



HAL
open science

Differentiable Optimisation: Theory and Algorithms – Part II: Algorithms

Andrea Simonetto

► **To cite this version:**

Andrea Simonetto. Differentiable Optimisation: Theory and Algorithms – Part II: Algorithms. Engineering school. ENSTA Paris, France. 2024. hal-04416966

HAL Id: hal-04416966

<https://hal.science/hal-04416966v1>

Submitted on 25 Jan 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Differentiable Optimisation: Theory and Algorithms

Part II: Algorithms

Lecture notes, version at December 1, 2023

Andrea Simonetto

Forewords

Context and Aim

This course follows naturally OPT201, which covers the theory part of continuous optimisation. OPT201 focuses on optimality conditions, convexity, and duality. In OPT202, we will look at how to use these notions to build algorithms that *solve* the problems.

In particular, **the aim of the course is to be able to answer the questions,**

1. Given an optimisation problem, which algorithm do I use to solve it?
2. Which properties and theoretical guarantees does the algorithm that I have chosen have?
3. Conversely, if I want to use a certain algorithm, which characteristics does the optimisation problem need to have?

In order to answer to these three questions, we will need to build a theory of algorithms, and ultimately understand what we really mean by *solving an optimisation problem*.

Note: The sections or subsections marked with * contain optional advanced material, which is not covered in class.

References

We will use standard references in convex optimisation and algorithms, which you can refer to.

[YN] Y. Nesterov, *Introductory Lectures on Convex Programming*, Kluiver, 2004

[BV] S. Boyd, L. Vandenberghe, *Convex Optimization*, Cambridge, 2004

I will also point out a few other ones, chapter by chapter.

Palaiseau, December 1, 2023.

Contents

1	First lecture	5
1.1	Introduction	5
1.2	Oracles, algorithms, and an impossibility theorem	5
1.3	Unconstrained optimisation	7
1.3.1	Setting	7
1.3.2	Lipschitz conditions	8
1.4	The gradient method	9
1.5	Newton's method	12
1.6	Damped Newton's method	14
1.7	Gradient vs. Newton's	15
1.8	Least-squares problems	15
1.9	References	15
1.10	Exercises	15
2	Second lecture	17
2.1	The convex problem class	17
2.1.1	Some necessary definitions	17
2.2	Standard gradient methods	19
2.2.1	Case I: Gradient for $f \in \mathcal{S}_{m,L}^{1,1}$	19
2.2.2	Case II: Gradient for $f \in \mathcal{S}_{0,L}^{1,1}$	20
2.3	Nesterov's accelerated gradient	21
2.3.1	Nesterov's alternative formulations*	22
2.4	The subgradient method	23
2.4.1	Adding strong convexity*	24
2.5	Damped Newton's method in the convex case	25
2.6	Main messages	25
2.7	References	25
2.8	Exercises	27
3	Third lecture	29
3.1	Setting	29
3.1.1	Setting: simplified	29
3.2	Splitting methods	30
3.2.1	Forward-backward splitting	30
3.3	The proximal gradient method	32
3.3.1	Prox-friendly functions	32
3.3.2	Proximal properties	33
3.3.3	Convergence of the proximal gradient method	33
3.3.4	Interpretation of the results	34
3.3.5	Numerical example	35
3.4	Duality	35
3.4.1	Equality constrained problems	35
3.4.2	The dual ascent method	36
3.4.3	Inequality and equality constrained problems	40
3.4.4	A primal-dual method	40
3.4.5	Convergence: an example	41

3.4.6	Summarising	42
3.5	References	42
3.6	Exercises	42
4	Fourth lecture	44
4.1	The full picture	44
4.1.1	Barrier functions and penalised problems	44
4.1.2	Interpretation of the penalised problem	45
4.2	A Newton's step for the penalised problem	46
4.2.1	Damping and the full method	46
4.3	Self concordance	47
4.3.1	Examples and properties	47
4.3.2	Matrix extensions	47
4.3.3	Newton's method for self-concordant functions	48
4.4	The interior-point method	49
4.4.1	Computational complexity analysis: outer	49
4.4.2	Computational complexity analysis: inner	49
4.4.3	Computational complexity analysis: complete	50
4.5	Classes of easy problems	50
4.5.1	Linear programs	51
4.5.2	Quadratic Programs: QPs	51
4.5.3	Second-Order Conic Programs: SOCPs	51
4.5.4	Semi-Definite Programs: SDPs	52
4.6	Summary	52
4.7	References	52
4.8	Exercises	53
A	Recap from OPT201	55
A.1	Optimality conditions: sets	55
A.2	Optimality conditions: equality constraints	55
A.3	Optimality conditions: complete case	56
A.4	Extension to matrix decisions	56
A.5	Useful reformulations: epigraph and Schur's complement	56
A.6	(Lagrangian) Duality	57
A.7	Determining the dual problem: examples	58
A.8	Conjugate function	59

Chapter 1

First lecture

Introduction and unconstrained optimisation (I)

1.1 Introduction

In this course, we look at generic optimisation problems of the form

$$(P) \quad \underset{\mathbf{x} \in X \subseteq \mathbf{R}^n}{\text{minimise}} \quad f(\mathbf{x}) \quad (1.1)$$

$$\text{subject to} \quad g(\mathbf{x}) \leq 0 \quad (1.2)$$

$$h(\mathbf{x}) = 0, \quad (1.3)$$

where the functions $f(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}$, $g(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}^p$, $h(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}^q$. In addition, X represents a generic closed set. As it appears, we look at continuous optimisation problems for which $\mathbf{x} \in \mathbf{R}^n$. We indicate with \mathbf{x}^* a (global) optimizer of the original problem, and by $f^* = f(\mathbf{x}^*)$ the unique (global) minimum. We also indicate the Euclidean inner product as $\langle v, u \rangle = v^\top u = u^\top v$ for two vectors $u, v \in \mathbf{R}^n$.

As we go along, we will answer the following questions:

- How do we solve the above? And what do we mean by solving?
- Can we set up **algorithms** to solve the above?
- Can we distinguish easy and hard problems in the algorithmic perspective, beyond the general divide between convex and nonconvex problems?

It is very non trivial to ask oneself the question: *what do we mean by solving optimisation problems?* Surely, on one side, you may remember that continuous optimisation problems may come with conditions that optimal points need to verify. These are the so-called KKT conditions. We will recall what they are in due course, but for now think of them as a system of non-linear equations and inequalities.

On the other side, even equipped with these KKT, how can we solve for the optimal points when we have thousands, or millions variables? What is the computational complexity associated with this solution?

As we can already see, we need to build some theory around the notion of algorithm to be able to characterise its complexity and the its inner workings. This is what we will do next.

1.2 Oracles, algorithms, and an impossibility theorem

We start by defining the main players of our new theory of algorithms. For a given optimisation problem, we denote its class as \mathcal{F} . The class is the collection of properties that that problem has. For example, \mathcal{F} can represent the class of convex problems.

We then define an oracle, denoted by \mathcal{O} , as a computational entity that can generate answers to our questions. An example is a gradient oracle, which can generate the value of the gradient of a function $\nabla_{\mathbf{x}} f(\mathbf{x})$, whenever we give it \mathbf{x} .

And finally, we denote by \mathcal{M} the method, that is the algorithm that we are going to employ to solve our optimisation problem.

Definition 1.1 (Performance) *We call performance or analytical complexity of a method \mathcal{M} for problem class \mathcal{F} using the oracle \mathcal{O} , the number of calls to the oracle to find an approximate solution to \mathcal{F} with an accuracy $\epsilon > 0$.*

The definition may be not very formal, but it will give us a clear direction for our theory. In particular we are interested in knowing how many iterations (i.e., how long do we have to wait) to obtain an approximate solution for a given optimisation problem. This question is both extremely pragmatic and practical: you want to know if you have to wait one minute, or three days. But the question is also very theoretical: we want to find theoretical results that are valid for whole problem classes.

So, how does a typical method look like? Well, a typical method is an iterative algorithm, which start with an initial guess \mathbf{x}_0 , and for all subsequent times $k \in \mathbf{N}$, does:

1. Calls the oracle \mathcal{O} at \mathbf{x}_k ;
2. Applies the method \mathcal{M} to the information up to time k and form the new test point \mathbf{x}_{k+1} ;
3. Stops if a stopping criterion is verified (the accuracy is reached), otherwise $k \leftarrow k + 1$.

Let's give some examples to fix the ideas on these new notions.

Example 1.1 *The quadratic problem*

$$\underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x}, \quad (1.4)$$

has optimality conditions: $\mathbf{Q} \mathbf{x} = -\mathbf{c}$. If \mathbf{Q} is positive semidefinite, these conditions are necessary and sufficient. Otherwise just necessary.

Suppose \mathbf{Q} is positive semidefinite. Then to find an optimal point \mathbf{x}^ , we could use a linear system solver. In this case we need an \mathcal{O} that delivers gradient and Hessian information, and we use the linear solver as a method \mathcal{M} . Which linear solvers do you know and which performance do they have?*

Example 1.2 *The generic differentiable problem*

$$\underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad f(\mathbf{x}), \quad (1.5)$$

has necessary optimality conditions: $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$.

We could use the gradient method with stepsize $\alpha_k > 0$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla_{\mathbf{x}} f(\mathbf{x}_k), \quad (1.6)$$

as a method \mathcal{M} , with a gradient oracle \mathcal{O} , and a stopping criterion $\|\nabla_{\mathbf{x}} f(\mathbf{x})\| \leq \epsilon$.

Let's get back to the development of our theory. We remind that the only information available for the method \mathcal{M} are the answers of the oracle \mathcal{O} . We distinguish, zero-order oracles, if they only answer the value of the cost $f(\mathbf{x})$; first-order oracle, if they only answer the value $f(\mathbf{x})$ and the gradient $\nabla_{\mathbf{x}} f(\mathbf{x})$; and second-order oracles, if they answer the value $f(\mathbf{x})$, the gradient $\nabla_{\mathbf{x}} f(\mathbf{x})$ and the Hessian $\nabla_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x})$.

In this context, the oracle of the first example was a second-order oracle, while the one of the second example was a first-order oracle.

Let us consider now a concrete optimisation problem and let us try to apply the formal language that we have introduced. Consider,

$$(P') \quad \underset{\mathbf{x} \in [0,1]^n}{\text{minimise}} \quad f(\mathbf{x}). \quad (1.7)$$

If $f(\mathbf{x})$ is defined by a complicated procedure, you may have only access to a zero-order oracle \mathcal{O} . This is typically the case if $f(\mathbf{x})$ is generated by a simulation.

Suppose $f(\mathbf{x})$ is of a particular Lipschitz class, for which there exists a scalar $L > 0$ that:

$$\forall \mathbf{x}, \mathbf{y} \in [0, 1]^n : |f(\mathbf{x}) - f(\mathbf{y})| \leq L \|\mathbf{x} - \mathbf{y}\|_\infty.$$

Lipschitz classes will play a big role in the following, but for now just accept the definition.

Remember that an accuracy ϵ means finding a point $\bar{\mathbf{x}}$ for which,

$$|f(\bar{\mathbf{x}}) - f(\mathbf{x}^*)| \leq \epsilon,$$

or, since we know that by optimality $f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \geq 0$, then we can also just focus on $f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) \leq \epsilon$.

So, let's try a zero-order oracle \mathcal{O} on problem (P') with functional class above.

The only method \mathcal{M} to obtain the accuracy we require is a gridding of the feasible space. In particular, a uniform gridding. It obtains an accuracy ϵ if $\|\mathbf{x} - \mathbf{y}\|_\infty \leq \frac{\epsilon}{L}$, since

$$f(\bar{\mathbf{x}}) - f(\mathbf{x}^*) = |f(\bar{\mathbf{x}}) - f(\mathbf{x}^*)| \leq L \|\mathbf{x} - \mathbf{y}\|_\infty,$$

which implies that we need $\lceil \frac{L}{\epsilon} \rceil^n$ grid points. This brings us to the first key message of this course.

Theorem 1.1 (Informal, adapted from Nesterov) *General optimisation problems are unsolvable.*

Proof. *The problem we just saw is an embodiment of a generic optimisation problem.*

To reach an accuracy ϵ , you need $\lceil \frac{L}{\epsilon} \rceil^n$ grid points.

For a very simple problem $L = 4, n = 10, \epsilon = 0.01$, you need $> 10^{26}$ calls of the oracle. Even with very performant machines you can do 10^{14} calls per second. Then you need 10^{12} seconds and therefore 32000 years. ♣

This is important, most of the problem you will write, they will be unsolvable. optimisation is hard.

In this course, we will see how to we solve problems with $n \sim 10^3 - 10^6$ and $\epsilon \sim 10^{-6}$. These problems will not be the most general ones, but they will cover a good portion of the relevant problems for society. We will see how the structure and class of the problem will play an essential role, and not just its convexity.

We start by looking at unconstrained optimisation.

1.3 Unconstrained optimisation

1.3.1 Setting

We consider now the unconstrained problem

$$(PU) \quad \underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} f(\mathbf{x}), \tag{1.8}$$

where the cost function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is not necessarily convex or differentiable, unless specified. We indicate by $C^k(\mathbf{R}^n)$ the class of functions in \mathbf{R}^n that are at least k -times differentiable.

For Problem (PU), we already know that,

- If $f \in C^1(\mathbf{R}^n)$, then $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$ is a necessary condition for optimality;
- If $f \in C^1(\mathbf{R}^n)$ and f is convex, then $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$ is a necessary and sufficient condition for optimality.

These two points are very important. Convexity defines the landscape but it does not necessarily says how hard is to solve for \mathbf{x}^* . In particular, when we have convexity, we have conditions that, if we are able to solve, deliver the optimal points. When we don't have convexity, we don't have even that!

In general, for nonconvex problems, we are happy to find stationary points.

1.3.2 Lipschitz conditions

It turns out that convexity is not the only important element for the theory of algorithms. Lipschitz continuity plays an equal (or even bigger) role too. Let us define it properly.

Definition 1.2 We label with $\mathcal{C}_L^{k,p}(\mathbf{R}^n)$, $p \leq k$, a k -differentiable function which has its p derivative L -Lipschitz continuous, i.e.,

$$\forall \mathbf{x}, \mathbf{y} : \|\nabla^p f(\mathbf{x}) - \nabla^p f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|.$$

Lipschitz continuity has a very intuitive geometric interpretation as a condition on growth, as we will see shortly.

An important class for us will be the class $\mathcal{C}_L^{1,1}$, that is the functions that are at least once differentiable and whose gradient is Lipschitz continuous. For these, we have the following result.

Lemma 1.2 (YN, Lemma 1.2.3) Let $f \in \mathcal{C}_L^{1,1}$. Then for any $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ we have:

$$|f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| \leq \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Proof. For all $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$, we have,

$$f(\mathbf{y}) = f(\mathbf{x}) + \int_0^1 \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle d\tau.$$

Hence, we can obtain the following chain of relations,

$$\begin{aligned} |f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| &= \left| \int_0^1 \langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\tau \right| \\ &\leq \int_0^1 |\langle \nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| d\tau \\ \text{by Cauchy-Schwarz} &\leq \int_0^1 \|\nabla f(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \|\mathbf{y} - \mathbf{x}\| d\tau \\ \text{by Lipschitz} &\leq \int_0^1 \tau d\tau L \|\mathbf{y} - \mathbf{x}\|^2 = \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2, \end{aligned}$$

from which the thesis follows. ♣

Geometrically, Lemma 1.2 means the following. Consider a function $f \in \mathcal{C}_L^{1,1}$. Let us fix some $\mathbf{x}_0 \in \mathbf{R}^n$ and form two quadratic functions

$$\varphi_1(\mathbf{x}) = f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 \quad (1.9)$$

$$\varphi_2(\mathbf{x}) = f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0 \rangle - \frac{L}{2} \|\mathbf{x} - \mathbf{x}_0\|^2 \quad (1.10)$$

Then, the graph of the function f is located between the graphs of $\varphi_1(\mathbf{x})$ and $\varphi_2(\mathbf{x})$.

$$\varphi_1(\mathbf{x}) \geq f(\mathbf{x}) \geq \varphi_2(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbf{R}^n.$$

This is exactly the condition of growth that we were talking about, which will make the algorithms easier or more difficult: we know upper and lower bounds on how the function changes. And these bounds are easy quadratic functions. This will be very handy.

Let's dig a bit deeper in Lipschitz conditions, since we are on the right track. We now look at the Hessian of a function $f \in \mathcal{C}_L^{2,1}(\mathbf{R}^n)$, which we indicated with $\nabla^2 f(\mathbf{x})$. For real-valued symmetric matrices like the Hessian, we can define an induced Euclidean norm, for which we know that,

$$\|\nabla^2 f(\mathbf{x})\| \leq L \iff \max_i |\text{eig}(\nabla^2 f(\mathbf{x}))| \leq L \iff -LI_n \leq \nabla^2 f(\mathbf{x}) \leq LI_n.$$

This chain of relations needs some explanations. Imposing a bound on the norm (left-side) is equivalent to impose bounds on the absolute values of the eigenvalues (eig) of the matrix (middle), which is equivalent to say that the matrix is bounded with conic inequalities \leq (right-side).

In particular, we say that $A \geq B$ for two squared symmetric matrices, if and only if $A - B \geq 0$, that is, the matrix $A - B$ is positive semidefinite.

We now have the following result.

Lemma 1.3 *Function f belongs to $\mathcal{C}_L^{2,1}(\mathbf{R}^n)$ if and only if*

$$\|\nabla^2 f(\mathbf{x})\| \leq L, \quad \forall \mathbf{x} \in \mathbf{R}^n,$$

where $\nabla^2 f(\mathbf{x})$ is the Hessian of f .

Proof. *Homework.* (Try also the generic $\mathcal{C}_L^{p+1,p}(\mathbf{R}^n)$). ♣

Lemma 1.3 is not very surprising and it is quite intuitive, given the quadratic growth condition we have established just before in Lemma 1.2.

1.4 The gradient method

We are ready to formulate and analyse our first algorithm: the gradient method. We recall that the gradient of a function is an ascent direction, so taking the opposite will bring us towards lower function values.

The **simplest gradient descent method** would use a stepsize $\alpha_k > 0$, and,

Gradient descent method

- Start with $\mathbf{x}_0 \in \mathbf{R}^n$
- Iterate $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$, $k = 0, 1, \dots$

There are different methods to choose the stepsize α_k (either a priori or online):

- Constant: $\alpha_k = \alpha$
- Vanishing, for example: $\alpha_k = \frac{\alpha}{\sqrt{k+1}}$
- Fully optimised: $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))$
- Backtracking: set $\alpha_k = 1$, $\tau \in (0, 1)$, $\beta \in (0, 1/2)$:

While $f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)) > f(\mathbf{x}_k) - \beta \alpha_k \|\nabla f(\mathbf{x}_k)\|^2$, set: $\alpha_k \leftarrow \tau \alpha_k$.

Different choices of stepsize will give rise to different properties of the gradient method. Some will not work in general, and some (like the fully optimised one) will be too expensive computationally. One that works well in most scenarios is the backtracking strategy. In this solution, we start with a large α and we refine it, so to guarantee at least a $\beta \alpha_k \|\nabla f(\mathbf{x}_k)\|^2$ reduction in cost.

What about convergence? And what is convergence anyway?

Now that we have defined an algorithm, it is important to talk about convergence, convergence rate, and ultimately performance. We know that performance means how many iterations we need to reach a certain accuracy. But what about convergence? An algorithm, like the gradient method, will deliver a sequence of points $\{\mathbf{x}_k\}_{k \in \mathbf{N}}$, and it is natural to ask whether this sequence

will converge to a limit point $\bar{\mathbf{x}}$ (possibly an optimal point). Then, it is also natural to ask about how fast this convergence goes (at least asymptotically), meaning evaluating the ratio,

$$\frac{\|\mathbf{x}_k - \bar{\mathbf{x}}\|}{\|\mathbf{x}_{k-1} - \bar{\mathbf{x}}\|},$$

or related metrics.

Let's characterise these notions for the gradient method.

Theorem 1.4 *Assuming f bounded below, the gradient descent method \mathcal{M} on unconstrained problems with cost $f \in \mathcal{C}^1(\mathbf{R}^n)$ with fully optimized α_k will generate a sequence $\{\mathbf{x}_k\}$ that converges to a stationary point $\bar{\mathbf{x}}$, so that $\nabla f(\bar{\mathbf{x}}) = 0$.*

Proof. We have $f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)) < f(\mathbf{x}_k)$, with equality only if $\nabla f(\mathbf{x}_k) = 0$, and since f is bounded below $f(\mathbf{x}_k)$ reaches a limit point for which $\nabla f(\mathbf{x}_k) = 0$. So the sequence $\{\mathbf{x}_k\} \rightarrow \bar{\mathbf{x}}$. \clubsuit

This is hardly satisfactory: typically the stepsize will vanish, i.e., $\alpha_k \rightarrow 0$, which makes the overall method very slow. Also it seems we don't have a notion of performance readily available, and not even a convergence rate result.

This is of course the worst scenario. You don't want to be in this scenario. However, in a lot of cases, for example, when dealing with training neural networks with non-differentiable activation functions (like ReLu), then you are in this case. Take it as warning: when training non-smooth machine learning models, you can wait for a long time for the result, and it can be just a stationary point, which may be quite bad.

Can we do something **much** better? Let's consider Lipschitz functions.

Theorem 1.5 *Assuming f bounded below, the gradient descent method \mathcal{M} on unconstrained problems with cost $f \in \mathcal{C}_L^{1,1}(\mathbf{R}^n)$ with backtracking α_k will generate a sequence $\{\mathbf{x}_k\}$ that converges to a stationary point as*

$$\min_{k=0, \dots, t} \|\nabla f(\mathbf{x}_k)\| \leq \frac{1}{\sqrt{t+1}} \sqrt{\frac{L}{\tau\beta} (f(\mathbf{x}_0) - f^*)} \quad . \quad (1.11)$$

Proof. Use Lipschitz and the fact that $\mathbf{x}_{k+1} - \mathbf{x}_k = -\alpha_k \nabla f(\mathbf{x}_k)$:

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}_{k+1} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \quad (1.12)$$

$$\leq f(\mathbf{x}_k) - \alpha_k \left(1 - \frac{\alpha_k L}{2}\right) \|\nabla f(\mathbf{x}_k)\|^2, \quad (1.13)$$

for which we know we need $\alpha_k < 2/L$. Which guarantees convergence to f^* .

We could stop here, but we want the rate. For backtracking line search, we have,

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \alpha_k \beta \|\nabla f(\mathbf{x}_k)\|^2. \quad (1.14)$$

Furthermore $\alpha_k = 1/L \leq 1$ is a valid choice since $\beta < 1/2$, so at least $\alpha_k \geq \frac{\tau}{L}$.

As such,

$$f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \alpha_k \beta \|\nabla f(\mathbf{x}_k)\|^2 \leq f(\mathbf{x}_k) - \frac{\tau}{L} \beta \|\nabla f(\mathbf{x}_k)\|^2. \quad (1.15)$$

Now, summing over $k = 0, \dots, t$:

$$\sum_{k=0}^t \frac{\tau}{L} \beta \|\nabla f(\mathbf{x}_k)\|^2 \leq \sum_{k=0}^t f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) = f(\mathbf{x}_0) - f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_0) - f^*, \quad (1.16)$$

and finally, since $\min_k \|\nabla f(\mathbf{x}_k)\|^2 \leq \frac{1}{t+1} \sum_{k=0}^t \|\nabla f(\mathbf{x}_k)\|^2$, the thesis follows. \clubsuit

What is the meaning of this result? We have obtained a global (starting from any point \mathbf{x}_0) convergence in terms of $\|\nabla f(\mathbf{x})\|$. The result implies that the sequence $\{\min \|\nabla f(\mathbf{x})\|\}$ converges to 0 with a certain rate.

In general for nonconvex problems it is hard to obtain global convergence in terms of objective $f(\mathbf{x}_k) - f^*$, or distance $\|\mathbf{x}_k - \mathbf{x}^*\|$.

In addition, by using an first-order oracle \mathcal{O} on method \mathcal{M} (gradient), for problem (PU) of class $\mathcal{C}_L^{1,1}$, we have obtain a **convergence rate** estimate of $O(\frac{1}{\sqrt{t+1}})$.

This means that to obtain an accuracy $\|\nabla f(\mathbf{x})\| \leq \epsilon$, we need

$$t \geq (f(\mathbf{x}_0) - f^*) \frac{L}{\tau\beta\epsilon^2} = O\left(\frac{1}{\epsilon^2}\right) \quad (1.17)$$

calls to our oracle. This is the result we wanted, a performance result, sometimes referred to as a performance certificate. This is also good: it does not depend on the problem dimension n . Look back at Theorem 1.1 and the impossibility result: that one depended on the problem dimension.

Before continuing, it is good to introduce some general notation and definitions. In general, we call results as in Eq. (1.11) as convergence certificates.

Formally, we define,

Definition 1.3 (Convergence certificate) *A convergence certificate is an inequality that upper bounds a convergence metric as a function of the number of iterations.*

Definition 1.4 (Local and global convergence) *A convergence is local if it is achieved starting close enough to the limit point. A convergence is global if it is achieved starting from any point.*

Definition 1.5 (Big-O, little-o notation) *Informally, the notation $O(\cdot)$ indicates the asymptotical behaviour (for big t 's) up to irrelevant constants. Formally, for two functions $f, g : \mathbf{R} \rightarrow \mathbf{R}$ with $g(t) > 0$ for all t , we say $f = O(g)$ iff, $\lim_{t \rightarrow \infty} \frac{|f(t)|}{g(t)} = \text{const.}$*

Similarly, we define the notation $o(\cdot)$, and we say that $f = o(g)$, if g grows much faster than f asymptotically, so that, for two functions $f, g : \mathbf{R} \rightarrow \mathbf{R}$ with $g(t) > 0$ for all t , we say $f = o(g)$ iff, $\lim_{t \rightarrow \infty} \frac{f(t)}{g(t)} = 0$.

So, in Eq. (1.11), the metric is $\min_k \|\nabla f(\mathbf{x}_k)\|$, the number of iterations is t , and the upper bound is $O(1/\sqrt{t+1})$.

Let us get back to the convergence certificate of Eq. (1.11). In general a rate of $O(\frac{1}{\sqrt{t}})$ is not great. We can do better if we look at functions, for which around an optimal point, the following extra property is verified,

$$mI_n \leq \nabla^2 f(\mathbf{x}) \leq LI_n, \quad \forall \mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq R, \quad (1.18)$$

for a given $R > 0$.

Note that the left inequality is true for strongly convex functions for all \mathbf{x} . The right inequality follows from $f \in \mathcal{C}_L^{2,1}(\mathbf{R}^n)$, so redundant. Here however we do not talk about convex functions, but only locally strongly convex ones.

We now have the result we wanted.

Theorem 1.6 *The gradient descent method \mathcal{M} on unconstrained problems, with constant step-size $\alpha < 2/L$, for cost functions $f \in \mathcal{C}_L^{2,1}(\mathbf{R}^n)$ verifying the property (1.18) above has a local (starting at $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$) convergence certificate of*

$$\|\mathbf{x}_t - \mathbf{x}^*\| \leq \rho^t \|\mathbf{x}_0 - \mathbf{x}^*\| \leq \rho^t R.$$

for $\rho = \max\{|1 - \alpha m|, |1 - \alpha L|\} < 1$, and t iterations.

Proof. Local means starting close enough to a local minimizer, so pick $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$. For this,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| = \|\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) - \mathbf{x}^*\| \quad (1.19)$$

$$= \|\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) - \mathbf{x}^* - \alpha \nabla f(\mathbf{x}^*)\| \quad (1.20)$$

Set $g(\mathbf{x}) := \mathbf{x} - \alpha \nabla f(\mathbf{x})$, for $\|\mathbf{x}_k - \mathbf{x}^*\| \leq R$ property (1.18) is valid and $\|\nabla g(\mathbf{x})\| \leq \rho$.

In fact,

$$\begin{aligned} \|\nabla g(\mathbf{x})\| &= \|I - \alpha \nabla^2 f(\mathbf{x})\| = \max |\text{eig}(I - \alpha \nabla^2 f(\mathbf{x}))| = \\ &= \max |1 - \alpha \text{eig}(\nabla^2 f(\mathbf{x}))| \leq \max\{|1 - \alpha m|, |1 - \alpha L|\}. \end{aligned} \quad (1.21)$$

For Lemma 1.3 we have also that,

$$\|g(\mathbf{x}) - g(\mathbf{y})\| \leq \rho \|\mathbf{x} - \mathbf{y}\|,$$

for which, $\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \rho \|\mathbf{x}_k - \mathbf{x}^*\|$. The choice $\alpha < 2/L$ guarantees that $\rho < 1$ and therefore, once in $\mathcal{R}(\mathbf{x}) := \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}^*\| \leq R\}$ we stay in $\mathcal{R}(\mathbf{x})$. So then, we can iterate the above inequality for t iterations, and obtain the claim. ♣

Let us analyse the result. We have a local convergence for special functions in terms of $\|\mathbf{x}_k - \mathbf{x}^*\|$ with constant stepsize. The convergence rate is exponential, or as we say in optimisation: it is a linear convergence. This nomenclature is due to the fact that in this case,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_k - \mathbf{x}^*\|}{\|\mathbf{x}_{k-1} - \mathbf{x}^*\|} = \rho,$$

so the ratio is constant.

In addition, for an accuracy $\|\mathbf{x}_t - \mathbf{x}^*\| \leq \epsilon$, we need

$$t \geq \frac{\log R + \log(1/\epsilon)}{\log(1/\rho)} = O\left(\log\left(\frac{1}{\epsilon}\right)\right)$$

calls of the oracle.

We can pick the best $\alpha = 2/(m + L)$, leading to the best $\rho = \frac{L-m}{L+m} = \frac{\kappa-1}{\kappa+1}$. The number κ is the condition number $\kappa := L/m$, and for badly conditioned functions (κ large), the gradient will not work well, since $\rho \approx 1$, and $t \rightarrow \infty$.

A linear convergence is much better than a $O(1/\sqrt{t})$ convergence. And it also turns out this is the fastest convergence result we can expect for gradient methods (we will look at these type of claims in the second class).

However, if we allow ourself the chance to use second-order oracles, we may do better.

1.5 Newton's method

We introduce now the Newton's method for functions at least in $\mathcal{C}^2(\mathbf{R}^n)$, which is a second-order (oracle) method.

Newton's method

- Start with $\mathbf{x}_0 \in \mathbf{R}^n$
- Iterate: solve the linear system $\nabla^2 f(\mathbf{x}_k) \mathbf{d}_k = \nabla f(\mathbf{x}_k)$, $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{d}_k$, $k = 0, 1, \dots$

The method involves solving a linear system of equations

$$\nabla^2 f(\mathbf{x}_k) \mathbf{d}_k = \nabla f(\mathbf{x}_k),$$

where the unknown is the Newton's step \mathbf{d}_k . Different methods will be generated by different choices of how you solve the linear system. We recall that the worst-case complexity of solving a quadratic linear system of n dimension is $O(n^3)$.

In the most naive (just in theory) version, we can set,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k).$$

The Newton's method is more computationally complex per iteration than the gradient method, which performs only vector additions, but what can we say in terms of convergence, and performance?

Consider functions $f \in \mathcal{C}_M^{2,2}(\mathbf{R}^n)$ with the additional property that

$$\nabla^2 f(\mathbf{x}) \geq mI_n, \quad \forall \mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq R, \quad (1.22)$$

for a local region $R > 0$. Then we have the following theorem.

Theorem 1.7 The Newton method applied to unconstrained problems, for cost functions $f \in \mathcal{C}_M^{2,2}(\mathbf{R}^n)$ verifying property (1.22) above, for any $R > 2m/M$, and starting in $\|\mathbf{x}_0 - \mathbf{x}^*\| < \min\{1, 2m/M\}$ has a local convergence certificate of

$$\frac{M}{2m} \|\mathbf{x}_t - \mathbf{x}^*\| \leq \left(\frac{M}{2m} \|\mathbf{x}_0 - \mathbf{x}^*\| \right)^{2^t} \leq C^{2^t},$$

when run for t iterations.

Proof. Consider the update,

$$\begin{aligned} r_{k+1} := \|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \|\mathbf{x}_k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)\| \\ &= \|\mathbf{x}_k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_k)]^{-1} \int_0^1 \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*)) (\mathbf{x}_k - \mathbf{x}^*) d\tau\| \\ &= \|[\nabla^2 f(\mathbf{x}_k)]^{-1} \times \int_0^1 [\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*))] d\tau (\mathbf{x}_k - \mathbf{x}^*)\| \end{aligned}$$

The first line is due to:

$$\nabla f(\mathbf{x}_k) = \nabla f(\mathbf{x}^*) + \int_0^1 \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*)) (\mathbf{x}_k - \mathbf{x}^*) d\tau.$$

The second line is due to:

$$\mathbf{x}_k - \mathbf{x}^* = [\nabla^2 f(\mathbf{x}_k)]^{-1} [\nabla^2 f(\mathbf{x}_k)] (\mathbf{x}_k - \mathbf{x}^*).$$

Use now Cauchy-Schwarz, the property $f \in \mathcal{C}_M^{2,2}(\mathbf{R}^n)$, and the local property $\nabla^2 f(\mathbf{x}) \geq mI_n$ to say,

$$\begin{aligned} r_{k+1} &\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1} \times \int_0^1 [\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*))] d\tau (\mathbf{x}_k - \mathbf{x}^*)\| \\ &\leq \frac{M}{m} \int_0^1 (1 - \tau) d\tau \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \frac{M}{2m} r_k^2. \end{aligned}$$

If $r_0 < 2m/M < R$, then $M/2m r_0 < 1$ and $r_1 < r_0 < 1$, so recursively $r_{k+1} < r_k$. The thesis follows considering the metric $\frac{M}{2m} r_k < 1$:

$$\frac{M}{2m} r_{k+1} \leq \left(\frac{M}{2m} r_k \right)^2 \leq \dots \leq \left(\frac{M}{2m} r_0 \right)^{2^{k+1}}.$$

♣

Let us, once again, analyse this result. First, the theorem offers a local convergence for special functions in terms of distance to the optimiser: $\|\mathbf{x}_k - \mathbf{x}^*\|$. The convergence is doubly exponential,

or as we say in optimisation: it is a quadratic convergence. In particular, to obtain an accuracy of $\|\mathbf{x}_t - \mathbf{x}^*\| \leq \epsilon$, we need,

$$t \gtrsim \log \log \frac{1}{\epsilon}.$$

calls to the oracle.

In practice this is very fast: to reach a very good accuracy $\epsilon \sim 10^{-20}$, typically 5-6 steps of Newton are enough; and this is basically independently of the problem conditioning κ . So the Newton's method works also for problems that are very badly conditioned. We will see how this is an essential feature for obtaining solvers to a large variety of problems.

As a corollary, for quadratic functions, the constant $M = 0$, so convergence is global, i.e., starting from anywhere, and it is achieved in exactly one iteration.

As for the disadvantages of the Newton's method we can cite two things. The first item is that the computational complexity per iteration is typically high. Solving a linear system involves $O(n^3)$ operations, and it may be not affordable in large-scale machine learning problem, or (in general) for any optimisation problems of reasonable size $n > 10^3$. This can be tackled in part by employing *quasi-Newton's* methods, which approximate the solution of the linear system or the Hessian. We will not discuss them here, but it is good to know that they exist.

An important message to remember here is that Newton's method is not the answer to all our problems, since it may be too computational complex, so we may need to use only first-order methods.

The second disadvantage of the Newton's method is that it is not a descent method, even in the convex case, and you have local convergence only. This is to say that, the method may diverge if not initialised properly, even in the convex case. This is however not very difficult to fix, as we see next.

1.6 Damped Newton's method

A way to fix the non-descending nature of the Newton's method is to add a stepsize to it, or as we say, we add damping. The method looks like the following,

Damped Newton's method

- Start with $\mathbf{x}_0 \in \mathbf{R}^n$
- Iterate:
 - Solve the linear system: $\nabla^2 f(\mathbf{x}_k) \mathbf{d}_k = \nabla f(\mathbf{x}_k)$,
 - Backtracking: set $\alpha_k = 1$, $\tau \in (0, 1)$, $\beta \in (0, 1/2)$:

While $f(\mathbf{x}_k - \alpha_k \mathbf{d}_k) > f(\mathbf{x}_k) - \beta \alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle$, *set:* $\alpha_k \leftarrow \tau \alpha_k$.
 - $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{d}_k$.

The backtracking step chooses a stepsize α_k to guarantee that the cost decreases of at least certain amount. In fact, when the backtracking condition is verified, then,

$$f(\mathbf{x}_k - \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) - \beta \alpha_k \langle \nabla f(\mathbf{x}_k), \mathbf{d}_k \rangle.$$

With this new relationship in place, and asking also for a Lipschitz gradient, we can establish convergence of the scheme and the fact that, after a few steps, $\alpha_k = 1$ and we get back to the standard Newton's method. This latter fact is important, and it means that damping does not alter the local quadratic convergence of the method.

By using damping, we render the Newton's method globally convergent (in the convex case we will show it in the next lesson), but it is still local in the nonconvex case and yet eventually quadratic.

While we will get back to these points in a few lessons, we close this lesson with some comparisons between gradient and Newton's.

1.7 Gradient vs. Newton's

We can summarise the methods as follows.

- For the gradient method we have derived sub-linear convergence rates $O(1/\sqrt{t})$ and linear ones $O(\rho^t)$ for different problem classes; for gradients the latter is the best you get, but you need only a first-order oracle: less computations per iteration!
- For the Newton's method we have derived quadratic convergence rates! Which is extremely fast. However Newton's require a second-order oracle and it is only a local method. We have shown how to overcome some of its drawbacks with damping.

For Newton's method, we have also seen that its convergence is nearly independent on the condition numbering κ , reaching good accuracy with a number of steps that is virtually independent from the problem. A way to peek into this property is to understand that Newton's method is **affine invariant**.

Let T be a non-singular squared matrix and define the change of coordinates $\tilde{f}(\mathbf{y}) = f(T\mathbf{y})$, $\mathbf{x} = T\mathbf{y}$. Then,

$$\nabla \tilde{f}(\mathbf{y}) = T^\top \nabla f(\mathbf{x}), \quad \nabla^2 \tilde{f}(\mathbf{y}) = T^\top \nabla^2 f(\mathbf{x}) T.$$

and the Newton's direction,

$$T\mathbf{d}_y = T[\nabla^2 \tilde{f}(\mathbf{y})]^{-1} \nabla \tilde{f}(\mathbf{y}) = \mathbf{d}_x,$$

so $T(\mathbf{y}_k - \mathbf{d}_y) = \mathbf{x}_k - \mathbf{d}_x$, meaning that "reconditioning" of the problem does not affect how Newton's behaves.

1.8 Least-squares problems

We close this lecture with an example of unconstrained problem which you will encounter over and over. For instance, when training a neural network.

Suppose you have some data that tells you the label $y_i \in \mathbf{R}$ of some feature $x_i \in \mathbf{R}^n$. You can imagine you have m of such data points (x_i, y_i) . You then build an neural network architecture, which is nothing but a series of interconnections and functions. At the end, the network is a map $h_i(\mathbf{x}; w) : x_i \mapsto y_i$, with some weights $w \in \mathbf{R}^p$ to train. The training problem is then,

$$\underset{w \in \mathbf{R}^p}{\text{minimise}} \frac{1}{2} \sum_{i=1}^m \|y_i - h_i(\mathbf{x}; w)\|_2^2,$$

and depending on the nature of h_i , the problem is easier or more complex. Typically, h_i is highly nonconvex and you have many local minima. Also, when the weights become big, your landscape can become very flat, making the progress of any method challenging.

The problem above is an example of a least-squares problem, first introduced by Gauss (aged 18) and then used by him to determine the position of the moon Ceres in ~ 1800 .

1.9 References

- YN: Introduction and Chapter 1;
- BV: Chapter 9.

◇◇◆

1.10 Exercises

Exercise 1.1 (Exam 2023) Consider the unconstrained problem,

$$\min_{x \in \mathbf{R}^n} f(x)$$

for a function $f : \mathbf{R}^n \rightarrow \mathbf{R}$, which is $\mathcal{C}^1(\mathbf{R}^n)$ and its gradient is L -Lipschitz continuous. Note that the function is **not** convex.

1. Consider the gradient method with variable stepsize $\alpha_k > 0$ to find the minimum f^* of f . Which convergence certificates can one expect (that is: which quantity converges and with which convergence rate)?
2. Let f now have the additional property that,

$$\frac{1}{2} \|\nabla f(x)\|^2 \geq \eta(f(x) - f^*), \quad \forall x \in \mathbf{R}^n,$$

for $L > \eta > 0$.

- (a) Prove that all the stationary points of $f(x)$, say \bar{x} , are global minimizers, i.e., $f(\bar{x}) = f^*$.
- (b) Prove that the gradient method with constant stepsize $\alpha = \frac{1}{L}$ offers a convergence certificate of the form,
$$f(x_k) - f^* \leq \rho^k (f(x_0) - f^*),$$
for $\rho = (1 - \eta/L) < 1$, and sequence $\{x_k\}_{k \in \mathbf{N}}$ generated by said gradient method.
- (c) How is this convergence called?

Exercise 1.2 Consider the problem

$$\min_{x \in \mathbf{R}^n} f(x),$$

for $f \in \mathcal{C}_L^{1,1}(\mathbf{R}^n)$. Consider the gradient method to find a local minimum.

Theorem. Assuming f bounded below, the gradient method on problem with $f \in \mathcal{C}_L^{1,1}(\mathbf{R}^n)$ with constant $\alpha < 2/L$ will converge to a stationary point as

$$\min_{k=0, \dots, t} \|\nabla f(x_k)\| \leq \frac{1}{\sqrt{t+1}} \sqrt{\frac{1}{\alpha(1-\alpha L/2)} (f(x_0) - f^*)}.$$

□

Prove the theorem.

Exercise 1.3 Consider the problem

$$\min_{x \in \mathbf{R}^n} f(x).$$

We study the convergence of Newton's method in the decrement,

$$\lambda(x) = (\nabla^\top f(x) \nabla^2 f(x)^{-1} \nabla f(x))^{1/2}$$

for special strongly convex functions f such that for $v = -\nabla^2 f(x)^{-1} \nabla f(x)$, $t \geq 0$:

$$(1 - t\lambda(x))^2 \nabla^2 f(x) \leq \nabla^2 f(x + tv) \leq \frac{1}{(1 - t\lambda(x))^2} \nabla^2 f(x).$$

(These functions are self-concordant, and they will be important in Class 4).

Suppose $\lambda(x_k) < 1$, and define $x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k) = x_k - v_k$. Prove that

$$\lambda(x_{k+1}) \leq \frac{\lambda^2(x_k)}{(1 - \lambda(x_k))^2}.$$

Furthermore if $\lambda(x_0) < 1/4$, prove that we can achieve quadratic convergence, with,

$$2\lambda(x_{k+1}) \leq (2\lambda(x_k))^2.$$

Hint: in both cases, assume without loss of generality that $\nabla^2 f(x_k) = I_n$.

Chapter 2

Second lecture

Unconstrained optimisation (II): convex case

2.1 The convex problem class

We now look at convex problems to understand whether they help to design algorithms that are better in terms of convergence, convergence rates, and generally computational complexity. In the interest of time, we only look here at simple methods, like the gradient (and small variations thereof), but our findings can be extended to fancier methods you will encounter in machine learning or in other fields.

We will see that convexity per se is not the only ingredient to obtain “better” results. Lipschitz continuity plays a hugely important part as well. Even though convexity will help in obtaining global optimality certificates, rather than stationarity.

Before starting, I recall that whenever we talk about convergence, we will need to be careful in saying what is converging to what exactly. In optimisation problems, we will distinguish at least three types of convergence results. The first being the fixed point residual, typically $\|\nabla_{\mathbf{x}} f(\mathbf{x}_k)\|$ converging to 0. The second is the objective convergence, $|f(\mathbf{x}_k) - f^*| \rightarrow 0$, and the third is the distance to the optimiser $\|\mathbf{x}_k - \mathbf{x}^*\| \rightarrow 0$. The former is weaker (it needs typically less assumptions and it is not a strong result), the latter is a stronger result, but harder to prove in general.

In fact, the convergence $\|\mathbf{x}_k - \mathbf{x}^*\| \rightarrow 0$ implies $|f(\mathbf{x}_k) - f^*| \rightarrow 0$, which in turn implies $\|\nabla_{\mathbf{x}} f(\mathbf{x}_k)\| \rightarrow 0$. But the converse is not true. For instance you may have multiple optimisers \mathbf{x}^* and, even though $|f(\mathbf{x}_k) - f^*| \rightarrow 0$, the sequence \mathbf{x}_k may not have a limit point. Also, you may have $\|\nabla_{\mathbf{x}} f(\mathbf{x}_k)\| \rightarrow 0$, but you are at a stationary point, not an optimiser, so $|f(\mathbf{x}_k) - f^*|$ does not converge to 0.

We are now ready to start.

2.1.1 Some necessary definitions

Definition 2.1 (Convex function) *A function $f : X \subseteq \mathbf{R}^n \rightarrow \mathbf{R}$ is convex iff X is convex and*

$$(C1) \quad \forall \mathbf{x}, \mathbf{y} \in X, \lambda \in [0, 1] : \quad f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

Multiple definitions exist if we have additional properties, for example:

$$(C1) + f \in \mathcal{C}^1(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle; \quad (2.1)$$

$$(C1) + f \in \mathcal{C}^2(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, \nabla^2 f(\mathbf{x}) \geq 0. \quad (2.2)$$

A proof of the relations above can be found in YN, by combining his Definition 2.1.1, with his Theorems 2.1.2 and 2.1.4.

Definition 2.2 (Strongly convex function) A convex function $f : X \subseteq \mathbf{R}^n \rightarrow \mathbf{R}$ is m -strongly convex iff

$$(SC) \quad f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2 \text{ is convex.}$$

We note here that f does not need to be differentiable to be strongly convex. Here too, multiple definitions exist if we have additional properties, for example:

$$(SC) + f \in \mathcal{C}^1(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq m\|\mathbf{x} - \mathbf{y}\|^2 \quad (2.3)$$

$$(SC) + f \in \mathcal{C}^2(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, \nabla^2 f(\mathbf{x}) \geq mI_n \quad (2.4)$$

A proof of the relations above can be found in YN, by using Theorem 2.1.8, 2.1.9, and 2.1.10.

Definition 2.3 (Smooth function) A convex function $f : X \subseteq \mathbf{R}^n \rightarrow \mathbf{R}$ is L -smooth iff

$$(LC) \quad \frac{L}{2}\|\mathbf{x}\|^2 - f(\mathbf{x}) \text{ is convex.}$$

Importantly $(LC) \implies f \in \mathcal{C}^1(X)$, that is a condition on the function implies its differentiability. Since this is rather important, we give a short proof here.

Proof. Define $g(x) = \frac{L}{2}\|x\|^2 - f(x)$, by the definition of subgradient of a convex function:

$$g(x) - g(x') \geq u^\top(x - x'), \quad \forall x, x', u \in \partial g,$$

where ∂g is the subdifferential of g . Furthermore, by definition of g and the convexity of f ,

$$\frac{L}{2}\|x\|^2 - \frac{L}{2}\|x'\|^2 \geq f(x) - f(x') + u^\top(x - x') \geq (u + v)^\top(x - x'), \quad \forall x, x', u \in \partial g, v \in \partial f,$$

where we have indicated with ∂f the subdifferential of f .

Since the function $\frac{L}{2}\|x\|^2$ is convex and differentiable, then we know that its subdifferential at x' is a singleton and it is equivalent to Lx' . Then, $u + v$ is a singleton and $u + v = Lx'$. Furthermore since u and v are elements of the subdifferentials and can be chosen independently, then u and v separately are also singletons. That is, f is differentiable. So $(LC) \implies f \in \mathcal{C}^1(\mathbf{R}^n)$. ♣

As before, depending on additional properties, multiple definitions of smoothness exist, for example:

$$(LC) \iff \forall \mathbf{x}, \mathbf{y} \in X, \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\| \quad (2.5)$$

$$(LC) \iff \forall \mathbf{x}, \mathbf{y} \in X, \frac{1}{L}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \quad (2.6)$$

$$(LC) + f \in \mathcal{C}^2(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, 0 \leq \nabla^2 f(\mathbf{x}) \leq LI_n \quad (2.7)$$

We see already that (2.5) is a Lipschitz condition on the gradient of f . Since our definition is equivalent to (2.5), often it is the latter than is used as a definition of smoothness.

The proof of the relations above can be found in YN, e.g., Theorem 2.1.5 and related.

We look here at the m -strongly convex L -smooth class, which is at least one time differentiable. To echo the class $\mathcal{C}_L^{p,q}$, we refer to the convex classes as

$$f \in \mathcal{S}_{m,L}^{p,q} \iff \{f \in \mathcal{C}_L^{p,q}\} \cap \{f \text{ is } m\text{-strongly convex}\}.$$

We are especially interested in the function class $f \in \mathcal{S}_{m,L}^{1,1}$. If $f \in \mathcal{S}_{0,L}^{1,1}$, the function is just convex. We also use a slight abuse of notation to indicate with $f \in \mathcal{S}_{m,+\infty}^{1,1}$ when the function is not Lipschitz and not differentiable. However, at least here, unless explicitly said, $0 < m \leq L < +\infty$. Finally, we remind that $m \leq L$ always and $\kappa := L/m \geq 1$.

Even though there are plenty of definitions, for this course, you will need to remember at least the following,

$$\forall \mathbf{x}, \mathbf{y} \in \mathbf{R}^n, \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq m \|\mathbf{x} - \mathbf{y}\|^2 \quad (2.8)$$

$$\mathbf{x}, \mathbf{y} \in \mathbf{R}^n, \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\|, \quad (2.9)$$

which will carry you a long way.

Geometrically, a function in the class $\mathcal{S}_{m,L}^{1,1}$ is a function that in all its domain grows faster than a quadratic function with concavity m but slower than a quadratic of concavity L . This is rather easy to see for doubly differentiable functions, for which $mI_n \leq \nabla^2 f(\mathbf{x}) \leq LI_n$, but it is also valid in general. In this context, the class $\mathcal{S}_{m,L}^{1,1}$ impose conditions on growth, which will be key in establishing efficient algorithms.

We have now all the tools to characterise convergence of our methods.

2.2 Standard gradient methods

2.2.1 Case I: Gradient for $f \in \mathcal{S}_{m,L}^{1,1}$

We look at the unconstrained problem

$$\underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} f(\mathbf{x}),$$

when $f \in \mathcal{S}_{m,L}^{1,1}$ ($0 < m \leq L < +\infty$), and we use the gradient method with constant stepsize $\alpha > 0$.

We already have a result for it (Theorem 1.6), and in particular:

Theorem 2.1 *The gradient descent method on unconstrained problems with constant stepsize $\alpha < 2/L$, for functions $f \in \mathcal{S}_{m,L}^{1,1}$ has a (global) convergence certificate of*

$$\|\mathbf{x}_t - \mathbf{x}^*\| \leq \rho^t \|\mathbf{x}_0 - \mathbf{x}^*\| \quad (2.10)$$

for $\rho = \max\{|1 - \alpha m|, |1 - \alpha L|\}$, and t iterations.

Proof. *The proof follows from the proof of Theorem 1.6, substituting the fact that the local property (1.18) is now valid for the whole space. Particular attention has to be put on the fact that f is now not doubly differentiable, so one has to prove $\|g(\mathbf{x}) - g(\mathbf{y})\| \leq \rho \|\mathbf{x} - \mathbf{y}\|$ by other means. How? ♣*

Can we do better in the convex case in terms of rate, more than just global? No we can't. The reasoning is rather technical, but there exist a function $g \in \mathcal{S}_{m,L}^{1,1}$ for which we have

$$\|\mathbf{x}_t - \mathbf{x}^*\| \geq \omega^t \|\mathbf{x}_0 - \mathbf{x}^*\|, \quad 0 < \omega \leq \rho.$$

This result is very strong: it tells you that in the worst case, linear convergence for $f \in \mathcal{S}_{m,L}^{1,1}$ is the best you can achieve. We will not prove this result here, but the interested reader can have a look at YN.

We can derive how other metrics (objective, fixed point) converge as well, and you will prove them in the exercises. In particular, for $f \in \mathcal{S}_{m,L}^{1,1}$, you will also derive similar linear converging conditions for $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ and $\|\nabla f(\mathbf{x}_t)\|$, completing the family picture for the gradient method.

Furthermore, one can show that there exist optimal stepsize and rates, which are, $\alpha^* = \frac{2}{m+L}$ and $\rho^* = \frac{\kappa-1}{\kappa+1}$.

2.2.2 Case II: Gradient for $f \in \mathcal{S}_{0,L}^{1,1}$

Let's move on to a non-strongly convex function class. In this case, there may be multiple solutions in terms of \mathbf{x}^* , all equally optimal.

Theorem 2.2 *The gradient descent method on unconstrained problems with constant stepsize $\alpha < 2/L$, for functions $f \in \mathcal{S}_{0,L}^{1,1}$ has an ergodic (global) convergence certificate of*

$$f(\bar{\mathbf{x}}_t) - f(\mathbf{x}^*) \leq \frac{C}{t+1} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad (2.11)$$

for $\bar{\mathbf{x}}_t = \frac{1}{t+1} \sum_{k=0}^t \mathbf{x}_k$, constant $C > 0$, and t iterations.

Proof. Start with the update rule,

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle + \|\alpha \nabla f(\mathbf{x}_k)\|^2 \end{aligned}$$

Use $\nabla f(\mathbf{x}^*) = 0$ to say that the right-hand side (RHS) is

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha \langle \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*), \mathbf{x}_k - \mathbf{x}^* \rangle + \alpha^2 \|\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*)\|^2.$$

Now use Eq. (2.6) to say,

$$RHS \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \alpha(2 - \alpha L) \langle \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*), \mathbf{x}_k - \mathbf{x}^* \rangle.$$

For $\alpha < 2/L$ the term $(2 - \alpha L)$ is positive and the error $\|\mathbf{x}_k - \mathbf{x}^*\|^2$ decreases.

This gives convergence, but not the rate.

Use convexity to say $f(\mathbf{x}^*) \geq f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x}^* - \mathbf{x}_k \rangle$ and hence: $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{x}^* \rangle$.

Therefore,

$$RHS \leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \alpha(2 - \alpha L)[f(\mathbf{x}_k) - f(\mathbf{x}^*)].$$

Sum over k to arrive at

$$\sum_{k=0}^t f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{1}{\alpha(2 - \alpha L)} \left(\sum_{k=0}^t -\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 + \|\mathbf{x}_k - \mathbf{x}^*\|^2 \right) \leq \underbrace{\frac{1}{\alpha(2 - \alpha L)}}_{=:C} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

Now we use Jensen's inequality valid for convex f : $f(\sum_i \beta_i \mathbf{x}_i) \leq \sum_i \beta_i f(\mathbf{x}_i)$ for $\sum_i \beta_i = 1$, to say,

$$f(\bar{\mathbf{x}}_t) - f(\mathbf{x}^*) \leq \frac{1}{t+1} \sum_{k=0}^t f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{C}{t+1} \|\mathbf{x}_0 - \mathbf{x}^*\|^2,$$

from which the thesis. ♣

Let us analyse the result. Theorem 2.2 gives us a global convergence result in terms of ergodic objective, with a rate of $O(1/t)$. An ergodic objective, or ergodic mean, is the fact that we are using $\bar{\mathbf{x}}_t$ instead of \mathbf{x}_t . Sometimes this helps to render the proofs easier. With a little more effort, we could obtain an objective convergence $f(\mathbf{x}_t) - f^* \leq O(1/t)$, but we do not look into that here.

An objective convergence of this type in the convex case is better than in the nonconvex case, for which we had a fixed point convergence to a stationary point with rate $O(1/\sqrt{t})$ (imposing Lipschitz).

Can we do better? Yes, actually the attainable complexity bound for the class $\mathcal{S}_{0,L}^{1,1}$ is $O(1/t^2)$. But how? The answer to this question will make us study one of the most important algorithms in optimisation of the past forty years: Nesterov's method.

2.3 Nesterov's accelerated gradient

Let's set the stage by defining some supporting sequences, $\{\lambda_k\}$ and $\{\gamma_k\}$ defined as,

$$\lambda_0 = 0, \lambda_k = \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2}, \text{ and } \gamma_k = \frac{1 - \lambda_k}{\lambda_{k+1}}. \quad (2.12)$$

Note already that $\gamma_k \leq 0$, $\lambda_1 = 1$, and $\lambda_k > 1$ for $k > 1$.

Now the algorithm is defined by the following equations,

Nesterov's accelerated gradient descent

- Start with $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbf{R}^n$ and the sequences (2.12)
- Iterate: $\mathbf{y}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)$, $\mathbf{x}_{k+1} = (1 - \gamma_k)\mathbf{y}_{k+1} + \gamma_k \mathbf{y}_k$, $k = 0, 1, \dots$

In other words, Nesterov's accelerated gradient descent performs a simple step of gradient descent to go from \mathbf{x}_k to \mathbf{y}_{k+1} , and then it 'slides' a little bit further than \mathbf{y}_{k+1} in the direction given by the previous point \mathbf{y}_k .

Theorem 2.3 *Nesterov's accelerated gradient method on unconstrained problems, with constant stepsize $\alpha = 1/L$, for functions $f \in \mathcal{S}_{0,L}^{1,1}$ has a (global) convergence certificate of*

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{2L}{t^2} \|\mathbf{x}_1 - \mathbf{x}^*\|^2, \quad (2.13)$$

for t iterations.

Proof. Start from convexity, Lipschitz, and $\alpha = 1/L$:

$$\begin{aligned} & f(\mathbf{x} - \alpha \nabla f(\mathbf{x})) - f(\mathbf{y}) \\ & \leq f(\mathbf{x} - \alpha \nabla f(\mathbf{x})) - f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{y}) \\ & \leq \nabla f(\mathbf{x})^\top (\mathbf{x} - \alpha \nabla f(\mathbf{x}) - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \alpha \nabla f(\mathbf{x}) - \mathbf{x}\|^2 + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{y}) \\ & = -\frac{1}{2L} \|\nabla f(\mathbf{x})\|^2 + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{y}). \end{aligned}$$

Now let us apply the inequality to $\mathbf{x} = \mathbf{x}_k$ and $\mathbf{y} = \mathbf{y}_k$, which gives

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

Similarly we apply it to $\mathbf{x} = \mathbf{x}_k$ and $\mathbf{y} = \mathbf{x}^*$ which gives

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

Now multiplying the first by $(\lambda_k - 1)$ which is positive for $k > 2$ and adding the result to the second, one obtains with $\delta_k = f(\mathbf{y}_k) - f(\mathbf{x}^*)$,

$$\lambda_k \delta_{k+1} - (\lambda_k - 1) \delta_k \leq -\frac{L}{2} \lambda_k \|\mathbf{y}_{k+1} - \mathbf{x}_k\|^2 - L(\mathbf{y}_{k+1} - \mathbf{x}_k)^\top (\lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*).$$

Multiplying the last inequality by $\lambda_k > 0$ and by using that by definition $\lambda_{k-1}^2 = \lambda_k^2 - \lambda_k$ one obtains

$$\begin{aligned} & \lambda_k^2 \delta_{k+1} - \lambda_{k-1}^2 \delta_k \\ & \leq -\frac{L}{2} \left(\|\lambda_k (\mathbf{y}_{k+1} - \mathbf{x}_k)\|^2 + 2\lambda_k (\mathbf{y}_{k+1} - \mathbf{x}_k)^\top (\lambda_k \mathbf{x}_k - (\lambda_k - 1) \mathbf{y}_k - \mathbf{x}^*) \right). \end{aligned}$$

We remark that $\lambda_{k-1}^2 = \lambda_k^2 - \lambda_k$ follows from $\lambda_k = \frac{1 + \sqrt{1 + 4\lambda_{k-1}^2}}{2}$, in fact the latter implies,

$$(2\lambda_k - 1)^2 = 1 + 4\lambda_{k-1}^2 \iff \lambda_{k-1}^2 = \lambda_k^2 - \lambda_k.$$

Now one can verify that

$$\begin{aligned} & \|\lambda_k(\mathbf{y}_{k+1} - \mathbf{x}_k)\|^2 + 2\lambda_k(\mathbf{y}_{k+1} - \mathbf{x}_k)^\top(\lambda_k\mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^*) \\ &= \|\lambda_k\mathbf{y}_{k+1} - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^*\|^2 - \|\lambda_k\mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^*\|^2. \end{aligned}$$

Next remark that, by definition of γ_k , one has

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{y}_{k+1} + \gamma_k(\mathbf{y}_k - \mathbf{y}_{k+1}) \\ &\iff \lambda_{k+1}\mathbf{x}_{k+1} = \lambda_{k+1}\mathbf{y}_{k+1} + (1 - \lambda_k)(\mathbf{y}_k - \mathbf{y}_{k+1}) \\ &\iff \lambda_{k+1}\mathbf{x}_{k+1} - (\lambda_{k+1} - 1)\mathbf{y}_{k+1} = \lambda_k\mathbf{y}_{k+1} - (\lambda_k - 1)\mathbf{y}_k. \end{aligned}$$

Putting together the previous relationships one gets with $u_k = \lambda_k\mathbf{x}_k - (\lambda_k - 1)\mathbf{y}_k - \mathbf{x}^*$,

$$\lambda_k^2\delta_{k+1} - \lambda_{k-1}^2\delta_k \leq \frac{L}{2} \left(\|u_k\|^2 - \|u_{k+1}\|^2 \right).$$

Summing these inequalities from $k = 1$ to $k = t - 1$ (rem: $\lambda_0 = 0, \lambda_1 = 1$) one obtains:

$$\delta_t \leq \frac{L}{2\lambda_{t-1}^2} \|u_1\|^2.$$

By induction: $\lambda_{t-1} \geq \frac{t}{2}$ which concludes the proof. ♣

What is the meaning of this convergence result? We have a global convergence result in terms of objective, with a rate of $O(1/t^2)$, which matches the lower bound. In this context, we say that Nesterov's accelerated method is an optimal method. This is way better than in the nonconvex case, for which we had a fixed point convergence to a stationary point with rate $O(1/\sqrt{t})$.

If one is attentive, one can see that we have proved a result for a specific $\alpha = 1/L$. However, a one-line proof can show that we can extend the theorem to any $\alpha \leq 1/L$, how?

Corollary 2.1 *Theorem 2.3 is valid for $\alpha \leq 1/L$.*

Proof. *Homework.* ♣

2.3.1 Nesterov's alternative formulations*

Nesterov's accelerated scheme is only one of the few accelerations that are possible for solving smooth convex problems. Moreover, Nesterov's acceleration may be presented in a few alternative (but almost equivalent) forms. In fact, in the proof of the method, we use the fact that the coefficients $\lambda_{k-1}^2 = \lambda_k^2 - \lambda_k$. However, for the rate to be proven, one would only need the looser condition, $-\lambda_{k-1}^2 \leq -(\lambda_k^2 - \lambda_k)$, which leads to $\lambda_{k-1}^2 \geq \lambda_k^2 - \lambda_k$. This yields a few different alternatives.

A popular choice is to set $\lambda_k = \frac{k}{2}$, this verifies,

$$\lambda_{k-1}^2 \geq \lambda_k^2 - \lambda_k \iff (k-1)^2 \geq k^2 - 2k \iff 1 \geq 0.$$

With this choice $\gamma_k = \frac{2-k}{k+1}$. So that the method is presented as follows.

Nesterov's accelerated gradient descent (variant)

- Start with $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbf{R}^n$
- Iterate: $\mathbf{y}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k), \quad \mathbf{x}_{k+1} = \mathbf{y}_{k+1} + \frac{k-2}{k+1}(\mathbf{y}_{k+1} - \mathbf{y}_k), \quad k = 0, 1, \dots$

2.4 The subgradient method

Let's now look at the case in which we lose differentiability. We let $f \in \mathcal{S}_{m,+\infty}^{1,1}$. In this case, since $\nabla f(\mathbf{x})$ does not exist, we need another notion, the one of subgradient $\partial f(\mathbf{x})$.

Definition 2.4 (Subdifferential) We define the subdifferential at \mathbf{x} as the set

$$\partial f(\mathbf{x}) := \{\mathbf{g} \in \mathbf{R}^n \mid f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle, \forall \mathbf{y} \in \mathbf{R}^n\}.$$

You have studied in details the subdifferential and the subgradient in OPT201. For the sake of this course, it is useful to see the subdifferential as a set of all the vectors that generate hyperplanes that stay “below” the function at any point. You can also see the parallel between the definition of subdifferential and the definition of convex function. In particular, if the function f is convex and differentiable, then, the subdifferential is a singleton: $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$.

You may encounter many definitions of subdifferential in the nonconvex domain. Luckily for us, in the convex domain they are equivalent to our definition. Finally, in many algorithms, we don't even need to get the whole subdifferential set, just one element $\mathbf{g} \in \partial f(\mathbf{x})$ will be enough!

With this new notion in place, we are ready for the subgradient method as follows.

Select a stepsize $\alpha_k > 0$:

Subgradient method

- Start with $\mathbf{x}_0 \in \mathbf{R}^n$
- Iterate: find a $\mathbf{g}_k \in \partial f(\mathbf{x}_k)$, compute $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{g}_k$, $k = 0, 1, \dots$

We are now also ready for its convergence certificate under the additional condition that the subgradient is bounded, as

$$f \in \mathcal{S}_{0,G}^{0,0}, \text{ or equivalently } \|\mathbf{g}\| \leq G, \text{ or equivalently } \|f(\mathbf{x}) - f(\mathbf{y})\| \leq G\|\mathbf{x} - \mathbf{y}\| \quad \forall \mathbf{x}, \mathbf{y}.$$

Theorem 2.4 The subgradient method on unconstrained problems with stepsize satisfying the conditions,

$$\lim_{k \rightarrow \infty} \alpha_k \rightarrow 0, \quad \sum_{k=0}^{\infty} \alpha_k = \infty,$$

for functions $f \in \mathcal{S}_{0,G}^{0,0}$ has a (global) convergence certificate of

$$\lim_{t \rightarrow \infty} \min_{k=0, \dots, t} f(\mathbf{x}_k) - f(\mathbf{x}^*) = 0. \quad (2.14)$$

Furthermore, at best, $\min_{k=0, \dots, t} f(\mathbf{x}_k) - f(\mathbf{x}^*) = O(1/\sqrt{t})$.

Valid conditions on α_k are, e.g., $\alpha_k = \frac{1}{k+1}$, $\alpha_k = \frac{1}{\sqrt{k+1}}$, ...

Proof. [for the choice of $\alpha_k = \frac{1}{k+1}$] We start from the algorithm definition and convexity,

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha_k \langle \mathbf{g}_k, \mathbf{x}_k - \mathbf{x}^* \rangle + \alpha_k^2 \|\mathbf{g}_k\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha_k (f(\mathbf{x}_k) - f(\mathbf{x}^*)) + \alpha_k^2 G^2. \end{aligned}$$

Summing over k , we obtain,

$$2 \sum_k \alpha_k (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \leq \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \sum_k \alpha_k^2 G^2.$$

Since $\sum_k \alpha_k (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \geq (\sum_k \alpha_k) \min_k (f(\mathbf{x}_k) - f(\mathbf{x}^*))$, then,

$$\min_k (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \leq \frac{1}{2 \sum_k \alpha_k} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{\sum_k \alpha_k^2}{2 \sum_k \alpha_k} G^2 \rightarrow 0.$$

The $\rightarrow 0$ is not immediately evident and requires some extra work. It is obvious if, e.g., $\alpha_k = \frac{1}{k+1}$.

In particular, the selection $\alpha_k = \frac{1}{k+1}$, we obtain the lowest rate bound,

$$\min_k (f(\mathbf{x}_k) - f(\mathbf{x}^*)) \sim O\left(\frac{1}{\log t}\right) \gtrsim O\left(\frac{1}{\sqrt{t}}\right)$$

where we recall that $O(\log t) < O(\sqrt{t})$. ♣

We have a global convergence result in terms of best objective, with a rate of $O(1/\sqrt{t})$. This is the rate of nonconvex fixed point residual in the Lipschitz case, so it's better in the convex case, but it's very slow compared to the convex Lipschitz case, for which we found $O(1/t^2)$.

In the next class, we will see how to avoid the use of the subgradient method (in some cases) by the use of splitting operator and restore reasonable convergence rates.

2.4.1 Adding strong convexity*

One may want to see if adding strong convexity to a non-differentiable cost f can change the convergence certificates we obtain. This is not the case, since adding strong convexity does not change the differentiability of the cost. So in general, you still have a $O(1/\sqrt{t})$ convergence guarantee.

However, adding strong convexity, that is, considering functions $f \in \mathcal{S}_{m,+\infty}^{1,1}$, can help the subgradient method to converge linearly to an error ball. In this sense, the subgradient method obtains two phases, a fast phase up to an error ball, and a slow phase, to the true optimizer.

Let's see how to show that. First, we need a improved strong convexity characterisation.

Lemma 2.5 *A function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is m -strongly convex iff $f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2$ is convex, or equivalently*

$$(\mathbf{v} - \mathbf{w})^\top (\mathbf{x} - \mathbf{y}) \geq m\|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{v} \in \partial f(\mathbf{x}), \mathbf{w} \in \partial f(\mathbf{y}),$$

for all $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$.

You can find the proofs in many standard references, such as E. K. Ryu and S. Boyd, *Primer on Monotone Operator Methods*, 2016.

Then, we can prove that,

Theorem 2.6 *The subgradient method on unconstrained problems with constant stepsize $\alpha < 1/2m$ for functions $f \in \mathcal{S}_{m,G}^{0,0}$ has a (global) convergence certificate of*

$$\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}^*\| = \sqrt{\frac{\alpha}{2m}} G. \quad (2.15)$$

Proof. *The function is strongly convex, so we have an unique optimiser \mathbf{x}^* . We start from the algorithm definition,*

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha \langle \mathbf{g}_k, \mathbf{x}_k - \mathbf{x}^* \rangle + \alpha^2 \|\mathbf{g}_k\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha \langle \mathbf{g}_k - \mathbf{g}^*, \mathbf{x}_k - \mathbf{x}^* \rangle + \alpha^2 G^2 \\ &\quad \text{since } \mathbf{g}^* \in \partial f(\mathbf{x}^*), \text{ and we can take } \mathbf{g}^* = 0 \\ &\leq (1 - 2m\alpha) \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \alpha^2 G^2 \quad \text{by strong convexity.} \end{aligned}$$

Then, since $\alpha < 1/2m$, by geometric recursion,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq (1 - 2m\alpha)^k \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{\alpha G^2}{2m},$$

from which the thesis follows. ♣

From the proof, we can see how the error $\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2$ decreases exponentially, up to an error floor $\frac{\alpha G^2}{2m}$. We can also see that if we want to reduce the error floor, we need to take α small, but this renders $(1 - 2m\alpha) \approx 1$, which damps convergence. This is a rather typical behaviour.

2.5 Damped Newton's method in the convex case

To close this chapter, we revisit the damped Newton's method we have seen in the previous class, and we provide a convergence certificate via the following theorem.

Theorem 2.7 *Apply the damped Newton's method on unconstrained problems, for cost functions $f \in \mathcal{S}_{m,M}^{2,2}$, and generate the sequence $\{\mathbf{x}_k\}$. Then, there exist two positive scalars η and γ with $0 < \eta \leq m^2/M$ and $\gamma > 0$ such that*

- If $\|\nabla f(\mathbf{x}_k)\| \geq \eta$, then

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\gamma.$$

- If $\|\nabla f(\mathbf{x}_k)\| < \eta$, then the backtracking line search selects $\alpha_k = 1$ and

$$\frac{M}{2m^2} \|\nabla f(\mathbf{x}_{k+1})\| \leq \left(\frac{M}{2m^2} \|\nabla f(\mathbf{x}_k)\| \right)^2.$$

As such, we obtain an accuracy of ϵ , i.e., $f(\mathbf{x}_k) - f^* \leq \epsilon$, in,

$$\frac{f(\mathbf{x}_0) - f^*}{\gamma} + \log_2 \log_2 \left(\frac{2m^3}{M^2\epsilon} \right)$$

iterations.

Proof. The proof can be found in [BV] Section 9.5.3. ♣

Theorem 2.7 tells us that the application of the damping strategy renders the Newton's method global. It can therefore converge from any initial points. The initial converge can be slow, but it is at least constant:

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\gamma,$$

for all k 's. Then, when we are close to an optimiser, the quadratic convergence phase kicks in, and we converge very fast.

2.6 Main messages

We finish this lesson with a recap table and a figure. The table captures all the results we have presented in this chapter. Your goal is to learn this table and understand what each entry means. In particular, for unconstrained problems, you should know which method to use in which situations, with which convergence certificate and performance.

For the sake of completeness, we present here also a figure to give more tangible intuitions on how slow or how fast certain methods can be.

Figure 2.1 depicts the different order of convergence rate if we could put them all on the same plot, for the same error measure and number of iterations. This graph tells you that $O(1/\sqrt{t})$ is unbelievably slow, reducing the error of one order of magnitude each *two* orders of magnitude iterations. Read it again. To get one more digit right, you need two order of magnitude more iterations!

On the other hand of the spectrum, the Newton's method is incredibly fast: the number of correct digits of your optimiser roughly *doubles* at each step! Needless to say: Go Newton!

2.7 References

- YN: Introduction and Chapter 2
- BV: Chapter 9
- S. Bubeck's blog *I'm a bandit — The complexities of optimization*, post 12, 13, 14, 15. [\[link to blog\]](#)

◇ ◆ ◇

Table 2.1. Summary of the certificates results obtained for nonconvex and convex unconstrained problems depending on the type of method and function class.

Function type	generic	$\mathcal{C}_L^{1,1}, \mathcal{S}_{0,L}^{1,1}$	$\mathcal{C}_{m,L}^{1,1}, \mathcal{S}_{m,L}^{1,1}$	$\mathcal{C}_{m,M}^{2,2}, \mathcal{S}_{m,M}^{2,2}$
		First-order		Second-order
nonconvex	convergence $f \in \mathcal{C}^1$ Theorem 1.4	$O(1/\sqrt{t})$ Theorem 1.5	locally $O(\rho^t)$ Theorem 1.6	locally $O(\rho^{2t})$ Theorem 1.7
Convex	$O(1/\sqrt{t})$ Theorem 2.4	$O(1/t), O(1/t^2)$ \rightsquigarrow Nesterov's Theorem 2.2-2.3	$O(\rho^t)$ Theorem 2.1 + Exercise 2.2	globally: const., locally $O(\rho_2^{2t})$ Theorem 2.7

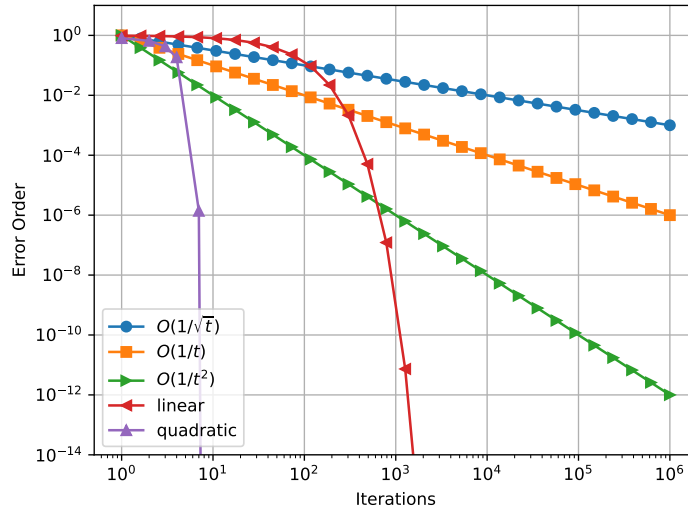


Figure 2.1. The intuition behind the convergence rates.

2.8 Exercises

Exercise 2.1 (Resit 2023) We aim at minimizing the following quadratic problem,

$$\min_{x \in \mathbf{R}^n} \frac{1}{2} x^\top A x + b^\top x + c,$$

for A symmetric and positive definite, b a vector, and c a scalar. Assume that $\ell I_n \leq A \leq L I_n$, for two scalars $0 < \ell \leq L < +\infty$.

1. Is the problem convex? Is the problem smooth? Is the problem strongly convex? Let x^* be an optimizer of the above problem. Prove that it is unique and find its value.
2. To solve the problem above, we look at the following variation of a gradient scheme. Start with random x_1, x_0 and define two stepsizes $\alpha, \beta > 0$, then for $k > 1$ do:

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

Prove that,

$$\begin{bmatrix} x_{k+1} - x^* \\ x_k - x^* \end{bmatrix} = \underbrace{\begin{bmatrix} (1 + \beta)I_n - \alpha A & -\beta I_n \\ I_n & 0 \end{bmatrix}}_{=:T} \begin{bmatrix} x_k - x^* \\ x_{k-1} - x^* \end{bmatrix}$$

3. Prove that the above algorithm converges (in the sense that $\|x_k - x^*\| \rightarrow 0$) if the eigenvalues of T are inside the unit ball, meaning they are all in modulus less than 1.
4. Let $n = 1$, $A = \ell$, $\beta = 0$. Find the allowed range for the stepsize α to ensure convergence. In this case, determine also the rate of convergence.

Exercise 2.2 The aim of this exercise is to derive linear convergence guarantees.

Consider the problem,

$$x^* \in \arg \min_{x \in \mathbf{R}^n} f(x),$$

for a μ -strongly convex and L -smooth function.

Consider the gradient method with constant stepsize $\alpha > 0$ to find the unique minimizer of f .

In class, we have seen that,

$$\|x_k - x^*\| \leq \rho^k \|x_0 - x^*\|,$$

for all $\alpha < 2/L$. Start from it and,

- Prove that,

$$|f(x_k) - f(x^*)| \leq \kappa \rho^{2k} |f(x_0) - f(x^*)|,$$

- Prove also that,

$$\|\nabla f(x_k)\| \leq \kappa \rho^k \|\nabla f(x_0)\|, \quad \kappa = L/m,$$

which shows the effect of the condition number explicitly.

Exercise 2.3 The aim of this exercise is to make you aware that different notion of convergences may not be necessarily better than others.

Consider the problem,

$$x^* \in \arg \min_{x \in \mathbf{R}^n} f(x),$$

for a μ -strongly convex and L -smooth function.

Consider the gradient method with constant stepsize $\alpha > 0$ to find the unique minimizer of f .

We have seen that the method converges as

$$\|x_k - x^*\| \leq \rho^k \|x_0 - x^*\|,$$

for all $\alpha < 2/L$. The convergence above is called linear convergence. Linear convergence of x_k is always better than $O(1/k)$ convergence of the same quantity. However, that may not be true for other quantities!

- Prove that the gradient method in this case is also converging in an ergodic sense as,

$$\|\bar{x}_k - x^*\| \leq \frac{1}{k+1} \frac{1}{1-\rho} \|x_0 - x^*\|, \quad \bar{x}_k = \frac{1}{k+1} \sum_{i=0}^k x_i$$

for all $\alpha < 2/L$. In particular, show that linear convergence of x_k implies ergodic convergence as $O(1/k)$.

Take-home message: when you have ergodic convergence as $O(1/k)$ it may not be so bad after all. And remember always to spell out convergence of what and with respect to what.

Chapter 3

Third lecture

Constrained optimisation (I): first-order methods

3.1 Setting

We look at generic convex optimisation problems of the form

$$(PC) \quad \underset{\mathbf{x} \in X \subseteq \mathbf{R}^n}{\text{minimise}} \quad f(\mathbf{x}) \quad (3.1)$$

$$\text{subject to} \quad g(\mathbf{x}) \leq 0 \quad (3.2)$$

$$A\mathbf{x} = b, \quad (3.3)$$

for which, we consider a convex cost function $f(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}$, convex inequality constraints $g(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}^l$, affine equality constraints, $h(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}^p \equiv A\mathbf{x} - b$, as well as a convex feasible set $X \subseteq \mathbf{R}^n$. From previous courses, see also the appendix, we know that (PC) is then a convex problem, for which necessary and sufficient conditions exist for finding its optimal points.

In this first class about the subject, we will focus on a specific version on (PC) and we will look at first-order methods (such as the gradient). In particular, we will look at proximal methods, projected gradient, dual ascent, and primal-dual methods.

I will give for granted some notions like KKT conditions and duality, but you can refer to the appendix for some recap.

3.1.1 Setting: simplified

We look at a specific convex optimisation problems of the form

$$(PC) \quad \underset{\mathbf{x} \in X \subseteq \mathbf{R}^n}{\text{minimise}} \quad f(\mathbf{x}), \quad (3.4)$$

with a convex cost $f(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}$ and a convex feasible set $X \subseteq \mathbf{R}^n$. Whereby here X is a “simple” set.

The notion of a “simple” set is rather vague, but we can loosely characterise it as follows: a set is simple if and only if one can easily project over it (e.g., a box constraint, the semidefinite cone). This will be clearer later. For now think of X as a subset of all the possible convex sets.

We talk here of projection, since a typical methods for our simply constrained optimisation problem is the projected gradient method, which performs the iterations,

$$\mathbf{x}_{k+1} = P_X[\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)], \quad P_X(\mathbf{v}) := \arg \min_{\mathbf{x} \in X} \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|^2.$$

Here $P_X(\mathbf{v})$ is the Euclidean projection of the vector \mathbf{v} onto the convex set X . And therefore you can already see why I want the set X to be simple: I want its projection to be carried over with as little complexity as possible, ideally in closed-form.

All these notions will be cleared later, when properly formalised. To do that, we can start with a little more general setting than (PC). We will start from

$$(PS) \quad \underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad f_1(\mathbf{x}) + f_2(\mathbf{x}), \quad (3.5)$$

for convex costs $f_1(\mathbf{x}), f_2(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}$. This seems rather surprising, but this setting will be very fertile in terms of results and reach. Let's see how.

3.2 Splitting methods

Let's go back to our new cost in terms of $f_1(\mathbf{x}) + f_2(\mathbf{x})$. We consider problems for which,

- f_1 is convex and at least $f_1 \in \mathcal{S}_{0,L}^{1,1}$. In this context, f_1 will carry the favourable aspects of the optimisation problem;
- f_2 is the "rest". In general f_2 is non-differentiable and "just" closed convex proper (CCP).

We recall that,

Definition 3.1 (Closed convex proper functions) *The extended real line is $\overline{\mathbf{R}} := \mathbf{R} \cup \{+\infty\}$. An extended real convex function $f : \mathbf{R}^n \rightarrow \overline{\mathbf{R}}$ is proper iff is $\forall \mathbf{x} : f(\mathbf{x}) > -\infty$, and $\exists \mathbf{x}_0 : f(\mathbf{x}_0) < +\infty$. A proper convex function is closed iff it is lower semi-continuous (l.s.c.). A proper, l.s.c., convex function is indicated with $f \in \Gamma(\mathbf{R}^n)$.*

Example 3.1 (Indicator function) *Define the indicator function of a set X as the function $\iota_X : \mathbf{R}^n \rightarrow \overline{\mathbf{R}}$,*

$$\iota_X(\mathbf{x}) : \begin{cases} 0 & \text{if } \mathbf{x} \in X, \\ +\infty & \text{otherwise.} \end{cases}$$

Then, the indicator function $\iota_X(\mathbf{x})$ of a closed and non-empty convex set X is $\Gamma(\mathbf{R}^n)$.

We know show that (PC) is a special case of (PS).

Proposition. (PC) and (PS) are have the same optimisers and optimal value iff f_2 is the indicator function for X .

Proof. *Consider*

$$(PC) \quad \min_{\mathbf{x} \in X} f(\mathbf{x}), \quad (PS') \quad \min_{\mathbf{x} \in \mathbf{R}^n} f(\mathbf{x}) + \iota_X(\mathbf{x})$$

and look at the optimality conditions

$$(PC-o) \quad \nabla f(\mathbf{x}) + N_X(\mathbf{x}) \ni 0,$$

$$(PS'-o) \quad \nabla f(\mathbf{x}) + \partial \iota_X(\mathbf{x}) \ni 0,$$

where $\partial \iota_X(\mathbf{x})$ is the subdifferential of $\iota_X(\mathbf{x})$. One can show (try it) that: $\partial \iota_X(\mathbf{x})$ is nothing else than the normal cone of X , i.e., $\partial \iota_X(\mathbf{x}) = N_X(\mathbf{x})$, for which the claim follows. ♣

Having established that our $f_1 + f_2$ is more general than the constrained setting, let us look at what we can do with it.

3.2.1 Forward-backward splitting

The first question we may ask is how do we find the optimal points of $f_1 + f_2$? If we studied well the second class, we know that if we were to use the subgradient method we would have a $O(1/\sqrt{t})$ convergence rate. This is not very enticing. On the other hand, since f_2 may be non-differentiable, we are not equipped with any better methods so far.

The following theorem may offer us some hints on how to proceed.

Theorem 3.1 *Let $\mathcal{A} = \nabla f_1$ and $\mathcal{B} = \partial f_2$. Then the following conditions are equivalent (for any constant $\alpha > 0$),*

$$(a) \quad (\mathcal{A} + \mathcal{B})(\mathbf{x}) \ni 0 \tag{3.6}$$

$$(b) \quad \mathbf{x} = \mathbf{x} - \alpha(\mathcal{A}(\mathbf{x}) + \mathbf{y}), \mathbf{y} \in \mathcal{B}(\mathbf{x}) \tag{3.7}$$

$$(c) \quad \mathbf{x} + \alpha\mathcal{B}(\mathbf{x}) \ni \mathbf{x} - \alpha\mathcal{A}(\mathbf{x}) \tag{3.8}$$

$$(d) \quad (I + \alpha\mathcal{B})(\mathbf{x}) \ni (I - \alpha\mathcal{A})(\mathbf{x}) . \tag{3.9}$$

Proof. *Immediate.* ♣

Let's look at these conditions very carefully. Condition (a) is just the optimality condition. Condition (b) describes a fixed point iteration: the subgradient method. When the subgradient method is converging, then (b) is verified and therefore we are at optimality. So far nothing new.

Conditions (c-d) describe a new fixed point iteration: the proximal method. The proximal method is an implicit method, in fact we can write:

$$\mathbf{x} + \alpha\mathcal{B}(\mathbf{x}) \ni \mathbf{x} - \alpha\mathcal{A}(\mathbf{x}) \iff \mathbf{x} \in (I + \alpha\mathcal{B})^{-1}(I - \alpha\mathcal{A})\mathbf{x}. \quad (3.10)$$

Let us stop for a moment. We have split the optimality condition into two parts \mathcal{A} and \mathcal{B} , then applied a forward pass $I - \alpha\mathcal{A}$ on \mathbf{x} , and subsequently a backward pass $(I + \alpha\mathcal{B})^{-1}$. This gives the name of what we are doing: a forward-backward splitting. If \mathcal{A} and \mathcal{B} are not just any operators, but (as in our case) gradient and subgradient of some functions, then the splitting is also named the proximal method.

How do we compute $(I + \alpha\mathcal{B})^{-1}$ and what is it? Let us first define the proximal operator.

Definition 3.2 (Proximal operator) *Given a function $\varphi \in \Gamma(\mathbf{R}^n)$ and a positive constant $\alpha > 0$, the proximal operator $\text{prox}_{\alpha\varphi} : \mathbf{R}^n \rightarrow \mathbf{R}^n$ is defined as,*

$$\text{prox}_{\alpha\varphi}(\mathbf{y}) := \arg \min_{\mathbf{x} \in \mathbf{R}^n} \left\{ \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{y}\|^2 + \varphi(\mathbf{x}) \right\}.$$

The proximal operator is the result of an optimisation problem. Since the cost function of the proximal operator is strongly convex for any positive α , then the solution is unique and the operator is well-defined. We will see later its properties, for now let's accept this definition and move on. How do we interpret $(I + \alpha\mathcal{B})^{-1}(I - \alpha\mathcal{A})\mathbf{x}$ in an algorithmic fashion?

Let us consider the step,

$$\mathbf{x} + \alpha\partial f_2(\mathbf{x}) \ni \mathbf{x} - \alpha\nabla f_1(\mathbf{x}), \quad (3.11)$$

which represents $\mathbf{x} + \alpha\mathcal{B}(\mathbf{x}) \ni \mathbf{x} - \alpha\mathcal{A}(\mathbf{x})$ with our choice of \mathcal{A} and \mathcal{B} . Consider now the following line of reasoning: let \mathbf{y} represent the right-hand side,

$$\begin{aligned} \mathbf{x} + \alpha\partial f_2(\mathbf{x}) &\ni \mathbf{x} - \alpha\nabla f_1(\mathbf{x}) := \mathbf{y} \\ \mathbf{x} + \alpha\partial f_2(\mathbf{x}) - \mathbf{y} &\ni 0. \end{aligned}$$

Now, the last inclusion can be interpreted as the optimality condition of an optimisation problem, which one? This one:

$$\begin{aligned} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \alpha f_2(\mathbf{x}) &\rightarrow \min \\ \mathbf{x} = \arg \min_{\mathbf{w} \in \mathbf{R}^n} &\left\{ \frac{1}{2\alpha} \|\mathbf{w} - \mathbf{y}\|^2 + f_2(\mathbf{w}) \right\}, \end{aligned}$$

which reminds us directly the definition of a proximal operator. In fact, the last equation is equivalent to,

$$\begin{aligned} \mathbf{x} &= \text{prox}_{\alpha f_2}(\mathbf{y}) \quad \text{by definition} \\ \mathbf{x} &= \text{prox}_{\alpha f_2}(\mathbf{x} - \alpha\nabla f_1(\mathbf{x})). \end{aligned}$$

So, finally, we have interpreted the forward-backward step as a gradient and proximal step respectively, i.e., $(I + \alpha\mathcal{B})^{-1}(I - \alpha\mathcal{A})\mathbf{x}$ is $\mathbf{x} = \text{prox}_{\alpha f_2}(\mathbf{x} - \alpha\nabla f_1(\mathbf{x}))$.

We are now fully ready for the proximal gradient method.

Remark 1 *Note, many splitting methods exist for different "structures". The proximal gradient method is one of the most used, e.g., in machine learning.*

3.3 The proximal gradient method

The method reads as follows, for a stepsize selection $\alpha_k > 0$

Proximal gradient method

- Start with $\mathbf{x}_0 \in \mathbf{R}^n$
- Iterate $\mathbf{x}_{k+1} = \text{prox}_{\alpha_k f_2}(\mathbf{x}_k - \alpha_k \nabla f_1(\mathbf{x}_k))$, $k = 0, 1, \dots$

The method is an upgrade of the gradient method, where we split the cost into a “nice” cost f_1 and the “rest”, and we take a forward step on the nice part and a backward step on the rest.

When f_2 is the indicator function of a closed convex set ι_X , then the proximal gradient is equivalent to a projected gradient method, in fact,

$$\text{prox}_{\alpha_k \iota_X}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathbf{R}^n} \left\{ \frac{1}{2\alpha} \|\mathbf{x} - \mathbf{y}\|^2 + \iota_X(\mathbf{x}) \right\} = \arg \min_{\mathbf{x} \in X} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2 \right\} = P_X[\mathbf{y}],$$

which is the Euclidean projection on X . So, once again, we are saying that the splitting strategy is a generalisation of the constrained problem (PC).

3.3.1 Prox-friendly functions

As you may already imagine, the proximal method is computationally reasonable when taking the proximal operator is computationally reasonable. To say it in another way, f_2 needs to be prox-friendly, otherwise the computations that we need to solve the optimisation problem in the definition of the proximal operator may exceed the computations to solve the whole problem.

If f_2 is prox-friendly, then the proximal gradient method can be applied. There are a few prox-friendly functions. Below two simple examples.

Example 3.2 (Prox-friendly functions 1) If $f_2 = \iota_X$, when X is a box constraint, then taking the proximal can be done in closed-form. For example, when $X = [0, 1]^n$, then,

$$\mathbf{y}_{k+1} = (\mathbf{x}_k - \alpha_k \nabla f_1(\mathbf{x}_k)), \quad \mathbf{x}_{k+1} = P_{[0,1]^n}(\mathbf{y}_{k+1}),$$

and, each component i of the projection can be computed as,

$$[P_{[0,1]^n}(\mathbf{y}_{k+1})]_i = \max\{0, \min\{1, [\mathbf{y}_{k+1}]_i\}\},$$

where $[\bullet]_i$ is the i -component of the vector.

Proof. The proximal gradient reads,

$$x_{k+1} = P_{[0,1]^n}(x_k - \alpha \nabla f(x_k))$$

Now the projection operator reads,

$$P_{[0,1]^n}(y_{k+1}) = \arg \min_{z \in [0,1]^n} \frac{1}{2} \|z - y_{k+1}\|^2 = \arg \min_{z \in [0,1]^n} \sum_{i=1}^n \frac{1}{2} ([z]_i - [y_{k+1}]_i)^2.$$

The above problem is separable into n optimisation problems, i.e.,

$$\arg \min_{z \in [0,1]^n} \sum_{i=1}^n \frac{1}{2} ([z]_i - [y_{k+1}]_i)^2 = \sum_{i=1}^n \frac{1}{2} \arg \min_{[z]_i \in [0,1]} ([z]_i - [y_{k+1}]_i)^2,$$

whose solution is

$$[x_{k+1}]_i = [z]_i^* = \max\{0, \min\{1, [y_{k+1}]_i\}\}.$$



Example 3.3 (Prox-friendly functions 2) For $f_2(\mathbf{x}) = \|\mathbf{x}\|_1$, then the proximal operator is the soft-thresholding operator:

$$[\text{prox}_{\alpha\|\bullet\|_1}(\mathbf{y})]_i \equiv \text{sign}([\mathbf{y}]_i)(|[\mathbf{y}]_i| - \alpha)_+,$$

where $[\bullet]_i$ is the i -th component of the vector, sign is the sign operator, and $(\bullet)_+$ is $\max\{\bullet, 0\}$.

Proof. The proximal gradient reads,

$$x_{k+1} = \text{prox}_{\alpha\|\bullet\|_1}(x_k - \alpha\nabla f(x_k)) = \arg \min_{z \in \mathbf{R}^n} (\|z\|_1 + \frac{1}{2\alpha}\|z - x_k + \alpha\nabla f(x_k)\|^2)$$

Solving the proximal step calling $y_{k+1} = x_k + \alpha\nabla f(x_k)$,

$$\arg \min_{z \in \mathbf{R}^n} (\|z\|_1 + \frac{1}{2\alpha}\|z - y_{k+1}\|^2) \iff \alpha\partial\|z\|_1 + (z - y_{k+1}) = 0$$

The last set of equations reads, for all i :

$$\alpha\partial[z]_i + [z]_i = [y_{k+1}]_i$$

We have now several possibilities,

$$\begin{aligned} [z]_i > 0 &\implies [z]_i = [y_{k+1}]_i - \alpha \quad (\text{which is valid then for } [y_{k+1}]_i > \alpha) \\ [z]_i < 0 &\implies [z]_i = [y_{k+1}]_i + \alpha \quad (\text{which is valid then for } [y_{k+1}]_i < -\alpha) \\ [z]_i = 0 &\implies \partial[z]_i \in [-1, 1] \text{ and } [y_{k+1}]_i \in [-\alpha, \alpha]. \end{aligned}$$

Reversing,

$$[x_{k+1}]_i = [z]_i^* = \begin{cases} [y_{k+1}]_i - \alpha & \text{for } [y_{k+1}]_i > \alpha \\ [y_{k+1}]_i + \alpha & \text{for } [y_{k+1}]_i < -\alpha \\ 0 & \text{otherwise} \end{cases}$$

and compactifying: $[x_{k+1}]_i = \text{sign}([y_{k+1}]_i)(|[y_{k+1}]_i| - \alpha)_+$. ♣

3.3.2 Proximal properties

The proximal operator has a number of useful properties that we can give and leverage to prove the convergence of the proximal gradient method. We focus on one in this course.

Theorem 3.2 The proximal operator of a function $f_2 \in \Gamma(\mathbf{R}^n)$ is non-expansive, meaning that:

$$\forall \alpha > 0, \forall \mathbf{y}', \mathbf{y} \in \mathbf{R}^n : \|\text{prox}_{\alpha f_2}(\mathbf{y}') - \text{prox}_{\alpha f_2}(\mathbf{y})\| \leq \|\mathbf{y}' - \mathbf{y}\|. \quad (3.12)$$

Proof. It follows from the implicit function theorem, and it is omitted here. ♣

We are now ready for the convergence results.

3.3.3 Convergence of the proximal gradient method

We have the following.

Theorem 3.3 Let $f_1 \in \mathcal{S}_{m,L}^{1,1}$ and $f_2 \in \Gamma(\mathbf{R}^n)$. Choose a constant stepsize $\alpha < 2/L$. Then the proximal gradient method on (PS) has the (global) convergence certificate of,

$$\|\mathbf{x}_t - \mathbf{x}^*\| \leq \rho^t \|\mathbf{x}_0 - \mathbf{x}^*\|, \quad (3.13)$$

for $\rho = \max\{|1 - \alpha m|, |1 - \alpha L|\}$, and t iterations.

Furthermore, let $f_1 \in \mathcal{S}_{0,L}^{1,1}$ and $f_2 \in \Gamma(\mathbf{R}^n)$. Choose a constant stepsize $\alpha < 2/L$. Then the proximal gradient method on (PS) has the (global) convergence certificate of,

$$|f_1(\mathbf{x}_t) + f_2(\mathbf{x}_t) - f_1(\mathbf{x}^*) - f_2(\mathbf{x}^*)| \leq \frac{\omega}{t+1} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad (3.14)$$

for a constant $\omega > 0$, and t iterations.

Consider now the following accelerated variant of the proximal gradient method. Define the supporting sequences, as before, in Equation (2.12), and the **proximal accelerated gradient method** as

Proximal accelerated gradient method

- Start with $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbf{R}^n$ and the sequences (2.12)
- Iterate: $\mathbf{y}_{k+1} = \text{prox}_{\alpha f_2}(\mathbf{x}_k - \alpha \nabla f_1(\mathbf{x}_k))$, $\mathbf{x}_{k+1} = (1 - \gamma_k)\mathbf{y}_{k+1} + \gamma_k \mathbf{y}_k$, $k = 0, 1, \dots$

Then,

Theorem 3.4 Let $f_1 \in \mathcal{S}_{0,L}^{1,1}$ and $f_2 \in \Gamma(\mathbf{R}^n)$. Choose a constant stepsize $\alpha < 1/L$. Then the proximal accelerated gradient method on (PS) has the (global) convergence certificate of,

$$|f_1(\mathbf{x}_t) + f_2(\mathbf{x}_t) - f_1(\mathbf{x}^*) - f_2(\mathbf{x}^*)| \leq \frac{\omega}{t^2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad (3.15)$$

for a constant $\omega > 0$, and t iterations.

Proof. To prove the first statement of Theorem 3.3, we use the non-expansivity property. Start from the algorithm update,

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \|\text{prox}_{\alpha f_2}(\mathbf{x}_k - \alpha \nabla f_1(\mathbf{x}_k)) - \text{prox}_{\alpha f_2}(\mathbf{x}^* - \alpha \nabla f_1(\mathbf{x}^*))\| \leq \\ &\leq \|\mathbf{x}_k - \alpha \nabla f_1(\mathbf{x}_k) - \mathbf{x}^* + \alpha \nabla f_1(\mathbf{x}^*)\|. \end{aligned} \quad (3.16)$$

Then, the rest follows as in the proof of Theorem 2.1.

The remaining results ($f_1 \in \mathcal{S}_{0,L}^{1,1}$ standard and accelerated) are more involved and omitted here. ♣

3.3.4 Interpretation of the results

The results in Theorem 3.3 and 3.4 are quite impressive. If compared to the results we obtained in Theorem 2.1 and 2.3 for the unconstrained case, they say that the proximal gradient method has exactly the same guarantees of an unconstrained algorithm that does not have the complicated f_2 part.

This is brilliant. Whenever f_2 is a prox-friendly function, we can just apply the proximal gradient method and its accelerated variant and forget it exists for the convergence analysis.

In particular for strongly convex and smooth f_1 costs, we obtain a linear convergence rate, while for just smooth f_1 costs, we obtain a $O(1/k)$ and $O(1/k^2)$ rate. The latter in case of a Nesterov's acceleration.

Incidentally, in the case $f_1 \equiv 0$, then we obtain a better method than the subgradient method, for the case of prox-friendly f_2 , called proximal point algorithm.

In this section, we have seen that (PC) can be generalized into (PS) and we have introduced your first splitting method to solve it (the proximal gradient algorithm). Albeit the method works best when f_2 and therefore the set X is prox-friendly, the proximal gradient is one of the workhorses of machine learning and modern optimisation as first-order algorithms are concerned.

In fact, researchers are very good at finding prox-friendly f_2 's.

The proximal gradient splits the objective in a nice part and in a not-so nice part, and then convergence goes as if only the nice part existed!

This brings us to the cases in which we can't avoid looking at non-prox-friendly constraints.

3.3.5 Numerical example

Before continuing, it is good to have a look at a specific numerical example, which you will explore in the second numerical project. We consider an image deblurring task: given an blurred image, our task is to retrieve the original, sharp one.

The task can be formulated in different ways. One possibility is to write an optimisation problem of the form,

$$x^* \in \arg \min_{x \in \mathbf{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \epsilon \|x\|_1, \quad (3.17)$$

for given matrix $A \in \mathbf{R}^{m \times n}$, vector $b \in \mathbf{R}^m$ and regularisation parameter $\epsilon > 0$. The regularisation is added, since typically the matrix A has very few rows $m \ll n$, and the problem can be ill-conditioned.

The cost function is convex, but non-differentiable, nor strongly convex. At face value, we should use the subgradient method. However, considering the splitting, $f_1(x) = \frac{1}{2} \|Ax - b\|_2^2$ and $f_2(x) = \epsilon \|x\|_1$, we can use the proximal gradient. Since f_1 is not just differentiable, but also smooth, we can also use the accelerated proximal gradient. The proximal gradient for this ℓ_1 -regularised problem is traditionally called ISTA (i.e., iterative soft-thresholding algorithm), while the accelerated variant is called FISTA (for fast ISTA), since we are using the soft-thresholding operator to solve the proximal step (see Example 3.3).

In Figure 3.1, you can see the convergence curves for the three methods, and appreciate the added value of the acceleration. Also, in Figure 3.2, you can see the end result. Note that in this example $m = 92416, n = 131072$.

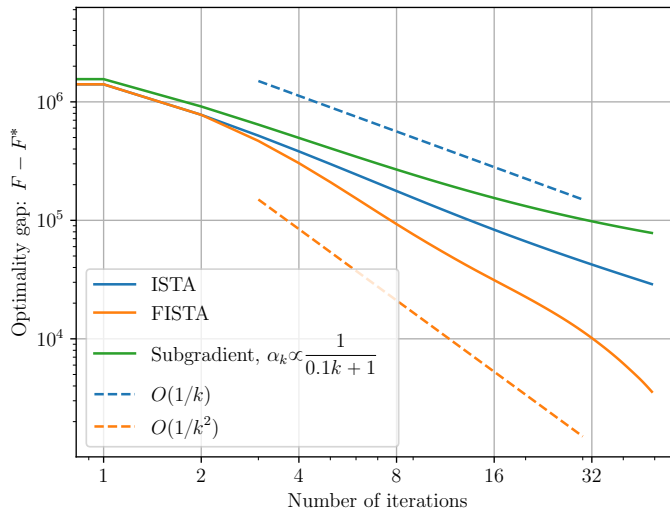


Figure 3.1. Convergence plots for Problem (3.17).

3.4 Duality

3.4.1 Equality constrained problems

We look now at a more explicit setting, for not so simple linear equality constraints,

$$(PE) \quad \underset{x \in \mathbf{R}^n}{\text{minimise}} \quad f(x) \quad (3.18)$$

$$\text{subject to} \quad Ax = b, \quad (3.19)$$

where A is a $\mathbf{R}^{p \times n}$ real-valued matrix, and $b \in \mathbf{R}^p$ with b in the image of A for feasibility. Here $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is a convex function, as usual.



Figure 3.2. Image results for Problem (3.17). Original file by Camille Enlart, reproduced under Creative Commons Attribution-Share Alike 4.0 International license.

We will look at first-order dual methods, so we will need the Lagrangian function for the problem above, which reads,

$$\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \langle \lambda, A\mathbf{x} - b \rangle = f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b).$$

We also need the dual problem,

$$\max_{\lambda \in \mathbf{R}^p} q(\lambda), \quad q(\lambda) := \inf_{\mathbf{x} \in \mathbf{R}^n} \mathcal{L}(\mathbf{x}, \lambda),$$

where $q(\lambda) : \mathbf{R}^p \rightarrow \mathbf{R}$ is the so-called dual function.

Dual problems are useful, since when there is no duality gap then the maximum of the dual problem coincides with the minimum of the primal problem; so if the former is easier to solve, we can just solve it instead of the latter.

Here, you should know that when f convex and constraints are linear (and feasible), constraint qualification holds, and there is no duality gap. If you don't remember the details, please refer to OPT201, or the appendix.

3.4.2 The dual ascent method

If we look at the dual problem, we see that it is an unconstrained problem in λ , so we could optimise it via a gradient method. Indeed, the method (renamed Dual Ascent/ Uzawa's method), look like the following:

- Start with $\lambda_0 \in \mathbf{R}^p$
- Iterate $\lambda_{k+1} = \lambda_k + \alpha_k \mathbf{y}$, $\mathbf{y} \in \partial q(\lambda_k)$ $k = 0, 1, \dots$

We already remark the following important point: in general we don't know if the dual function q is differentiable, so we may need to use the subgradient of the dual function.

You may have already noticed, the properties of convergence and rate of the dual ascent method are linked to the properties of the dual function q . In particular, most of the theorems we have already proved in the second lesson will apply here, once we can determine the properties of q in terms of properties of the original cost f and constraints A .

Remark 2 We give for granted some standard linear algebra results, like the notion of singular values and how to compute them. If you are not familiar with that, please revise Appendix A.5 of [BV].

We have the following key result.

Theorem 3.5 Consider the problem (PE). The dual function $q(\lambda)$ is always concave. Furthermore, if we label the singular values of A as $\sigma_{\min} \leq \dots \leq \sigma_{\max}$, then the following statements hold.

- (a) If $f \in \mathcal{S}_{m,+\infty}^{1,1}$, $\sigma_{\max} > 0$, then $-q \in \mathcal{S}_{0,L'}^{1,1}$ with $L' = \sigma_{\max}^2/m$.
- (b) If $f \in \mathcal{S}_{0,L}^{1,1}$, $\sigma_{\min} > 0$, then $-q \in \mathcal{S}_{m',+\infty}^{1,1}$ with $m' = \sigma_{\min}^2/L$.
- (c) The converse is true if f is convex.

Proof. We prove only the second part of the theorem, since you should know that the dual function is always concave from OPT201.

We start by (a). If $f \in \mathcal{S}_{m,+\infty}^{1,1}$, $\sigma_{\max} > 0$, then the function,

$$\mathbf{x}^*(\lambda) = \arg \min_{\mathbf{x}} f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b)$$

is single-valued, since the solution is unique (due to strong convexity of the cost). Optimality conditions state:

$$\nabla f(\mathbf{x}^*(\lambda)) = -A^\top \lambda.$$

Strong convexity of $f(\mathbf{x})$ also implies,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \|\mathbf{x} - \mathbf{y}\| \geq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq m \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y},$$

hence:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \geq m \|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y},$$

and setting $\mathbf{x} = \mathbf{x}^*(\lambda)$, $\mathbf{y} = \mathbf{x}^*(\lambda')$, we obtain

$$\frac{1}{m} \|A^\top \lambda' - A^\top \lambda\| \geq \|\mathbf{x}^*(\lambda) - \mathbf{x}^*(\lambda')\|. \quad (3.20)$$

Consider now the dual function's gradient along any direction $e \in \mathbf{R}^p$,

$$\begin{aligned} e^\top \nabla q(\lambda) &= \lim_{\epsilon \rightarrow 0} \frac{\min_{\mathbf{x}} (f(\mathbf{x}) + (\lambda^\top + \epsilon e^\top)(A\mathbf{x} - b)) - \min_{\mathbf{x}} (f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b))}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} \frac{(f(\mathbf{x}^*(\lambda + \epsilon e)) + (\lambda^\top + \epsilon e^\top)(A\mathbf{x}^*(\lambda + \epsilon e) - b)) - (f(\mathbf{x}^*(\lambda)) + \lambda^\top (A\mathbf{x}^*(\lambda) - b))}{\epsilon} \end{aligned}$$

Use then Taylor,

$$\begin{aligned} f(\mathbf{x}^*(\lambda + \epsilon e)) &= f(\mathbf{x}^*(\lambda)) + \nabla f(\mathbf{x}^*(\lambda))(\mathbf{x}^*(\lambda + \epsilon e) - \mathbf{x}^*(\lambda)) + O(\|\epsilon\|^2) = \\ &= f(\mathbf{x}^*(\lambda)) - \lambda^\top A(\mathbf{x}^*(\lambda + \epsilon e) - \mathbf{x}^*(\lambda)) + O(\|\epsilon\|^2), \end{aligned}$$

where we have used the Lipschitz property of $\mathbf{x}^*(\lambda)$ for $O(\cdot)$ (Eq. (3.20)) and the optimality condition.

Substituting the latter into the gradient of the dual,

$$\begin{aligned} e^\top \nabla q(\lambda) &= \lim_{\epsilon \rightarrow 0} \frac{-\lambda^\top A(\mathbf{x}^*(\lambda + \epsilon e) - \mathbf{x}^*(\lambda)) + (\lambda^\top + \epsilon e^\top)(A\mathbf{x}^*(\lambda + \epsilon e) - b) - \lambda^\top (A\mathbf{x}^*(\lambda) - b)}{\epsilon} \\ &= e^\top (A\mathbf{x}^*(\lambda) - b) \end{aligned}$$

and therefore $\nabla q(\lambda) = A\mathbf{x}^*(\lambda) - b$.

Lipschitz continuity of the gradient means analyzing:

$$\|\nabla q(\lambda) - \nabla q(\lambda')\| = \|A\mathbf{x}^*(\lambda) - A\mathbf{x}^*(\lambda')\| \leq \frac{\sigma_{\max}^2}{m} \|\lambda - \lambda'\|,$$

where we have used the Lipschitz property of $\mathbf{x}^*(\lambda)$ (Eq. (3.20)), which proves the claim.

We continue by proving (b). If $f \in \mathcal{S}_{0,L}^{1,1}$, $\sigma_{\min} > 0$, we need to prove strong convexity of $-q$. Here we cheat a bit by allowing f to be strictly convex, so q is differentiable, but the proof can be

extended to the general case with a little more work. So, with the cheat, we need a lower bound as,

$$\langle -\nabla q(\lambda) + \nabla q(\lambda'), \lambda - \lambda' \rangle \geq \kappa \|\lambda - \lambda'\|^2,$$

but,

$$\begin{aligned} \langle -\nabla q(\lambda) + \nabla q(\lambda'), \lambda - \lambda' \rangle &= \langle -A\mathbf{x}^*(\lambda) + A\mathbf{x}^*(\lambda'), \lambda - \lambda' \rangle \\ &= \langle -\mathbf{x}^*(\lambda) + \mathbf{x}^*(\lambda'), A^\top \lambda - A^\top \lambda' \rangle \\ &= \langle -\mathbf{x}^*(\lambda) + \mathbf{x}^*(\lambda'), \nabla f(\mathbf{x}^*(\lambda')) - \nabla f(\mathbf{x}^*(\lambda)) \rangle \end{aligned}$$

For L -smoothness of f ,

$$\langle -\mathbf{x}^*(\lambda) + \mathbf{x}^*(\lambda'), \nabla f(\mathbf{x}^*(\lambda')) - \nabla f(\mathbf{x}^*(\lambda)) \rangle \geq \frac{1}{L} \|\nabla f(\mathbf{x}^*(\lambda')) - \nabla f(\mathbf{x}^*(\lambda))\|^2,$$

and using the optimality conditions,

$$\langle -\nabla q(\lambda) + \nabla q(\lambda'), \lambda - \lambda' \rangle \geq \frac{\sigma_{\min}^2}{L} \|\lambda - \lambda'\|^2,$$

as required.

We leave (c) for homework. ♣

This theorem is key: it relates the properties of f to the properties (i.e., the functional class) of q . In particular,

- If f is strongly convex and A is not identically zero, then the negative of the dual function $-q$ is smooth. I can then apply a gradient method, and its accelerated variant with a convergence rate of $O(1/t)$ and $O(1/t^2)$, respectively. I can also compute the Lipschitz constant, and therefore I know the stepsize upper bound.
- If f is smooth and A is full row rank, then the negative of the dual function $-q$ is strongly convex, but not necessarily differentiable. I will have to apply a subgradient method, with a convergence rate of $O(1/\sqrt{t})$.
- If f is both strongly convex and smooth, and A is full row rank, then $-q$ is also strongly convex and smooth, and I can apply a gradient method with linear convergence guarantees.

This is all quite good. However, we need to compute the gradient or the subgradient of $-q$. How do we proceed?

Theorem 3.6 Consider (PE) and $f \in \Gamma(\mathbf{R}^n)$. Indicate with conv the convex hull, then:

$$\partial q(\lambda) = \text{conv}(A\mathbf{x}^*(\lambda) - b),$$

where, $\mathbf{x}^*(\lambda) \in \arg \min_{\mathbf{x} \in \mathbf{R}^n} \mathcal{L}(\mathbf{x}, \lambda)$, which is a set in general.

If f is strictly convex, then $\mathbf{x}^*(\lambda)$ is a singleton for every λ , and $\partial q(\lambda) \equiv \nabla q(\lambda)$, i.e., q is differentiable.

Proof. The proof follows from the proof of Theorem 3.5 and the definition of subdifferential. It is also known as Danskin's theorem, and an alternate proof can be found in Bestekas, *Nonlinear Programming*, Section B.25. ♣

Armed with our gradient definition, we can rewrite the dual ascent method as,

Dual (sub)-gradient ascent method

- Start with $\lambda_0 \in \mathbf{R}^p$
- Iterate:
 - Solve the primal optimisation problem: $\mathbf{x}^*(\lambda_k) \in \arg \min_{\mathbf{x} \in \mathbf{R}^n} \mathcal{L}(\mathbf{x}, \lambda_k)$,
 - Update: $\lambda_{k+1} = \lambda_k + \alpha_k (A\mathbf{x}^*(\lambda_k) - b)$ $k = 0, 1, \dots$

In the algorithm, we have already indicated the gradient of the dual function as $A\mathbf{x}^*(\lambda_k) - b$. When the dual function is not differentiable, the algorithm does not change and we take an

element of the subdifferential of q , i.e., an optimiser $\mathbf{x}^*(\lambda_k)$. We recall that we do not need to take the whole subdifferential set.

If q is smooth, then we can accelerate the gradient via a Nesterov's acceleration in pretty much the same way we discussed in the previous lectures.

The convergence and rates of the (sub)-gradient ascent method and its accelerated variant, in terms of the sequence $\{\lambda_k\}$ as well as the cost $q(\lambda_k)$ follow readily from the discussion we had above, and in particular, let's make it a theorem.

Theorem 3.7 *Consider problem (PE). By using the (sub)-gradient dual ascent method we generate a sequence $\{\lambda_k\}$, which has the following convergence certificates,*

- If $f \in \mathcal{S}_{m,+\infty}^{1,1}$, $\sigma_{\max} > 0$, and the stepsize is chosen constant as $\alpha < 2m/\sigma_{\max}^2$ then,

$$q^* - q(\lambda_k) \leq O(1/t),$$

for a number of iterations t .

- If $f \in \mathcal{S}_{0,L}^{1,1}$, $\sigma_{\min} > 0$, then choosing the stepsize as in Theorem 2.4, at best,

$$\min_{k=0,\dots,t} q^* - q(\lambda_k) \leq O(1/\sqrt{t}),$$

for a number of iterations t .

- If $f \in \mathcal{S}_{m,L}^{1,1}$, $\sigma_{\min} > 0$, then we can choose the stepsize as $\alpha < 2m/\sigma_{\max}^2$ to obtain,

$$\|\lambda_k - \lambda^*\| \leq O(\rho^t),$$

for a number of iterations t and $\rho = \max\{|1 - \alpha\sigma_{\min}^2/L|, |1 - \alpha\sigma_{\max}^2/m|\} < 1$.

By using a Nesterov's accelerated version, then, if $f \in \mathcal{S}_{m,+\infty}^{1,1}$, $\sigma_{\max} > 0$, and the stepsize is chosen constant as $\alpha < m/\sigma_{\max}^2$ then,

$$q^* - q(\lambda_k) \leq O(1/t^2),$$

for a number of iterations t .

Proof. *The proof follows from the discussion above.* ♣

The theorem tells us how to construct dual sequences that convergence in a pertinent sense. We recall that strong duality holds, and as such $f^* = q^*$, so we are converging to the minimum of our problem. However, no mention is given to the primal variable \mathbf{x}^* . That is, given a dual sequence, how can we recover the corresponding primal variable? This is important, right? At the end of the day, you want to solve (PE) and you don't necessarily care about its dual variant.

Theorem 3.8 (Primal recovery) *Consider problem (PE). If $f \in \mathcal{S}_{m,L}^{1,1}$, $L \leq +\infty$, then we can set*

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbf{R}^n} \mathcal{L}(\mathbf{x}, \lambda^*).$$

Otherwise, construct the auxiliary ergodic variable,

$$\bar{\mathbf{x}}_k = \frac{1}{k+1} \sum_{i=0}^k \mathbf{x}_i, \quad \mathbf{x}_i := \mathbf{x}^*(\lambda_i).$$

Then, $\|\bar{\mathbf{x}}_k - \mathbf{x}^*\| \rightarrow 0$.

Proof. *The first part of the proof is obvious, since \mathbf{x}^* is unique. The second part is more involved and left here.* ♣

Summarising: We have transformed (PE) into an unconstrained dual problem, which we can fully characterise and solve as we did in the second class.

Dual ascent is only one of the possible dual algorithms, and not even the best performing. Many algorithms out there, but all of them follow the basic principles that we have looked at here. In particular, you use duality to derive properties of the dual problem; then you use your basic unconstrained theory to derive convergence certificates. Finally, you can say something about primal recovery in some more complex cases.

Homework. Can you think of a dual Newton's method? Can you think of a dual proximal method? When would you apply these methods and with which guarantees?

Before moving on, stop a bit and reflect. In the proximal setup, convergence certificates were very easy to obtain, and the disadvantage was to have prox-friendly costs. In the dual setup, convergence certificates require a bit more work, and now we need that the gradient of the dual function is easy to compute. This latter requirement means that we can compute

$$\mathbf{x}^*(\lambda_k) = \arg \min_{\mathbf{x} \in \mathbf{R}^n} \mathcal{L}(\mathbf{x}, \lambda_k)$$

in an easier way than solving the full, constrained, problem. What if it is not the case? Let's look at primal-dual methods.

3.4.3 Inequality and equality constrained problems

When we play with dual ascent methods, you may be tempted to try out different modifications. As we have just discussed, you may find the optimisation $\arg \min_{\mathbf{x} \in \mathbf{R}^n} \mathcal{L}(\mathbf{x}, \lambda_k)$ awkward, and tempted to substitute it with one or a few steps a gradient method,

$$\mathbf{x}_{\kappa+1} = \mathbf{x}_\kappa - \alpha_\kappa \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_\kappa, \lambda_k), \quad \kappa = 0, 1, \dots, T,$$

where you could start with $\mathbf{x}_\kappa = \mathbf{x}_k$, and let $\mathbf{x}^*(\lambda_k) \approx \mathbf{x}_{T+1}$. This is of course a naive modification that may (or not) work. This type of modifications goes under the name of primal-dual methods, since we both change the primal and the dual at the same time. Let's look at a more general setting and then specify our results for a few algorithms.

We look now at the complete setting

$$(PI) \quad \underset{\mathbf{x} \in X \subseteq \mathbf{R}^n}{\text{minimise}} \quad f(\mathbf{x}) \tag{3.21}$$

$$\text{subject to} \quad A\mathbf{x} = b, g(\mathbf{x}) \leq 0, \tag{3.22}$$

where A is a $\mathbf{R}^{p \times n}$ real-valued matrix, and $b \in \mathbf{R}^p$ with b in the image of A for feasibility. The functions $f : \mathbf{R}^n \rightarrow \mathbf{R}$ and $g : \mathbf{R}^n \rightarrow \mathbf{R}^l$ are convex functions, and X is a convex set. Assume that Slater's constraint qualification holds (i.e. $\exists \bar{\mathbf{x}} : \bar{\mathbf{x}} \in \text{relint}(X), g(\bar{\mathbf{x}}) < 0$). In this context, problem (PI) is convex and there is no duality gap.

Notice that (PE) \subset (PI).

We look here at primal-dual methods, whose convergence is harder and more technical to characterise, and mostly open in many settings. But, we will see what we can say, at least in a simple case.

3.4.4 A primal-dual method

Start by considering f, g to be differentiable. A simple (perhaps naive) primal gradient descent - dual ascent method is the following.

Form the Lagrangian

$$\mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b) + \nu^\top g(\mathbf{x}),$$

then look at the saddle-point problem,

$$\min_{\mathbf{x} \in X} \max_{\lambda \in \mathbf{R}^p, \nu \in \mathbf{R}_+^l} \mathcal{L}(\mathbf{x}, \lambda, \nu). \tag{3.23}$$

We recall here that the dual variable associated to the inequality constraint $g(\mathbf{x}) \leq 0$ has to be constrained in the positive orthant, i.e., $\nu \in \mathbf{R}_+^l$. This comes from the KKT conditions.

Then, one can ask how to solve (3.23), for instance iterating as follows (where we indicate with \mathbf{P} the projection operator).

(Basic) Primal-dual descent-ascent method

- Start with $\mathbf{x}_0, \lambda_0, \nu_0$ and stepsizes $\alpha, \beta, \gamma > 0$,
- Iterate for $k = 0, 1, \dots$:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{P}_X[\mathbf{x}_k - \alpha \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}_k, \lambda_k, \nu_k)] \\ \lambda_{k+1} &= \lambda_k + \beta \nabla_{\lambda} \mathcal{L}(\mathbf{x}_k, \lambda_k, \nu_k) \\ \nu_{k+1} &= \mathbf{P}_{\mathbf{R}_{\geq 0}^l}[\nu_k + \gamma \nabla_{\nu} \mathcal{L}(\mathbf{x}_k, \lambda_k, \nu_k)] \end{aligned}$$

This primal-dual method (a method which alternates between a primal update and a dual update) is rather basic. Its convergence is hard to characterise in general (since \mathcal{L} is not strictly concave in λ, ν). Better and more involved methods exist that tune the stepsize or take more or less steps in the primal and in the dual. Also, you can find methods that add penalisation or regularisation terms to help convergence. We do not explore this here, but it is useful to understand to which we are converging to, in general.

Convergence is to a saddle point, i.e., a point $(\mathbf{x}^*, \lambda^*, \nu^*)$ such that:

$$\forall \mathbf{x} \in X, \lambda, \nu \in \mathbf{R}_{\geq 0}^l \quad \mathcal{L}(\mathbf{x}^*, \lambda, \nu) \leq \mathcal{L}(\mathbf{x}^*, \lambda^*, \nu^*) \leq \mathcal{L}(\mathbf{x}, \lambda^*, \nu^*).$$

This convergence is hard to characterise since you are not descending, not ascending with a certain rate, but you are alternating. So convergence involves determining how far you are from the saddle point, and how much you are violating your constraints.

I stop here, but let's have a look at one interesting example.

3.4.5 Convergence: an example

Let's go back to our (PE) and try to solve it with our basic primal-dual method. In this case $g \equiv 0$, and we don't need ν . We have the following theorem.

Theorem 3.9 *Let $f \in S_{m,L}^{1,1}$ and $g \equiv 0$. Let A be full row rank. Then, the primal-dual algorithm with $\alpha = \beta < 1/L$ delivers a sequence $\{\mathbf{x}_k, \lambda_k\}_{k \in \mathbf{N}}$ such that,*

$$\{\mathbf{x}_k, \lambda_k\}_{k \in \mathbf{N}} \rightarrow (\mathbf{x}^*, \lambda^*)$$

linearly.

Proof. *The proof is interesting since we are going to use a different tool: Energy decrease! By the assumptions on the function and problem, we know that the saddle point $(\mathbf{x}^*, \lambda^*)$ exists and is unique.*

Consider then the variable: $\mathbf{z}_k = \text{col}(\mathbf{x}_k - \mathbf{x}^, \lambda_k - \lambda^*)$, and the energy function $V_k = \frac{1}{2} \|\mathbf{z}_k\|^2$. If we show that the energy decreases with k , then $\|\mathbf{z}_k\| \rightarrow 0$. The energy function V_k is the so-called Lyapunov function, for the students who know what I am talking about.*

First, we can write

$$\mathbf{z}_{k+1} = \mathbf{z}_k + \alpha \begin{bmatrix} -\nabla_{\mathbf{x}} f(\mathbf{x}_k) - A^T \lambda_k + \nabla_{\mathbf{x}} f(\mathbf{x}^*) + A^T \lambda^* \\ A \mathbf{x}_k - b - A \mathbf{x}^* + b \end{bmatrix}.$$

Then, by using Taylor expansion and the mean value theorem,

$$\nabla_{\mathbf{x}} f(\mathbf{x}_k) - \nabla_{\mathbf{x}} f(\mathbf{x}^*) = \nabla_{\mathbf{x}\mathbf{x}} f(\boldsymbol{\xi}_k)(\mathbf{x}_k - \mathbf{x}^*),$$

for a certain $\boldsymbol{\xi}_k$ on the line with end point \mathbf{x}_k and \mathbf{x}^* . In particular, let $\mathbf{Q}_k = \nabla_{\mathbf{x}\mathbf{x}} f(\boldsymbol{\xi}_k)$. Then,

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \alpha \begin{bmatrix} \mathbf{Q}_k & A^T \\ -A & \mathbf{0} \end{bmatrix} \mathbf{z}_k = \begin{bmatrix} \mathbf{I} - \alpha \mathbf{Q}_k & -\alpha A^T \\ \alpha A & \mathbf{I} \end{bmatrix} \mathbf{z}_k.$$

This is a dynamical system. Convergence to zero requires the system to be stable, therefore the system matrix to have eigenvalues in the unit circle. If the matrix $[\mathbf{Q}_k, A^\top; -A, \mathbf{0}]$ has all the eigenvalues in the positive half plane and its maximum eigenvalue is bounded, then we can make the system matrix to have eigenvalues in the unit circle by choosing a small enough α .

In order to check for the eigenvalues of $[\mathbf{Q}_k, A^\top; -A, \mathbf{0}]$ to be positive, we need to verify that

$$\mathbf{v}^\top \begin{bmatrix} \mathbf{Q}_k & A^\top \\ -A & \mathbf{0} \end{bmatrix} \mathbf{v} > 0, \quad \forall \mathbf{v},$$

which is trivially verified. In particular $\sigma_{\max}([\mathbf{Q}_k, A^\top; -A, \mathbf{0}]) = \sigma_{\max}(\mathbf{Q}_k) = \|\mathbf{Q}_k\| \leq L$. Thus, choosing $\alpha < 1/L$ leads stability of the dynamical system and convergence of the primal-dual method. In particular, $\mathbf{z}_k \rightarrow 0$, linearly.

♣

3.4.6 Summarising

We have transformed (PI) into a simply-constrained saddle-point problem, which however it is more complex to characterise and solve than what we have seen for (PE) with purely dual methods. In fact, the research in primal-dual methods is very active and we do not have a complete picture yet.

As we have seen, when moving from simply-constrained problems to general problems, we find more and more general algorithms, which are more and more difficult to handle in theory. While the scope here is just to give you the tools to understand the methods, don't forget that there are many many methods out there, even if the basic ingredients are all the same!

3.5 References

- E. K. Ryu and S. Boyd, *Primer on Monotone Operator Methods*, 2016
- N. Parikh and S. Boyd, *Proximal Algorithms*, 2013
- J. Eckstein, *Splitting methods for monotone operators with applications to parallel optimisation*, 1989

◆◆◆

3.6 Exercises

Exercise 3.1 (Exam 2023) Consider the following quadratic optimisation problem in two scalar variables,

$$\min_{x_1 \in \mathbf{R}, x_2 \in \mathbf{R}} \frac{1}{2}(x_1 - x_2)^2 + 4(x_1 - 1) + a(x_1^2 + 2x_2^2), \quad \text{subject to } (x_1, x_2) \in X,$$

where a is a real scalar, and X a generic set.

1. Prove that it is sufficient that $a \geq 0$ and X convex for the problem above to be convex.
2. Consider now $a = 0$ and $X = \{x_1, x_2 \mid x_1 \in [0, 1], x_2 \in [0, 1]\}$, give an example of a first-order algorithm that achieves the optimal convergence rate for this class of problems. Write the algorithm explicitly (that is, write every steps of the algorithm and detail all the operations involved including the allowed numerical values for the stepsize) and provide its optimal convergence certificate in terms of distance to f^* .
3. Consider now $a = 1$ and $X = \{x_1, x_2 \mid x_1 - x_2 = 1\}$:
 - (a) Derive the dual function. Which functional properties does this dual function have? (Be as precise as you can be: for example if the function is L -smooth, derive the numerical value for the constant L , and so forth).

- (b) Derive the dual problem. Which first-order algorithm would solve the dual problem with the fastest convergence rate? Write the algorithm explicitly (that is, write every steps of the algorithm and detail all the operations involved including the allowed numerical values for the stepsize), and provide a convergence certificate for it.

Exercise 3.2 (Resit 2023) Consider the problem,

$$(P) \quad \min_{x \in \mathbf{R}} f(x) + \varphi(x),$$

where, $f(x) = \frac{1}{2}px^2$, with $p > 0$ and,

$$\varphi(x) = \begin{cases} \frac{1}{2}x^2 & \text{if } x \in [-1, 1] \\ |x| - \frac{1}{2} & \text{otherwise} \end{cases}.$$

1. Is the problem convex? Strongly convex?
2. Prove that

$$\mathbf{prox}_{\alpha\varphi}(x) = \begin{cases} \frac{x}{1+\alpha} & \text{if } |x| \leq 1 + \alpha \\ x - \alpha \mathbf{sign}(x) & \text{otherwise} \end{cases}.$$

3. Write the proximal gradient method and its accelerated variant to solve (P). What are the convergence guarantees that you can give? What are the conditions on α to ensure convergence?
4. Considering the convergence rate alone. Is the accelerated variant better?

Exercise 3.3 (Exam 2023) Consider the set of real symmetric matrices of dimension $n \times n$, and indicate it as $\mathbb{S}(\mathbf{R}^n)$. Consider their eigenvalue decomposition as $X = U\Lambda U^\top$, where $X \in \mathbb{S}(\mathbf{R}^n)$, and Λ is the matrix collecting on the diagonal the eigenvalues of X . Let the eigenvalues be $\lambda_1, \dots, \lambda_n$, and collect them in a vector $\lambda \in \mathbf{R}^n$, such that, $\Lambda = \mathbf{diag}(\lambda)$.

Consider now special convex functions $g : \mathbb{S}(\mathbf{R}^n) \rightarrow \mathbf{R} \cup \{+\infty\}$, the so-called spectral functions, which can be defined by considering eigenvalue functions alone. For example, the trace is a spectral function, since $\mathbf{trace}(X) = \mathbf{1}^\top \lambda$. Other such functions are the indicator function of the sets $\mathcal{X}_1 = \{X \in \mathbb{S} | X \geq 0 \equiv \lambda \geq 0\}$, or $\mathcal{X}_2 = \{X \in \mathbb{S} | \|X\|_F^2 \equiv \mathbf{trace}(X^\top X) \leq 1 \equiv \|\lambda\|_2^2 \leq 1\}$.

For these spectral functions, then the proximal map simplifies as

$$\mathbf{prox}_g(X) = U \mathbf{diag}[\mathbf{prox}_g(\lambda)] U^\top,$$

where with a slight abuse of notation we use the function g both for matrix X and its eigenvalues λ .

Now, consider the problem,

$$(\star) \quad \min_{X \in \mathcal{X} \subseteq \mathbb{S}(\mathbf{R}^n)} f(X) + f_2(X),$$

for a differentiable convex function $f : \mathbb{S}(\mathbf{R}^n) \rightarrow \mathbf{R}$, such that $\nabla f(X) \in \mathbb{S}(\mathbf{R}^n)$, and a generic convex function $f_2 : \mathbb{S}(\mathbf{R}^n) \rightarrow \mathbf{R} \cup \{+\infty\}$.

1. Write the standard proximal gradient method for the problem (\star) for the set $\mathcal{X} = \mathcal{X}_1$ and $f_2 \equiv 0$. Give the proximal operator in closed-form.
2. Write the standard proximal gradient method for the problem (\star) for the set $\mathcal{X} = \mathcal{X}_2$ and $f_2 \equiv 0$. Give the proximal operator in closed-form.
3. Write the standard proximal gradient method for the problem (\star) for the set $\mathcal{X} = \mathbb{S}(\mathbf{R}^n)$ and $f_2 = \mathbf{trace}(X)$. Give the proximal operator in closed-form.
4. Consider problem (\star) for $f(X) = \exp(X)$, $\mathcal{X} = \mathcal{X}_1$, and $f_2 \equiv 0$. With the theory we have studied in class, which convergence guarantees has the standard proximal gradient method applied to (\star) for $f(X) = \exp(X)$? (You could consider $n = 1$ for simplicity, if you need).

Chapter 4

Fourth lecture

Constrained optimisation (II): second-order methods, i.e., interior-point method

4.1 The full picture

We look at the complete explicit convex setting,

$$(PI) \quad \underset{\mathbf{x} \in X \subseteq \mathbf{R}^n}{\text{minimise}} \quad f(\mathbf{x}) \tag{4.1}$$

$$\text{subject to} \quad A\mathbf{x} = \mathbf{b}, g(\mathbf{x}) \leq 0. \tag{4.2}$$

where A is a $\mathbf{R}^{p \times n}$ real-valued matrix, and $\mathbf{b} \in \mathbf{R}^p$ with \mathbf{b} in the image of A for feasibility. Furthermore, both $f(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}$ and $g : \mathbf{R}^n \rightarrow \mathbf{R}^l$ are convex functions, and X is a convex set. Assume that Slater's constraint qualification holds (i.e. $\exists \bar{\mathbf{x}} : \bar{\mathbf{x}} \in \text{relint}(X), g(\bar{\mathbf{x}}) < 0$). Remark that (PE) \subset (PI).

In this lecture, we are interested in second-order algorithms that can deliver very fast convergence, since they will use variants of the Newton's algorithm in some inner loops. We are mainly concerned with interior-point methods, which are the state of the art, if one can afford the computations to run them. For example, you can see them implemented in off-the-shelf solvers like `cvxpy` in python.

To develop the methods, we will assume that both the cost f and the constraint function g are at least double differentiable, so that $f, g \in \mathcal{C}^2$. This is reasonable, since we will need their Hessians for constructing a second-order oracle.

We will then make use of some properties of the Newton's algorithm (namely its affine invariance) to build polynomial-time algorithms to a large class of optimisation problems. As far as computational complexity and time go, polynomial-time algorithms are the "good guys", meaning that they allow to solve optimisation problems in a time which is bounded by a polynomial function of the inputs. This is in contrast with exponential-time algorithms, which are in general the "bad guys". We will discuss this a bit further in the lecture, but keep in mind that the polynomial-time algorithm that we will study made the front page of the New York Times in 1984, when it was discovered! Polynomial-time algorithms are the stars in optimisation.

And, finally, we will use our star method to gauge which convex problems are easy (i.e., can be solved in polynomial time), and which not. Please, remember that convex problems are still in general very hard to solve, even though their optimality conditions can be written down.

4.1.1 Barrier functions and penalised problems

We are now ready to introduce some constructions and notions to help us develop the algorithm. Assume without loss of generality that $X \equiv \mathbf{R}^n$. This is without loss of generality, since you can load g with the constraints represented by X , as well.

Suppose now that we can construct a convex l.s.c. barrier function

$$\Phi_g(\mathbf{x}) : \begin{cases} \approx 0 & \forall \mathbf{x} : g(\mathbf{x}) \leq 0 \\ +\infty & \text{otherwise.} \end{cases} \tag{4.3}$$

As a typical example, consider the constraint $x \leq 0$, and the barrier

$$\Phi_{x \leq 0}(x) = -\log(-x).$$

This logarithmic barrier is $+\infty$ if $x \rightarrow 0$, and stays small for small negative x . As a generalisation, for the convex function $g(\mathbf{x}) : \mathbf{R}^n \rightarrow \mathbf{R}^l$ and the constraint $g(\mathbf{x}) \leq 0$, we can consider each constraint g_i for $i = 1, \dots, l$, and then build a barrier as,

$$\Phi_g(\mathbf{x}) = -\sum_{i=1}^l \log(-g_i(\mathbf{x})), \quad (4.4)$$

both barrier functions extended outside the domain by $+\infty$.

Having now constructed a barrier, instead of the original problem, consider the following penalised version

$$\text{(PIP)} \quad \underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad tf(\mathbf{x}) + \Phi_g(\mathbf{x}) \quad (4.5)$$

$$\text{subject to} \quad A\mathbf{x} = b, \quad (4.6)$$

for $t \geq 0$.

Since the barrier function $\Phi_g(\mathbf{x})$ is convex, and $t \geq 0$, then the penalised problem is still convex. We define as $\mathbf{x}^*(t)$ a solution of this penalised problem and we further call the curve $\{\mathbf{x}^*(t)\}_{t \geq 0}$ the central path. The central path is then the curve generated by a solution of (PIP) by varying t . Moreover, we notice that $\mathbf{x}^*(0)$ is a particular feasible point, called the analytical center of the constraint function $g(\mathbf{x})$.

By simple inspection, we can immediately see that if we let $\lim_{t \rightarrow +\infty} \mathbf{x}^*(t)$ we obtain \mathbf{x}^* , a minimiser of the original problem PI, that is the central path converges to the solution of our original problem.

The idea is to move along the central path, using some basic algorithm to do small steps along the central path. The “basic algorithm” that is used in interior-point methods is the Newton’s method, because of its very fast local convergence.

The technical details are for which classes of problems the (PIP) problems do not become harder and impossibly hard as $t \rightarrow \infty$, and how to select the increment of t such that we stay in the local convergence zone.

Since gradient methods depend on the condition number and they would become very slow for large values of t , we can safely say that they are useless in this lecture.

4.1.2 Interpretation of the penalised problem

To understand a bit better the idea behind the penalised problem, let us have a look at the KKT conditions for it. In particular, consider a logarithmic barrier function, then the KKT conditions for (PIP) are

$$t\nabla f(\mathbf{x}) + A^T\lambda + \sum_{i=1}^l \frac{1}{-g_i(\mathbf{x})} \nabla g_i(\mathbf{x}) = 0 \quad (4.7)$$

$$A\mathbf{x} = b, \quad g_i(\mathbf{x}) < 0 \quad i = 1, \dots, l. \quad (4.8)$$

Or equivalently, introducing the new variables $\nu_i = -1/(tg_i(\mathbf{x}))$,

$$\nabla f(\mathbf{x}) + A^T\lambda + \sum_{i=1}^l \nu_i \nabla g_i(\mathbf{x}) = 0 \quad (4.9)$$

$$A\mathbf{x} = b, \quad g_i(\mathbf{x}) < 0 \quad \nu_i > 0 \quad i = 1, \dots, l \quad (4.10)$$

$$-\nu_i g_i(\mathbf{x}) = 1/t \quad i = 1, \dots, l. \quad (4.11)$$

By the second set of equation, we can see that the (PIP) approximates the complementary slackness condition, and as $t \rightarrow \infty$, we regain the original problem. Excellent.

4.2 A Newton's step for the penalised problem

We are now ready to develop a Newton's step/method for the penalised problem, which is an equality-constraint convex optimisation problem. The developments here are rather general and can be applied at any equality-constraint convex optimisation problems.

Let's start by constructing the Lagrangian function for (PIP) as,

$$F_t(\mathbf{x}, \lambda) = tf(\mathbf{x}) + \lambda^T(A\mathbf{x} - b) - \sum_{i=1}^l \log(-g_i(\mathbf{x})), \quad (4.12)$$

where we have indicated with λ the dual variable associated to the constraint $A\mathbf{x} = b$.

We can further write the optimality conditions for (PIP) as,

$$\nabla_{\mathbf{x}} F_t(\mathbf{x}, \lambda) = 0 \quad A\mathbf{x} = b. \quad (4.13)$$

To build a Newton's method, we need to expand the optimality conditions by using a Taylor approximation around a point $(\mathbf{x}_k, \lambda_k)$, i.e.,

$$\nabla_{\mathbf{x}} F_t(\mathbf{x}_k, \lambda_k) + \nabla_{\mathbf{x}\mathbf{x}}^2 F_t(\mathbf{x}_k, \lambda_k) \Delta \mathbf{x} + \nabla_{\lambda \mathbf{x}}^2 F_t(\mathbf{x}_k, \lambda_k) \Delta \lambda = 0 \quad (4.14)$$

$$A\mathbf{x}_k - b + A\Delta \mathbf{x} = 0. \quad (4.15)$$

We can now build the linear system for the Newton's increment $(\Delta \mathbf{x}, \Delta \lambda)$ as,

$$\underbrace{\begin{bmatrix} \nabla_{\mathbf{x}\mathbf{x}}^2 F_t(\mathbf{x}_k, \lambda_k) & A^T \\ A & 0 \end{bmatrix}}_{=: H_k} \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \lambda \end{bmatrix} = \underbrace{\begin{bmatrix} -\nabla_{\mathbf{x}} F_t(\mathbf{x}_k, \lambda_k) \\ -A\mathbf{x}_k + b \end{bmatrix}}_{=: r_k}, \quad (4.16)$$

and solve for $(\Delta \mathbf{x}, \Delta \lambda)$.

Finally, we can use a damped strategy for the update, $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \Delta \mathbf{x}$, $\lambda_{k+1} = \lambda_k + \alpha_k \Delta \lambda$, where $\alpha_k > 0$ is the stepsize and it is based on the chosen damped strategy.

4.2.1 Damping and the full method

We focus here on a special choice of damping to fix the ideas, namely the backtracking strategy. Since now the problem is constrained, we cannot use a strategy based on the cost, as done in previous lectures. It is much better to use instead a strategy based on the residual $\|r_k\|$, meaning that you would like to reduce it at least by a given amount per iteration.

The full method then looks like the following.

Damped Newton's method for equality-constrained problems

- Start with $\mathbf{x}_0, \lambda_0 \in \mathbf{R}^n$ (Need: $\mathbf{x}_0 : g_i(\mathbf{x}_0) < 0$ if a penalised problem)
- Iterate:
 - Solve the linear system: $H_k \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \lambda \end{bmatrix} = r_k$,
 - Backtracking: set $\alpha_k = 1$, $\tau \in (0, 1)$, $\beta \in (0, 1/2)$:
 - While $\|r(\mathbf{x}_k + \alpha_k \Delta \mathbf{x}, \lambda_k + \alpha_k \Delta \lambda)\| > (1 - \beta \alpha_k) \|r_k\|$, set: $\alpha_k \leftarrow \tau \alpha_k$.
 - $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \Delta \mathbf{x}$, $\lambda_{k+1} = \lambda_k + \alpha_k \Delta \lambda$.

As long as the generalised function F_t is strongly convex and had a Lipschitz Hessian in \mathbf{x} , i.e., it is $\in \mathcal{S}_{m,M}^{2,2}$, then convergence goes as in Theorem 2.7, and in particular it is locally quadratic. The assumptions on F_t may be reasonable in some situations, but in our case, they may be awkward. In fact the logarithmic barrier is not strongly convex, and it may even have problems to have a Lipschitz Hessian.

In the next section, we will see how to modify the analysis and to change the assumption for a more general one, which will help us with the convergence proof of the interior-point method, independently of the conditioning parameter t .

4.3 Self concordance

We start with defining a new function class: the self-concordance class as follows.

Definition 4.1 (Self-concordance on \mathbf{R}) A convex function $f : \mathbf{R} \rightarrow \mathbf{R}$ is self-concordant iff

$$|f'''(x)| \leq 2f''(x)^{3/2}$$

for all $x \in \text{dom}(f)$. Where $'''$ indicated the third order derivative and $''$ the second order one.

Definition 4.2 (Self-concordance on \mathbf{R}^n) A convex function $f : \mathbf{R}^n \rightarrow \mathbf{R}$ is self-concordant iff it is self-concordant along every line in its domain, i.e., if the function $\tilde{f}(\mu) = f(\mathbf{x} + \mu v)$ is a self-concordant function of μ for all $\mathbf{x} \in \text{dom}(f)$ and for all v .

Self concordance will be indicated as SC . At first, this definition may sound a mathematical weirdness, but we will see how it can help us a great deal, and it encompass a good deal of functions!

Let's look at some examples and properties.

4.3.1 Examples and properties

Let us look at $|f'''(x)| \leq 2f''(x)^{3/2}$ in the scalar case. We notice immediately that, linear and positive quadratic functions are self-concordant since convex and $f'''(x) = 0$ and $f''(x) \geq 0$.

We also notice that, $-\log(x)$ is self-concordant on its domain: this key for us, since it will be the basis for our barrier functions, which are not strongly convex. Since this is the key step, let's make it a theorem.

Theorem 4.1 (Self concordance of the logarithmic function) The function $-\log(x)$ is self-concordant on its domain.

Proof. *Homework.* ♣

Some useful properties of the SC class are the following ones, whose proofs are immediate.

- If $f_1 \in SC, f_2 \in SC$ then, $f_1 + f_2 \in SC$
- If $f \in SC$ and $a \geq 1$, then, $af \in SC$
- If $f \in SC$ then $f(Ax + b) \in SC$

4.3.2 Matrix extensions

In some cases, we also look at problems defined over the semi-definite cone. In particular, we look at problems defined over symmetric real-valued matrix variables $X \in \mathbb{S}(\mathbf{R}^n)$, such as,

$$\underset{X \in \mathbb{S}(\mathbf{R}^n)}{\text{minimise}} \quad F(X) \tag{4.17}$$

$$\text{subject to} \quad X \geq 0, \tag{4.18}$$

where $F : \mathbb{S}(\mathbf{R}^n) \rightarrow \mathbf{R}$ is a convex function. These types of problems are often referred to as semi-definite programs, and they are useful in many domains.

First of all, it should be immediate to prove that the above is a convex problem, since $X \geq 0$ is a convex set (try it out).

Then, we need two things. First an inner product, which is the trace. Second, a self-concordant barrier.

Theorem 4.2 An allowed barrier function for the constraint $X \geq 0$ is $\phi(X) = -\log \det(X)$.

Proof. We start by proving that $\phi(X)$ is convex.

For the function $\phi(X)$, $X > 0$, we can verify convexity by considering an arbitrary line, given by $X = Z + tV$, where Z, V are symmetric matrices of dimension n . We define $g(t) = \phi(Z + tV)$, and restrict g to the interval of values of t for which $Z + tV > 0$. Without loss of generality, we can assume that $t = 0$ is inside this interval, i.e., $Z > 0$.

We now show that $g''(t) \geq 0$ (g'' is the Hessian), and conclude that $\phi(X)$ is a convex function.

For the defined $g(t)$,

$$\begin{aligned} g(t) &= -\log \det(Z + tV) \\ &= -\log \det(Z^{1/2}(I + tZ^{-1/2}VZ^{-1/2})Z^{1/2}) \\ &= -\sum_{i=1}^n \log(1 + t\lambda_i) - \log \det Z, \end{aligned}$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $Z^{-1/2}VZ^{-1/2}$, therefore for the scalar function $g(t)$

$$g'(t) = \frac{dg}{dt} = -\sum_{i=1}^n \frac{\lambda_i}{1 + t\lambda_i}, \quad g''(t) = \frac{d^2g}{dt^2} = \sum_{i=1}^n \frac{\lambda_i^2}{(1 + t\lambda_i)^2} \geq 0.$$

Note, we used the fact that $\det(A) = \prod_{i=1}^n \lambda_i(A)$, that $\det(AB) = \det(A)\det(B)$, and the usual $\log(ab) = \log(a) + \log(b)$.

We continue by proving the self-concordance of $\phi(X)$. Let us compute the derivatives along a line $X + tV$, for $X > 0$ and V symmetric.

$$\begin{aligned} \left. \frac{d\phi(X + tV)}{dt} \right|_{t=0} &= -\left. \frac{d}{dt} \right|_{t=0} \{\log \det(X) + \log \det(I + tX^{-1}V)\} = -\mathbf{trace}(X^{-1}V), \\ \left. \frac{d^2\phi(X + tV)}{dt^2} \right|_{t=0} &= -\left. \frac{d}{dt} \right|_{t=0} \mathbf{trace}((X + tV)^{-1}V) = \mathbf{trace}(X^{-1}VX^{-1}V), \\ \left. \frac{d^3\phi(X + tV)}{dt^3} \right|_{t=0} &= \left. \frac{d}{dt} \right|_{t=0} \mathbf{trace}((X + tV)^{-1}V(X + tV)^{-1}V) = -2\mathbf{trace}(X^{-1}VX^{-1}VX^{-1}V). \end{aligned}$$

Let $\tilde{V} = X^{-1/2}VX^{-1/2}$, then,

$$\left. \frac{d^2\phi(X + tV)}{dt^2} \right|_{t=0} = \mathbf{trace}(\tilde{V}^2) > 0, \quad \left. \frac{d^3\phi(X + tV)}{dt^3} \right|_{t=0} = -2\mathbf{trace}(\tilde{V}^3).$$

Hence,

$$\left| \left. \frac{d^3\phi(X + tV)}{dt^3} \right|_{t=0} \right| = \left| -2\mathbf{trace}(\tilde{V}^3) \right| \leq 2\mathbf{trace}(\tilde{V}^2)^{3/2} = 2 \left(\left. \frac{d^2\phi(X + tV)}{dt^2} \right|_{t=0} \right)^{3/2},$$

from which the thesis. ♣

4.3.3 Newton's method for self-concordant functions

We are now ready for the analysis of the damped Newton's method for self-concordant functions.

Theorem 4.3 *Damped Newton's method convergence globally for strictly convex and self-concordant functions on the problem $\min_{\mathbf{x} \in \mathbf{R}^n} f(\mathbf{x})$. In particular, there exists a $\gamma > 0$ depending on the backtracking method, for which it obtains a solution \mathbf{x}_k for which $f(\mathbf{x}_k) - f(\mathbf{x}^*) < \epsilon$ in*

$$\frac{f(\mathbf{x}_0) - f(\mathbf{x}^*)}{\gamma} + \log \log(1/\epsilon)$$

iterations.

Typically $1/\gamma \sim 100$, while $\log \log(1/\epsilon) \approx 6$. Note that the accuracy does not depend on the conditioning of the function or on its strong convexity.

Proof. *The proof is rather involved and left to the reader to check out in [BV] book.* ♣

The theorem can be extended to equality-constrained Newton's methods and the result stay the same. With this result in place, let's look at interior-point methods.

4.4 The interior-point method

Recall the damped Newton's method for equality-constrained problems and consider the following method to solve (PI) with a sequence of (PIP) problems,

Interior-point method

- Start with \mathbf{x}_0 such that $g(\mathbf{x}_0) < 0$, $t > 0$, $\mu > 1$ and required accuracy ϵ . Rem: $g : \mathbf{R}^n \rightarrow \mathbf{R}^l$
- Iterate $\kappa = 0, 1 \dots$
 - Solve (PIP) starting from \mathbf{x}_κ with a damped Newton's method. Let the solution be $\mathbf{x}_{\kappa+1}$
 - Quit if $l/t < \epsilon$ (remember complementary slackness: $-\nu_i g_i(\mathbf{x}) = 1/t$)
 - Update: $t \leftarrow \mu t$

Remark 3 *You can't start with t satisfying $l/t < \epsilon$, because t would be too big and we may be far from the local convergence zone.*

Our aim is now to show that if we use damped Newton's method for each (PIP) the number of required iterations is bounded by a polynomial function of the inputs (e.g., number of constraints), hence, the interior-point method for self-concordant functions and barriers has a polynomial complexity and can solve (some) convex problems efficiently.

4.4.1 Computational complexity analysis: outer

First, we look at the outer steps (updates on t). Starting for a certain t , say t_0 for short, each steps we obtain an accuracy of

$$\frac{l}{\mu^\kappa t_0}.$$

Therefore for a desired accuracy ϵ , we need

$$\left\lceil \frac{\log(l/\epsilon t_0)}{\log \mu} \right\rceil$$

steps.

4.4.2 Computational complexity analysis: inner

Theorem 4.4 *Consider the (PIP)s generated by an interior-point method. Assume the technical requirements that,*

- *The function $tf(\mathbf{x}) - \sum_{i=1}^l \log(-g_i(\mathbf{x}))$ is self-concordant for all $t > 0$;*
- *The (PIP)s have bounded sub-level sets (i.e., they are solvable).*

Then, the number of required inner iterations of the damped Newton's method to achieve an accuracy ϵ_{nt} can be upper bounded as,

$$\frac{l(\mu - 1 - \log(\mu))}{\gamma} + \log \log(1/\epsilon_{\text{nt}}) \approx \frac{l(\mu - 1 - \log(\mu))}{\gamma} + 6.$$

with γ depending only on the backtracking parameters.

The proof is based on Theorem 4.3 on the convergence of the damped Newton’s method with self-concordance assumption.

We see immediately that the number of iterations are linear in l . Each iteration requires a Newton’s step computation, which is in the worst case a matrix inversion, which is of the order of $O((n+m)^3)$, so also polynomial. We also see that the bound does not depend on the penalisation term t .

4.4.3 Computational complexity analysis: complete

Now we can put things together

Theorem 4.5 *Assume (PI) is convex and enjoys a self-concordant barrier function and each of the (PIP) is solvable. Then employing an interior-point method with parameters (t_0, μ, γ) delivers a solution with accuracy ϵ in*

$$\left\lceil \frac{\log(l/\epsilon t_0)}{\log \mu} \right\rceil \times \left\lceil \frac{l(\mu - 1 - \log(\mu))}{\gamma} + 6 \right\rceil = o(l^2)$$

iterations. Each of them requiring the solution of a matrix inversion with complexity $< O((n+m)^3)$.

As such the interior-point method with self-concordant barrier is a polynomial-time method, with overall complexity of $< O(l^2(n+m)^3)$.

Theorem 4.5 delivers the convergence certificate guarantees we need: it tells us how many iterations we need to reach a certain accuracy, in terms of complementary slackness. It also tells us that each iteration can be carried out in polynomial time, making our life easy.

We have therefore developed a polynomial-time method for generic convex optimisation problems. This is a huge feat: we can actually solve some of the problems (using a second-order oracle) with very high accuracy!

I want to recall that first-order methods can offer better computational complexity *per iteration*, but they may take many many iteration (not polynomially bounded) to converge. Sometimes however, that’s the only thing you have.

The next natural question is: which optimisation problems can be solve then in polynomial time? All convex ones?

4.5 Classes of easy problems

We now look at the classes of easy problems, problems that can be solved in polynomial time by using the interior-point method that we have developed. There are four of such classes, which can model a widespread array of applications. This is what makes (a sub-set of convex) optimisation a useful technology. In particular, this sub-set of convex optimisation is a “complete” domain:

- We have an amazing theory and we understand it inside out;
- We have efficient algorithms and we can get their complexity;
- The problems we can solve efficiently *matter* for society!

Frankly, you will not see much of these domains in your professional life. For example, machine learning matter for society but the theory behind is still not well understood. Binary optimisation matter for society, but we do not have efficient algorithms.

Our four musketeers are:

- Linear programs (LPs);
- Quadratic programs (QPs);
- Second-order conic programs (SOCPs);

- Semidefinite programs (SDPs).

And for all the above you have a self-concordant barrier, so that you can apply Theorem 4.5.

4.5.1 Linear programs

Linear programs (or LPs) are convex optimisation problems with linear cost and linear constraints. Linear programs in continuous optimisation do not have to be mixed with linear programs in integer programming, where the decision variables are integer. The latter are not linear at all!

A typical linear program is the following,

$$\underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad c^\top \mathbf{x}, \quad (4.19)$$

$$\text{subject to} \quad A\mathbf{x} = b, D\mathbf{x} \leq e. \quad (4.20)$$

For example:

$$\underset{x_1 \in \mathbf{R}, x_2 \in \mathbf{R}}{\text{minimise}} \quad 2x_1 + x_2$$

$$\text{subject to} \quad x_1 \geq 0, x_2 \geq 0$$

$$x_1 + x_2 \leq 0.$$

The solution is on the boundary (vertex), from which the simplex method can be developed (which is a combinatorial method with worst-case exponential-time complexity). The interior-point method was applied to LPs successfully in 1984 with a complexity of $O(n^2m)$, $m \geq n$ for n variables and m constraints. Typically one can solve instances with 10k variables and constraints. Furthermore, LPs can be used to approximate more complex optimisation problems by linearisation.

4.5.2 Quadratic Programs: QPs

Quadratic Programs or QPs are convex optimisation problems with convex quadratic cost and linear constraints,

$$\underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + c^\top \mathbf{x}, \quad (4.21)$$

$$\text{subject to} \quad A\mathbf{x} = b, D\mathbf{x} \leq e, \quad (4.22)$$

with $Q \geq 0$.

QPs have better properties than LP and their solution is not on the boundary. The level sets are ellipsoids, and in general, they are not more complicated to solve than LPs. In fact, writing the Lagrangian and the KKT conditions (here the constraint qualifications are satisfied)

$$S : \begin{cases} Q\mathbf{x} + c + A^\top \lambda + D^\top \nu = 0 \\ A\mathbf{x} - b = 0 \\ D\mathbf{x} - e \leq 0 \\ \nu \geq 0 \\ (D\mathbf{x} - e)_i \nu_i = 0 \end{cases} \rightarrow \begin{cases} \text{find} & \mathbf{x}, \lambda, \nu \\ \text{s.t.} & \mathbf{x}, \lambda, \nu \in S \end{cases} = \text{LP} + \text{sth. else}$$

we see that QPs are similar to LPs.

Here the computational complexity with interior-point is limited by matrix inversions. We can cite a few notable examples of QPs, such as least-squares problems, regression problems, optimal control, and so forth. QPs have been, and still are, the cornerstone of nonconvex optimisation with sequential quadratic programming (SQP).

4.5.3 Second-Order Conic Programs: SOCPs

Second-Order Conic Programs or SOCPs are convex optimisation problems with quadratic cost and quadratic constraints and a bit more. Compactly, they can be written as,

$$\underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad c^\top \mathbf{x}$$

$$\text{subject to} \quad \|D_i x + e_i\|_2 \leq h_i^\top x + d_i, \quad i = 1, \dots, m$$

$$A\mathbf{x} = b \quad (4.23)$$

4.5.4 Semi-Definite Programs: SDPs

Semi-Definite Programs or SDPs are convex optimisation problems with linear cost and with semidefinite constraints.

Let's us define a Linear matrix $F : \mathbf{R}^n \rightarrow \mathbf{S}^n$:

$$F(x) = x_1 F_1 + x_2 F_2 + \dots + x_n F_n + G \quad \mathbf{x} \in \mathbf{R}^n, F_1, \dots, F_n, G \in \mathbf{S}^n.$$

and, the SDP:

$$\begin{array}{ll} \underset{\mathbf{x}}{\text{minimise}} & c^\top \mathbf{x} \\ \text{subject to} & F(\mathbf{x}) \geq 0 \quad (\text{SDP constraint, LMI..}) \\ & A\mathbf{x} = b. \end{array}$$

4.6 Summary

In these four classes we have developed a theory of algorithms for optimisation. Without it, we wouldn't be able to solve anything. Without it, we wouldn't be able to know how many iterations we need to run to obtain a result with a certain accuracy. Without it, we could wait forever to get nothing in return.

When talking about algorithms, we care about complexity, and which type of oracles we can afford to use: zero-order, first-order, second-order, and so forth.

And, we also care about the properties (convexity, smoothness) of the problems, and we use them to decide which algorithm to employ.

And, finally, by using the algorithmic lens, we can decide how to model our problem, so that we are able to solve it (we are guided in terms of assumptions). To further hammer this last point, consider the following last example.

Example 4.1 Consider the problem,

$$(Pex) \quad \underset{\mathbf{x} \in \mathbf{R}^n}{\text{minimise}} \quad \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + c^\top \mathbf{x} + \lambda \|\mathbf{x}\|_\infty,$$

$\mathbf{Q} \geq 0$. How would you solve the problem?

At face value: since the cost is not differentiable, we could employ the sub-gradient method. This would converge very slowly.

However, isn't the term $\|\cdot\|_\infty$ prox-friendly? Then the proximal method can be employed instead. Here we could even use its accelerated variant.

Is the dimension not-too big? Say $n \sim 1000$? Then, transform the problem into a QP and solve it via interior-point methods (which are much faster), how? By using the epigraph form:

$$(Pex) \iff \begin{array}{ll} \underset{\mathbf{x} \in \mathbf{R}^n, t > 0}{\text{minimise}} & \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + c^\top \mathbf{x} + \lambda t, \\ \text{subject to} & \|\mathbf{x}\|_\infty \leq t \iff |[\mathbf{x}]_i| \leq t \forall i \\ & \iff [\mathbf{x}]_i \leq t, [\mathbf{x}]_i \geq -t, \forall i \end{array}$$

Optimisation problems can be transformed into other ones.

4.7 References

- BV: Chapters 10, 11.

◆◆◆

4.8 Exercises

Exercise 4.1 Consider the problem,

$$\min_{x_1 \in \mathbf{R}, x_2 \in \mathbf{R}} \frac{1}{2}(x_1^2 + x_2^2), \quad \text{subject to } x_2 \geq 1 - x_1.$$

1. Introduce a logarithmic barrier for the constraint, write the penalised problem, and find the analytical centre.
2. Plot in a 2D sketch the qualitative level curves of the penalised problem for $t = 1$ and $t = 10$.
3. Argue that the problem is equivalent to the one dimensional problem,

$$\min_{y \in \mathbf{R}} \frac{1}{2}y^2, \quad \text{subject to } y \geq 1/2,$$

in the sense that $x_1^* = x_2^* = y^*$.

4. Write the penalised problem for y and solve it analytically for $t = 1, 4, 10$.
5. What is the value of y as $t \rightarrow \infty$? Could you have guessed that?

Exercise 4.2 (Exam 2023) Consider the problem,

$$\min_{x \in \mathbf{R}^n} \frac{1}{2}x^\top Qx + c^\top x + \|x\|_1,$$

with $Q \geq 0$.

1. Prove that the problem above can be formulated as a quadratic program as,

$$\begin{aligned} \min_{x \in \mathbf{R}^n, y \in \mathbf{R}^n} \quad & \frac{1}{2}x^\top Qx + c^\top x + \mathbf{1}^\top y, \\ \text{subject to} \quad & x_i - y_i \leq 0, \quad x_i + y_i \geq 0, \quad y_i \geq 0, \quad \forall i, \end{aligned}$$

where \bullet_i represents the i -th component of a given vector \bullet .

2. Apply the interior-point method to the quadratic program above using a logarithmic barrier, in particular write:
 - (a) The modified cost function with the logarithmic barrier;
 - (b) The Newton's step for the inner problems, that is: the linear system;
 - (c) A damped strategy to solve each inner problems;
 - (d) How many Newton iterations do you need overall? What is the total computational complexity (i.e., number of iterations times the complexity per iterations)?
 - (e) Finally, focus on the constraint $y_i \geq 0$. Could one use $\log^2(y_i)$ as a barrier?

Exercise 4.3 (Resit 2023) Consider the problem

$$\min_{x \in \mathbf{R}^n} x^\top Wx, \quad \text{subject to } x_i^2 = 1, \quad i = 1, \dots, n.$$

Assume that W is symmetric and its diagonal is 0.

1. Is this problem convex?
2. Prove that the dual of the above problem is,

$$\min_{\nu \in \mathbf{R}^n} \mathbf{1}^\top \nu, \quad \text{subject to } W + \mathbf{diag}(\nu) \geq 0,$$

where $\mathbf{diag}(\nu)$ is a $n \times n$ matrix, with diagonal equal to ν . Prove that the dual problem is convex and it is in fact an SDP.

3. We are going to solve this SDP via an interior-point method.

- (a) Write the modified cost function with the barrier function $\phi(X) : \mathbb{S}_+^n \rightarrow \mathbf{R}$, defined as $\phi(X) = \log \det(X^{-1}) = -\log \det(X)$ for the set $X \geq 0$. Recall that the space \mathbb{S}_+^n is the space of symmetric matrices in n dimension.
- (b) Write the Newton's step, recalling that

$$\nabla_{\nu} \log \det(W + \mathbf{diag}(\nu)) = \mathbf{Diag}[(W + \mathbf{diag}(\nu))^{-1}],$$

where $\mathbf{Diag}[\cdot]$ is a vector, whose entries are the diagonal elements of (\cdot) . And,

$$\nabla_{\nu\nu} \log \det(W + \mathbf{diag}(\nu)) = -(W + \mathbf{diag}(\nu))^{-2}.$$

Appendix A

Recap from OPT201

A.1 Optimality conditions: sets

Consider the convex problem

$$\underset{\mathbf{x} \in X \subseteq \mathbf{R}^n}{\text{minimise}} f(\mathbf{x}),$$

with $f : \mathbf{R}^n \rightarrow \mathbf{R}$ convex, and X convex. Then, necessary and sufficient conditions for optimality are,

$$\partial f(\mathbf{x}^*) + N_X(\mathbf{x}^*) \ni \mathbf{0},$$

where ∂f is the sub-differential of f and $N_X(\mathbf{x})$ is the normal cone operator of X at \mathbf{x} .

If f is differentiable then $\partial f \equiv \nabla f$ and we can write,

$$\nabla f(\mathbf{x}^*) + N_X(\mathbf{x}^*) \ni \mathbf{0}.$$

The above is equivalent (in the convex case, in which we are) to

$$\nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \forall \mathbf{x} \in X$$

All of these are set conditions, and it is hard to compute \mathbf{x}^* from them.

A.2 Optimality conditions: equality constraints

Consider now the case $X = \{\mathbf{x} | A\mathbf{x} = b\}$, $A \in \mathbf{R}^{q \times n}$, and for simplicity the case in which $f \in \mathcal{C}(\mathbf{R}^n)$. The problem is still convex. We have the following implication.

Theorem A.1 *Consider the problem, $\min_{\mathbf{x} \in \mathbf{R}^n} f(\mathbf{x})$, subject to $A\mathbf{x} = b$, if the constraints are qualified, then the optimality conditions $\nabla f(\mathbf{x}^*) + N_X(\mathbf{x}^*) \ni \mathbf{0}$ can be equivalently written as*

$$\begin{cases} \nabla f(\mathbf{x}^*) + A^\top \lambda^* = 0 \\ A\mathbf{x}^* = b \end{cases}$$

- The variables $\lambda \in \mathbf{R}^q$ are the Lagrangian multipliers;
- We can write the optimality conditions as the saddle-points of the Lagrangian function,

$$\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b).$$

- These are non-linear equations: easier to solve them with first-order methods or Newton's.
- In the convex case as we are, the constraint $A\mathbf{x} = b$ is always qualified, provided that a solution exist! (I.e., b is in the image of A).

A.3 Optimality conditions: complete case

Consider now the case $X = \{\mathbf{x} | A\mathbf{x} = b, g(\mathbf{x}) \leq 0\}$, $A \in \mathbf{R}^{q \times n}$, $g : \mathbf{R}^n \rightarrow \mathbf{R}^p$, convex and differentiable, and for simplicity the case in which $f \in \mathcal{C}(\mathbf{R}^n)$. The problem is still convex. We have the following implication.

Theorem A.2 Consider the problem, $\min_{\mathbf{x} \in \mathbf{R}^n} f(\mathbf{x})$, subject to $A\mathbf{x} = b, g(\mathbf{x}) \leq 0$ if the constraints are qualified, then the optimality conditions $\nabla f(\mathbf{x}^*) + N_X(\mathbf{x}^*) \ni \mathbf{0}$ can be equivalently written as

$$\begin{cases} \nabla f(\mathbf{x}^*) + A^\top \lambda^* + \nabla^\top g(\mathbf{x}^*) \nu^* = 0 \\ A\mathbf{x}^* = b, \quad g(\mathbf{x}^*) \leq 0 \\ \nu^* \geq 0, \quad [\nu^*]_i [g(\mathbf{x}^*)]_i = 0, \forall i \in \{1, p\} \end{cases}$$

- The variables $\lambda \in \mathbf{R}^q, \nu^* \in \mathbf{R}^p$ are the Lagrangian multipliers or dual variables;
- We can write the optimality conditions as the saddle-points of the Lagrangian function,

$$\mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b) + \nu^\top g(\mathbf{x}),$$

with the restriction that $\nu \geq 0$.

- These are non-linear inequalities: not so easy to solve.
- In the convex case as we are, the constraint $g(\mathbf{x}) \leq 0$ is qualified for example if there exists at least a feasible point \mathbf{x}' for which $A\mathbf{x}' = b$ and $g(\mathbf{x}') < 0$. We call this condition: Slater's constraint qualification.

A.4 Extension to matrix decisions

It is often useful to look at problems whose decision variables are symmetric matrices. We can extend most of the results of these classes to this setting (and beyond to infinite dimensional Hilbert spaces), by considering the proper inner product. For squared symmetric matrices the inner product is the trace, as such,

$$\underset{X \in \mathbb{S}^n}{\text{minimise}} f(X), \quad \text{subject to } X \geq 0$$

we can form the Lagrangian with dual matrix variable Z and derive the optimality conditions,

$$\mathcal{L}(X, Z) = f(X) - \text{trace}(ZX) \quad [Z \geq 0] \quad \begin{cases} \nabla_X \mathcal{L}(X, Z) = 0 \rightarrow \nabla f(X) - Z = 0 \\ X \geq 0, Z \geq 0, \text{trace}(ZX) = 0. \end{cases}$$

A.5 Useful reformulations: epigraph and Schur's complement

Optimisation problems can be transformed into equivalent ones. By equivalent we mean different problems from which we can derive the solution of each other, for example:

$$\min_{x \in \mathbf{R}} f(x), \text{ subject to } |x| \leq 1 \quad \longleftrightarrow \quad \min_{x \in \mathbf{R}} f(x), \text{ subject to } x \leq 1, x \geq -1.$$

We are going to recap a few useful rules.

(A) Epigraph. (i.e., all the opt. problems can have a linear cost).

Consider the optimisation problem,

$$\min_{x \in \mathbf{R}^n} f(x), \text{ subject to } x \in X.$$

This problem is equivalent to the optimisation problem in $n + 1$ variables, with linear cost,

$$\min_{x \in \mathbf{R}^n, t \in \mathbf{R}} t, \text{ subject to } f(x) \leq t, x \in X.$$

(B) Substitutions.

$$\min_{x \in \mathbf{R}^n, y \in \mathbf{R}^m} f(x) + g(y), \text{ subject to } y = Ax + b \quad \longleftrightarrow \quad \min_{x \in \mathbf{R}^n} f(x) + g(Ax + b).$$

(C) Partial Min.

$$\min_{x \in X, y \in Y} f(x, y), \quad \longleftrightarrow \quad \min_{y \in Y} \underbrace{\left(\min_{x \in X} f(x, y) \right)}_{=\tilde{f}(y)}.$$

(D) Monotone transfer. Ψ_0 is a monotone increasing function, $\min f(x) \rightarrow \min \Psi_0 \circ f(x)$, ex:

$$\underbrace{\min_{x \in \mathbf{R}^n} \|Ax - b\|_2}_{\text{Non-diff, hard to solve}}, \quad \longleftrightarrow \quad \underbrace{\min_{x \in \mathbf{R}^n} \|Ax - b\|_2^2}_{\text{Super-easy to solve}}$$

(E) Schur's complement (very useful in SDPs)

Consider the squared symmetric matrix X , and partition it as follows,

$$X = \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}.$$

Then, the convex semidefinite constraint $X \geq 0$ can be equivalently formulated in two ways,

$$\text{First possibility: } \{A > 0, \quad S_1 = C - B^\top A^{-1} B \geq 0\};$$

$$\text{Second possibility: } \{C > 0, \quad S_2 = A - BC^{-1}B^\top \geq 0\};$$

and conversely.

A.6 (Lagrangian) Duality

Consider the problem,

$$\underset{\mathbf{x} \in X}{\text{minimise}} f(\mathbf{x}), \quad \text{subject to } A\mathbf{x} = b, g(\mathbf{x}) \leq 0$$

Define the Lagrangian,

$$\mathcal{L}(\mathbf{x}, \lambda, \nu) = f(\mathbf{x}) + \lambda^\top (A\mathbf{x} - b) + \nu^\top g(\mathbf{x}),$$

We call primal problem:

$$(P) \quad \inf_{\mathbf{x} \in X} \left(\sup_{\lambda, \nu \geq 0} \mathcal{L}(\mathbf{x}, \lambda, \nu) \right) = \begin{cases} \inf_{\mathbf{x} \in X} f(\mathbf{x}) & \text{if } A\mathbf{x} = b, g(\mathbf{x}) \leq 0, \\ +\infty & \text{otherwise.} \end{cases}$$

We define the dual function:

$$q(\lambda, \nu) := \inf_{\mathbf{x} \in X} \mathcal{L}(\mathbf{x}, \lambda, \nu),$$

and call the dual problem:

$$(D) \quad \sup_{\lambda, \nu \geq 0} \left(\inf_{\mathbf{x} \in X} \mathcal{L}(\mathbf{x}, \lambda, \nu) \right) = \sup_{\lambda, \nu \geq 0} q(\lambda, \nu).$$

The dual function $q(\lambda, \nu)$ is always concave, so D is always convex! The proof of it is that q is the point-wise infimum of affine functions. Rem: $\max_i \{a_i^\top \mathbf{x} + b_i\}$ is convex.

Furthermore, $d^* = \text{val}(D) \leq \text{val}(P) = p^*$ or otherwise stated

$$\sup_{\lambda, \nu \geq 0} \inf_{\mathbf{x} \in X} \mathcal{L}(\mathbf{x}, \lambda, \nu) \leq \inf_{\mathbf{x} \in X} \sup_{\lambda, \nu \geq 0} \mathcal{L}(\mathbf{x}, \lambda, \nu).$$

For convex problems that have a solution, and under constraint qualification (e.g., under Slater's condition), then the KKT point is a saddle-point and we don't have duality gap (i.e., $d^* = p^*$). So for most convex problems solving the dual problem or the primal is the same.

A.7 Determining the dual problem: examples

Example A.1 (Dual of an LP) *Let's look at an LP in the (standard) form:*

$$\underset{x \in \mathbf{R}^n}{\text{minimise}} \quad c^\top x, \quad \text{subject to } Ax = b, \quad x \geq 0.$$

for matrix $A \in \mathbf{R}^{m \times n}$, and $b \in \mathbf{R}^m$. Introduce the Lagrangian multipliers $\lambda \in \mathbf{R}^m$ for $Ax - b = 0$ and $\nu \in \mathbf{R}^n$ for $-x \leq 0$, and write the Lagrangian function:

$$\mathcal{L}(x, \lambda, \nu) = c^\top x + \lambda^\top (Ax - b) - \nu^\top x.$$

The dual function is given by

$$q(\lambda, \nu) = \inf_x \mathcal{L}(x, \lambda, \nu) = -\lambda^\top b + \inf_x (c^\top x + \lambda^\top Ax - \nu^\top x) = -b^\top \lambda + \inf_x (x^\top (A^\top \lambda - \nu + c)),$$

which yields,

$$q(\lambda, \nu) = \begin{cases} -b^\top \lambda & \text{if } A^\top \lambda - \nu + c = 0 \\ -\infty & \text{otherwise.} \end{cases}$$

The dual problem is then given as,

$$\sup_{\lambda, \nu \geq 0} q(\lambda, \nu),$$

so,

$$\underset{\lambda, \nu \geq 0}{\text{minimise}} \quad b^\top \lambda \quad \text{subject to: } A^\top \lambda - \nu + c = 0, \quad \nu \geq 0.$$

In fact $\nu \geq 0$ is redundant and can be eliminated as

$$\underset{\lambda \in \mathbf{R}^m}{\text{minimise}} \quad b^\top \lambda \quad \text{subject to: } A^\top \lambda + c \geq 0.$$

Note the duality between primal LP and dual LP, which is also a consequence of Farkas' lemma.

Example A.2 (Dual of a QCQP) *Consider the convex QCQP*

$$\begin{aligned} \underset{x \in \mathbf{R}^n}{\text{minimise}} \quad & (1/2)x^\top P_0 x + c_0^\top x + r_0 \\ \text{subject to} \quad & (1/2)x^\top P_i x + c_i^\top x + r_i \leq 0, \quad i = 1, \dots, m, \end{aligned}$$

with P_0 positive definite and P_i positive semi-definite.

The Lagrangian is

$$\mathcal{L}(x, \nu) = (1/2)x^\top P(\nu)x + s(\nu)^\top x + r(\nu),$$

with

$$P(\nu) = P_0 + \sum_{i=1}^m \nu_i P_i, \quad s(\nu) = s_0 + \sum_{i=1}^m \nu_i s_i, \quad r(\nu) = r_0 + \sum_{i=1}^m \nu_i r_i.$$

This gives a dual function of the form,

$$q(\nu) = \inf_x ((1/2)x^\top P(\nu)x + s(\nu)^\top x + r(\nu))$$

which can be computed analytically for $\nu \geq 0$, since $P(\nu)$ is positive definite in this case. By optimality conditions,

$$x^*(\nu) = -P(\nu)^{-1}s(\nu),$$

and therefore

$$q(\nu) = -(1/2)s(\nu)^\top P(\nu)^{-1}s(\nu) + r(\nu).$$

The dual problem is therefore,

$$\begin{aligned} \underset{\nu \in \mathbf{R}^m}{\text{minimise}} \quad & (1/2)s(\nu)^\top P(\nu)^{-1}s(\nu) - r(\nu) \\ \text{subject to} \quad & \nu \geq 0. \end{aligned}$$

Strong duality holds under Slater's condition, i.e., if there exists a x such that $(1/2)x^\top P_i x + c_i^\top x + r_i < 0$, $i = 1, \dots, m$.

For $P_i = 0$ (the QP case), then $P(\nu)^{-1} = P_0^{-1}$ and

$$s(\nu)^\top P(\nu)^{-1} s(\nu) = (s_0 + S^\top \nu)^\top P_0^{-1} (s_0 + S^\top \nu) = \nu^\top S P_0^{-1} S^\top \nu + 2s_0^\top P_0^{-1} \nu + s_0^\top P_0^{-1} s_0$$

where $S \in \mathbf{R}^{n \times m}$ is the matrix such that $S^\top = [s_1 | s_2 | \dots | s_m]$. As such, with some re-labelling, the dual is

$$\begin{aligned} & \underset{\nu \in \mathbf{R}^m}{\text{minimise}} && (1/2)\nu^\top Q\nu + w^\top \nu + z \\ & \text{subject to} && \nu \geq 0, \end{aligned}$$

which is a convex QP.

Let's get back to the general case. The dual problem is

$$\begin{aligned} & \underset{\nu \in \mathbf{R}^m}{\text{minimise}} && (1/2)s(\nu)^\top P(\nu)^{-1} s(\nu) - r(\nu) \\ & \text{subject to} && \nu \geq 0. \end{aligned}$$

Use the epigraph form,

$$\begin{aligned} & \underset{t, \nu \in \mathbf{R}^m}{\text{minimise}} && t - r(\nu) \\ & \text{subject to} && (1/2)s(\nu)^\top P(\nu)^{-1} s(\nu) \leq t \\ & && \nu \geq 0, t \geq 0, \end{aligned}$$

and then Schur's complement (since $P(\nu) > 0$),

$$\begin{aligned} & \underset{t, \nu \in \mathbf{R}^m}{\text{minimise}} && t - r(\nu) \\ & \text{subject to} && \begin{bmatrix} t & s(\nu)^\top \\ s(\nu) & 2P(\nu) \end{bmatrix} \geq 0 \\ & && \nu \geq 0, t \geq 0, \end{aligned}$$

Since $r(\nu)$, $s(\nu)$ and $P(\nu)$ are affine functions of ν , then the above is a SDP.

A.8 Conjugate function

Consider a function $f : \mathbf{R}^n \rightarrow \mathbf{R}$. We define its conjugate function as,

$$f_\star(\mathbf{y}) := \sup_{\mathbf{x} \in \mathbf{R}^n} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x}))$$

When f is closed convex and proper, then

$$(\partial f)^{-1} = \partial f_\star$$

where ∂f is the subdifferential set.