

## Lagrange duality

Another way to arrive at the KKT conditions, and one which gives us some insight on solving constrained optimization problems, is through the *Lagrange dual*. The dual is a maximization program in  $\boldsymbol{\lambda}, \boldsymbol{\nu}$  — it is always concave (even when the original program is not convex), and gives us a systematic way to lower bound the optimal value.

### The Lagrangian

We consider an optimization program of the form

$$\begin{aligned} \underset{\boldsymbol{x} \in \mathbb{R}^N}{\text{minimize}} \quad & f_0(\boldsymbol{x}) \quad f_m(\boldsymbol{x}) \leq 0, \quad m = 1, \dots, M \\ & h_p(\boldsymbol{x}) = 0, \quad p = 1, \dots, P. \end{aligned} \tag{1}$$

Much of what we will say below applies equally well to nonconvex programs as well as convex programs, so we will make it clear when we are taking the  $f_m$  to be convex and the  $h_p$  to be affine. We will take the domain of all of the  $f_m$  and  $h_p$  to be all of  $\mathbb{R}^N$  below; this just simplifies the exposition, we can easily replace this with the intersections of the  $\text{dom } f_m$  and  $\text{dom } h_p$ . We will assume that the feasible set

$$\mathcal{C} = \{\boldsymbol{x} : f_m(\boldsymbol{x}) \leq 0, h_p(\boldsymbol{x}) = 0, m = 1, \dots, M, p = 1, \dots, P\}$$

is non-empty and a subset  $\mathbb{R}^N$ .

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

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\mathbf{x}) + \sum_{m=1}^M \lambda_m f_m(\mathbf{x}) + \sum_{p=1}^P \nu_p h_p(\mathbf{x})$$

The  $\mathbf{x}$  above are referred to as **primal variables**, and the  $\boldsymbol{\lambda}, \boldsymbol{\nu}$  as either **dual variables** or **Lagrange multipliers**.

The **Lagrange dual function**  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) : \mathbb{R}^M \times \mathbb{R}^P \rightarrow \mathbb{R}$  is the minimum of the Lagrangian over all values of  $\mathbf{x}$ :

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x} \in \mathbb{R}^N} \left( f_0(\mathbf{x}) + \sum_{m=1}^M \lambda_m f_m(\mathbf{x}) + \sum_{p=1}^P \nu_p h_p(\mathbf{x}) \right).$$

Since the dual is a pointwise infimum of a family of affine functions in  $\boldsymbol{\lambda}, \boldsymbol{\nu}$ ,  $g$  is **concave** regardless of whether or not the  $f_m, h_p$  are convex.

The key fact about the dual function is that it is everywhere a lower bound on the optimal value of the original program. If  $p^\star$  is the optimal value for (1), then

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq p^\star, \quad \text{for all } \boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\nu} \in \mathbb{R}^P.$$

This is (almost too) easy to see. For any feasible point  $\mathbf{x}_0$ ,

$$\sum_{m=1}^M \lambda_m f_m(\mathbf{x}_0) + \sum_{p=1}^P \nu_p h_p(\mathbf{x}_0) \leq 0,$$

and so

$$L(\mathbf{x}_0, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f_0(\mathbf{x}_0), \quad \text{for all } \boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\nu} \in \mathbb{R}^P,$$

meaning

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{x} \in \mathbb{R}^N} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq L(\mathbf{x}_0, \boldsymbol{\lambda}, \boldsymbol{\nu}) \leq f_0(\mathbf{x}_0).$$

Since this holds for all feasible  $\mathbf{x}_0$ ,  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq \inf_{\mathbf{x} \in \mathcal{C}} f_0(\mathbf{x}) = p^*$ .

The (Lagrange) dual to the optimization program (1) is

$$\underset{\boldsymbol{\lambda} \in \mathbb{R}^M, \boldsymbol{\nu} \in \mathbb{R}^P}{\text{maximize}} \quad g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \quad \text{subject to} \quad \boldsymbol{\lambda} \geq \mathbf{0}. \quad (2)$$

The dual optimal value  $d^*$  is

$$d^* = \sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\nu}} g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \sup_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\nu}} \inf_{\mathbf{x} \in \mathbb{R}^N} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}).$$

Since  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq p^*$ , we know that

$$d^* \leq p^*.$$

The quantity  $p^* - d^*$  is called the **duality gap**. If  $p^* = d^*$ , then we say that (1) and (2) exhibit **strong duality**.

## Certificates of (sub)optimality

Any dual feasible<sup>1</sup>  $(\boldsymbol{\lambda}, \boldsymbol{\nu})$  gives us a lower bound on  $p^*$ , since  $g(\boldsymbol{\lambda}, \boldsymbol{\nu}) \leq p^*$ . If we have a primal feasible  $\boldsymbol{x}$ , then we know that

$$f_0(\boldsymbol{x}) - p^* \leq f_0(\boldsymbol{x}) - g(\boldsymbol{\lambda}, \boldsymbol{\nu}).$$

We will refer to  $f_0(\boldsymbol{x}) - g(\boldsymbol{\lambda}, \boldsymbol{\nu})$  as the **duality gap** for the primal/dual (feasible) variables  $\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\nu}$ . We know that

$$p^* \in [g(\boldsymbol{\lambda}, \boldsymbol{\nu}), f_0(\boldsymbol{x})], \quad \text{and likewise} \quad d^* \in [g(\boldsymbol{\lambda}, \boldsymbol{\nu}), f_0(\boldsymbol{x})].$$

If we are ever able to reduce this gap to zero, then we know that  $\boldsymbol{x}$  is primal optimal, and  $\boldsymbol{\lambda}, \boldsymbol{\nu}$  are dual optimal.

There are certain kinds of “primal-dual” algorithms that produce a series of (feasible) points  $\boldsymbol{x}^{(k)}, \boldsymbol{\lambda}^{(k)}, \boldsymbol{\nu}^{(k)}$  at every iteration. We can then use

$$f_0(\boldsymbol{x}^{(k)}) - g(\boldsymbol{\lambda}^{(k)}, \boldsymbol{\nu}^{(k)}) \leq \epsilon,$$

as a stopping criteria, and know that our answer would yield an objective value no further than  $\epsilon$  from optimal.

## Strong duality and the KKT conditions

Suppose that for a convex program, the primal optimal value  $p^*$  and the dual optimal value  $d^*$  are equal

$$p^* = d^*.$$

---

<sup>1</sup>We simply need  $\boldsymbol{\lambda} \geq \mathbf{0}$  for  $(\boldsymbol{\lambda}, \boldsymbol{\nu})$  to be dual feasible.

If  $\mathbf{x}^*$  is a primal optimal point and  $\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$  is a dual optimal point, then we must have

$$\begin{aligned}
f_0(\mathbf{x}^*) &= g(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \\
&= \inf_{\mathbf{x} \in \mathbb{R}^N} \left( f_0(\mathbf{x}) + \sum_{m=1}^M \lambda_m^* f_m(\mathbf{x}) + \sum_{p=1}^P \nu_p^* h_p(\mathbf{x}) \right) \\
&\leq f_0(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m^* f_m(\mathbf{x}^*) + \sum_{p=1}^P \nu_p^* h_p(\mathbf{x}^*) \\
&\leq f_0(\mathbf{x}^*).
\end{aligned}$$

The last inequality follows from the fact that  $\lambda_m^* \geq 0$  (dual feasibility),  $f_m(\mathbf{x}^*) \leq 0$ , and  $h_p(\mathbf{x}^*) = 0$  (primal feasibility). Since we started out and ended up with the same thing, all of the things above must be equal, and so

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

Also, since we know  $\mathbf{x}^*$  is a minimizer of  $L(\mathbf{x}, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*)$  (second equality above), which is an unconstrained convex function (with  $\boldsymbol{\lambda}, \boldsymbol{\nu}$  fixed), the gradient with respect to  $\mathbf{x}$  must be zero:

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) = \nabla f_0(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m^* \nabla f_m(\mathbf{x}^*) + \sum_{p=1}^P \nu_p^* \nabla h_p(\mathbf{x}^*) = \mathbf{0}.$$

Thus strong duality immediately leads to the KKT conditions holding at the solution.

Also, if you can find  $\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$  that obey the KKT conditions, not only do you know that you have a primal optimal point on your hands, but also we have strong duality (and  $\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*$  are dual optimal). For if KKT holds,

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

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

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

thus

$$\begin{aligned} g(\boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) &= L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\nu}^*) \\ &= f_0(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m^* f_m(\mathbf{x}^*) + \sum_{p=1}^P \nu_p^* h_p(\mathbf{x}^*) \\ &= f_0(\mathbf{x}^*), \quad (\text{by KKT}), \end{aligned}$$

and we have strong duality.

The upshot of this is that the conditions for strong duality are essentially the same as those under which KKT is necessary.

The program (1) and its dual (2) have strong duality if the  $f_m$  are affine inequality constraints, or there is an  $\mathbf{x} \in \mathbb{R}^N$  such that for all the  $f_i$  which are not affine we have  $f_i(\mathbf{x}) < 0$ .

## Examples

1. **Inequality LP.** Calculate the dual of

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

**Answer:** 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 functional in  $\mathbf{x}$  — it is unbounded below unless

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

Thus

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

So the Lagrange dual program is

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

2. **Least-squares.** Calculate the dual of

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad \mathbf{Ax} = \mathbf{b}.$$

Check that the duality gap is zero.

**Answer:** The Lagrangian is

$$L(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{x}^T \mathbf{x} - \boldsymbol{\nu}^T \mathbf{b} + \boldsymbol{\nu}^T \mathbf{Ax}.$$

This is quadratic in  $\mathbf{x}$  and will attain its minimum for

$$\mathbf{x} = -\frac{1}{2} \mathbf{A}^T \boldsymbol{\nu}.$$

Thus

$$\begin{aligned} g(\boldsymbol{\nu}) &= \frac{1}{4} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} - \mathbf{b}^T \boldsymbol{\nu} - \frac{1}{2} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} \\ &= -\frac{1}{4} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} - \mathbf{b}^T \boldsymbol{\nu}, \end{aligned}$$

and the Lagrange dual problem is

$$\underset{\boldsymbol{\nu} \in \mathbb{R}^M}{\text{maximize}} \quad -\frac{1}{4} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} - \mathbf{b}^T \boldsymbol{\nu}.$$

Note that this will be maximized when  $-\frac{1}{2} \mathbf{AA}^T \boldsymbol{\nu} = \mathbf{b}$ , which, when substituted into the dual problem yields

$$-\frac{1}{4} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} + \frac{1}{2} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} = \frac{1}{4} \boldsymbol{\nu}^T \mathbf{AA}^T \boldsymbol{\nu} = \left\| -\frac{1}{2} \mathbf{A}^T \boldsymbol{\nu} \right\|_2^2,$$

which shows that strong duality holds.



3. **Minimum norm.** Calculate the dual of

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{x}\| \quad \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b},$$

where  $\|\cdot\|$  is a general valid norm.

**Answer:** Use  $f_0(\mathbf{x}) = \|\mathbf{x}\|$  to ease notation below. We start with the Lagrangian:

$$\begin{aligned} L(\mathbf{x}, \boldsymbol{\nu}) &= f_0(\mathbf{x}) + \sum_{p=1}^P \nu_p (\langle \mathbf{x}, \mathbf{a}_p \rangle - b_p) \\ &= f_0(\mathbf{x}) - \langle \boldsymbol{\nu}, \mathbf{b} \rangle + (\mathbf{A}^T \boldsymbol{\nu})^T \mathbf{x} \end{aligned}$$

and so

$$\begin{aligned} g(\boldsymbol{\nu}) &= -\langle \boldsymbol{\nu}, \mathbf{b} \rangle + \inf_{\mathbf{x}} \left( f_0(\mathbf{x}) + (\mathbf{A}^T \boldsymbol{\nu})^T \mathbf{x} \right) \\ &= -\langle \boldsymbol{\nu}, \mathbf{b} \rangle - \sup_{\mathbf{x}} \left( -f_0(\mathbf{x}) - (\mathbf{A}^T \boldsymbol{\nu})^T \mathbf{x} \right) \\ &= -\langle \boldsymbol{\nu}, \mathbf{b} \rangle - f_0^*(-\mathbf{A}^T \boldsymbol{\nu}), \end{aligned}$$

where  $f_0^*$  is the Fenchel dual of  $f_0$ :

$$f_0^*(\mathbf{y}) = \sup_{\mathbf{x}} (\langle \mathbf{x}, \mathbf{y} \rangle - f_0(\mathbf{x})).$$

With  $f_0 = \|\cdot\|$ , we know already that

$$f_0^*(\mathbf{y}) = \begin{cases} 0, & \|\mathbf{y}\|_* \leq 1, \\ \infty, & \text{otherwise} \end{cases},$$

so

$$g(\boldsymbol{\nu}) = \begin{cases} -\langle \boldsymbol{\nu}, \mathbf{b} \rangle, & \|\mathbf{A}^T \boldsymbol{\nu}\|_* \leq 1 \\ -\infty, & \text{otherwise} \end{cases}.$$

Thus the dual program is

$$\underset{\boldsymbol{\nu} \in \mathbb{R}^P}{\text{maximize}} \quad -\langle \boldsymbol{\nu}, \mathbf{b} \rangle \quad \text{subject to} \quad \|\mathbf{A}^T \boldsymbol{\nu}\|_* \leq 1,$$

where  $\|\cdot\|_*$  is the dual norm of  $\|\cdot\|$ . Note that since  $\boldsymbol{\nu}$  is unconstrained and  $\|-\mathbf{A}^T \boldsymbol{\nu}\|_* = \|\mathbf{A}^T \boldsymbol{\nu}\|_*$ , we can equivalently express this as

$$\underset{\boldsymbol{\nu} \in \mathbb{R}^P}{\text{maximize}} \quad \langle \boldsymbol{\nu}, \mathbf{b} \rangle \quad \text{subject to} \quad \|\mathbf{A}^T \boldsymbol{\nu}\|_* \leq 1,$$

which is identical to the Fenchel dual derived earlier in this course.