

# III. Constrained Optimization

So far we have focused exclusively on *unconstrained* optimization problems. In such a setting, our goal is typically clear: find a point where the gradient (or subgradient) is equal to zero. All of the algorithms we have explored so far were different strategies for finding such a point. Once we add constraints, however, things get a bit more complicated. In particular, there may no longer be any points that satisfy the constraints we are imposing where the gradient vanishes. Showing that we have found an optimal point will now involve a more complicated relationship between the gradient of the function we are minimizing together with the constraints.

In these notes we will look at a specific class of constrained optimization problems of the form

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) \quad g_m(\mathbf{x}) \leq 0, \quad m = 1, \dots, M. \quad (1)$$

Here, we represent the constraints as functions  $g_1, \dots, g_M$ , which by convention we define so that we are always imposing  $g_m(\mathbf{x}) \leq 0$ . Note that if we had a constraint of the form  $h(\mathbf{x}) = 0$  we could write this as  $h(\mathbf{x}) \leq 0$  combined with  $-h(\mathbf{x}) \leq 0$ , so equality constraints can also be handled, although we will encounter these less frequently in this course.<sup>1</sup>

While some of what we will say actually applies to the case where the  $g_m$  are nonconvex, we will mostly only be interested in the case where the  $g_m$  are convex functions.

---

<sup>1</sup>In practice, you would want to handle equality constraints more explicitly, but focusing only on inequality constraints will make the exposition quite a bit cleaner without really sacrificing any intuition.

An important consideration in constrained optimization problems is the concept of *feasibility*. A vector  $\mathbf{x}$  is **feasible** if it satisfies the constraints of (1). Specifically, a feasible  $\mathbf{x}$  must satisfy  $g_m(\mathbf{x}) \leq 0$  for  $m = 1, \dots, M$ . It is not a given that any feasible  $\mathbf{x}$  exists. For instance, I might demand that  $\sum_{n=1}^N x_n > 1$  and also  $\sum_{n=1}^N x_n < -1$ . Clearly, no  $\mathbf{x}$  that simultaneously satisfies both of these constraints can exist. In our discussion below, we will assume that the **feasible set**

$$\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^N : g_m(\mathbf{x}) \leq 0 \ m = 1, \dots, M, \}$$

is non-empty.

## The Lagrangian

The **Lagrangian** takes the constraints in the program above and integrates them into the objective function. Specifically, the Lagrangian  $\mathcal{L} : \mathbb{R}^N \times \mathbb{R}^M \rightarrow \mathbb{R}$  associated with this optimization program is

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) := f(\mathbf{x}) + \sum_{m=1}^M \lambda_m g_m(\mathbf{x}).$$

For reasons that will become clearer below, the  $\mathbf{x}$  above are referred to as **primal variables**, and the  $\boldsymbol{\lambda}$  as either **dual variables** or **Lagrange multipliers**.

The Lagrangian allows us to transform the *constrained* optimization problem in (1) into an *unconstrained* one. Specifically, consider the problem given by

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) + \sum_{m=1}^M \lambda_m g_m(\mathbf{x}). \quad (2)$$

To get some intuition, suppose that we set the  $\lambda_1, \dots, \lambda_M$  to be very large (positive) numbers. In this case, violating any of the constraints (allowing  $g_m(\mathbf{x}) > 0$ ) will result in a very large penalty being added to the objective function, so that by setting the corresponding  $\lambda_m$  to be large we will eventually guarantee that the resulting solution will satisfy the desired constraints. The problem here is that large values of  $\lambda_m$  not only avoid the setting where  $g_m(\mathbf{x}) > 0$ , but actually encourages  $g_m(\mathbf{x}) \ll 0$  (since we can potentially benefit by not just satisfying the constraints but by exceeding them by a large margin).

This raises a natural question: can we set  $\boldsymbol{\lambda}$  so that the solution to the unconstrained problem (2) is the same as the constrained problem (1)? Here we will provide an answer in the case where the objective function  $f$  and the constraints  $g_1, \dots, g_M$  are both convex and differentiable.

Suppose that  $\mathbf{x}^*$  is a solution to the constrained problem (1). If we want  $\mathbf{x}^*$  to be a solution to (2), then we need

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}) = \nabla_{\mathbf{x}} f(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m \nabla_{\mathbf{x}} g_m(\mathbf{x}^*) = \mathbf{0}. \quad (3)$$

If we knew  $\mathbf{x}^*$  already, finding a  $\boldsymbol{\lambda}$  that would make the unconstrained and constrained problems equivalent (meaning that they both have the same solution  $\mathbf{x}^*$ ) would just amount to finding a  $\boldsymbol{\lambda}$  such that (3) holds. Unfortunately, this is not particularly useful since  $\mathbf{x}^*$  is what we are trying to find to begin with.

To see how we might compute a  $\boldsymbol{\lambda}$  that makes the unconstrained and constrained problems equivalent, we will need to take a brief detour to discuss the notion of **duality**.

# Lagrangian duality

## The Lagrange dual function

We can think of the unconstrained optimization problem (2) as actually representing a family of different optimization problems (depending on  $\boldsymbol{\lambda}$ ). For any fixed  $\boldsymbol{\lambda}$ , imagine solving (2) and computing the minimal value of the objective function – we can think of this as actually defining a function that maps  $\boldsymbol{\lambda} \in \mathbb{R}^M$  to  $\mathbb{R}$ . We call this the **Lagrange dual function**  $d(\boldsymbol{\lambda})$ , which is defined as

$$d(\boldsymbol{\lambda}) = \min_{\boldsymbol{x} \in \mathbb{R}^N} \left( f(\boldsymbol{x}) + \sum_{m=1}^M \lambda_m g_m(\boldsymbol{x}) \right).$$

Note that, using a result from a previous homework, since the dual is the minimum of a family of affine functions in  $\boldsymbol{\lambda}$ , the Lagrange dual function **always concave**.<sup>2</sup>

A key fact about the dual function is that it can provide a lower bound on the optimal value of the original program. In the discussion below, we assume throughout that  $\boldsymbol{\lambda} \geq 0$ , meaning that  $\lambda_m \geq 0$  for all  $m$ , but can otherwise be arbitrary. The main claim is that if  $f^*$  is the optimal value for (1), then we have

$$d(\boldsymbol{\lambda}) \leq f^*.$$

This is very easy to show. Specifically, for any feasible point  $\boldsymbol{x}'$ , we

---

<sup>2</sup>While it is not really significant in the context of this class, since we are focusing on convex optimization problems, a remarkable fact is that the dual is concave regardless of whether or not the  $g_m$  are convex. This can be very useful when dealing with nonconvex problems, but we will not explore this here.

must have  $g_m(\mathbf{x}') \leq 0$  for all  $m$  and hence

$$\sum_{m=1}^M \lambda_m g_m(\mathbf{x}') \leq 0.$$

From this we have that

$$\mathcal{L}(\mathbf{x}', \boldsymbol{\lambda}) \leq f(\mathbf{x}'),$$

meaning

$$d(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \leq \mathcal{L}(\mathbf{x}', \boldsymbol{\lambda}) \leq f(\mathbf{x}').$$

Since this holds for all feasible  $\mathbf{x}'$ , including the minimizer of (1), we have  $d(\boldsymbol{\lambda}) \leq f^*$ .

## The (Lagrange) dual problem

Given that  $d(\boldsymbol{\lambda})$  provides a lower bound on  $f^*$ , if you wanted to get an idea of what  $f^*$  looks like (for example, to see if you are close to convergence), it is natural to see how large you can make this lower bound. This gives rise to what we call the **dual problem** of (1):

$$\text{maximize}_{\boldsymbol{\lambda} \in \mathbb{R}^M} d(\boldsymbol{\lambda}) \quad \text{subject to} \quad \boldsymbol{\lambda} \geq \mathbf{0}. \quad (4)$$

The dual optimal value  $d^*$  is

$$d^* = \max_{\boldsymbol{\lambda} \geq \mathbf{0}} d(\boldsymbol{\lambda}) = \max_{\boldsymbol{\lambda} \geq \mathbf{0}} \min_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}).$$

Since  $d(\boldsymbol{\lambda}) \leq f^*$ , we know that

$$d^* \leq f^*.$$

The quantity  $f^* - d^*$  is called the **duality gap**. If  $f^* = d^*$ , then we say that (1) and (4) exhibit **strong duality**.

We will soon discuss when strong duality holds, but first, why is it important? Suppose that  $\mathbf{x}^*$  is a solution to the original constrained problem (1) – which we will call the **primal problem** to distinguish it from the dual problem – and suppose that  $\boldsymbol{\lambda}^*$  is a solution to the dual problem (4). It turns out that if we have strong duality, then  $\boldsymbol{\lambda}^*$  is exactly what we need to make  $\mathbf{x}^*$  the solution to the unconstrained problem (2).

To see why, note that

$$\begin{aligned}
 f(\mathbf{x}^*) &= d(\boldsymbol{\lambda}^*) \\
 &= \min_{\mathbf{x} \in \mathbb{R}^N} \left( f(\mathbf{x}) + \sum_{m=1}^M \lambda_m^* g_m(\mathbf{x}) \right) \\
 &\leq f(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m^* g_m(\mathbf{x}^*) \\
 &\leq f(\mathbf{x}^*).
 \end{aligned} \tag{5}$$

where the last inequality follows from the fact that we must have  $\lambda_m^* \geq 0$  and  $g_m(\mathbf{x}^*) \leq 0$ . Looking at this entire chain of inequalities, where the first and last term are both  $f(\mathbf{x}^*)$ , means that

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*) = \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*).$$

In words, the solution to the primal problem  $\mathbf{x}^*$  is also a minimizer of  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$ .

## Example

Consider the optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \langle \mathbf{x}, \mathbf{c} \rangle \quad \text{subject to} \quad \mathbf{Ax} \leq \mathbf{b}.$$

This is called a **linear program** (since the objective function is just a linear function of  $\mathbf{x}$ ). The Lagrangian is

$$\begin{aligned} L(\mathbf{x}, \boldsymbol{\lambda}) &= \langle \mathbf{x}, \mathbf{c} \rangle + \sum_{m=1}^M \lambda_m (\langle \mathbf{x}, \mathbf{a}_m \rangle - b_m) \\ &= \mathbf{c}^T \mathbf{x} - \boldsymbol{\lambda}^T \mathbf{b} + \boldsymbol{\lambda}^T \mathbf{Ax}. \end{aligned}$$

This is a linear function of  $\mathbf{x}$ . Note that it is unbounded below unless

$$\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0}.$$

Thus

$$\begin{aligned} g(\boldsymbol{\lambda}) &= \min_{\mathbf{x}} (\mathbf{c}^T \mathbf{x} - \boldsymbol{\lambda}^T \mathbf{b} + \boldsymbol{\lambda}^T \mathbf{Ax}) \\ &= \begin{cases} -\langle \boldsymbol{\lambda}, \mathbf{b} \rangle, & \mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} = \mathbf{0} \\ -\infty, & \text{otherwise.} \end{cases} \end{aligned}$$

So the Lagrange dual program is

$$\begin{aligned} \underset{\boldsymbol{\lambda} \in \mathbb{R}^M}{\text{maximize}} \quad & -\langle \boldsymbol{\lambda}, \mathbf{b} \rangle \quad \text{subject to} \quad \mathbf{A}^T \boldsymbol{\lambda} = -\mathbf{c} \\ & \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

Note that the dual is another linear program.



## Strong duality and the KKT conditions

So when does strong duality hold? The answer can be pretty complicated and depends on the structure of the constraints. There are a variety of so-called “constraint qualifications” that serve as sufficient conditions to guarantee strong duality.

Probably the simplest and most widely applicable is known as **Slater’s condition**, which is essentially that the  $g_m$  are affine inequality constraints (i.e., they can be expressed as  $\mathbf{Ax} \leq \mathbf{b}$ ), and that there is an  $\mathbf{x}$  that is strictly feasible for the remaining constraints (i.e., an  $\mathbf{x}$  such that for all the  $g_m$  which are not affine we have  $g_m(\mathbf{x}) < 0$ ). Actually proving that this condition implies strong duality is somewhat involved. We will not worry too much about this, in all of the problems that we will encounter in this course, strong duality will hold.

Above we argued that when strong duality holds, if  $\mathbf{x}^*$  is a solution to the primal problem and  $\boldsymbol{\lambda}^*$  is a solution to the dual problem, then  $\mathbf{x}^*$  is also a minimizer of  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$ , or equivalently, that the condition (3) holds for  $\boldsymbol{\lambda} = \boldsymbol{\lambda}^*$ , i.e.,

$$\nabla_{\mathbf{x}} f(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m^* \nabla_{\mathbf{x}} g_m(\mathbf{x}^*) = \mathbf{0}.$$

We can also say something else about  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  from the analysis on the previous page. Specifically, the final inequality in (5) implies that

$$\sum_{m=1}^M \lambda_m^* g_m(\mathbf{x}^*) = 0,$$

but since each term in this sum has to be less than or equal to zero,

this actually implies that

$$\lambda_m^* g_m(\mathbf{x}^*) = 0, \quad m = 1, \dots, M.$$

If we combine these facts with the fact that if  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  must be feasible in order to be solutions to the primal/dual problems, we arrive at a set of conditions that solutions  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  to the primal and dual problems must satisfy. These are known as the Karush-Kuhn-Tucker (KKT) conditions.

### **KKT**

The KKT conditions for for an  $\mathbf{x} \in \mathbb{R}^N$  and  $\boldsymbol{\lambda} \in \mathbb{R}^M$  are

$$g_m(\mathbf{x}) \leq 0, \quad m = 1, \dots, M, \quad (\text{K1})$$

$$\lambda_m \geq 0, \quad m = 1, \dots, M, \quad (\text{K2})$$

$$\lambda_m g_m(\mathbf{x}) = 0, \quad m = 1, \dots, M, \quad (\text{K3})$$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) + \sum_{m=1}^M \lambda_m \nabla_{\mathbf{x}} g_m(\mathbf{x}) = \mathbf{0}. \quad (\text{K4})$$

We have already shown (if strong duality holds) that if  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  are primal/dual optimal, then  $\mathbf{x}^*, \boldsymbol{\lambda}^*$  must satisfy the KKT conditions. It is also the case that if you can find  $\mathbf{x}^*, \boldsymbol{\lambda}^*$  that obey the KKT conditions, then you necessarily have strong duality and  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  are primal/dual optimal. This is easy to show. Note that if KKT4 holds,

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0},$$

meaning that  $\mathbf{x}^*$  is a minimizer of  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$ , i.e.,

$$\mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) \leq \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*),$$

thus

$$\begin{aligned}d(\boldsymbol{\lambda}^*) &= \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \\ &= f(\boldsymbol{x}^*) + \sum_{m=1}^M \lambda_m^* g_m(\boldsymbol{x}^*) \\ &= f(\boldsymbol{x}^*), \quad (\text{by KKT3}),\end{aligned}$$

and we have strong duality. Since the dual  $d(\boldsymbol{\lambda})$  always provides a lower bound to the primal  $f(\boldsymbol{x})$  (for any feasible  $\boldsymbol{x}$  and  $\boldsymbol{\lambda}$ ), if  $\boldsymbol{x}^*$  and  $\boldsymbol{\lambda}^*$  satisfy  $d(\boldsymbol{\lambda}^*) = f(\boldsymbol{x}^*)$ , we know that  $\boldsymbol{\lambda}^*$  and  $\boldsymbol{x}^*$  are optimal since we clearly cannot further decrease  $f$  to be smaller than  $d(\boldsymbol{\lambda}^*)$  or increase  $d$  to be larger than  $f(\boldsymbol{x}^*)$ .

The KKT conditions are a very useful tool in optimization. As we will soon see, one algorithmic approach to constrained optimization is to simply find  $\boldsymbol{x}^*$  and  $\boldsymbol{\lambda}^*$  satisfying these conditions using an iterative method. The KKT conditions also allow us to easily “translate” a solution to the primal problem into a solution to the dual, or a solution to the dual problem into a solution to the primal. This can be useful because sometimes, for example, the dual might be much easier to solve than the primal.