

14

Karush-Kuhn-Tucker Conditions

For convex problems of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{w}) \\ & \text{subject to} && \mathbf{f}(\mathbf{w}) \leq \mathbf{0} \\ & && \mathbf{A}\mathbf{w} = \mathbf{b}, \end{aligned} \tag{P}$$

we introduced the Lagrangian $\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ and defined the Lagrange dual as

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{w} \in \mathcal{D}} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

We saw that $g(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is a lower bound on the optimal value of (P). Note that we wrote the linear equality conditions $h_j(\mathbf{w}) = 0$ as system of linear equations $\mathbf{A}\mathbf{w} = \mathbf{b}$. We will derive conditions under which the lower bound provided by the Lagrange dual matches the upper bound, and derive a system of equations and inequalities that certify optimality, the Karush-Kuhn-Tucker (KKT) conditions. These conditions can be seen as generalizations of the first-order optimality conditions to the setting when equality and inequality constraints are present.

Constraint qualification

Let p^* and d^* denote the primal and dual optimal values, so that

$$d^* = \sup_{\boldsymbol{\lambda} \geq 0, \boldsymbol{\mu}} g(\boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \inf_{\mathbf{w} \in \mathcal{D}} \{f(\mathbf{w}) \mid f_i(\mathbf{w}) \leq 0, \mathbf{A}\mathbf{w} = \mathbf{b}\} = p^*.$$

Once certain conditions, called **constraint qualifications**, hold, we can ensure that **strong duality** holds, which means $d^* = p^*$. One particular such constraint qualification is Slater's Theorem.

Theorem 14.1. (*Slater conditions*) Assume that the interior of the domain \mathcal{D} of (P) is non-empty, that the problem (P) is convex, and that there exists a point $\mathbf{w} \in \mathcal{D}$ such that

$$f_i(\mathbf{w}) < 0, \quad 1 \leq i \leq m, \quad \mathbf{A}\mathbf{w} = \mathbf{b}, \quad 1 \leq j \leq p.$$

Assume that \mathbf{A} has maximal rank. Then $d^* = p^*$.

The proof makes use of a fundamental result on convex sets, the **separating hyperplane theorem**. For an affine hyperplane $H = \{\mathbf{w} : \mathbf{a}^T \mathbf{w} + b = 0\}$, we denote by $H_+ = \{\mathbf{w} : \mathbf{a}^T \mathbf{w} + b \geq 0\}$ one of the two closed halfspaces defined by the hyperplane, and by H_- the other.

Theorem 14.2 (Separating Hyperplane Theorem). *Let C and D be disjoint, nonempty convex subsets of \mathbb{R}^d . Then there exists an affine hyperplane H such that $C \subset H_+$ and $D \subset H_-$.*

A hyperplane with the properties of Theorem 14.2 is called a **separating hyperplane**. The hyperplane separates the two sets strictly, if each of the sets is contained in the interior of the respective halfspaces. It can be shown that for *compact* convex sets, there exists a hyperplane separating them strictly.

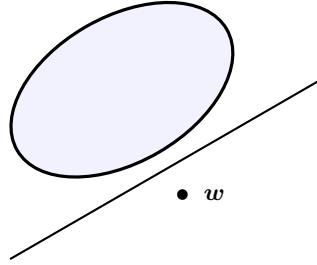


Figure 14.1: A hyperplane separating a convex sets and a point

Proof of Theorem 14.1. Assume \mathbf{A} has rank p (the number of rows). Assume moreover that p^* is finite, since if $p^* = -\infty$, then by weak duality we already have $d^* = p^*$. Define the convex set

$$\mathcal{A} = \{(\mathbf{u}, \mathbf{v}, t) \in \mathbb{R}^m \times \mathbb{R}^p \times \mathbb{R} \mid \forall \mathbf{w} \in \mathcal{D}, f_i(\mathbf{w}) \leq u_i, \mathbf{A}\mathbf{w} - \mathbf{b} = \mathbf{v}, f(\mathbf{w}) \leq t\}.$$

Then the optimal value of (P) is

$$p^* = \inf\{t \mid (\mathbf{0}, \mathbf{0}, t) \in \mathcal{A}\}.$$

Define the convex set \mathcal{B} as

$$\mathcal{B} = \{(\mathbf{0}, \mathbf{0}, s) \in \mathbb{R}^m \times \mathbb{R}^p \times \mathbb{R} \mid s < p^*\}.$$

The sets \mathcal{A} and \mathcal{B} are disjoint. To see this, assume to the contrary that there is a point $\mathbf{z} \in \mathcal{A} \cap \mathcal{B}$. Then since $\mathbf{z} \in \mathcal{B}$, $\mathbf{z} = (\mathbf{0}, \mathbf{0}, s)$ with $s < p^*$, but also since $\mathbf{z} \in \mathcal{A}$, there exists an $\mathbf{w} \in \mathcal{D}$ with $f_i(\mathbf{w}) \leq 0$, $\mathbf{A}\mathbf{w} - \mathbf{b} = \mathbf{0}$, and $f(\mathbf{w}) \leq s < p^*$, in contradiction to the optimality of p^* .

By the separating hyperplane theorem, there exists a hyperplane separating \mathcal{A} and \mathcal{B} (but not necessarily strictly!), defined by a vector $(\tilde{\lambda}, \tilde{\mu}, \nu) \neq \mathbf{0}$ and $\alpha \neq 0$ with the property that

$$(\mathbf{u}, \mathbf{v}, t) \in \mathcal{A} \implies \tilde{\lambda}^\top \mathbf{u} + \tilde{\mu}^\top \mathbf{v} + \nu t \geq \alpha \quad (14.1)$$

and

$$(\mathbf{u}, \mathbf{v}, t) \in \mathcal{B} \implies \tilde{\lambda}^\top \mathbf{u} + \tilde{\mu}^\top \mathbf{v} + \nu t \leq \alpha. \quad (14.2)$$

If $\tilde{\lambda} < \mathbf{0}$ or $\nu < 0$, we could find values of \mathbf{u} and t for which $(\mathbf{u}, \mathbf{v}, t) \in \mathcal{A}$, but that make the right-hand side of (14.1) arbitrary small, contradicting the bound by α . It follows that $\tilde{\lambda} \geq \mathbf{0}$ and $\nu \geq 0$. Condition (14.2) simply means that $\nu t \leq \alpha$ for all $t < p^*$, so that $\nu p^* \leq \alpha$. Combining this bound with (14.1), we get the two inequalities, valid for any $\mathbf{w} \in \mathcal{D}$,

$$\sum_{i=1}^m \tilde{\lambda}_i f_i(\mathbf{w}) + \tilde{\mu}^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) + \nu f(\mathbf{w}) \geq \alpha \geq \nu p^*. \quad (14.3)$$

Note that the left-hand side has the form of a Lagrangian function scaled by ν . If $\nu > 0$ we can divide by ν and set $\lambda_i = \tilde{\lambda}_i/\nu$, $\mu_i = \tilde{\mu}_i/\nu$, to obtain

$$\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \geq p^*.$$

By weak duality we have $p^* \geq g(\boldsymbol{\lambda}, \boldsymbol{\mu})$, so that we get strong duality if $\nu > 0$. If $\nu = 0$, then (14.3) implies

$$\sum_{i=1}^m \tilde{\lambda}_i f_i(\mathbf{w}) + \tilde{\mu}^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) \geq 0, \quad (14.4)$$

and for a point $\tilde{\mathbf{w}}$ satisfying the conditions $\mathbf{A}\mathbf{w} = \mathbf{b}$ and $f_i(\mathbf{w}) < 0$ for $1 \leq i \leq m$, this means that $\tilde{\lambda} = \mathbf{0}$. As $(\tilde{\lambda}, \tilde{\mu}, \nu) \neq \mathbf{0}$, we have $\boldsymbol{\mu} \neq \mathbf{0}$. Since $\tilde{\mu}^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) \geq 0$ and $= 0$ for some $\tilde{\mathbf{w}}$ in the interior of \mathcal{D} , there must be a \mathbf{w} in the interior such that $\tilde{\mu}^\top (\mathbf{A}\mathbf{w} - \mathbf{b}) < 0$, in contradiction to (14.4) (unless $\tilde{\mu}^\top \mathbf{A} = \mathbf{0}$, which would contradict the condition that \mathbf{A} has maximal rank p). \square

Example 14.3. Consider a linear programming problem of the form

$$\begin{aligned} &\text{minimize} && \mathbf{c}^\top \mathbf{w} \\ &\text{subject to} && \mathbf{A}\mathbf{w} = \mathbf{b} \\ &&& \mathbf{w} \geq \mathbf{0}. \end{aligned}$$

The inequality constraints are $-w_i \leq 0$, while the equality constraints are $\mathbf{a}_i^\top \mathbf{w} = b_i$. The Lagrangian has the form

$$\begin{aligned} \mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\mu}) &= \mathbf{c}^\top \mathbf{w} - \sum_{i=1}^m \lambda_i w_i + \sum_{j=1}^p \mu_j (\mathbf{a}_j^\top \mathbf{w} - b_j) \\ &= (\mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu})^\top \mathbf{w} - \mathbf{b}^\top \boldsymbol{\mu}. \end{aligned}$$

The infimum over \mathbf{w} of this function is $-\infty$ unless $\mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0}$. The Lagrange dual is therefore

$$g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \begin{cases} -\boldsymbol{\mu}^\top \mathbf{b} & \text{if } \mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0} \\ -\infty & \text{else.} \end{cases}$$

From the previous chapter we conclude that

$$\max\{-\mathbf{b}^\top \boldsymbol{\mu} \mid \mathbf{c} - \boldsymbol{\lambda} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\lambda} \geq \mathbf{0}\} \leq \min\{\mathbf{c}^\top \mathbf{w} \mid \mathbf{A}\mathbf{w} = \mathbf{b}, \mathbf{w} \geq \mathbf{0}\}.$$

Note that if we write $\mathbf{v} = -\boldsymbol{\mu}$ and $\mathbf{s} = \boldsymbol{\lambda}$, then we get the dual version of the linear programming problem we started out with, and in this case it is known that

$$\max_{\boldsymbol{\lambda} \geq \mathbf{0}, \boldsymbol{\mu}} g(\boldsymbol{\lambda}, \boldsymbol{\mu}) = p^*.$$

A more common way to express **linear programming duality** is as follows: the optimal value of

$$\text{maximize } \mathbf{b}^\top \mathbf{v} \quad \text{subject to } \mathbf{A}^\top \mathbf{v} \leq \mathbf{c} \tag{D}$$

coincides with the optimal value of

$$\text{minimize } \mathbf{c}^\top \mathbf{w} \quad \text{subject to } \mathbf{A}\mathbf{w} = \mathbf{b}, \mathbf{w} \geq \mathbf{0}, \tag{P}$$

provided both (D) and (P) have a finite solution (here, we set $\mathbf{v} = -\boldsymbol{\mu}$).

Example 14.4. Consider the problem

$$\text{minimize } \|\mathbf{w}\|^2 \quad \text{subject to } \mathbf{A}\mathbf{w} = \mathbf{b}.$$

Note that $\|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{w}$. The Lagrangian is $\mathcal{L}(\mathbf{w}, \boldsymbol{\mu}) = \|\mathbf{w}\|^2 + \boldsymbol{\mu}^\top (\mathbf{A}\mathbf{w} - \mathbf{b})$. For any $\boldsymbol{\mu}$, we can find the infimum

$$g(\boldsymbol{\mu}) = \inf_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\mu})$$

by setting the derivative of the Lagrangian to \mathbf{w} to zero:

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \boldsymbol{\mu}) = 2\mathbf{w} + \mathbf{A}^\top \boldsymbol{\mu} = \mathbf{0},$$

which gives the solution

$$\mathbf{w} = -\frac{1}{2} \mathbf{A}^\top \boldsymbol{\mu}.$$

The dual function is therefore

$$g(\boldsymbol{\mu}) = -\frac{1}{4} \boldsymbol{\mu}^\top \mathbf{A}^\top \mathbf{A} \boldsymbol{\mu} - \mathbf{b}^\top \boldsymbol{\mu}.$$

As the negative of a positive semidefinite quadratic function, it is concave. Moreover, we get the lower bound

$$-\frac{1}{4} \boldsymbol{\mu}^\top \mathbf{A}^\top \mathbf{A} \boldsymbol{\mu} - \mathbf{b}^\top \boldsymbol{\mu} \leq \inf\{\|\mathbf{w}\|^2 \mid \mathbf{A}\mathbf{w} = \mathbf{b}\}.$$

The problem we started out with is convex, and if we assume that there exists a feasible primal point, then the above inequality is in fact an equality by Slater's conditions.

Karush-Kuhn-Tucker optimality conditions

Consider now a not necessarily convex problem of the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{w}) \\ & \text{subject to} && \mathbf{f}(\mathbf{w}) \leq \mathbf{0} \\ & && \mathbf{h}(\mathbf{w}) = \mathbf{0}, \end{aligned} \tag{14.5}$$

If p^* is the optimal solution of (14.5) and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ dual variables, then we have seen that (this holds even in the non-convex case)

$$p^* \geq g(\boldsymbol{\lambda}, \boldsymbol{\mu}).$$

From this follows that for any primal feasible point \mathbf{w} ,

$$f(\mathbf{w}) - p^* \leq f(\mathbf{w}) - g(\boldsymbol{\lambda}, \boldsymbol{\mu}).$$

The difference $f(\mathbf{w}) - g(\boldsymbol{\lambda}, \boldsymbol{\mu})$ between the primal objective function at a primal feasible point and the dual objective function at a dual feasible point is called the **duality gap** at \mathbf{w} and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$. For any such points we know that

$$p^*, q^* \in [g(\boldsymbol{\lambda}, \boldsymbol{\mu}), f(\mathbf{w})],$$

and if the gap is small we have a good approximation of the primal and dual optimal values. The duality gap can be used in iterative algorithms to define stopping criteria: if the algorithm generates a sequence of primal-dual variables $(\mathbf{w}^k, \boldsymbol{\lambda}^k, \boldsymbol{\mu}^k)$, then we can stop if the duality gap is less than, say, a predefined tolerance ε .

Now suppose that $(\mathbf{w}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ is such that the duality gap is zero. Then

$$\begin{aligned} f(\mathbf{w}^*) &= g(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \inf_{\mathbf{w}} \left(f(\mathbf{w}) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{w}) + \sum_{j=1}^p \mu_j^* h_j(\mathbf{w}) \right) \\ &\leq f(\mathbf{w}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{w}^*) + \sum_{j=1}^p \mu_j^* h_j(\mathbf{w}^*) \\ &\leq f(\mathbf{w}^*), \end{aligned}$$

where the last inequality follows from the fact that $h_j(\mathbf{w}^*) = 0$ and $\lambda_i^* f_i(\mathbf{w}^*) \leq 0$ for $1 \leq j \leq p$ and $1 \leq i \leq m$. It follows that the inequalities are in fact equalities. From the identity

$$f(\mathbf{w}^*) = f(\mathbf{w}^*) + \sum_{i=1}^m \lambda_i^* f_i(\mathbf{w}^*)$$

and $\lambda_i^* \geq 0$ and $f_i(\mathbf{w}^*) \leq 0$ we also conclude that at such optimal points,

$$\lambda_i^* f_i(\mathbf{w}^*) = 0, \quad 1 \leq i \leq m.$$

This condition is known as **complementary slackness**. From the above we also see that \mathbf{w}^* minimizes the Lagrangian $\mathcal{L}(\mathbf{w}, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, so that

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

Collecting these conditions (primal and dual feasibility, complementary slackness, vanishing gradient), we arrive at a set of optimality conditions known as the Karush-Kuhn-Tucker (KKT) conditions.

Theorem 14.5. (KKT conditions) Let \mathbf{w}^* and $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ be primal and dual optimal solutions of (14.5) with zero duality gap. Then the following conditions are satisfied:

$$\begin{aligned} \mathbf{f}(\mathbf{w}^*) &\leq \mathbf{0} \\ \mathbf{h}(\mathbf{w}^*) &= \mathbf{0} \\ \boldsymbol{\lambda}^* &\geq \mathbf{0} \\ \lambda_i^* f_i(\mathbf{w}^*) &= 0, \quad 1 \leq i \leq m \\ \nabla_{\mathbf{w}} f(\mathbf{w}^*) + \sum_{i=1}^m \lambda_i^* \nabla_{\mathbf{w}} f_i(\mathbf{w}^*) + \sum_{j=1}^p \mu_j^* \nabla_{\mathbf{w}} h_j(\mathbf{w}^*) &= \mathbf{0}. \end{aligned}$$

Moreover, if the problem is convex and the Slater Conditions (Theorem 14.1) are satisfied, then any points satisfying the KKT conditions have zero duality gap.

Notes

The Karush-Kuhn-Tucker conditions were introduced by Kuhn and Tucker [1], and the necessity was shown by William Karush in his 1939 MSc thesis at the University of Chicago. These conditions generalize and unify different previously known results: the optimality conditions for linear programming based on linear programming duality, and Lagrange duality in the presence of inequalities. The KKT conditions also form the basis of algorithms for non-linear optimization, such as interior-point methods. A detailed treatment can be found in [2]. The separating hyperplane theorem is a central result in convexity, and lies at the heart of many duality results. In particular, it is the basis of a classic “theorem of the alternative” known as Farkas’ Lemma, which states that given a matrix $\mathbf{A} \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$, there exists a vector \mathbf{w} such that

$$\mathbf{A}\mathbf{w} = \mathbf{b}, \quad \mathbf{w} \geq \mathbf{0}$$

if and only if there is no $\mathbf{v} \in \mathbb{R}^m$ such that

$$\mathbf{A}^\top \mathbf{v} \geq \mathbf{0}, \quad \mathbf{v}^\top \mathbf{b} < 0.$$

This result, in turn, is an ingredient for deriving linear programming duality.

- [1] Harold W Kuhn and Albert W Tucker. Nonlinear programming. In *Traces and emergence of nonlinear programming*, pages 247–258. Springer, 2014.
- [2] J. Nocedal and S. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York, 2006.