_k) + \beta^{-1} \alpha_k \nabla f(x_k)^\top d_k + \frac{L}{2} \beta^{-2} \alpha_k^2 \|d_k\|^2, \quad (3.49)$$

which implies:

$$\alpha \ge -(1-c_A)\frac{2\beta}{L}\frac{\nabla f(x_k)^\top d_k}{\left\|d_k\right\|^2}$$
(3.50)

Now, let us plug (3.50) into (3.47b):

$$f(x_{k+1}) \leq f(x_k) - c_A \left( (1 - c_A) \frac{2\beta}{L} \frac{\nabla f(x_k)^\top d_k}{\|d_k\|^2} \right) \nabla f(x_k)^\top d_k$$
  

$$= f(x_k) - c_A (1 - c_A) \frac{2\beta}{L} \left( \frac{\nabla f(x_k)^\top d_k}{\|d_k\|} \right)^2,$$
  

$$\leq f(x_k) - c_A (1 - c_A) \frac{2\beta\varepsilon^2}{L} \left( \frac{\|\nabla f(x_k)\|^2}{\|d_k\|} \right)^2 \quad \text{using (3.44a)},$$
  

$$\leq f(x_k) - c_A (1 - c_A) \frac{2\beta\varepsilon^2}{L\gamma^2} \left( \frac{\|\nabla f(x_k)\|^2}{\|\nabla f(x_k)\|} \right)^2 \quad \text{using (3.44b)},$$
  

$$= f(x_k) - c_A (1 - c_A) \frac{2\beta\varepsilon^2}{L\gamma^2} \|\nabla f(x_k)\|^2,$$
  

$$\leq f(x_k) - C \|\nabla f(x_k)\|^2,$$

#### **Corollary 3.4**

Thanks to Theorem 3.4, the convergence theorems 3.1, 3.3, 3.2 and Corollary 3.3 also apply when using the backtracking line search method.

### 3.5 Other examples of descent methods

#### 3.5.1 Quasi-Newton methods

Withouth going to much into the details of these methods, a popular class of methods is called second-order methods. There, the idea is to minimize a convex quadratic approximation of the function f at each iteration. Such quadratic approximation takes the following form:

$$f(x) \approx f(x_k) + \nabla f(x_k)^{\top} (x - x_k) + \frac{1}{2} (x - x_k)^{\top} H(x_k) (x - x_k)$$
(3.51)

for some matrix  $H(x_k) \succ 0$ .

When minimizing the right hand-side of (3.51), we obtain  $x = x_k + d_k$  where  $d_k$  is defined as follows:

$$d_k = -H(x_k)^{-1} \nabla f(x_k) \eqqcolon \phi(x_k) \tag{3.52}$$

When  $H(x_k) = \nabla^2 f(x_k)$ , the method is called the *exact Newton method*. Note that this method is well defined only when  $H(x_k) \succ 0$ .

If we define  $M(x) \coloneqq H(x)^{-1}$ , we get the following class of methods.

#### **Definition 3.5: Quasi-Newton method**

The quasi-Newton methods are descent methods where the direction  $d_k$  is computed as follows:

$$d_k = -M(x_k)\nabla f(x_k) \quad (=: \phi(x_k)) \tag{3.53}$$

for some matrices  $M(x_k) \geq 0$ .

#### **Proposition 3.8**

Assume that M(x) is such that:

$$\forall x \in \mathbb{R}^n, \quad \varepsilon I_n \preccurlyeq M(x) \preccurlyeq \gamma I_n, \tag{3.54}$$

with some  $\varepsilon > 0$  and  $\gamma > 0$ .

Then, the direction  $d_k$  defined by (3.53) fulfills Assumption 3.2 (with the same  $\varepsilon$  and  $\gamma$ ).

*Proof.* Writting Assumption 3.2 for  $d_k = -M(x_k)\nabla f(x_k)$  reads:

$$\forall x \in \mathbb{R}^n \quad g^\top M(x)g \ge \varepsilon \|g\|^2 \quad \text{with } g = \nabla f(x) \tag{3.55a}$$

$$\forall x \in \mathbb{R}^n \quad \|M(x)g\| \le \gamma \|g\| \quad \text{with } g = \nabla f(x) \tag{3.55b}$$

Equation (3.55a) is directly derived from  $M(x) \succeq \varepsilon I_n$  Regarding equation (3.55b), we can rearrange it at follows:

$$\forall x \in \mathbb{R}^n g^\top M(x)^2 g \le \gamma^2 \|g\|^2 \quad \text{with } g = \nabla f(x), \tag{3.56}$$

which is verified by remarking that  $0 \preccurlyeq M(x) \preccurlyeq \gamma I_n \implies M(x)^2 \preccurlyeq \gamma^2 I_n$ .

#### **Corollary 3.5**

The quasi-Newton method combined with a backtracking line search for the globalization strategy inherits from the convergence results from theorems 3.1, 3.3, 3.2 and the corollary 3.3.

*Proof.* Direct consequence of Proposition 3.8.

The Gauss-Newton method For solving the non-linear least squares problem (1.10), we discussed about the Gauss-Newton method in Section 2.4, which consists in solving iteratively the optimization problem:

$$\underset{d \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \left\| y_j - \left( \varphi(a_j; x_k) + \nabla \varphi(a_j; x_k)^\top d \right) \right\|^2.$$
(3.57)

Note that this method is an example of quasi-Newton methods, with the choice:

$$M(x) = \left(\frac{1}{m}\sum_{j=1}^{m} \nabla \varphi(a_j; x) \nabla \varphi(a_j; x)^{\top}\right)^{-1}$$

54

#### 3.5.2 Stochastic gradient methods

In Chapter ??, we will discuss in detail the stochastic gradient descent methods. Here, we can at least mention the setup: instead of having access to the true gradient  $\nabla f(x_k)$ , we only have access to a random variable  $g_k$  which is such that  $\mathbb{E}[g_k] = \nabla f(x_k)$ . In this case, the direction  $d_k = -g_k$  does not exactly verify the descent conditions (3.44a) and (3.44b), but it will be verified in average.

#### 3.5.3 Coordinate descent methods

In coordinate descent methods, at each iteration, only one coordinate of the solution  $x_k$  is updated. More precisely, let  $e_i$  be the vector with a 1 at the *i*-th position and 0 elsewhere. The coordinate descent methods take the form:

$$x_{k+1} = x_k - \alpha_k \, (\nabla f(x))_{i_k} \, e_{i_k} \tag{3.58}$$

for some index  $i_k$  and some step-size  $\alpha_k$ .

A typical choice is to choose the index  $i_k$  randomly at each iteration, or to cycle through the indices.

For these choices, the descent conditions discussed in the previous sections are not exactly verified, but similar results can be obtained.

There exists, however, one variation where the Assumption 3.2 is verified: the Gauss-Southwell method. In this specific method, we choose the index  $i_k$  as follows:

$$i_k = \arg\max_i |(\nabla f(x_k))_i| \tag{3.59}$$

Regarding the motivation behind this choice, one could be reluctant: why would we use only one coordinate of the gradient if we have to access to all of them anyway? The answer is that, in some specific cases, the gradient can be updated very efficiently when only one coordinate is updated. This is the case, for example, when the objective function is a sum of many functions, each depending only on a limited number of variables.

## Chapter 4

## **Descent** Methods with Momentum

In the previous methods that we saw, at each step, we only keep track of the current estimate of the solution. However, it could be that by using more information from the previous steps, the convergence is improved. In descent methods, the algorithm typically takes the following form:

> Initialize  $x_0, y_0$  to some initial guess For k = 0, ..., T or until some convergence criterion is satisfied compute the gradient  $g_k = \nabla f(x_k)$  (4.1) compute  $d_k$  according to some rule  $d_k = \Phi_k(g_1, ..., g_k, x_1, ..., x_k)$ update  $x_k x_{k+1} = x_k + \alpha_k d_k$

In the sketch of algorithm (4.1), the rule is very general, and will see some more specificity later on.

## 4.1 Derivation of descent methods with momentum

#### 4.1.1 Motivation from differential equations

To guide the choice of the update rule, one analogy is meaningful: seeing the gradient descent method as an approximation of the following differential equation:

$$\dot{x}(t) = -\nabla f(x) \tag{4.2}$$

The differential equation (4.2) is called the gradient flow method and has very strong properties, except that it is not numerically implementable. Let us note that if one would approximate (4.2) with an Euler scheme:

$$x(t + \Delta t) = x(t) - \Delta t \nabla f(x(t)), \qquad (4.3)$$

one would recover the gradient descent method, with  $\alpha = \Delta t$  and  $x_k = x(k \cdot \Delta t)$ .

Now, if one replaces the equation (4.2) by a physics-inspired equation, where f(x) represents the energy associated with the position x, one would get the following second-order differential equation:

$$\ddot{x}(t) = -\nabla f(x) - \nu \dot{x}(t) \tag{4.4}$$

where  $\nu$  is a friction coefficient. In (4.4), the vector  $d = -\nabla f(x)$  does not directly influence the velocity of x(t) anymore, but rather the acceleration.

One of the motivations behind (4.4) is that "small local minim" might be skipped thanks to the inertia of the particle x(t). On the other hand, the friction term  $-\nu \dot{x}(t)$  ensures that in the long term, the particle x(t) will stabilize at a local minimum of its energy f(x).

#### 4.1.2 Derivation of the heavy-ball method

The heavy-ball method is derived from a discretization of the differential equation (4.4). Using finite differences on (4.4), we get the following discretization:

$$\frac{x(t+\Delta t) - 2x(t) + x(t-\Delta t)}{(\Delta t)^2} = -\nabla f(x(t)) - \nu \frac{x(t) - x(t-\Delta t)}{\Delta t}$$

$$\tag{4.5}$$

Like for the gradient flow, we define  $x_k = x(k \cdot \Delta t)$ , and rearrange equation (4.5):

$$x_{k+1} - x_k = -(\Delta t)^2 \nabla f(x(t)) + (1 - \nu \Delta t) (x_k - x_{k-1})$$
(4.6)

Finally, after defining  $\alpha = (\Delta t)^2$  and  $\beta = 1 - \nu \Delta t$ , we can see that (4.6) is equivalent to the heavy-ball method defined as follows.

#### Definition 4.1: The heavy-ball method

The heavy-ball method is the following iterative algorithm:

$$x_{k+1} = x_k + \beta \ (x_k - x_{k-1}) - \alpha \nabla f(x_k) \tag{4.7}$$

where  $\alpha$  and  $\beta$  are two parameters hyperparameters.

#### **Proposition 4.1: Alternative formulation**

The heavy-ball method (4.7) can also be formulated as follows:

$$x_{k+1} = x_k + \alpha p_k, d_{k+1} = -\nabla f(x_{k+1}), p_{k+1} = \gamma p_k + d_{k+1}$$
(4.8)

*Proof.* Define  $p_k := \frac{x_{k+1}-x_k}{\alpha}$ , and  $\gamma = \beta$  permits to rewrite the heavy-ball method (4.7) as in (4.8).

*Remark.* By setting  $r = 1 - \gamma$ ,  $\tilde{\alpha} = \frac{\alpha}{r}$  and  $\tilde{p}_k = rp_k$ , we can rewrite (4.8) as follows:

$$x_{k+1} = x_k + \tilde{\alpha} \tilde{p}_k, d_{k+1} = -\nabla f(x_{k+1}), \tilde{p}_{k+1} = (1-r)\tilde{p}_k + rd_{k+1}$$
(4.9)

From this formulation, we can see that in the heavy-ball method, the descent direction  $\tilde{p}_k$  is a weighted average of the previous direction and the current negative gradient. By rearranging the equations (4.9), we can see that that the descent direction is a weighted sum of all the previous negative gradients, with a forgetting factor r:

$$x_{k+1} = x_k + \tilde{\alpha} \frac{\sum_{i=0}^{k} (1-r)^i d_{k-i}}{\sum_{i=0}^{k} (1-r)^i}$$
(4.10)

From this perspective, this method has another advantage: when one has only access to a noisy estimate of the gradient, averaging the gradient over time reduces the noise for the descent direction.

#### 4.1.3 Nesterov's accelerated gradient method

Now that we have seen the heavy-ball method, we will see a small variation of it.

#### Definition 4.2: Nesterov's accelerated gradient method

Nesterov's accelerated gradient method is the following iterative method:

$$x_{k+1} = x_k + \beta(x_k - x_{k-1}) - \alpha \nabla f(x_k + \beta(x_k - x_{k-1}))$$
(4.11)

where  $\alpha$  and  $\beta$  are two parameters hyperparameters.

*Remark.* The present method is very similar to the heavy-ball method, except that the gradient is computed after adding the inertia term  $\beta(x_k - x_{k-1})$  instead of before.

*Remark.* For a fine choice of the parameters  $\alpha$  and  $\beta$ , the Nesterov's accelerated gradient method actually converges faster than the heavy-ball method. This is the reason why this method is also called *The Nesterov's Optimal Method*. In the next section, we will dive into the convergence rate of this method.

Note that one might prefer to use the following notation for Nesterov's method:

Definition 4.3: Nesterov's accelerated gradient method (alternative formulation)

Nesterov's accelerated gradient method can also be written as follows

$$y_{k} = x_{k} + \beta(x_{k} - x_{k-1}),$$
  

$$x_{k+1} = y_{k} - \alpha \nabla f(y_{k}),$$
(4.12)

### 4.2 Convergence analysis of Nesterov's accelerated gradient method

In this section, we will analyze the convergence rate of Nesterov's method, for  $\mu$ -strongly convex and *L*-smooth functions. We recall that, in the case of twice continuously differentiable functions, these are the functions that verifies:

$$\forall x \in \mathbb{R}^n, \quad \mu I_n \preccurlyeq \nabla^2 f(x) \preccurlyeq L I_n \tag{4.13}$$

We will also define the following quantity:

$$c \coloneqq \sqrt{\frac{\mu}{L}} \in (0, 1) \tag{4.14}$$

*Remark.* In the case where f is quadratic,  $c = \sqrt{\text{cond}(Q^{-1})}$  where cond(P) is the condition number of a matrix P: the ration between its largest and smallest eigenvalues.

We also recall that for such functions, we found in Chapter 3 that the steepest gradient descent method (i.e. the gradient descent method with  $\alpha = \frac{1}{L}$ ) converged exponentially fast, with the following rate (cf. (3.28) in Corollary 3.1):

$$f(x_k) - f(x^*) \le (1 - c^2)^k \left( f(x_0) - f(x^*) \right)$$
(4.15)

In the case of Nesterov's method, a better convergence rate can be obtained, given that both  $\mu$  and L are known a priori, as we see in the following theorem.

#### Theorem 4.1: Convergence of Nesterov's method for strongly convex functions

Assume that f is  $\mu$ -strongly convex and L-smooth. Consider Nesterov's method (4.11) with the following choice of parameters:

$$\alpha = \frac{1}{L}, \qquad \beta = \frac{\sqrt{L} - \sqrt{\mu}}{\sqrt{L} + \sqrt{\mu}} = \frac{1 - c}{1 + c}$$

$$(4.16)$$

Let  $x^*$  be the solution of the optimization problem (3.1). Then, the sequence  $x_k$  converges to  $x^*$ , with the following rate:

$$f(x_k) - f(x^*) \le (1 - c)^k M(x_0) \tag{4.17}$$

with 
$$M(x_0) \coloneqq f(x_0) - f(x^*) + \frac{\mu^2}{2L} ||x_0 - x^*||^2$$

*Remark.* Since  $1 - c \le 1 - c^2$ , we can see that the Nesterov's method converges faster than the steepest gradient descent method.

*Proof.* First of all, let us introduce some notations:

$$g_k \coloneqq \alpha \nabla f(y_k) = \frac{1}{L} \nabla f(y_k)$$

$$e_k \coloneqq y_k - x^*$$

$$v_k \coloneqq \beta(x_k - x_{k-1}) = \frac{1-c}{1+c} (x_k - x_{k-1})$$
(4.18)

The proof will be divided in several paragraphs because it is quite long.

Bound on the decrease of the objective function Let us derive three inequalities based on the assumptions on f.

• Since f is L-smooth, we can use the inequality (3.7) from Proposition 3.2:

$$f(x_{k+1}) = f(y_k - g_k) \le f(y_k) - \nabla f(y_k)^\top g_k + \frac{L}{2} \|g_k\|^2$$
  
=  $f(y_k) - Lg_k^\top g_k + \frac{L}{2} \|g_k\|^2$   
=  $f(y_k) - \frac{L}{2} \|g_k\|^2$  (4.19)

• Then, since f is convex, we can use the inequality (2.5b) from Proposition 2.9:

$$f(y_k) + \nabla f(y_k)^{\top} (x_k - y_k) \le f(x_k)$$

which directly implies:

$$f(y_k) \le f(x_k) + \nabla f(y_k)^\top v_k = f(x_k) + Lg_k^\top v_k$$
 (4.20)

• Finally, since f is  $\mu$ -strongly convex, we can use the inequality (2.16b) from Proposition 2.11:

$$f(x^{\star}) \ge f(y_k) + \nabla f(y_k)^{\top} (x^{\star} - y_k) + \frac{\mu}{2} \|x^{\star} - y_k\|^2$$

which directly implies:

$$f(y_k) \le f(x^*) + \nabla f(y_k)^\top e_k - \frac{\mu}{2} \|e_k\|^2 = f(x^*) + Lg_k^\top e_k - \frac{\mu}{2} \|e_k\|^2$$
(4.21)

Now, using inequalities (4.19), (4.20) and (4.21), we can derive the following inequality:

$$f(x_{k+1}) - f(x^*) \leq f(y_k) - f(x^*) - \frac{L}{2} ||g_k||^2 \quad \text{from (4.19)},$$

$$= (1 - c) (f(y_k) - f(x^*)) + c (f(y_k) - f(x^*)) - \frac{L}{2} ||g_k||^2,$$

$$\leq (1 - c) \left( f(x_k) + Lg_k^\top v_k - f(x^*) \right) + c (f(y_k) - f(x^*)) - \frac{L}{2} ||g_k||^2 \quad \text{from (4.20)},$$

$$,$$

$$\leq (1 - c) \left( f(x_k) + Lg_k^\top v_k - f(x^*) \right) + c \left( Lg_k^\top e_k - \frac{\mu}{2} ||e_k||^2 \right) - \frac{L}{2} ||g_k||^2 \quad \text{from (4.21)},$$

$$,$$

$$= (1 - c) (f(x_k) - f(x^*)) + \underbrace{(1 - c) Lg_k^\top v_k + c Lg_k^\top e_k - c\frac{\mu}{2} ||e_k||^2 - \frac{L}{2} ||g_k||^2}_{=:r_k}.$$

This brings the conclusion that:

$$f(x_{k+1}) - f(x^*) \le (1 - c) \left( f(x_k) - f(x^*) \right) + r_k, \tag{4.22}$$

#### Simplification of $r_k$ :

The quantity  $r_k$  simplifies as follows:

$$\begin{split} r_k &\coloneqq (1-c) L g_k^\top v_k + c L g_k^\top e_k - c \frac{\mu}{2} \|e_k\|^2 - \frac{L}{2} \|g_k\|^2 \,, \\ &= \frac{L}{2} \left[ 2(1-c) g_k^\top v_k + 2c g_k^\top e_k - c \frac{\mu}{L} \|e_k\|^2 - \|g_k\|^2 \right] \,, \\ &= \frac{L}{2} \left[ 2g_k^\top \left( (1-c) v_k + c e_k \right) - c^3 \|e_k\|^2 - \|g_k\|^2 \right] \,, \\ &= \frac{L}{2} \left[ \|(1-c) v_k + c e_k\|^2 - \|(1-c) v_k + c e_k - g_k\|^2 - c \|c e_k\|^2 \right] \,, \end{split}$$

## Express $r_k$ as the difference of two squares :

By expanding both sides, one can proves the following equality:

$$\forall a, b \in \mathbb{R}^n, c \in \mathbb{R}, \quad \|(1-c)a + b\|^2 = c \|b\|^2 + (1-c) \|a + b\|^2 - c(1-c) \|a\|^2 \tag{4.23}$$

Let us apply (4.23) to  $a = v_k$  and  $b = ce_k$ :

$$\|(1-c)v_k + ce_k\|^2 = c \|ce_k\|^2 + (1-c) \|v_k + ce_k\|^2 - c(1-c) \|v_k\|^2$$
(4.24)

Now, plugging the inequality (4.24) into the expression of  $r_k$ , we get:

$$r_{k} = \frac{L}{2} \left[ \|(1-c)v_{k} + ce_{k}\|^{2} - \|(1-c)v_{k} + ce_{k} - g_{k}\|^{2} - c \|ce_{k}\|^{2} \right],$$
  

$$= \frac{L}{2} \left[ (1-c) \|v_{k} + ce_{k}\|^{2} - c(1-c) \|v_{k}\|^{2} - \|(1-c)v_{k} + ce_{k} - g_{k}\|^{2} \right],$$
  

$$\leq \frac{L}{2} \left[ (1-c) \|v_{k} + ce_{k}\|^{2} - \|(1-c)v_{k} + ce_{k} - g_{k}\|^{2} \right]$$
(4.25)

Express  $r_k$  in the form  $-\frac{L}{2}\left(\|z_{k+1}\|^2 - \|z_k\|^2\right)$  :

Let us now define  $z_k := v_k + ce_k$ . In order to express  $z_{k+1}$ , let us first express  $v_{k+1}$  and  $e_{k+1}$  using the definition of Nesterov's method (4.11):

• Expression of  $e_{k+1}$ :

$$e_{k+1} = y_{k+1} - x^*$$
  
=  $x_{k+1} + v_{k+1} - x^*$   
=  $y_k - g_k + v_{k+1} - x^*$   
=  $e_k - g_k + v_{k+1}$ 

• Expression of  $v_{k+1}$ :

$$v_{k+1} = \beta(x_{k+1} - x_k)$$
$$= \beta(y_k - g_k - x_k)$$
$$= \beta(v_k - g_k)$$

Using the two equations above, we can express  $z_{k+1}$  as follows:

$$z_{k+1} = v_{k+1} + ce_{k+1}$$
  
=  $v_{k+1} + c\left(e_k - g_k + v_{k+1}\right)$   
=  $(1 + c)v_{k+1} + ce_k - cg_k$   
=  $(1 + c)\beta(v_k - g_k) + ce_k - cg_k$   
=  $(1 + c)\frac{1 - c}{1 + c}(v_k - g_k) + ce_k - cg_k$   
=  $(1 - c)(v_k - g_k) + ce_k - cg_k$   
=  $(1 - c)v_k + ce_k - g_k$ 

Using the expressions of  $z_k$  and  $z_{k+1}$ , and the bound derived in (4.25), we can draw the following conclusion:

$$r_k \le \frac{L}{2} \left[ (1-c) \|z_k\|^2 - \|z_{k+1}\|^2 \right]$$
(4.26)

#### **Conclusion step** :

The inequality (4.22) combined with the bound (4.26), gives:

$$f(x_{k+1}) - f(x^{\star}) \le (1-c) \left( f(x_k) - f(x^{\star}) \right) + \frac{L}{2} \left[ (1-c) \|z_k\|^2 - \|z_{k+1}\|^2 \right], \quad (4.27)$$

which can be rewritten as:

$$\underbrace{f(x_{k+1}) - f(x^*) + \frac{L}{2} \|z_{k+1}\|^2}_{=V_{k+1}} \le (1 - c) \left(\underbrace{f(x_k) - f(x^*) + \frac{L}{2} \|z_k\|^2}_{=V_k}\right)$$
(4.28)

where we defined  $V_k \coloneqq f(x_k) - f(x^*) + \frac{L}{2} ||z_k||^2$ .

To conclude, applying inequality (4.28) recursively concludes the proof:

$$f(x_k) - f(x^*) \le V_k \le (1 - c)^k V_0 \tag{4.29}$$

Finally, a little rearrangement can be made to get to the final result (4.17):

$$V_0 = f(x_0) - f(x^*) + \frac{L}{2} ||z_0||^2$$
  
=  $f(x_0) - f(x^*) + \frac{Lc^2}{2} ||e_0||^2$   
=  $f(x_0) - f(x^*) + \frac{\mu^2}{2L} ||x_0 - x^*||^2$ 

This directly gives the desired result (4.17):

$$f(x_k) - f(x^*) \le (1 - c)^k \left( \underbrace{f(x_0) - f(x^*) + \frac{\mu^2}{2L} \|x_0 - x^*\|^2}_{=:M(x_0)} \right)$$
(4.30)

*Remark.* It is classical in a lot of proof to find a quantity  $V_k$  that decreases at each iteration, as implied by (4.28). This is called a Lyapunov function.

## 4.3 The conjugate gradient method

Now, we will see another gradient method with momentum called the *Conjugate Gradient* (CG) method. This method is specific to the case of unconstrained (strictly) convex quadratic programming, i.e., it is for solving the following kinds of problems:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} x^\top Q x - c^\top x + r \rightleftharpoons f(x), \tag{4.31}$$

whith  $Q \succ 0$ .

#### 4.3.1 Motivation

As we saw in the previous chapters, the problem (4.31) has a closed-form solution given by  $x^* = Q^{-1}c$ . A legitimate question is: Why would we derive an iterative method to solve (4.31) if we already have a closed-form solution?

The answer is that when n >> 1, the computation of  $Q^{-1}$  might be costly. The CG method is a good alternative to solve the problem in this case. More precisely, the computation cost of computing  $x^*$  is on the order of  $\frac{n^3}{3}$  (with the Cholesky factorization, which is out of the scope here).

Instead, as we will see in this section, the CG method is very fast (maximum of *n* iterations). On each iteration, the computational cost is majorated by the computation of matrix-vector product on the form "Qp" (where  $Q \in \mathbb{R}^{n \times n}$  and  $p \in \mathbb{R}^n$ ).

An example of how the CG method is useful is when only an approximate solution is needed. In that case, one could run less than n iteration. Therefore, the CG might be cheeper that  $\frac{n^3}{3}$  operations (what is needed to compute  $Q^{-1}$ ).

Another case where the CG method is useful is when the products Qp can be computed efficiently. More precisely, if Q is a sparse matrix, or a sum-product of a few sparse matrices, one can compute the product Qp with much less computational effort than  $n^2$  operations. More precisely, one could have: complexity  $(Qp) << n^2$ . In that case, running n steps of the CG method would yield the solution  $x^*$  with a computational cost of  $\approx n$  complexity  $(Qp) << n^3$ , which would be much more efficient than the closed-form expression, which requires  $\approx \frac{n^3}{3}$  operations. In fact, in the case where Qp can be computed efficiently, the CG method can even be used as a linear solver for the system Qx = c.

#### Example 4.1: Ridge regression with $m \ll n$

Consider the Ridge regression example (1.12) with  $\dim(y_j) = 1$  and  $m \ll n$ :

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2m} \sum_{j=1}^m \left\| y_j - a_j^\top x \right\|^2 + \frac{\lambda}{2} \|x\|^2.$$
(4.32)

Then, the optimization problem (4.32) can be written as a quadratic program (4.31) with  $Q = \lambda I_n + \frac{1}{m} \sum_{j=1}^m a_j a_j^{\top}$ .

In that case, the matrix-vector product Qp can be computed with  $\approx 2nm \ll n^2$  operations:

$$Qp = \lambda p + \frac{1}{m} \sum_{j=1}^{m} a_j (a_j^{\top} p).$$
(4.33)

#### 4.3.2 Construction of the CG method

Let us dive into the construction of the CG method. The idea is to take the heavy-ball formulation (4.8), but to let the parameters  $\gamma$  and  $\alpha$  vary over time:

$$x_{k+1} = x_k + \alpha_k p_k, d_{k+1} = -\nabla f(x_{k+1}), p_{k+1} = \gamma_k p_k + d_{k+1}$$
(4.34)

Here is how the parameters  $\gamma_k$  and  $\alpha_k$  are chosen:

• The parameters  $\alpha_k$  are found via exact line search:

$$\alpha_k = \arg\min_{\alpha} f(x_k + \alpha p_k)$$
  
=  $\arg\min_{\alpha} \frac{1}{2} (\alpha p_k)^\top \nabla^2 f(x_k) (\alpha p_k) + \nabla f(x_k)^\top (\alpha p_k) + f(x_k)$  (4.35)

Solving equation (4.35) yields:

$$\alpha_k = -\frac{\nabla f(x_k)^\top p_k}{p_k^\top \nabla^2 f(x_k) p_k} = \frac{d_k^\top p_k}{p_k^\top Q p_k}$$
(4.36)

• The parameters  $\gamma_k$  are chosen such that  $p_{k+1}$  and  $p_k$  are orthogonal with respect to the scalar product  $x \cdot y := x^\top Q y$ :

$$p_k^{\top} Q p_{k+1} = p_k^{\top} Q \left( \gamma_k p_k + d_{k+1} \right) = 0$$
(4.37)

Solving equation (4.37) yields:

$$\gamma_k = -\frac{p_k^{\top} Q d_{k+1}}{p_k^{\top} Q p_k} \tag{4.38}$$

An important remark is that if  $p_k = 0$ , then  $\alpha_k$  and  $\gamma_k$  can not be computed. In that case, the algorithm terminates.

Now, using the definition of the method (4.34), and the definitions of  $\gamma_k$  and  $\alpha_k$  in (4.38) and (4.36), we can write the CG method explicitly.

#### **Definition 4.4: The Conjugate Gradient Method**

The Conjugate Gradient (CG) method is the following iterative algorithm:

While  $p_k \neq 0$ :

$$\alpha_{k} = \frac{d_{k}^{\top} p_{k}}{p_{k}^{\top} Q p_{k}},$$

$$x_{k+1} = x_{k} + \alpha_{k} p_{k},$$

$$d_{k+1} = -\nabla f \left( x_{k+1} \right) \qquad (= c - Q x_{k+1})$$

$$\gamma_{k} = -\frac{p_{k}^{\top} Q d_{k+1}}{p_{k}^{\top} Q p_{k}},$$

$$p_{k+1} = \gamma_{k} p_{k} + d_{k+1}$$

$$(4.39)$$

The quantity  $x_0$  is initialized with some initial guess, and  $p_0 = -\nabla f(x_0) = c - Qx_0$ .

#### 4.3.3 Some properties of the CG method

Now, let us prove some of the properties of the CG method.

#### **Proposition 4.2**

At each iteration, the following equality holds:

$$d_{k+1} \, \, ^{\prime} p_k = 0 \tag{4.40}$$

*Proof.* The first-order optimality condition for the line search (4.35) at iteration k yields:

$$\nabla f(x_{k+1})^{\top} p_k = 0 \tag{4.41}$$

Using the definition  $d_{k+1} = -\nabla f(x_{k+1})$ , the property (4.40) follows.

### Proposition 4.3

When applying the CG method, the following holds for each k:

$$||d_k||^2 \le ||p_k||^2, \tag{4.42}$$

*Proof.* From the definition of  $p_k$ :

$$p_k = \gamma_{k-1} p_{k-1} + d_k \tag{4.43}$$

(note that we define  $\gamma_{-1} = 1$  and  $p_{-1} = 0$  for (4.43) to be also true for k = 0). Now, combining equations (4.40) at step k - 1 and (4.43), we get:

$$\|p_k\|^2 = \|\gamma_{k-1}p_{k-1} + d_k\|^2 = \gamma_{k-1}^2 \|p_{k-1}\|^2 + \|d_k\|^2 \ge \|d_k\|^2$$
(4.44)

Now, we are able to prove that if the CG method terminates, then the solution is found.

#### Proposition 4.4: Success of CG in case of termination

If the CG method terminates at iteration k, then  $x_k$  is  $x^*$ , the global minimizer.

*Proof.* If the CF method terminates at iteration k, it means that  $p_k = 0$ . Using Proposition 4.3, this implies that  $d_k = -\nabla f(x_k) = 0$ , hence  $x_k$  is a stationary point of f. Since f is strongly convex, the unique stationary point is the global minimizer:  $x_k = x^*$ .

Proposition 4.5: Termination of CG when a solution is found

If  $x_k = x^*$ , then the CG method terminates at iteration k.

*Proof.* If  $x_k = x^*$ , then  $d_k = -\nabla f(x_k) = 0$ . From equation (4.38), this implies  $\gamma_{k-1} = 0$ . Therefore,  $p_k = \gamma_{k-1}p_{k-1} + d_k = d_k = 0$ . Hence, the CG method terminates at iteration k.

**Proposition 4.6:**  $\alpha_k \neq 0$ 

At every step k (where the algorithm has not terminated),  $\alpha_k \neq 0$ 

*Proof.* By contradiction, assume that  $\alpha_k = 0$ . Then  $x_{k+1} = x_k$ . From equation (4.40), we have  $d_k^{\top} p_k = 0$ . Now using the definition  $p_k = \gamma_{k-1} p_{k-1} d_k$ , the following holds:

$$0 = d_k^{\top} p_k$$
  
=  $d_k^{\top} (\gamma_{k-1} p_{k-1} + d_k)$   
=  $\gamma_{k-1} d_k^{\top} p_{k-1} + ||d_k||^2$ 

Moreover, from equation (4.40) applied at iteration k - 1, we have  $d_k^{\top} p_{k-1} = 0$ . This implies that  $||d_k||^2 = 0$ .

Therefore,  $\nabla f(x_k) = -d_k = 0$ , so  $x_k = x^*$ , which implies that the algorithm should have terminated, which is a contradiction.

#### 4.3.4 Termination of the CG method

The CG method has a very nice convergence property: it converges in at most n steps.

In fact, after k steps, the iterate  $x_k$  is the optimal value of f over the affine space spanned by  $\{p_0, \ldots, p_{k-1}\}$  and shifted by  $x_0$ .

More precisely, let us define the vectorial subspace  $S_k$  spanned by  $\{p_0, \ldots, p_{k-1}\}$ :

$$S_k \coloneqq \operatorname{span}\{p_0, \dots, p_{k-1}\} = \left\{ \sum_{i=0}^{k-1} \lambda_i p_i \mid \lambda_i \in \mathbb{R} \right\}.$$
(4.45)

By convention,  $S_0 = \{0\}$  to keep the property  $S_{k+1} = \operatorname{span}(p_k) + S_k$  consistent.

In the following theorem, we will show that the iterate  $x_k$  is the optimal value of f over the affine space  $S_k + x_0$ .

Furthermore, the family of vectors  $\{p_0, \ldots, p_k\}$  is an orthogonal family for the scalar product  $x \cdot y \coloneqq x^\top Q y$ .

Note that by construction,  $p_k$  is already orthogonal to  $p_{k-1}$ . The surprise here, is that it is also orthogonal to all  $p_s$  with s < k.

Before we dive into the theorem, we need one important lemma.

#### Lemma 4.1: An important property

The following holds for all t (where the algorithm has not terminated):

$$QS_t \subset S_{t+1} \tag{4.46}$$

*Proof.* For all  $k \leq t - 1$ , we have:

$$p_{k+1} - \gamma_k p_k = d_{k+1}, = d_k - Q(x_{k+1} - x_k), = d_k - \alpha_k Q p_k, = p_k - \gamma_{k-1} p_{k-1} - \alpha_k Q p_k.$$
(4.47)

This can be rearranged as follows, using  $\alpha_k \neq 0$  from Proposition 4.6:

$$Qp_k = \frac{1}{\alpha_k} \left( \gamma_k p_k - p_{k+1} - \gamma_{k-1} p_{k-1} + p_k \right) \in S_{t+1}.$$
(4.48)

Since this holds for all  $k \leq t - 1$ , we can conclude the following:

$$QS_t = Q \operatorname{span}\{p_0, \dots, p_{t-1}\} = \operatorname{span}\{Qp_0, \dots, Qp_{t-1}\} \subset S_{t+1}.$$
(4.49)

#### Theorem 4.2: The key property of CG

For all k, the two following properties hold:

$$x_k = \operatorname*{arg\,min}_{x \in S_k + x_0} f(x),\tag{4.50a}$$

$$p_k \in S_k^{\perp_Q} \coloneqq \{ p \in \mathbb{R}^n \mid \forall z \in S_k, \ z^\top Q p = 0 \}$$

$$(4.50b)$$

*Proof.* We will prove this theorem by induction.

The properties (4.50a) and (4.50b) are trivial for k = 0 since  $S_0 = \{0\}$ .

Now, assume that the properties (4.50a) and (4.50b) hold for  $k \leq t$ . We will show that they also hold for k = t + 1.

• **Property** (4.50a) for k = t + 1:

Let x be an element of  $\{x_0\} + S_t$ , i.e.  $x - x_0 \in S_t$ . Using the induction hypothesis for (4.50b) at k = t, we have  $(x - x_0)Qp_t = 0$ . This implies:

$$\nabla f(x)^{\top} p_t = (\nabla f(x_0) + Q(x - x_0))^{\top} p_t$$
$$= \nabla f(x_0)^{\top} p_t + \underbrace{(x - x_0)^{\top} Q p_t}_{=0}$$
$$= \nabla f(x_0)^{\top} p_t$$

It follows that:

$$\forall x \in (\{x_0\} + S_t), \ f(x + \alpha p_t) = f(x) + \alpha^2 p_t^\top \nabla^2 f(x) p_t + \alpha \nabla f(x)^\top p_t$$
$$= f(x) + \underbrace{\alpha^2 p_t^\top Q p_t + \alpha \nabla f(x_0)^\top p_t}_{=:h_t(\alpha)},$$

which is summarized by:

$$\forall x \in (\{x_0\} + S_t), \ f(x + \alpha p_t) = f(x) + h_t(\alpha)$$
(4.51)

Now let z be an element of  $\{x_0\} + S_{t+1}$ . Since  $\{x_0\} + S_{t+1} = (\{x_0\} + S_t) + p_t$ , we can write  $z = x + \alpha p_t$  for some  $x \in (\{x_0\} + S_t)$  and  $\alpha \in \mathbb{R}$ .

Now, we use (4.51), the induction hypothesis for (4.50a) at k = t, and the definition of  $\alpha_t$ 

in (4.36), to derive the following inequalities:

$$f(z) = f(x + \alpha p_t) = f(x) + h_t(\alpha)$$
  

$$\geq \min_{x \in (\{x_0\} + S_t)} f(x) + h_t(\alpha)$$
  

$$= f(x_t) + h_t(\alpha)$$
  

$$= f(x_t + \alpha p_t)$$
  

$$\geq \min_{\alpha \in \mathbb{R}} f(x_t + \alpha p_t)$$
  

$$= f(x_{t+1})$$

It follows that  $f(z) \ge f(x_{t+1})$ , and this is true for any  $z \in \{x_0\} + S_{t+1}$ . On the other hand,  $x_{t+1} = x_t + \alpha_t p_t \in (\{x_0\} + S_t) + p_t = \{x_0\} + S_{t+1}$ . This proves:

$$x_{t+1} = \underset{x \in \{x_0\} + S_{t+1}}{\arg\min} f(x), \qquad (4.52)$$

which proves the property (4.50a) for k = t + 1.

#### • Intermediate step :

Let z be an element of  $S_t$ . Using (4.46), we have  $Qz \in QS_t \subset S_{t+1}$ . Therefore, the following holds:

$$\forall \tau \in \mathbb{R}, \ x_{t+1} + \tau Qz \in \{x_0\} + S_{t+1} \tag{4.53}$$

Using (4.52) and (4.53), we can write:

$$0 = \arg\min_{\tau} f(x_{t+1} + \tau Qz),$$
(4.54)

for which the first-order optimality condition is:

$$\nabla f(x_{t+1})^{\top} Q z = 0 \tag{4.55}$$

This leads to the conclusion:

$$d_{k+1} = \nabla f(x_{t+1}) \in (S_t)^{\perp_Q}$$
(4.56)

• **Property** (4.50b) for k = t + 1:

The induction hypothesis for (4.50b) at k = t implies

$$p_t \in (S_t)^{\perp_Q} \tag{4.57}$$

Furthermore, we have already shown that  $\nabla f(x_{t+1}) \in (S_t)^{\perp_Q}$ . Therefore, using the definition  $p_{t+1} = \gamma_t p_t - \nabla f(x_{t+1})$ , we have:

$$p_{t+1} = \gamma_t \underbrace{p_t}_{\in (S_t)^{\perp_Q} \text{ from } (4.57)} + \underbrace{d_{k+1}}_{\in (S_t)^{\perp_Q} \text{ from } (4.56)} \in (S_t)^{\perp_Q}$$
(4.58)

Moreover, from the construction of  $\gamma_t$  in (4.37), we have:

$$p_t \,^{\top} Q p_{t+1} = 0 \tag{4.59}$$

Every element of z in  $S_{t+1}$  can be written as  $z = \tilde{z} + \lambda p_t$  for some  $\lambda \in \mathbb{R}$ , and  $\tilde{z} \in S_t$ . Therefore:

$$z^{\top}Qp_{t+1} = \underbrace{z^{\top}Qp_{t+1}}_{=0 \text{ by } (4.58)} + \lambda \underbrace{p_t^{\top}Qp_{t+1}}_{=0 \text{ by } (4.59)} = 0$$
(4.60)

This implies that  $p_{t+1} \in (S_{t+1})^{\perp_Q}$ , which proves the property (4.50b) for k = t + 1.

The induction is complete; hence, the properties (4.50a) and (4.50b) hold for all k.  $\Box$ 

#### Corollary 4.1: Linear independence of $p_k$

If the algorithm has not terminated at step t, then the vectors  $p_0, \ldots, p_t$  are linearly independent.

*Proof.* This comes directly from the fact that the vectors  $p_0, \ldots, p_t$  are non-zero and orthogonal with respect to the scalar product  $x \cdot y \coloneqq x^\top Q y$ .

The theorem being proven, we can now state the following corollary.

Corollary 4.2: Termination of CG

The CG method terminates in at most n steps and yields the global minimizer  $x^*$ .

*Proof.* We prove this corollary by contradiction. Assume that after n steps, the CG method has not terminated. Using Corollary 4.1, this implies that the vectors  $p_0, \ldots, p_n$  are linearly independent. But since dim $(\mathbb{R}^n) = n$ , there can not be n + 1 linearly independent vectors in  $\mathbb{R}^n$ . Since this is a contradiction, the CG method must terminate in at most n steps.

Finally, using Proposition 4.4, we can conclude that the solution  $x^*$  is found at the termination of the algorithm.

#### 4.3.5 More efficient formulation

After some algebraic manipulations, one can rearrange the equations (4.39) to get a more efficient formulation of the CG method.

Proposition 4.7: The Conjugate Gradient Method (more efficient formulation)

The CG method defined in (4.39) can be rewritten as follows:

$$q_{k} = Qp_{k},$$

$$\alpha_{k} = \frac{d_{k}^{\top} p_{k}}{p_{k}^{\top} q_{k}},$$

$$x_{k+1} = x_{k} + \alpha_{k} p_{k},$$

$$d_{k+1} = d_{k} - \alpha_{k} q_{k}$$

$$\gamma_{k} = -\frac{q_{k}^{\top} d_{k+1}}{p_{k}^{\top} q_{k}},$$

$$p_{k+1} = \gamma_{k} p_{k} - d_{k+1}$$

$$(4.61)$$

*Proof.* The proof is pretty straightforward, by remarking the following fact:

$$d_{k+1} = -\nabla f(x_{k+1})$$
  
=  $c - Qx_{k+1}$   
=  $c - Q(x_k + \alpha_k p_k)$   
=  $d_k - \alpha_k Qp_k$   
=  $d_k - \alpha_k q_k$ 

*Remark.* From equations (4.61), we see that at each iteration, the only computational cost that may exceed  $\mathcal{O}(n)$  is the computation of  $Qp_k$ . This means that the complexity of a CG iteration is  $\approx$  complexity("Qp"), as we have mentionned in part 4.3.1.

## 4.3.6 Computing $Q^{-1}$ with the CG method

As we have mentioned in part 4.3.1, the CG method can be used as a linear solver for the system Qx = c. In fact, the CG method can even be used to compute  $Q^{-1}$ .

**Proposition 4.8: Computation of**  $Q^{-1}$ Assume that the CG method has terminated at iteration n (i.e., the worst case). Then the following holds: $Q^{-1} = \sum_{k=0}^{n-1} \frac{1}{p_k^{\top} q_k} p_k p_k^{\top}$ (4.62) *Proof.* Let us define the matrix J as follows:

$$J \coloneqq \left(\sum_{k=0}^{n-1} \frac{1}{p_k^{\top} q_k} p_k p_k^{\top}\right) Q = \sum_{k=0}^{n-1} \frac{1}{p_k^{\top} Q p_k} p_k p_k^{\top} Q$$
(4.63)

We recall that for  $t \neq k$ ,  $p_k^{\top} Q p_t = 0$  (from (4.50b) in Theorem 4.2). Hence, for t = 0, ..., n - 1:

$$Jp_{t} = \sum_{k=0}^{n-1} \frac{1}{p_{k}^{\top} Q p_{k}} p_{k} p_{k}^{\top} Q p_{t}$$

$$= \sum_{k=0}^{n-1} \frac{1}{p_{k}^{\top} Q p_{k}} (p_{k}^{\top} Q p_{t}) p_{k}$$

$$= p_{t} + \sum_{k \neq p_{t}} \frac{1}{p_{k}^{\top} Q p_{k}} \underbrace{p_{k}^{\top} Q p_{t}}_{=0} p_{k}$$

$$= p_{t}$$

$$(4.64)$$

This implies that  $(J - I_n)p_t = 0$  for t = 0, ..., n - 1. Furthermore, the vectors  $p_0, ..., p_{n-1}$  are linearly independent (cf. Corollary 4.1), therefore, they form a basis of  $\mathbb{R}^n$ .

This implies that  $J - I_n$  is the zero matrix, i.e.  $J = I_n$ . From the definition of J in (4.63),  $J = I_n$  implies:

$$\left(\sum_{k=0}^{n-1} \frac{1}{p_k^{\top} q_k} p_k p_k^{\top}\right) Q = I_n, \tag{4.65}$$

which proves the property (4.62).

## Appendix A

## Very basics of mathematics

In this small chapter, we provide a few basic mathematical concepts that are used throughout the book. These are basic definitions and well-known properties/theorems that will be stated without proof.

The goal is to have written support that can be referred to when needed. Note that the definitions, theorems, and properties mentioned here are not necessarily used in the rest of the script, but it is good to have them in mind when reading the script or attending the course.

## A.1 Basics of linear algebra

#### Definition A.1: Euclidean norm

The norm  $\|\cdot\|$  denotes the L2-norm, also called the *Euclidean norm* defined as follows:

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

Proposition A.1: The Cauchy-Schwarz inequality

If x and y are vectors in  $\mathbb{R}^n$ , then:

$$\left|x^{\top}y\right| \le \|x\| \,\|y\| \tag{A.2}$$

#### **Definition A.2: Eigenvalues of a matrix**

Let  $A \in \mathbb{R}^{n \times n}$  be a square matrix. An eigenvalue of A is a scalar  $\lambda \in \mathbb{C}$  such that there exists a non-zero vector  $x \in \mathbb{C}^n$  satisfying:

$$Ax = \lambda x \tag{A.3}$$

We write sp(A) the set of eigenvalues of A.

#### Theorem A.1: Spectral theorem

Let  $S \in \mathbb{R}^{n \times n}$  be a symmetric matrix (i.e.  $S^{\top} = S$ ). Then  $\operatorname{sp}(S) \subset \mathbb{R}$ , i.e. the eigenvalues of S are real. Furthermore, S is orthogonally diagonalizable:

$$S = P\Lambda P^{\top} \tag{A.4}$$

where  $P \in \mathbb{R}^{n \times n}$  is an orthogonal matrix, i.e. such that  $P^{\top}P = I$ , and  $\Lambda$  is a diagonal matrix whose diagonal elements are the eigenvalues of S (possibly repeated).

#### **Definition A.3: Positiveness of symmetric matrices**

Let  $S \in \mathbb{R}^{n \times n}$  be a symmetric matrix. We say that S is positive semi-definite (resp. positive definite) if for all  $x \in \mathbb{R}^n \setminus \{0\}$ ,  $x^{\top}Sx \ge 0$  (resp.  $x^{\top}Sx > 0$ ). We denote this property by  $S \succeq 0$  (resp.  $S \succ 0$ ).

For  $S_1, S_2 \in \mathbb{R}^{n \times n}$  two symmetric matrices, we write  $S_1 \succeq S_2$  (resp.  $S_1 \succeq S_2$ ) when  $S_1 - S_2 \succeq 0$  (resp.  $S_1 - S_2 \succeq 0$ ).

#### Proposition A.2: Characterization of positive symmetric matrices

Let  $S \in \mathbb{R}^{n \times n}$  be a symmetric matrix. The two following properties are equivalent:

$$S \succcurlyeq 0 \text{ (resp. } S \succ 0) \tag{A.5a}$$

$$\lambda \in \operatorname{sp}(S), \ \lambda \ge 0 \ (\text{resp. } \lambda > 0) \tag{A.5b}$$

#### **Proposition A.3: Inequalities for symmetric matrices**

Let  $S \in \mathbb{R}^{n \times n}$  be a symmetric matrix. Let  $\lambda_1, \lambda_2$  be two scalars. Then the three following properties are equivalent:

$$\lambda_1 I_n \preccurlyeq S \preccurlyeq \lambda_2 I_n \tag{A.6a}$$

$$\langle x \in \mathbb{R}^n, \quad \lambda_1 \| x \|^2 \le x^\top S x \le \lambda_2 dx^2$$
 (A.6b)

 $\forall \lambda \in \operatorname{sp}(S), \qquad \lambda_1 \le \lambda \le \lambda_2 \tag{A.6c}$ 

#### **Definition A.4: Orthogonal sets**

Let  $E \subset \mathbb{R}^n$  be a vectorial subspace of  $\mathbb{R}^n$ . Then, we define the orthogonal of E as:

$$E^{\perp} = \{ x \in \mathbb{R}^n \mid \forall y \in E, \ x^{\perp} y = 0 \}$$
(A.7)

Note that  $E^{\perp}$  is also a vectorial subspace of  $\mathbb{R}^n$ .

#### **Proposition A.4: Inclusion of orthogonal sets**

Then  $F^{\perp} \subset E^{\perp}$ . Let  $E, F \subset \mathbb{R}^n$  be two vectorial subspaces of  $\mathbb{R}^n$  such that  $E \subset F$ .

#### Proposition A.5: Orthogonal sets of images and kernels

Let  $A \in \mathbb{R}^{n \times m}$  be a matrix. Then the following properties hold:

$$\operatorname{Im}(A)^{\perp} = \operatorname{Ker}(A^{\top}) \tag{A.8}$$

$$\operatorname{Ker}(A)^{\perp} = \operatorname{Im}(A^{\top}) \tag{A.9}$$

## A.2 Basics of differential calculus

#### **Definition A.5: Differentiability**

Let  $\mathcal{X} \subset \mathbb{R}^n$ . Let  $f : \mathcal{X} \to \mathbb{R}$  be a function. We say that f is differentiable at  $x \in \mathbb{R}^n$  if there exists a vector  $\nabla f(x) \in \mathbb{R}^n$  such that:

$$\forall d \in \mathbb{R}^n, \quad \frac{f\left(x + \varepsilon d\right) - f\left(x\right)}{\varepsilon} \xrightarrow[\varepsilon \to 0]{} \nabla f(x)^\top d \tag{A.10}$$

The vector  $\nabla f(x)$  is called the gradient of f at x. We say that f is differentiable on  $\mathcal{X}$  if it is differentiable at all points of  $\mathcal{X}$ . If  $\nabla f(x)$  is a continuous function, we say that f is continuously differentiable.

#### **Definition A.6: Twice differentiable functions**

We say that f is twice differentiable if it is differentiable, and if there exists a matrix  $\nabla^2 f(x) \in \mathbb{R}^{n \times n}$  such that:

$$\forall x \in \mathcal{X}, \ \forall d \in \mathbb{R}^n, \quad \frac{\nabla f\left(x + \varepsilon d\right) - \nabla f\left(x\right)}{\varepsilon} \xrightarrow[\varepsilon \to 0]{\varepsilon \to 0} \nabla^2 f(x)d \tag{A.11}$$

The matrix  $\nabla^2 f(x)$  is called the Hessian of f at x. We say that f is twice continuously differentiable if  $\nabla^2 f(x)$  is continuous.

#### Theorem A.2: Schwarz's theorem

If f is twice continuously differentiable, then  $\nabla^2 f(x)$  is symmetric for all  $x \in \mathcal{X}$ .

#### Proposition A.6: First order Taylor's approximation

Assume that f is continuously differentiable. Then the following formula holds:

$$f(x+d) = f(x) + \nabla f(x)^{\top} d + r(x,d)$$
 (A.12)

where r(x, d) is defined as follows:

$$r(x,d) = \left(\int_0^1 \nabla f(x+sd) - \nabla f(x) \mathrm{d}s\right)^\top d \tag{A.13}$$

#### Proposition A.7: First-order Taylor's approximation with limits

The following holds:

$$\frac{r(x,d)}{\|d\|} \xrightarrow[d \to 0]{} 0 \tag{A.14}$$

Proposition A.8: First order Taylor's approximation for a twice differentiable function

Furthermore, if f is twice continuously differentiable, then:

$$r(x,d) = d^{\top} \left( \int_0^1 s \nabla^2 f(x+sd) \mathrm{d}s \right) d \tag{A.15}$$

### Proposition A.9: Inequality for first order Taylor's approximation

If f is twice continuously differentiable, and:

$$\forall s \in [0,1], \quad \lambda_{\min} I_n \preccurlyeq \nabla^2 f(x+sd) \preccurlyeq \lambda_{\max} I_n, \tag{A.16}$$

then:

$$\lambda_{\min} \|d\|^2 \le r(x, d) \le \lambda_{\max} \|d\|^2 \tag{A.17}$$

## A.3 Basics of topology

#### **Definition A.7: Neighborhoods**

Let  $x \in \mathbb{R}^n$ . We say that  $\mathcal{N} \subset \mathbb{R}^n$  is a neighborhood of x if there exists an  $\varepsilon > 0$  such that:

$$\forall y \in \mathbb{R}^n, \quad \|x - y\| < \varepsilon \Rightarrow y \in \mathcal{N}. \tag{A.18}$$

#### **Definition A.8: Open sets**

We say that x is in the interior of a set  $\mathcal{X}$  when it admits a neighborhood  $\mathcal{N}$  such that  $\mathcal{N} \subset \mathcal{X}$ .

#### **Definition A.9: Open sets**

We say that a set  $\mathcal{O} \subset \mathbb{R}^n$  is open when it is its own interior.

#### **Definition A.10: Closed-set**

We say that a set  $\mathcal{C} \subset \mathbb{R}^n$  is closed when its complement  $\mathbb{R}^n \setminus \mathcal{C}$  is open.

*Remark.* Interestingly, the empty set and  $\mathbb{R}^n$  are both open and closed.

#### **Proposition A.10: Characterization of closed-sets**

A set  $\mathcal{C}$  is closed if and only if when  $(x_k)_{k \in \mathbb{N}}$  in  $\mathcal{C}$  converges to  $x \in \mathbb{R}^n$ ,  $x \in \mathcal{C}$ .

**Definition A.11: Compact sets** 

We say that  $\mathcal{K} \subset \mathbb{R}^n$  is compact when is closed and bounded, i.e.

 $\exists M > 0, \quad \forall x \in \mathcal{K}, \quad \|x\| \le M \tag{A.19}$ 

#### **Definition A.12: Accumulation points**

Let  $(x_k)_{k\in\mathbb{N}}$  be a sequence in  $\mathbb{R}^n$ . We say that  $\bar{x} \in \mathbb{R}^n$  is an accumulation point of the sequence  $(x_k)_{k\in\mathbb{N}}$  if there exists an increasing sequence of integers  $k_j$  such that  $x_{k_j} \xrightarrow{j \to \infty} \bar{x}$ .

#### Theorem A.3: Bolzano-Weierstrass theorem

If  $\mathcal{K} \subset \mathbb{R}^n$  is a compact set, then any sequence  $(x_k)_{k \in \mathbb{N}}$  in  $\mathcal{K}$  has at least an accumulation point in  $\mathcal{K}$ .

# Bibliography

[1] Stephen J Wright and Benjamin Recht. *Optimization for data analysis*. Cambridge University Press, 2022.