

Duality and augmented Lagrangian

Lecture notes

TAF MCE - UE Advanced optimization for IA

thierry.chonavel@imt-atlantique.fr

2021

Contents

1	Introduction	1
2	Dual problem and duality gap	1
2.1	Definition	1
2.2	Weak duality	2
2.3	Example	2
2.4	Concavity of the dual function	2
2.5	Strong duality and saddle point interpretation	2
3	Uzawa and Arrow-Hurwicz algorithms	3
4	Augmented Lagrangian	4
4.1	Equality constraints	4
4.2	Dual approach for augmented Lagrangian	4
4.3	Augmented Lagrangian with equality and inequality constraints	6
5	Example: sphere packing	6
6	Exercise: SVM and duality	7

1 Introduction

In these notes we introduce duality principle that is often useful in optimization. In constrained optimization problems, it can be noted that if Lagrange multipliers were known, finding solution often would be much simpler. Thus, focussing on the search for Lagrange multipliers can be of particular interest. The problem of finding them is known as the **dual problem**. In the case of convex problems, we will see that that there is some symmetry between the initial (also called **primal**) problem and the dual problem, while in the general case, the difference between the solutions of these two problems is measured by the **duality gap**. Let us consider the standard constrained optimization problem

$$\begin{cases} \min_{\mathbf{x}} f(\mathbf{x}) \\ \mathbf{g}(\mathbf{x}) \leq \mathbf{0} \\ \mathbf{h}(\mathbf{x}) = \mathbf{0} \end{cases} \quad (1)$$

with $\mathbf{h} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^p$. Problem (1) is often referred to as the **primal problem**. Alternatively, the

dual problem looks for optimization of Lagrange multipliers. Before going into details, let us give a few reasons why duality is worth being considered.

1. **the dual problem of (1) is a concave optimization problem**, regardless problem (1) is convex or not. Thus the dual problem is often simpler than the initial problem.
2. If problem (1) has a global solution given by $f(\mathbf{x}^*) = p^*$ and d^* denotes the optimum of the dual problem, we have $p^* \geq d^*$. In other words **solving the dual problem yields a lower bound for the solution of problem (1)**, what can be helpful to decide whether an approximate solution $\mathbf{x} \in \mathcal{D} = \text{dom}(\mathbf{g}) \cap \text{dom}(\mathbf{h})$ of (1) could be accepted by comparing $f(\mathbf{x})$ to d^* .
3. **The initial problem involves n variables and $m+p$ constraints while its dual problem involves $m+p$ variables and p simple constraints**. What makes a problem difficult is often handling a large number of complex inequality constraints rather than a large number of variables, making thus the dual problem often appealing when p is large.
4. Interesting algorithms can be derived from the study of duality principle.

In the following we are going to define the dual problem and justify these arguments.

2 Dual problem and duality gap

2.1 Definition

The Lagrangian of problem (1) writes

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}).$$

Clearly, for $\boldsymbol{\mu} \geq \mathbf{0}$, we have

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\mathbf{x}) + \mathcal{I}_{\{0\}}(\mathbf{h}(\mathbf{x})) + \mathcal{I}_{\mathbb{R}_-^m}(\mathbf{g}(\mathbf{x})), \quad (2)$$

where $\mathcal{I}_A(\mathbf{u}) = 0$ if all entries of \mathbf{u} are in A and $\mathcal{I}_A(\mathbf{x}) = +\infty$ otherwise. Thus, $L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ is a continuous underestimate of

the right hand term which is an objective equivalent to the initial constrained problem. The **Dual function** is defined by

$$\phi(\lambda, \mu) = \inf_{\mathbf{x}} L(\mathbf{x}, \lambda, \mu). \quad (3)$$

2.2 Weak duality

Letting

$$d^* \doteq \sup_{\lambda; \mu \geq \mathbf{0}} \phi(\lambda, \mu),$$

the quantity $p^* - d^*$ is called the **duality gap**:

Theorem 1 (weak duality) *The duality gap $p^* - d^*$ is positive:*

$$p^* \geq d^*.$$

Proof Relations (2) and (3) yield

$$\phi(\lambda, \mu) \leq \min_{\{\mathbf{x}; \mathbf{h}(\mathbf{x})=0, \mathbf{g}(\mathbf{x}) \leq 0\}} f(\mathbf{x}) = p^*.$$

Then,

$$d^* = \max_{\lambda; \mu \geq \mathbf{0}} \phi(\lambda, \mu) \leq p^*$$

□

In fact, this result is a particular case of the **max-min inequality** [4] which states that for any function $k(\mathbf{w}, \mathbf{z})$ the **max-min inequality** [4] is satisfied:

$$\sup_{\mathbf{z}} \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}) \leq \inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z}). \quad (4)$$

Indeed, since $k(\mathbf{w}, \mathbf{z}) \leq \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z})$, we have $\inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}) \leq \inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z})$ and as the right hand side of this inequality does not depend on \mathbf{z} , considering the supremum of both sides with respect to \mathbf{z} yields the max-min inequality. Theorem 1 is obtained letting $k(\mathbf{w}, \mathbf{z}) = L(\mathbf{w}, \mathbf{z})$, with $\mathbf{w} = \mathbf{x}$ and $\mathbf{z} = (\mu, \lambda)$.

2.3 Example

As an example, let us consider the problem $\min \frac{1}{2} \|\mathbf{x}\|^2$, s.t. (subject to) $\mathbf{A}\mathbf{x} = \mathbf{b}$. Then, $\nabla L(\mathbf{x}, \lambda) = \mathbf{x} + \mathbf{A}^T \lambda$ and for fixed λ the minimum of L is reached for $\mathbf{x} = -\mathbf{A}^T \lambda$. This leads to

$$\begin{aligned} \phi(\lambda) &= L(-\mathbf{A}^T \lambda, \lambda) \\ &= (1/2) \lambda^T \mathbf{A} \mathbf{A}^T \lambda + \lambda^T (-\mathbf{A} \mathbf{A}^T \lambda - \mathbf{b}) \\ &= -(1/2) \lambda^T \mathbf{A} \mathbf{A}^T \lambda - \lambda^T \mathbf{b}. \end{aligned}$$

2.4 Concavity of the dual function

Now, for the sake of simplicity, we can assume inequality constraints only since $\mathbf{h}(\mathbf{x}) = 0 \Leftrightarrow [\mathbf{h}(\mathbf{x}) \leq 0 \text{ and } \mathbf{h}(\mathbf{x}) \geq 0]$. In the previous example ϕ is clearly a concave function. In fact, this is a general property:

Theorem 2 *The set of \mathbb{R}_+^p where $\phi(\mu)$ is lower bounded is convex and in this set ϕ is concave.*

Proof For $0 \leq \alpha \leq 1$, and $\mu_1, \mu_2 \in \mathbb{R}_+^p$ where $\phi(\mu_1), \phi(\mu_2) > -\infty$

$$\begin{aligned} \phi(\alpha \mu_1 + (1-\alpha) \mu_2) &= \inf_{\mathbf{x}} \{f(\mathbf{x}) + (\alpha \mu_1 + (1-\alpha) \mu_2)^T \mathbf{g}(\mathbf{x})\} \\ &\geq \alpha \inf_{\mathbf{x}} \{f(\mathbf{x}) + \mu_1^T \mathbf{g}(\mathbf{x})\} \\ &\quad + (1-\alpha) \inf_{\mathbf{x}} \{f(\mathbf{x}) + \mu_2^T \mathbf{g}(\mathbf{x})\} \\ &\geq \alpha \phi(\mu_1) + (1-\alpha) \phi(\mu_2). \end{aligned}$$

□

2.5 Strong duality and saddle point interpretation

Let us consider again the example in subsection 2.3 for which first order necessary optimality conditions (NC1) clearly write

$$(\mathbf{x}^*, \lambda^*) = (\mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{b}, -(\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{b}).$$

We observe that

$$\begin{aligned} d^* &= \max_{\lambda} \phi(\lambda) \\ &= \phi(\lambda^*) \\ &= (1/2) \mathbf{b}^T (\mathbf{A} \mathbf{A}^T)^{-1} \mathbf{b} \\ &= (1/2) \|\mathbf{x}^*\|^2 \\ &= p^* \end{aligned}$$

Thus, it appears that there are cases where maximizing ϕ yields the multipliers that solve NC1 optimality conditions and the duality gap is 0. When it is satisfied, this property is called **strong duality**.

This leads us to a geometric interpretation of strong duality. As discussed at the end of section 2.2, for any multivariate function $k(\mathbf{w}, \mathbf{z})$ the max-min inequality is always satisfied: $\sup_{\mathbf{z}} \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}) \leq \inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z})$. On another hand, a **saddle point** for k is defined as a point $(\mathbf{w}^*, \mathbf{z}^*)$, such that

$$\forall \mathbf{w}, \mathbf{z}, k(\mathbf{w}^*, \mathbf{z}) \leq k(\mathbf{w}^*, \mathbf{z}^*) \leq k(\mathbf{w}, \mathbf{z}^*),$$

that is, at $(\mathbf{w}^*, \mathbf{z}^*)$ k has a minimum w.r.t. (with respect to) \mathbf{w} and a maximum w.r.t. \mathbf{z} . The next theorem shows that if there is a saddle point, at this point the max-min inequality for k can be replaced by an equality.

Theorem 3 *If $(\mathbf{w}^*, \mathbf{z}^*)$ is a saddle point of $k(\mathbf{w}, \mathbf{z})$, then*

$$\sup_{\mathbf{z}} \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}) = k(\mathbf{w}^*, \mathbf{z}^*) = \inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z})$$

Proof As $(\mathbf{w}^*, \mathbf{z}^*)$ is a saddle point,

$$\begin{aligned} \inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z}) &\leq \sup_{\mathbf{z}} k(\mathbf{w}^*, \mathbf{z}) \leq k(\mathbf{w}^*, \mathbf{z}^*) \\ k(\mathbf{w}^*, \mathbf{z}^*) &\leq \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}^*) \leq \sup_{\mathbf{z}} \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}). \end{aligned}$$

Thus,

$$\inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z}) \leq k(\mathbf{w}^*, \mathbf{z}^*) \leq \sup_{\mathbf{z}} \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}), \quad (5)$$

and from the max-min inequality we also have $\sup_{\mathbf{z}} \inf_{\mathbf{w}} k(\mathbf{w}, \mathbf{z}) \leq \inf_{\mathbf{w}} \sup_{\mathbf{z}} k(\mathbf{w}, \mathbf{z})$ so that all inequalities in Eq. (5) become equalities. □

Then, it appears that strong duality holds when (\mathbf{x}^*, μ^*) represents a saddle point of the Lagrangian:

$$d^* = \sup_{\mu \geq \mathbf{0}} \inf_{\mathbf{x}} L(\mathbf{x}, \mu) = \inf_{\mathbf{x}} \sup_{\mu \geq \mathbf{0}} L(\mathbf{x}, \mu) = p^*, \quad (6)$$

The connections between local optima and saddle points are summarized in the following theorem.

Theorem 4 *If (\mathbf{x}^*, μ^*) , with $\mu^* \in \mathbb{R}_+^p$ is a saddle point of $L(\mathbf{x}, \mu)$, then \mathbf{x}^* is a solution of the primal problem. Conversely, if the primal problem is convex with solution \mathbf{x}^* , there exists $\mu^* \in \mathbb{R}_+^p$ such that (\mathbf{x}^*, μ^*) , is a saddle point of $L(\mathbf{x}, \mu)$.*

Proof If $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a saddle point of $L(\mathbf{x}, \boldsymbol{\mu})$, the condition $L(\mathbf{x}^*, \boldsymbol{\mu}) \leq L(\mathbf{x}^*, \boldsymbol{\mu}^*)$ writes $(\boldsymbol{\mu} - \boldsymbol{\mu}^*)^T \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}$, what can occur for any $\boldsymbol{\mu} \in \mathbb{R}_+^p$ only if $\mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}$, showing thus that \mathbf{x}^* is feasible. Then $\boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}$, since $\boldsymbol{\mu}^* \in \mathbb{R}_+^p$. But letting $\boldsymbol{\mu} = \mathbf{0}$, we also get $(\mathbf{0} - \boldsymbol{\mu}^*)^T \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}$, that is, $\boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}^*) \geq \mathbf{0}$. Finally, $\boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}^*) = \sum_i \mu_i^* \mathbf{g}_i(\mathbf{x}^*) = \mathbf{0}$. Then, for any $\mathbf{x} \in \mathcal{D}$, inequality $L(\mathbf{x}^*, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*)$ yields

$$L(\mathbf{x}^*, \boldsymbol{\mu}^*) = f(\mathbf{x}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*) = f(\mathbf{x}) + \boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}) \leq f(\mathbf{x}),$$

which completes the proof of the first part of the theorem. Conversely, from KKT conditions, if \mathbf{x}^* is a solution of the primal problem, there exists $\boldsymbol{\mu}^* \in \mathbb{R}_+^p$ with $\boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}^*) = \mathbf{0}$ and $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\mu}^*) = \mathbf{0}$. Consider any $\boldsymbol{\mu} \in \mathbb{R}_+^p$ and $\mathbf{x} \in \mathcal{D}$. First note that

$$L(\mathbf{x}^*, \boldsymbol{\mu}) \leq f(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\mu}^*).$$

In addition, for a convex problem, $\mathbf{x} \rightarrow L(\mathbf{x}, \boldsymbol{\mu}^*)$ is convex and the necessary condition $\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\mu}^*) = \mathbf{0}$ is also sufficient to guarantee that $L(\mathbf{x}^*, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*)$. Finally, we get $L(\mathbf{x}^*, \boldsymbol{\mu}) \leq L(\mathbf{x}^*, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*)$ which proves that $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a saddle point of L . \square

This result suggests solving the primal problem by looking for a saddle point of the Lagrangian. A natural idea that will be developed in the next section consists in performing iteratively moves of \mathbf{x} that decrease the Lagrangian and moves of $\boldsymbol{\mu}$ (and $\boldsymbol{\lambda}$ when there are equality constraints), that increase the Lagrangian. This can be done via gradient steps. Recall the primal and dual problems:

$$\begin{aligned} (P) : & \text{ solve } \inf_{\mathbf{x}} \sup_{\boldsymbol{\mu}} L(\mathbf{x}, \boldsymbol{\mu}) = p^* \\ (Q) : & \text{ solve } \sup_{\boldsymbol{\mu} \geq \mathbf{0}} \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu}) = d^* \end{aligned} \quad (7)$$

Consider also the problems $(P_{\boldsymbol{\mu}}) : \text{ solve } \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu})$. From theorem 4, for a saddle point $(\mathbf{x}^*, \boldsymbol{\mu}^*)$, \mathbf{x}^* is a solution of (P) . In addition, since $\inf_{\mathbf{x}} \sup_{\boldsymbol{\mu}} L(\mathbf{x}, \boldsymbol{\mu}) = L(\mathbf{x}^*, \boldsymbol{\mu}^*) = \sup_{\boldsymbol{\mu}} \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu})$, $\boldsymbol{\mu}^*$ is a solution of the dual problem (Q) . Conversely, assuming that we can define uniquely $\mathbf{x}_{\boldsymbol{\mu}} = \arg \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu})$ (then $\phi(\boldsymbol{\mu}) = L(\mathbf{x}_{\boldsymbol{\mu}}, \boldsymbol{\mu})$), we would like the solution $\mathbf{x}_{\boldsymbol{\mu}^*}$ of $(P_{\boldsymbol{\mu}^*})$ where $\boldsymbol{\mu}^*$ denotes a solution of (Q) to be a solution of (P) . This is guaranteed by the following result:

Theorem 5 *If \mathbf{g} is continuous and for any $\boldsymbol{\mu} \in \mathbb{R}_+^p$ there is a unique solution $\mathbf{x}_{\boldsymbol{\mu}}$ of $(P_{\boldsymbol{\mu}})$, where $\boldsymbol{\mu} \rightarrow \mathbf{x}_{\boldsymbol{\mu}}$ is continuous, then $\phi(\boldsymbol{\mu})$ is derivable with $\nabla_{\boldsymbol{\mu}} \phi(\boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}})$. In addition, if $\boldsymbol{\mu}^*$ denotes any solution of (Q) , $\mathbf{x}_{\boldsymbol{\mu}^*}$ is a solution of (P) .*

Proof First, let us show that $\phi(\boldsymbol{\mu})$ is derivable. Let $\boldsymbol{\mu}, \boldsymbol{\mu} + \boldsymbol{\delta} \in \mathbb{R}_+^p$. Then, since $\phi(\boldsymbol{\mu}) = L(\mathbf{x}_{\boldsymbol{\mu}}, \boldsymbol{\mu}) \leq L(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}, \boldsymbol{\mu})$ and $\phi(\boldsymbol{\mu} + \boldsymbol{\delta}) = L(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}, \boldsymbol{\mu} + \boldsymbol{\delta}) \leq L(\mathbf{x}_{\boldsymbol{\mu}}, \boldsymbol{\mu} + \boldsymbol{\delta})$, we have

$$\begin{aligned} \phi(\boldsymbol{\mu} + \boldsymbol{\delta}) - \phi(\boldsymbol{\mu}) & \leq L(\mathbf{x}_{\boldsymbol{\mu}}, \boldsymbol{\mu} + \boldsymbol{\delta}) - L(\mathbf{x}_{\boldsymbol{\mu}}, \boldsymbol{\mu}) \\ & \leq \boldsymbol{\delta}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}}) \\ \phi(\boldsymbol{\mu} + \boldsymbol{\delta}) - \phi(\boldsymbol{\mu}) & \geq L(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}, \boldsymbol{\mu} + \boldsymbol{\delta}) - L(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}, \boldsymbol{\mu}) \\ & \geq \boldsymbol{\delta}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}). \end{aligned}$$

Then, $\phi(\boldsymbol{\mu} + \boldsymbol{\delta}) - \phi(\boldsymbol{\mu}) \in [\boldsymbol{\delta}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}), \boldsymbol{\delta}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}})]$ and there exists $\alpha \in [0, 1]$ such that

$$\begin{aligned} \phi(\boldsymbol{\mu} + \boldsymbol{\delta}) - \phi(\boldsymbol{\mu}) & = \boldsymbol{\delta}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}}) + \alpha \boldsymbol{\delta}^T [\mathbf{g}(\mathbf{x}_{\boldsymbol{\mu} + \boldsymbol{\delta}}) - \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}})] \\ & = \boldsymbol{\delta}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}}) + \|\boldsymbol{\delta}\| \epsilon(\boldsymbol{\delta}) \end{aligned}$$

with $\lim_{\boldsymbol{\delta} \rightarrow \mathbf{0}} \epsilon(\boldsymbol{\delta}) = 0$, from the continuity of \mathbf{g} . This shows that $\phi(\boldsymbol{\mu})$ is derivable, with

$$\nabla_{\boldsymbol{\mu}} \phi(\boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}}). \quad (8)$$

Now, as $\boldsymbol{\mu}^*$ maximizes ϕ over \mathbb{R}_+^p , Euler inequality writes $(\boldsymbol{\mu} - \boldsymbol{\mu}^*)^T \nabla_{\boldsymbol{\mu}} \phi(\boldsymbol{\mu}^*) \leq 0$ for any $\boldsymbol{\mu} \in \mathbb{R}_+^p$, that is, $\boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}^*}) \leq \boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}^*})$. This leads to

$$\begin{aligned} L(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}) & = f(\mathbf{x}_{\boldsymbol{\mu}^*}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}^*}) \\ & \leq f(\mathbf{x}_{\boldsymbol{\mu}^*}) + \boldsymbol{\mu}^{*T} \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}^*}) \\ & \leq L(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}^*) \end{aligned} \quad (9)$$

Then we clearly have $L(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}) \leq L(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*)$ which shows that $(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}^*)$ is a saddle point. Thus, from the first part of theorem 4, $\mathbf{x}_{\boldsymbol{\mu}^*}$ is a solution of (P) . \square

Remarks

1. It is interesting to note that theorem 5 does not even require the constraints \mathbf{g} to be derivable (one can check that this applies to the first part of theorem 4 it resorts to).
2. The property $\nabla_{\boldsymbol{\mu}} \phi(\boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}_{\boldsymbol{\mu}})$ is very interesting for using gradient ascent to increase L w.r.t. $\boldsymbol{\mu}$ from point $(\mathbf{x}_{\boldsymbol{\mu}}, \boldsymbol{\mu})$ as will be done in the next section.
3. Consider $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu}) : \mathbb{R}^{n+p} \rightarrow \mathbb{R}^n$. Then, using arguments similar to those used in the sensitivity analysis ([5] section 2.1) that involves implicit function theorem, it can be shown that a sufficient condition for the function $\boldsymbol{\mu} \rightarrow \mathbf{x}_{\boldsymbol{\mu}}$ to be uniquely defined about $(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}^*)$ is that $\mathbf{x}_{\boldsymbol{\mu}^*}$ is a regular point and $\nabla^2 L(\mathbf{x}_{\boldsymbol{\mu}^*}, \boldsymbol{\mu}^*) \succ \mathbf{0}$.

3 Uzawa and Arrow-Hurwicz algorithms

Considering saddle point interpretation of strong duality suggests the **Uzawa algorithm** or the **Arrow-Hurwicz algorithm** (algorithm 1) for a problem in the form $\min_{\mathbf{x}} f(\mathbf{x})$ s.t. $\mathbf{h}(\mathbf{x}) = \mathbf{0}, \mathbf{g}(\mathbf{x}) \leq \mathbf{0}$

Algorithm 1 Uzawa-Arrow-Hurwicz algorithms

- 1: **init** \mathbf{x}_0 with $\mathbf{g}(\mathbf{x}_0) < \mathbf{0}, \boldsymbol{\mu}_0 > \mathbf{0}$
 - 2: **while** stopping condition \neq true **do**
 - 3: $\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu}_k)$ ▷ Uzawa
or

$$\begin{cases} \alpha_k & \approx \arg \min_r L(\mathbf{x}_k - r \nabla_{\mathbf{x}} L(\mathbf{x}_k, \boldsymbol{\mu}_k)) \\ \mathbf{x}_{k+1} & = \mathbf{x}_k - \alpha_k \nabla_{\mathbf{x}} L(\mathbf{x}_k, \boldsymbol{\mu}_k) \end{cases}$$
▷ Arrow-Hurwicz
 - 4: $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \beta_k \mathbf{h}(\mathbf{x}_{k+1})$
 - 5: $\boldsymbol{\mu}_{k+1} = [\boldsymbol{\mu}_k + \gamma_k \mathbf{g}(\mathbf{x}_{k+1})]_+$
 - 6: **end while**
-

Arrow-Hurwicz algorithm replaces the minimization of the Lagrangian w.r.t \mathbf{x} in Uzawa algorithm by a gradient step. In algorithm 1, $[\mathbf{z}]_+$ denotes the vector with i -th entry $[\mathbf{z}]_{+,i} = \max(\mathbf{z}_i, 0)$. Note also that a gradient step projected on \mathbb{R}_+^p is considered for the update of $\boldsymbol{\mu}$.

Convergence guarantees can be obtained for Uzawa algorithm. In particular, restricting our interest to the case of linear inequality constraints, we have

Theorem 6 *For a problem in the form $\min_{\mathbf{x}} f(\mathbf{x})$ s.t. $\mathbf{A}\mathbf{x} \leq \mathbf{b}$, where f is α -strongly convex [6], variables \mathbf{x}_k generated by Uzawa algorithm converge to a solution \mathbf{x}^* of the primal problem for fixed step size $\gamma_k = \gamma$ provided*

$$0 < \gamma < \frac{2\alpha}{\|\mathbf{A}\|^2},$$

with $\| \cdot \|$ the operator norm [7]. In addition, if \mathbf{A} is full rank, $\boldsymbol{\mu}_k$ converges to the corresponding Lagrange multiplier $\boldsymbol{\mu}^*$ solution of KKT conditions.

Proof Uzawa iterations write

$$\begin{cases} \mathbf{x}_{k+1} &= \arg \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu}_k) \\ \boldsymbol{\mu}_{k+1} &= [\boldsymbol{\mu}_k + \gamma \mathbf{g}(\mathbf{x}_k)]_+ \end{cases} \quad (10)$$

Then, $\nabla_{\mathbf{x}} f(\mathbf{x}_k) + \boldsymbol{\mu}_k^T \nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x}_k) = \mathbf{0}$. On another hand, a fixed point solution $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ should satisfy $\mathbf{x}^* = \arg \min_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\mu}^*)$ and $\boldsymbol{\mu}^* = [\boldsymbol{\mu}^* + \gamma \mathbf{g}(\mathbf{x}^*)]_+$. Indeed, the necessary Euler maximality condition for ϕ at $\boldsymbol{\mu}^*$ writes

$$\nabla_{\boldsymbol{\mu}} \phi(\mathbf{x}^*)^T (\boldsymbol{\mu} - \boldsymbol{\mu}^*) \leq 0, \quad \forall \boldsymbol{\mu} \in \mathbb{R}_+^p.$$

But since $\nabla_{\boldsymbol{\mu}} \phi(\boldsymbol{\mu}^*) = \mathbf{g}(\mathbf{x}^*)$ (Eq. 8), for $\gamma > 0$, we also get

$$(\boldsymbol{\mu}^* + \gamma \nabla \phi(\mathbf{x}^*) - \boldsymbol{\mu}^*)^T (\boldsymbol{\mu} - \boldsymbol{\mu}^*) \leq 0, \quad \forall \boldsymbol{\mu} \in \mathbb{R}_+^p.$$

Then, from the characterization of the projection on a closed convex set (here \mathbb{R}_+^p) [8], it appears that $\boldsymbol{\mu}^* = [\boldsymbol{\mu}^* + \mathbf{g}(\mathbf{x}^*)]_+$. Then, since $\boldsymbol{\mu}_{k+1} = [\boldsymbol{\mu}_k + \gamma \mathbf{g}(\mathbf{x}_{k+1})]_+$, we get

$$\begin{aligned} & \| \boldsymbol{\mu}_{k+1} - \boldsymbol{\mu}^* \|^2 \\ &= \| [\boldsymbol{\mu}_k + \gamma \mathbf{g}(\mathbf{x}_{k+1})]_+ - [\boldsymbol{\mu}^* + \gamma \mathbf{g}(\mathbf{x}^*)]_+ \|^2 \\ &\leq \| \boldsymbol{\mu}_k + \gamma \mathbf{g}(\mathbf{x}_{k+1}) - \boldsymbol{\mu}^* + \gamma \mathbf{g}(\mathbf{x}^*) \|^2 \\ &\leq \| \boldsymbol{\mu}_k - \boldsymbol{\mu}^* + \gamma \mathbf{A}(\mathbf{x}_{k+1} - \mathbf{x}^*) \|^2 \\ &\leq \| \boldsymbol{\mu}_k - \boldsymbol{\mu}^* \|^2 + \gamma^2 \| \mathbf{A} \|^2 \| \mathbf{x}_{k+1} - \mathbf{x}^* \|^2 \\ &\quad + 2\gamma (\boldsymbol{\mu}_k - \boldsymbol{\mu}^*)^T \mathbf{A}(\mathbf{x}_{k+1} - \mathbf{x}^*) \\ &\leq \| \boldsymbol{\mu}_k - \boldsymbol{\mu}^* \|^2 + \gamma^2 \| \mathbf{A} \|^2 \| \mathbf{x}_{k+1} - \mathbf{x}^* \|^2 \\ &\quad - 2\gamma (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}^*))^T (\mathbf{x}_{k+1} - \mathbf{x}^*) \\ &\leq \| \boldsymbol{\mu}_k - \boldsymbol{\mu}^* \|^2 - \gamma (2\alpha - \gamma \| \mathbf{A} \|^2) \| \mathbf{x}_{k+1} - \mathbf{x}^* \|^2 \end{aligned}$$

where the first inequality holds because projection onto a convex set is contractive [9], the second to last from necessary optimality conditions $\nabla_{\mathbf{x}} f(\mathbf{x}_{k+1}) + \mathbf{A}^T \boldsymbol{\mu}_k = \mathbf{0}$ and $\nabla_{\mathbf{x}} f(\mathbf{x}^*) + \mathbf{A}^T \boldsymbol{\mu}^* = \mathbf{0}$. and the last one from strong convexity characterization:

$$(\nabla_{\mathbf{x}} f(\mathbf{x}_{k+1}) - \nabla_{\mathbf{x}} f(\mathbf{x}^*))^T (\mathbf{x}_{k+1} - \mathbf{x}^*) \geq \alpha \| \mathbf{x}_{k+1} - \mathbf{x}^* \|^2.$$

Then, from hypothesis $2\alpha - \gamma \| \mathbf{A} \|^2 > 0$, $\| \boldsymbol{\mu}_k - \boldsymbol{\mu}^* \|$ decreases to a limit and we must have $\| \mathbf{x}_{k+1} - \mathbf{x}^* \| \rightarrow 0$. In addition, if \mathbf{A} is full rank, relations $\nabla_{\mathbf{x}} f(\mathbf{x}_{k+1}) + \mathbf{A}^T \boldsymbol{\mu}_k = \mathbf{0}$ yield $\boldsymbol{\mu}_k = (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{A} \nabla_{\mathbf{x}} f(\mathbf{x}_k)$ and $\boldsymbol{\mu}_k$ converges to $\boldsymbol{\mu}^* = (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{A} \nabla_{\mathbf{x}} f(\mathbf{x}^*)$. \square

4 Augmented Lagrangian

If we look at the Lagrangian, for a linear program $\min \mathbf{f}^T \mathbf{x}$ s.t. $\mathbf{A}\mathbf{x} \leq \mathbf{b}$, for fixed $\boldsymbol{\mu}$ the Lagrangian

$$L(\mathbf{x}, \boldsymbol{\mu}) = (\mathbf{f}^T + \boldsymbol{\mu}^T \mathbf{A})\mathbf{x} - \boldsymbol{\mu}^T \mathbf{b}$$

has no lower bound. Thus, theorem 5 does not apply and we cannot recover \mathbf{x}^* by solving the dual problem. **Augmented Lagrangian** techniques permit to combat this kind of difficulty by adding a penalty term to the objective.

In the case of equality constraints $\mathbf{h}(\mathbf{x}) = \mathbf{0}$ the penalty can be $\frac{c}{2} \| \mathbf{h}(\mathbf{x}) \|^2$ while for inequality constraints $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$ it can be taken as $\frac{c}{2} \sum_j \max(0, \mathbf{g}_j(\mathbf{x}))^2$ since condition $\mathbf{g}_j(\mathbf{x}) \leq 0$ also writes $\max(0, \mathbf{g}_j(\mathbf{x})) = 0$. Thus in what follows we simply assume equality constraints. We will come back to the case of inequality constraints at the end of the section.

4.1 Equality constraints

For the problem

$$(P) \begin{cases} \min f(\mathbf{x}) \\ \mathbf{h}(\mathbf{x}) = 0, \end{cases}$$

the augmented Lagrangian method combines ideas from dual algorithms described in the previous section and penalty techniques, simply replacing the initial problem (P) by problem

$$(P_c) \begin{cases} \min f(\mathbf{x}) + \frac{c}{2} \| \mathbf{h}(\mathbf{x}) \|^2 \\ \mathbf{h}(\mathbf{x}) = 0. \end{cases}$$

It is clear that problems (P) and (P_c) share the same minima. But we are going to see that introducing the penalty term $\frac{c}{2} \| \mathbf{h}(\mathbf{x}) \|^2$ can significantly improve performance.

For the problem (P) the augmented Lagrangian writes

$$L_c(\mathbf{x}, \lambda) = f(\mathbf{x}) + \frac{c}{2} \| \mathbf{h}(\mathbf{x}) \|^2 + \lambda^T \mathbf{h}(\mathbf{x})$$

with $c > 0$. In other words, L_c is the Lagrangian of the problem (P_c).

Note that

$$\nabla_{\mathbf{x}} (f(\mathbf{x}) + \frac{c}{2} \| \mathbf{h}(\mathbf{x}) \|^2) = \nabla_{\mathbf{x}} f(\mathbf{x}) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}) [\mathbf{c}\mathbf{h}(\mathbf{x})].$$

Assume that \mathbf{x}_c minimizes the unconstrained penalized objective $f_c(\mathbf{x}) = f(\mathbf{x}) + \frac{c}{2} \| \mathbf{h}(\mathbf{x}) \|^2$ and $(\mathbf{x}^*, \lambda^*)$ is the solution of (P). Then,

$$\begin{aligned} \nabla_{\mathbf{x}} f_c(\mathbf{x}_c) &= \nabla_{\mathbf{x}} f(\mathbf{x}_c) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_c) [\mathbf{c}\mathbf{h}(\mathbf{x}_c)] \\ &= \nabla_{\mathbf{x}} f(\mathbf{x}_c) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_c) \boldsymbol{\lambda}_c \\ &= 0. \end{aligned}$$

where $\boldsymbol{\lambda}_c = \mathbf{c}\mathbf{h}(\mathbf{x}_c)$. As c becomes large, the minimum of f_c tends to come closer to set $\{\mathbf{x}; \mathbf{h}(\mathbf{x}) = \mathbf{0}\}$. Assuming that $(\mathbf{x}^*, \lambda^*)$ is the solution of (P), \mathbf{x}_c approaches \mathbf{x}^* while $\boldsymbol{\lambda}_c$ approaches λ^* . An interesting aspect of problem (P_c) compared to the unconstrained minimization of f_c with increasing c is that it can benefit from improvement of the behavior of the objective brought by the penalty for sufficiently large c without having to let it tend to infinity since problems (P) and (P_c) are equivalent. In addition, adapting the dual algorithms considered in the previous section to augmented Lagrangian problems (P_c) yields several benefits that are discussed in the next section.

4.2 Dual approach for augmented Lagrangian

Convergence speed of gradient algorithms With a view to study the benefits of dual approaches for augmented Lagrangian let us first recall the following result related to the convergence of gradient algorithms

Theorem 7 *Let f a convex function with minimum f^* such that $m\mathbf{I} \leq \nabla^2 f(x) \leq M\mathbf{I}$. Then, the gradient algorithm with optimal step-size yields a sequence $(\mathbf{x}_k)_k$ such that*

$$f(\mathbf{x}_{k+1}) - f^* \leq (1 - m/M)(f(\mathbf{x}_k) - f^*).$$

Proof Note first that Taylor-Lagrange formula yields

$$f(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k) \leq f(\mathbf{x}_k) - \alpha \| \nabla f(\mathbf{x}_k) \|^2 + \frac{\alpha^2 M}{2} \| \nabla f(\mathbf{x}_k) \|^2.$$

The minimum of the right hand side is reached for $\alpha = 1/M$, leading to $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) - \frac{1}{2M} \| \nabla f(\mathbf{x}_k) \|^2$. Then,

$$f(\mathbf{x}_{k+1}) - f^* \leq f(\mathbf{x}_k) - f^* - \frac{1}{2M} \| \nabla f(\mathbf{x}_k) \|^2. \quad (11)$$

We also have

$$f(\mathbf{x}) \geq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{m}{2} \|\mathbf{x} - \mathbf{x}_k\|^2.$$

and the minimum of the right hand side is obtained for $\mathbf{x} = \mathbf{x}_k - \frac{1}{m} \nabla f(\mathbf{x}_k)$. thus,

$$f^* \geq f(\mathbf{x}_k) - \frac{1}{2m} \|\nabla f(\mathbf{x}_k)\|^2,$$

leading to $-\|\nabla f(\mathbf{x}_k)\|^2 \leq 2m(f^* - f(\mathbf{x}_k))$. Inserting this inequality in Eq. (11) yields

$$f(\mathbf{x}_{k+1}) - f^* \leq (1 - \frac{m}{M})(f(\mathbf{x}_k) - f^*). \quad (12)$$

□

This result shows that the convergence speed of gradient descent algorithms is related to the condition number $K = M/m$ ($K \geq 1$) of the Hessian matrix of the objective f : the lower K the better for convergence speed.

Hessian matrices of Lagrangians and dual functions
Now, we are going to calculate the Hessian matrices of Lagrangians L and L_c and of the dual functions $\phi(\boldsymbol{\lambda})$ and $\phi_c(\boldsymbol{\lambda})$ related to problems (P) and (P_c) respectively.

Letting $\mathbf{x}_\lambda = \arg \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})$, we have the following equations:

$$\begin{aligned} \nabla_{\mathbf{x}} L(\mathbf{x}_\lambda, \boldsymbol{\lambda}) &= \nabla_{\mathbf{x}} f(\mathbf{x}_\lambda) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda) \boldsymbol{\lambda} \\ &= 0 \\ \nabla_{\boldsymbol{\lambda}} [\nabla_{\mathbf{x}} L(\mathbf{x}_\lambda, \boldsymbol{\lambda})] &= \nabla_{\boldsymbol{\lambda}} \mathbf{x}_\lambda (\nabla_{\mathbf{x}}^2 f(\mathbf{x}_\lambda) + \nabla_{\mathbf{x}}^2 \mathbf{h}(\mathbf{x}_\lambda) \boldsymbol{\lambda}) \\ &\quad + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda)^T \\ &= \nabla_{\boldsymbol{\lambda}} \mathbf{x}_\lambda \nabla_{\mathbf{x}}^2 L(\mathbf{x}_\lambda, \boldsymbol{\lambda}) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda)^T \\ &= 0. \end{aligned} \quad (13)$$

On another hand, as discussed at the end of section 2.5, $\nabla_{\boldsymbol{\lambda}} \phi(\boldsymbol{\lambda}) = \mathbf{h}(\mathbf{x}_\lambda)$ and

$$\nabla_{\boldsymbol{\lambda}}^2 \phi(\boldsymbol{\lambda}) = \nabla_{\boldsymbol{\lambda}} \mathbf{x}_\lambda \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda) \quad (14)$$

The last equation of (13) yields

$$\nabla_{\boldsymbol{\lambda}} \mathbf{x}_\lambda = -\nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda)^T [\nabla_{\mathbf{x}} L(\mathbf{x}_\lambda, \boldsymbol{\lambda})]^{-1}, \quad (15)$$

and replacing (15) inside (14), we get

$$\nabla_{\boldsymbol{\lambda}}^2 \phi(\boldsymbol{\lambda}) = -\nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda)^T [\nabla_{\mathbf{x}} L(\mathbf{x}_\lambda, \boldsymbol{\lambda})]^{-1} \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda). \quad (16)$$

Now, let us consider the same calculations for ϕ_c : omitting variables to simplify notations, at $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ we get

$$\begin{aligned} \nabla_{\mathbf{x}} L_c &= \nabla_{\mathbf{x}} L + c \nabla_{\mathbf{x}} \mathbf{h} \times \mathbf{h} \\ &= 0 \\ \nabla_{\boldsymbol{\lambda}} \nabla_{\mathbf{x}} L_c &= \nabla_{\boldsymbol{\lambda}} \mathbf{x} \nabla_{\mathbf{x}}^2 L + \nabla_{\mathbf{x}} \mathbf{h}^T \\ &\quad + c \nabla_{\boldsymbol{\lambda}} \mathbf{x} (\nabla_{\mathbf{x}}^2 \mathbf{h} \times \mathbf{h} + \nabla_{\mathbf{x}} \mathbf{h} \nabla_{\mathbf{x}} \mathbf{h}^T) \\ &= \nabla_{\boldsymbol{\lambda}} \mathbf{x} \nabla_{\mathbf{x}}^2 L_c + \nabla_{\mathbf{x}} \mathbf{h}^T \\ &= 0, \end{aligned} \quad (17)$$

where

$$\nabla_{\mathbf{x}}^2 L_c = \nabla_{\mathbf{x}}^2 L + c \nabla_{\mathbf{x}} \mathbf{h} \nabla_{\mathbf{x}} \mathbf{h}^T + c \nabla_{\mathbf{x}}^2 \mathbf{h} \times \mathbf{h}. \quad (18)$$

Then putting all pieces together just like for the study of $\nabla_{\boldsymbol{\lambda}}^2 \phi$, we get

$$\nabla_{\boldsymbol{\lambda}}^2 \phi_c(\boldsymbol{\lambda}) = -\nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda)^T [\nabla_{\mathbf{x}} L_c(\mathbf{x}_\lambda, \boldsymbol{\lambda})]^{-1} \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}_\lambda). \quad (19)$$

Note that for \mathbf{x}_λ feasible, $\mathbf{h}(\mathbf{x}_\lambda) = \mathbf{0}$ and the last term of (18) cancels: $\nabla_{\mathbf{x}}^2 L_c = \nabla_{\mathbf{x}}^2 L + c \nabla_{\mathbf{x}} \mathbf{h} \nabla_{\mathbf{x}} \mathbf{h}^T$. Close to the optimum, $\mathbf{x}_\lambda \approx \mathbf{x}_\lambda^* = \mathbf{x}^*$ and $\mathbf{h}(\mathbf{x}_\lambda) \approx 0$, leading to

$$\nabla_{\boldsymbol{\lambda}}^2 \phi_c(\boldsymbol{\lambda}) \approx -\nabla_{\mathbf{x}} \mathbf{h}^T [\nabla_{\mathbf{x}}^2 L + c \nabla_{\mathbf{x}} \mathbf{h} \nabla_{\mathbf{x}} \mathbf{h}^T]^{-1} \nabla_{\mathbf{x}} \mathbf{h}, \quad (20)$$

with equality when $\mathbf{h}(\mathbf{x}_\lambda) = 0$, and in particular at $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$.

Local strict convexity of the augmented Lagrangian

When second order minimality conditions are satisfied for problem (P) at $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$, it can be shown that for sufficiently large c the augmented Lagrangian L_c is strictly convex locally, w.r.t. \mathbf{x} :

Theorem 8 *Assume $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ satisfies second order sufficient optimality conditions for problem (P) . Then, there exists c^* such that for $c > c^*$, the function $\mathbf{x} \rightarrow L_c(\mathbf{x}, \boldsymbol{\lambda}^*)$ has a strict local minimum at \mathbf{x}^* .*

Proof First, let us show that there exists $c^* > 0$ such that $c > c^* \Rightarrow \nabla_{\mathbf{x}}^2 L_c(\mathbf{x}^*, \boldsymbol{\lambda}^*) \succ \mathbf{0}$. Indeed, if this was untrue, there would exist a sequence $(\mathbf{x}_n)_{n \geq 0}$ with $\|\mathbf{x}_n\| = 1$ such that $\mathbf{x}_n^T \nabla_{\mathbf{x}}^2 L_c(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x}_n \leq 0$, that is,

$$\mathbf{x}_n^T \nabla_{\mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x}_n + n \|\nabla \mathbf{h}(\mathbf{x}^*) \mathbf{x}_n\|^2 \leq 0.$$

\mathbf{x}_n admits a converging subsequence with limit \mathbf{x} such that $\|\nabla \mathbf{h}(\mathbf{x}^*) \mathbf{x}_n\|^2 = 0$ since $\mathbf{x}_n^T \nabla_{\mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x}_n$ is bounded. Thus, \mathbf{x} belongs to the tangent space at \mathbf{x}^* . But then we should also have $\mathbf{x}^T \nabla_{\mathbf{x}}^2 L(\mathbf{x}^*, \boldsymbol{\lambda}^*) \mathbf{x} \leq 0$ what would be in contradiction with optimality second order sufficient conditions.

Thus $\nabla_{\mathbf{x}}^2 L_c(\mathbf{x}^*, \boldsymbol{\lambda}^*) \succ \mathbf{0}$ for c large enough. In addition, $\nabla_{\mathbf{x}} L_c(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) + c \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x}^*) \times \mathbf{h}(\mathbf{x}^*) = 0$ since $\mathbf{h}(\mathbf{x}^*) = 0$. Thus second order sufficient conditions are also satisfied for (P_c) and $\mathbf{x} \rightarrow L_c(\mathbf{x}, \boldsymbol{\lambda}^*)$ has a strict local minimum at \mathbf{x}^* . □

The proof of theorem 8 shows that second order necessary conditions yield strict convexity property $\nabla^2 L_c \succ \mathbf{0}$ at $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ when c is large enough. Slightly moving about $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ maintains the strict convexity of L_c what guarantees that the Lagrangian can be minimized, hence the existence of the dual function.

Example For the problem (P) : $\min f(x, y) = x^2 + y$ s.t. $y = 0$, the dual cannot be calculated while it can be checked that for any $c > 0$, $\nabla^2 L_c \succ \mathbf{0}$ and the dual of (P_c) is well defined and dual optimization techniques can be considered.

Behavior of the dual Hessian After considering the Hessian of the Lagrangian, we now study the eigenvalues of the Hessian of the dual function for the augmented problem (P_c) . Recall that at $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ we have

$$\nabla_{\boldsymbol{\lambda}}^2 \phi_c = -\nabla_{\mathbf{x}} \mathbf{h}^T [\nabla_{\mathbf{x}}^2 L + c \nabla_{\mathbf{x}} \mathbf{h} \nabla_{\mathbf{x}} \mathbf{h}^T]^{-1} \nabla_{\mathbf{x}} \mathbf{h}.$$

Letting $\mathbf{A} = \nabla_{\mathbf{x}}^2 L$, $\mathbf{B} = \nabla_{\mathbf{x}} \mathbf{h}$ and $\mathbf{M} = -\nabla_{\boldsymbol{\lambda}}^2 \phi = \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}$, Woodbury formula (a.k.a. the matrix inversion lemma) [10] yields

$$\begin{aligned} -\nabla_{\boldsymbol{\lambda}}^2 \phi_c &= \mathbf{B}^T (\mathbf{A} + c \mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B} \\ &= \mathbf{B}^T [\mathbf{A}^{-1} - c \mathbf{A}^{-1} \mathbf{B} (\mathbf{I} + c \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{A}] \mathbf{B} \\ &= \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B} \\ &\quad - c \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B} (\mathbf{I} + c \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B} \\ &= \mathbf{M} - c \mathbf{M} (\mathbf{I} + c \mathbf{M})^{-1} \mathbf{M} \\ &= \mathbf{M} (\mathbf{I} + c \mathbf{M})^{-1} (\mathbf{I} + c \mathbf{M} - c \mathbf{M}) \\ &= \mathbf{M} (\mathbf{I} + c \mathbf{M})^{-1} \\ &= \mathbf{U} \mathbf{D} (\mathbf{I} + c \mathbf{D})^{-1} \mathbf{U}^T \end{aligned}$$

where $\mathbf{M} = \mathbf{U} \mathbf{D} \mathbf{U}^T$ denotes the eigen-decomposition of \mathbf{M} . Then, if \mathbf{u} is an eigenvector of $\mathbf{M} = -\nabla_{\boldsymbol{\lambda}}^2 \phi$ with corresponding eigenvalue α , it is also an eigenvector of $-\nabla_{\boldsymbol{\lambda}}^2 \phi_c$ with corresponding eigenvalue

$$\alpha_c = \frac{\alpha}{1 + c\alpha} = \left(\frac{1}{\alpha} + c\right)^{-1}$$

Letting α_m and α_M denote the smallest and largest eigenvalues of $-\nabla_{\lambda}^2 \phi$, the condition number of $-\nabla_{\lambda}^2 \phi_c$ is given by

$$K_c = \frac{\frac{1}{\alpha_m} + c}{\frac{1}{\alpha_M} + c}.$$

K_c decreases from α_M/α_m to 1 as c grows. This makes the steepest ascent for the maximization of ϕ_c interesting due to the faster convergence when the condition number is close to 1, as shown in theorem 7 about convergence speed of gradient algorithms.

Now, note that for large c the eigenvalues of $-\nabla_{\lambda}^2 \phi_c$ are close to c^{-1} so that $\nabla_{\lambda}^2 \phi_c \approx -c^{-1}\mathbf{I}$. Thus, updating λ via Newton algorithm for ϕ_c maximization simply writes

$$\lambda_{k+1} = \lambda_k + \mathbf{c}\mathbf{h}(\mathbf{x}_{\lambda_k}). \quad (21)$$

4.3 Augmented Lagrangian with equality and inequality constraints

In order to complete the discussion about the augmented Lagrangian approach, we now briefly discuss one popular formulation for the standard problem (1) with equality and inequality constraints. Clearly, such a problem can be written in the form

$$\begin{cases} \min f(\mathbf{x}) \\ \mathbf{h}(\mathbf{x}) = \mathbf{0} \\ \mathbf{g}(\mathbf{x}) + \mathbf{z} = \mathbf{0}, \end{cases} \quad (22)$$

with $\mathbf{z}^2 = [\mathbf{z}_1^2, \dots, \mathbf{z}_p^2]^T$. Then, the augmented Lagrangian writes

$$L_c(\mathbf{x}, \mathbf{z}, \lambda, \mu) = f(\mathbf{x}) + \lambda^T \mathbf{h}(\mathbf{x}) + \mu^T (\mathbf{g}(\mathbf{x}) + \mathbf{z}^2) + \frac{c}{2} (\|\mathbf{h}(\mathbf{x})\|^2 + \|\mathbf{g}(\mathbf{x}) + \mathbf{z}^2\|^2). \quad (23)$$

Then, letting $\mathbf{z}_i^2 = \mathbf{s}_i$, the dual writes

$$\phi(\lambda, \mu) = \min_{\mathbf{x}, \mathbf{s} \geq \mathbf{0}} f(\mathbf{x}) + \lambda^T \mathbf{h}(\mathbf{x}) + \mu^T (\mathbf{g}(\mathbf{x}) + \mathbf{s}) + \frac{c}{2} (\|\mathbf{h}(\mathbf{x})\|^2 + \|\mathbf{g}(\mathbf{x}) + \mathbf{s}\|^2). \quad (24)$$

Note that L_c depends on \mathbf{s}_i only via the term

$$\xi(\mathbf{s}_i) = (\mu_i + c\mathbf{g}_i(\mathbf{x}))\mathbf{s}_i + (c/2)\mathbf{s}_i^2.$$

and the minimum of $\xi(\mathbf{s}_i)$ is achieved for $\mathbf{s}_i = -\mathbf{g}_i(\mathbf{x}) - \mu_i/c$. Then, accounting for the positivity of \mathbf{s}_i , we get

$$\mathbf{s}_i = \max[0, -\mathbf{g}_i(\mathbf{x}) - \mu_i/c].$$

So, if $\mathbf{s}_i = 0$, $\mu_i(\mathbf{g}_i(\mathbf{x}) + \mathbf{s}_i) + \frac{c}{2}(\mathbf{g}_i(\mathbf{x}) + \mathbf{s}_i)^2$ is equal to

$$\mu_i \mathbf{g}_i(\mathbf{x}) + \frac{c}{2} \mathbf{g}_i(\mathbf{x})^2 = \frac{1}{2c} [(c\mathbf{g}_i(\mathbf{x}) + \mu_i)^2 - \mu_i^2],$$

and if $\mathbf{s}_i > 0$, it is equal to $\frac{-\mu_i^2}{2c}$. Then, letting

$$\psi_c(v, \mu) = \frac{1}{2c} [\max(0, \mu + cv)^2 - \mu^2], \quad (25)$$

we have

$$\mu_i(\mathbf{g}_i(\mathbf{x}) + \mathbf{s}_i) + \frac{c}{2}(\mathbf{g}_i(\mathbf{x}) + \mathbf{s}_i)^2 = \psi_c(\mathbf{g}_i(\mathbf{x})\mu_i)$$

and after minimization w.r.t. \mathbf{z} , L_c finally writes

$$L_c(\mathbf{x}, \lambda, \mu) = f(\mathbf{x}) + \lambda^T \mathbf{h}(\mathbf{x}) + \frac{c}{2} \|\mathbf{h}(\mathbf{x})\|^2 + \sum_i \psi_c(\mathbf{g}_i(\mathbf{x}), \mu_i). \quad (26)$$

Considering again the definition (23) of L_c , we get

$$\nabla_{\mathbf{x}} L_c(\mathbf{x}, \mathbf{z}, \lambda, \mu) = \nabla_{\mathbf{x}} f(\mathbf{x}) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x})\lambda + \mathbf{g}(\mathbf{x})\mu + c([\nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x})]\mathbf{h}(\mathbf{x}) + [\nabla_{\mathbf{x}} \mathbf{g}(\mathbf{x})](\mathbf{g}(\mathbf{x}) + \mathbf{z}^2)). \quad (27)$$

Then, solutions of necessary optimality conditions for problem (22) must satisfy $\mathbf{h}(\mathbf{x}) = \mathbf{0}$, $\mathbf{g}(\mathbf{x}) + \mathbf{z}^2 = \mathbf{0}$, and

$$\nabla_{\mathbf{x}} f(\mathbf{x}) + \nabla_{\mathbf{x}} \mathbf{h}(\mathbf{x})\lambda + \mathbf{g}(\mathbf{x})\mu = \mathbf{0}.$$

For such a solution, $(\mathbf{x}, \lambda, \mu)$ is also solution of problem (1) if we add the constraint $\mu \in \mathbb{R}_+^p$. This constraint is also necessary to be able to apply second order sufficient conditions in the presence of inequality constraints and extend theorem (8) to this situation.

To express the gradient ascent iterations for μ , let us remark that

$$\begin{aligned} \frac{\partial L_c}{\partial \mu_i} &= \frac{\partial}{\partial \mu_i} \psi_c(\mathbf{g}_i(\mathbf{x}), \mu_i) \\ &= \frac{1}{2c} [\mathbb{I}_{\{\mu_i + c\mathbf{g}_i \geq 0\}} \frac{\partial}{\partial \mu_i} ((\mu_i + c\mathbf{g}_i(\mathbf{x}))^2 - \mu_i^2) \\ &\quad + \mathbb{I}_{\{\mu_i + c\mathbf{g}_i < 0\}} \frac{\partial}{\partial \mu_i} (-\mu_i^2)] \\ &= \mathbb{I}_{\{\mu_i + c\mathbf{g}_i \geq 0\}} \mathbf{g}_i(\mathbf{x}) - \mathbb{I}_{\{\mu_i + c\mathbf{g}_i < 0\}} (\mu_i/c) \end{aligned} \quad (28)$$

Then, from the discussion about the behavior of the Hessian of the dual function as c grows, it is clear that at a minimum of (1), for active constraints the curvature of the dual function is close to $-c$. Then the following Newton step $\mu \leftarrow \mu + c\nabla_{\mu} L_c$ writes

$$\mu \leftarrow \mu + c[\mathbb{I}_{\{\mu_i + c\mathbf{g}_i \geq 0\}} \mathbf{g}_i(\mathbf{x}) - \mathbb{I}_{\{\mu_i + c\mathbf{g}_i < 0\}} (\mu_i/c)],$$

that is, $\mu \leftarrow [\mu + c\mathbf{g}]_+$, a form that clearly accounts for positivity constraints $\mu \geq \mathbf{0}$. In particular, it is clear that in the neighborhood of a solution, μ_i remains equal to 0 for an inactive constraint $\mathbf{g}_i(\mathbf{x}) < 0$.

Finally, the corresponding complete augmented Lagrangian method is given in algorithm 2

Algorithm 2 Augmented Lagrangian algorithm

1: **define**

$$\begin{aligned} \psi_c(v, \mu) &= \frac{1}{2c} [\max(0, \mu + cv)^2 - \mu^2], \\ L_c(\mathbf{x}, \lambda, \mu) &= f(\mathbf{x}) + \lambda^T \mathbf{h}(\mathbf{x}) + \frac{c}{2} \|\mathbf{h}(\mathbf{x})\|^2 + \sum_i \psi_c(\mathbf{g}_i(\mathbf{x}), \mu_i). \end{aligned}$$

2: **init** $\mathbf{x}_0, \lambda_0, \mu_0, c > 0$

3: **while** $c < c_{\max}$ **do**

4: **while** stopping condition \neq true **do**

5: $\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} L_c(\mathbf{x}, \lambda_k, \mu_k)$

or

$$\begin{cases} \alpha_k & \approx \arg \min_r L_c(\mathbf{x}_k - r\nabla_{\mathbf{x}} L(\mathbf{x}_k, \mu_k)) \\ \mathbf{x}_{k+1} & = \mathbf{x}_k - \alpha_k \nabla_{\mathbf{x}} L(\mathbf{x}_k, \mu_k) \end{cases}$$

6: $\lambda_{k+1} = \lambda_k + \mathbf{c}\mathbf{h}(\mathbf{x}_{k+1})$

7: $\mu_{k+1} = [\mu_k + c\mathbf{g}(\mathbf{x}_{k+1})]_+$

8: **end while**

9: $c = \rho c$

10: **end while**

5 Example: sphere packing

Let us consider the disk packing problem that consists in packing n disks of radius $r = 1$ within a minimum area. Here we assume in addition some shape constraint for the

enclosing surface. For instance, assuming a circular boundary with radius z for the enclosing surface, we want to solve

$$(P) \begin{cases} \min_{\mathbf{x}, \mathbf{y}, z, \mathbf{p}} z \\ (\mathbf{x}_i - \mathbf{x}_j)^2 + (\mathbf{y}_i - \mathbf{y}_j)^2 > (2r)^2 & 1 \leq i < j \leq n \\ (\mathbf{x}_i - \mathbf{p}_x)^2 + (\mathbf{y}_i - \mathbf{p}_y)^2 \leq (z - r)^2 & 1 \leq i \leq n \end{cases}$$

The objective z is the radius that must be minimized, \mathbf{x} and \mathbf{y} are the vectors of x and y coordinates for the centers of the disks, The first set of constraints, $(\mathbf{x}_i - \mathbf{x}_j)^2 + (\mathbf{y}_i - \mathbf{y}_j)^2 > (2r)^2$, are the non overlap constraints for disks. Constraints $(\mathbf{x}_i - \mathbf{p}_x)^2 + (\mathbf{y}_i - \mathbf{p}_y)^2 \leq (z - r)^2$ represent circular boundary constraints, where \mathbf{p} denotes the center of the enclosing circle. Note that the problem is not convex and using convex optimization approaches can lead to local minima. Letting the center \mathbf{p} vary instead of setting for instance $\mathbf{p} = \mathbf{0}$ supplies additional degrees of freedom to the problem that can be helpful to avoid some local minima.

Figure 1 shows examples of this approach for different boundary conditions using augmented Lagrangian.

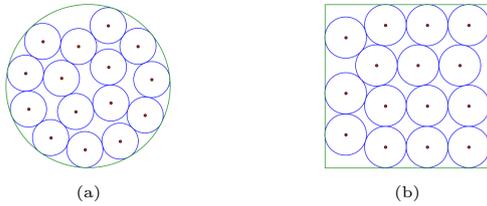


Figure 1: Packing 15 spheres (a) in a disk and (b) in a box.

6 Exercise: SVM and duality

In this exercise, we recall first some elements of SVM classification. For additional information about SVM principles see for instance the corresponding Wikipedia page.

SVM techniques rely on two main ideas: classify data by (i) searching for a separation boundary between two classes such that labelled data [1] exhibit largest minimum distance to the boundary and (ii) transform data into a space of larger, possibly infinite, dimension where linear separation of classes works better (but then in the original classes separation becomes nonlinear).

When possible, linear separation between 2 classes, say $c \in \{-1, 1\}$, can simply be achieved by looking for an hyperplane with parameters (\mathbf{w}, b) such that in the learning set $\{(\mathbf{x}_k, c_k)\}_{k=1:K}$ elements satisfy $c_k(\mathbf{w}^T \mathbf{x}_k - b) \geq 0$. We can also search for two hyperplanes such that $\mathbf{w}^T \mathbf{x} - b \leq -1$ for class -1 and $\mathbf{w}^T \mathbf{x} - b \geq 1$ for class 1. We would like that they represent limiting hyperplanes for both classes with maximum separation. Then, $c_k(\mathbf{w}^T \mathbf{x}_k - b) \geq 1$ with equality for at least one point in each class.

1) - Assuming the hyperplanes $\mathbf{w}^T \mathbf{x}_k - b \leq -1$ and $\mathbf{w}^T \mathbf{x}_k - b \geq 1$ are limiting hyperplanes for classes $c = -1$ and $c = 1$ respectively, prove that the maximum distance between these class limiting hyperplanes is $2 \|\mathbf{w}\|^{-1}$, that is,

$$\max_{c_k=1, c_l=-1} \|\mathbf{x}_k - \mathbf{x}_l\|^2 = \frac{2}{\|\mathbf{w}\|^2}.$$

Then, check that we can get (\mathbf{w}, b) by solving

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \\ c_k(\mathbf{w}^T \mathbf{x}_k - b) \geq 1, & k = 1 : K. \end{cases}$$

This way of determining \mathbf{w} enables maximum margin between class limiting hyperplanes $\mathbf{w}^T \mathbf{x}_k - b \leq -1$ and

$$\mathbf{w}^T \mathbf{x}_k - b \geq 1.$$

2) - When both classes are not exactly linearly separable, a possible linear approximation is obtained by weighting the objective $\|\mathbf{w}\|^2$ by the sum of distances to the class boundary for wrongly classified data. Explain how this can be expressed by the following problem by discussing the meaning of variables ζ_k and parameter λ :

$$\begin{cases} \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + \lambda \sum_{k=1:K} \zeta_k \\ c_k(\mathbf{w}^T \mathbf{x}_k - b) \geq 1 - \zeta_k, & \zeta_k \geq 0, \quad k = 1 : K. \end{cases}$$

3) Write the Lagrangian of the problem and express KKT conditions.

4) From KKT conditions, show that the dual of the problem writes

$$\begin{cases} \max_{\alpha} f(\alpha) = \sum_{k=1:K} \alpha_k - \frac{1}{2} \sum_{i,j=1:K} \alpha_i \alpha_j c_i c_j (\mathbf{x}_j^T \mathbf{x}_i) \\ \sum_{k=1:K} \alpha_k c_k = 0 \\ 0 \leq \alpha_k \leq \lambda, & k = 1 : K. \end{cases}$$

What do variables $(\alpha_k)_{k=1:K}$ represent ?

Note that, just as for the primal problem, this is a standard quadratic optimization problem. However, it appears that the dual formulation presents the benefit of decoupled variables in the formulation of inequality constraints, what can help algorithms to enforce these constraints. Indeed, for instance we can combine a penalty approach to handle the equality constraint and projection on the cube satisfying inequality constraints. Then, the gradient is calculated on the penalized convex criterion $f_t(\alpha) = -f(\alpha) + t(\alpha^T \mathbf{c})^2$ and the updated α is projected:

$$\alpha^{(n+1)} = \Pi_{[0, \lambda]^K} (\alpha^{(n)} - \mu_n \nabla f_t(\alpha^{(n)})),$$

where $\Pi_{[0, \lambda]^K}$ represents the projection onto the cube $[0, \lambda]^K$ of \mathbb{R}^K . t is slightly increased at each iteration ($t \rightarrow \beta t$, with e.g. $\beta = 1.1$).

5) Check that the solution of the primal for \mathbf{w} is then given by $\mathbf{w} = \sum_k \alpha_k c_k \mathbf{x}_k$. To find b , check first that if $0 < \alpha_k < \lambda$ then \mathbf{x}_k lies on the margin. Then, $c_k(\mathbf{w}^T \mathbf{x}_k - b) = 1$, that is, $b = \mathbf{w}^T \mathbf{x}_k - c_k$.

Finally, the decision variable writes $\text{sign}((\mathbf{w}^T \mathbf{x} - b)) = \text{sign}(\sum_k \alpha_k c_k \mathbf{x}_k^T \mathbf{x} - b)$.

In order to improve classes separation, we now consider possibly nonlinear separation. A nice way to achieve nonlinear separation of classes is to apply a nonlinear transform of \mathbf{x}_k into a space of larger (possibly infinite) dimension, say $\mathbf{x}_k \rightarrow \varphi(\mathbf{x}_k)$ and consider the linear separation of transformed variables $\varphi(\mathbf{x}_k)$. Note that the decision variable is now in the form $\text{sign}(\sum_k \alpha_k c_k \varphi(\mathbf{x}_k)^T \varphi(\mathbf{x}) - b) = \text{sign}(\sum_k \alpha_k c_k k(\mathbf{x}_k, \mathbf{x}) - b)$ and φ needs not to be known: only the bilinear function $k(\cdot, \cdot)$ must be set. In fact $k(\mathbf{x}, \mathbf{x}')$ represents a scalar product $\varphi(\mathbf{x}')^T \varphi(\mathbf{x})$ if and only if it is positive definite. This result is known as Mercer's theorem. The Gaussian kernel $k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / h)$ ($h > 0$) is a classical example of such a similarity function. The dual of the optimisation problem for data $(\varphi(\mathbf{x}_k), c_k)_k$ directly formulates in terms of $k(\cdot, \cdot)$ and does not require φ .

With all the calculations above you should be able to implement the dual approach for SVM proposed for the laboratory session.

References

- [1] S. Boyd, L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
<https://web.stanford.edu/boyd/cvxbook/>
- [2] D.G. Luenberger, Y. Ye, *Linear and Nonlinear Programming*, 3rd Ed., Springer, 2008.
- [3] P.G. Ciarlet, *Introduction à l'Analyse Numérique Matricielle et à l'Optimisation*, Masson, 1982.
- [4] Wikipedia: Max–min inequality
https://en.wikipedia.org/wiki/Max-min_inequality
- [5] T. Chonavel, *Notes on unconstrained and constrained optimization algorithms*, IMT Atlantique, TAF MCE, UE NUMMET, 2021
<https://moodle.imt-atlantique.fr/mod/resource/view.php?id=43984>
- [6] Wikipedia: *Strong convexity*
https://en.wikipedia.org/wiki/Convex_function#Strongly_convex_functions
- [7] Wikipedia: *Operator norm*
https://en.wikipedia.org/wiki/Operator_norm
- [8] Stackexchange: *projection on closed convex sets*
<https://math.stackexchange.com/questions/2900841/orthogonal-projection-of-a-vector-onto-convex-set>
- [9] Stackexchange: *projection on closed convex sets are contractive*
<https://math.stackexchange.com/questions/3809431/geometric-proof-that-projections-on-convex-sets-are-contractive> <https://jump.dev/JuMP.jl/stable>
- [10] Wikipedia: *Woodbury matrix identity*
https://en.wikipedia.org/wiki/Woodbury_matrix_identity
- [11] SVM, https://en.wikipedia.org/wiki/Support_vector_machine