# Lectures 5-7: KKT and Lagrange Duality

## Magnus Önnheim

## November 20, 2014

[Note: these notes are prelimiary and subject to change.]

We are now going to turn our attention to constrained optimization problems, that is, problems of the form

$$\min f(\mathbf{x}), \tag{1a}$$
$$\text{subject to } \mathbf{x} \in S, \tag{1b}$$

where $f : \mathbb{R}^n \to \mathbb{R}$ and $S \subset \mathbb{R}^n$. In the case where $S$ is convex and $f \in C^1$, we have already worked out an optimality condition for this type of problem, i.e., we have a theorem of the form

$$\mathbf{x}^* \text{is a local minimum} \implies \mathbf{x}^* \text{is a stationary point}$$

We formulated stationarity in several different ways, one of which was

$$-\nabla f(\mathbf{x}^*) \in N_S(\mathbf{x}^*),$$

where $N_S(\mathbf{x}^*)$ is the normal cone of $S$ at $\mathbf{x}^*$, i.e.,

$$N_S(\mathbf{x}^*) := \{\mathbf{p} \in \mathbb{R}^n \mid \mathbf{p}^{\mathrm{T}}(\mathbf{y} - \mathbf{x}^*) \leq 0, \forall \mathbf{y} \in S\}.$$

Let us pause for a minute to reflect on this; according to what we know about unconstrained optimization the condition that $-\nabla f(\mathbf{x}^*) \in N_S(\mathbf{x}^*)$ simply means that for all $y \in S$, we have $(\nabla f(\mathbf{x}^*))^{\mathrm{T}}(\mathbf{y} - \mathbf{x}^*) \geq 0$, which tells us that the vector pointing from our locally optimal point $\mathbf{x}^*$ to the feasible point $\mathbf{y} \in S$ does not look like a descent direction for $f$. The optimality condition $-\nabla f(\mathbf{x}^*)$ thus says nothing else than that it should not be possible to move from $\mathbf{x}^*$ in a direction allowed by $S$, such that $f$ decreases.

This is also the approach we will take to develop optimality conditions for more general non-linearly constrained problems. We first formalize the notion of what a "direction allowed by $S$" is going, and then require that these allowed directions do not contain any descent directions for $f$. It will however turn out that formulating a good notion of "allowed direction" is quite possibly the technically most challenging part of this course!

# 1 Geometric optimality conditions

First we introduce the most natural way of measuring allowed direction.

**Definition 1** (cone of feasible direction, the radial cone). *The cone of feasible directions $R_S(\mathbf{x})$ for $S$ at $\mathbf{x} \in S$ is defined as*

$$\{\mathbf{p} \in \mathbb{R}^n \mid \exists \delta > 0, \mathbf{x} + \alpha\mathbf{p} \in S, \forall 0 \leq \alpha \leq \delta\}. \tag{2}$$

An element $\mathbf{p} \in R_s(\mathbf{x})$ is thus simply a vector such that the feasible set $S$ contains a non-trivial part of the half-line $\mathbf{x} + \alpha\mathbf{p}$, $\alpha \geq 0$. Although natural this cone is too small to use for optimality conditions for non-linearly constrained programs[1].

**Example 1.** *Let $S := \{x_2 = x_1^2\}$. Then $R_S(\mathbf{x}) = \emptyset$ for all $\mathbf{x} \in S$, simply because the feasible set is a curved line in $\mathbb{R}^n$.*

The perhaps most widely used object in the literature to develop optimality conditions is therefore a significantly more complicated object.

**Definition 2** (The tangent cone $T_S(\mathbf{x})$). *The tangent cone for $S$ at $\mathbf{x} \in S$ is defined as*

$$\begin{aligned}
T_S(\mathbf{x}) := \big\{\mathbf{p} \mid &\exists \{\mathbf{x}_k\}_{k=1}^{\infty} \subset S, \{\lambda_k\}_{k=1}^{\infty}, \text{ such that} \\
&\lim_{k \to \infty} \mathbf{x}_k = \mathbf{x}, \\
&\lim_{k \to \infty} \lambda_k(\mathbf{x}_k - \mathbf{x}) = \mathbf{p}\big\}.
\end{aligned} \tag{3}$$

It may not be obvious to the reader in what way this horrenduous looking definition of a cone $T_S(\mathbf{x})$ actually measures some kind of allowed direction of $S$. However, in words, the above definition tells us that to check whether a vector $\mathbf{p} \in T_S(\mathbf{x})$ we should check whether there is a *feasible* sequence of points $\mathbf{x}_k \in S$ that approaches $\mathbf{x}$, such that $\mathbf{p}$ is the asymptotic direction from which $\mathbf{x}_k$ approaches $\mathbf{x}$. Seen this way, we can convince ourselves that $T_S(\mathbf{x})$ consists precisely of all the possible directions in which $\mathbf{x}$ can be asymptotically approached through $S$. In other words, $T_S(\mathbf{x})$ consists of all vectors pointing 'along' $S$ from $\mathbf{x} \in S$, that is all vectors that are *tangent* to $S$ at $\mathbf{x}$.

**Example 2.** *Let again $S := \{x_2 = x_1^2\}$. Then $T_S(0) = \{\mathbf{p} \mid p_2 = 0\}$. Note that in this example, we can identify the tangent cone with the ordinary tangent line that we learn how to compute in multivariable courses.*

**Example 3.** *Suppose that we have a smooth curve in $S$ starting at $\mathbf{x} \in S$, that is, we have a $C^1$ map $\gamma : [0, T] \to S$ for some $T > 0$. Then $\gamma'(0) \in T_S(\mathbf{x})$. To see this, note that the very definition of (one-sided) derivative is just that*

$$\gamma'(0) = \lim_{t \to 0} \frac{\gamma(t) - \gamma(0)}{t}, \tag{4}$$

---

[1]it will, however, work perfectly for *linear* programs!

*so that if we fix any sequence $t_k \to 0$, and let $\mathbf{x}_k := \gamma(t_k)$, $\lambda_k = 1/t_k$, we have defined the sequences required in the definition of $T_S(\mathbf{x})$.*

It remains to formulate a notion of descent directions to $f$, fortunately we can use the same characerization as in the unconstrained case

**Definition 3** (descent cone). $\overset{\circ}{F}(\mathbf{x}) := \{\mathbf{p} \in \mathbb{R}^n \mid \nabla f(\mathbf{x})^\mathrm{T} \mathbf{p} < 0\}$.

The above examples should then make the following theorem intuitively obvious.

**Theorem 1** (geometric optimality conditions). *Consider the problem* (1), *where $f \in C^1$. Then*

$$\mathbf{x}^* \text{ is a local minimum } \implies \overset{\circ}{F}(\mathbf{x}^*) \cap T_S(\mathbf{x}^*) = \emptyset. \tag{5}$$

*Proof.* See theorem 5.10 in the book. □

**Example 4.** *If we again return to our example with smooth curves, we showed that for any smooth curve $\gamma$ thorugh $S$ starting at $\mathbf{x}^*$, we had $\gamma'(0) \in T_S(\mathbf{x}^*)$. Try to convince yourself that the geometric optimality condition reduces to the statement that $\dfrac{d}{dt}|_{t=0} f(\gamma(t)) \geq 0$ when applied to this tangent vector.*

# 2 Going from geometric to useful

Now we have develeped a quite elegant optimality condition, however there is a huge catch. There is no practical way to compute $T_S(\mathbf{x})$ directly from its definition! There are two ways out of this dilemma. The first (which lead to the Fritz John conditions) is to simply replace the cone $T_S(\mathbf{x})$ by smaller cones.

**Lemma 1.** *If the family of cones $C(\mathbf{x}) \subseteq T_S(\mathbf{x})$ for all $\mathbf{x} \in S$, then $\overset{\circ}{F}(\mathbf{x}^*) \cap C(\mathbf{x}^*) = \emptyset$ is a neccessary optimality condition.*

*Proof.* Using the geometric optimality condition we have for any locally optimal $\mathbf{x}^* \in S$,
$$\overset{\circ}{F}(\mathbf{x}^*) \cap C(\mathbf{x}^*) \subseteq \overset{\circ}{F}(\mathbf{x}^*) \cap T_S(\mathbf{x}^*) = \emptyset.$$
□

The obvious danger of introducing smaller cones is that the optimality conditions we get are *weaker* than the geometric optimality conditions.

**Example 5.** *Let $C(\mathbf{x}) = R_S(\mathbf{x})$ and consider again the example $S := \{x_2 = x_1^2\}$. Since $R_S(\mathbf{x}) = \emptyset$, the optimality condition $\overset{\circ}{F}(\mathbf{x}) \cap R_S(\mathbf{x}) = \emptyset$ holds for any feasible $\mathbf{x} \in S$, which is obviously a totally useless optimality condition*

The second way out is to introduce regularity conditions, or *constraint qualifications*, which will allow us to actually compute the tangent cone $T_S(\mathbf{x})$ by other means. This approach leads to the Karush-Kuhn-Tucker (KKT) conditions. The obvious drawback of this approach is that, although the KKT conditions are equally strong as the geometric conditions, they are *less general*, i.e., they do not apply for irregular problems.

From now on we will consider a problem of the form

$$\min f(\mathbf{x}), \tag{6a}$$

$$\text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m \tag{6b}$$

where now $f : \mathbb{R}^n \to \mathbb{R}$, and $g_i : \mathbb{R}^n \to \mathbb{R}$, $i = 1, \ldots, m$ are all $C^1$, i.e., we take $S$ to be of the form $S := \{\mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0, i = 1, \ldots, m\}$. This allows us to define (even more) cones that are related to $T_S(\mathbf{x})$. From now on we will also use the notation $\mathcal{I}(\mathbf{x})$ to denote the *active set of constraints* at $\mathbf{x}$, that is,

$$\mathcal{I}(\mathbf{x}) := \{i \in \{1, \ldots, m\} \mid g_i(\mathbf{x}) = 0\}. \tag{7}$$

**Definition 4** (gradient cones)**.** *We define the inner gradient cone* $\overset{\circ}{G}(\mathbf{x})$ *as*

$$\overset{\circ}{G}(\mathbf{x}) := \{\mathbf{p} \in \mathbb{R}^n \mid \nabla g_i(\mathbf{x})^{\mathrm{T}}\mathbf{p} < 0, \forall i \in \mathcal{I}(\mathbf{x})\}. \tag{8}$$

*Similarly we define the gradient cone* $G(\mathbf{x})$ *as*

$$G(\mathbf{x}) := \{\mathbf{p} \in \mathbb{R}^n \mid \nabla g_i(\mathbf{x})^{\mathrm{T}}\mathbf{p} \leq 0, \forall i \in \mathcal{I}(\mathbf{x})\}. \tag{9}$$

Note that in the inner gradient cone $\overset{\circ}{G}(\mathbf{x})$ consists of all vectors $\mathbf{p}$ that can be guaranteed to be descent directions of all defining functions for the active constraints, while the gradient cone $G(\mathbf{x})$ consists of all directions that can be guaranteed not to be ascent directions for the active constraints. Since ascent/descent of the active constraints captures the intuition behing what a 'feasible movement' from $\mathbf{x}$ through $S$ is, the following theorem should not come as a surprise.

**Theorem 2** (Relations between cones)**.** *For the problem* (6) *it holds that*

$$\mathrm{cl}\, \overset{\circ}{G}(\mathbf{x}) \subseteq \mathrm{cl}\, R_S(\mathbf{x}) \subseteq T_S(\mathbf{x}) \subseteq G(\mathbf{x}) \tag{10}$$

*Proof.* See the book for a complete proof. The moral of the story is that $\overset{\circ}{G}(\mathbf{x})$ consists of all directions $\mathbf{p}$ that are descent directions for all active constraints. Thus the active constraints must decrease along the direction $\mathbf{p}$, thus defining a feasible direction, i.e., $\overset{\circ}{G}(\mathbf{x}) \subseteq R_S(\mathbf{x})$.

The inclusion $R_S(\mathbf{x}) \subseteq T_S(\mathbf{x})$ follows the above example of smooth curves: a straight line in $S$ is most definitely a smooth curve.

The final inclusion remains. Take an arbitrary $\mathbf{p} \in T_S(\mathbf{x})$, and let $\{\mathbf{x}_k\}$, $\{\lambda_k\}$ be the sequences as in the definition of $T_S(\mathbf{x})$. Then for any $i \in \mathcal{I}(\mathbf{x})$

$$\nabla g_i(\mathbf{x})^{\mathrm{T}} \mathbf{p} = \lim \frac{g_i(\mathbf{x}_k) - g_i(\mathbf{x})}{\|\mathbf{x}_k - \mathbf{x}\|} \leq 0 \tag{11}$$

Since $\mathbf{x}_k \in S$ for all $k$, so that $g_i(\mathbf{x}_k) \leq 0$, and $g_i(\mathbf{x}) = 0$ since $i \in \mathcal{I}(\mathbf{x})$.

Finally taking closures the theorem follows as $T_S(\mathbf{x})$ is a closed set (the proof of which can be found in the book). $\qquad\square$

Please note that the above inclusions in general are strict.

# 3 The Fritz John conditions

The Fritz John conditions are what we when we replace the tangent $T_S$ in the geometric optimality condition by $\overset{\circ}{G}(\mathbf{x})$. According to the above discussion the Fritz John conditions are *weaker* than the geometric optimality conditions.

$$\mathbf{x}^* \text{is locally optimal in } (6) \implies \overset{\circ}{G}(\mathbf{x}) \cap \overset{\circ}{F}(\mathbf{x}) = \emptyset. \tag{12}$$

Again, this condition looks fairly abstract, however it is actually quite easy to turn the above equation into something more practically viable. The moral of the story is that the above equation is for a fix $\mathbf{x}$ just the statement that a linear system of inequalities does not have solution. Fortunately we have a tool at our disposal for turning an inconsistent set of linear inequalities into a consistent set of inequalities, namely Farkas' Lemma.

**Theorem 3** (The Fritz John conditions). *If $\mathbf{x}^*$ is locally optimal in* (6), *then the system*

$$\mu_0 \nabla f(\mathbf{x}^*) + \sum_{i=1}^{m} \mu_i \nabla g_i(\mathbf{x}^*) = 0, \tag{13}$$

$$\mu_i g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m, \tag{14}$$

$$(\mu_0, \mu_i) \geq 0, \tag{15}$$

$$(\mu_0, \mu_i) \neq 0.. \tag{16}$$

*has a solution $\boldsymbol{\mu}$.*

*Proof.* For full details, see the book. The core of the argument is to use Farkas' Lemma to convert an inconsistent system to a consistent one. That is, we formulate the local optimality condition $\overset{\circ}{G}(\mathbf{x}) \cap \overset{\circ}{F}(\mathbf{x}) = \emptyset$ as the unsolvability of the system of linear inequalities

$$\nabla g_i(\mathbf{x}^*)^\mathrm{T}\mathbf{p} + \alpha \leq 0, \quad i \in \mathcal{I}(\mathbf{x}^*) \tag{17}$$

$$\nabla f(\mathbf{x}^*)^\mathrm{T}\mathbf{p} + \alpha \leq 0, \tag{18}$$

$$\alpha > 0. \tag{19}$$

This looks just like a system in Farkas Lemma, so the unsolvability of the above is equivalent the solvability of

$$\mu_0 \nabla f(\mathbf{x}^*) + \sum_{i \in \mathcal{I}(\mathbf{x}^*)} \mu_i \nabla g_i(\mathbf{x}^*) = 0, \tag{20}$$

$$\mu_0 + \sum_{i \in \mathcal{I}(\mathbf{x}^*)} \mu_i = 1. \tag{21}$$

To turn this system into the one claimed in the theorem we add the multipliers $\mu_i = 0$ for $i \notin \mathcal{I}(\mathbf{x}^*)$, this yields the complementarity conditions (14), and by a simple scaling argument the second row above is equivalent to (16).

$\square$

It might be somewhat illuminating to try and use the above reasoning to get optimality conditions for unconstrained problems. The logic is exactly the same: we have a geometric condition (namely 'no descent directions') which we formulate as the unsolvability of the 'system' $\nabla f(\mathbf{x}^*)^\mathrm{T} p < 0$. This can be turned into a solvability of another system, namely $\nabla f(\mathbf{x}^*) = 0$.

The main drawback of the Fritz-John conditions is that they are too weak, consider the example from before where $R_S(\mathbf{x}) = \emptyset$. We can also see this in the Fritz-John system; there is a multiplier in front of the objective function term. If there is a solution to the Fritz-John system where the multiplier $\mu_0 = 0$, in effect the objective function does not play any role whatsoever in the system, which is indeed a very weak optimality condition. This insight gives us at least one way to think of regularity conditions (constraint qualifications); it is conditions that guarantee that any solution of the Fritz-John system must satisfy $\mu_0 \neq 0$.

## 4 KKT conditions

We begin by developing the KKT conditions from a standpoint where we simply assume some regularity of the problem at hand; that is we simply assume that the problem we consider is nicely behaved, and postpone the issue of whether any given problem is indeed well behaved until later.

**Definition 5** (Abadies constraint qualification)**.** *We say that the problem* (6) *satifies Abadies constraint qualification if* $T_S(\mathbf{x}) = G(\mathbf{x})$ *for all* $\mathbf{x} \in S$.

*Remark:* Abadies constraint qualification should be viewed as an abstract condition, it really is just the mathematical way of expressing "(6) is well-behaved"

This allows us to the following (important!) theorem

**Theorem 4.** *Assume that the problem* (6) *satisfies Abadies CQ, then at any locally optimal point* $\mathbf{x}^*$ *the system*

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^{m} \mu_i \nabla g(\mathbf{x}^*) = 0, \tag{22}$$

$$\mu_i g_i(\mathbf{x}^*) = 0, \tag{23}$$

$$\mu_i \geq 0, \quad i = 1, \ldots, m. \tag{24}$$

*has a solution* $\boldsymbol{\mu}$.

*Proof.* The geometric optimality and the CQ yields that $\emptyset = T_S(\mathbf{x}^*) \cap \overset{\circ}{F}(\mathbf{x}^*) = G(\mathbf{x}^*) \cap \overset{\circ}{F}(\mathbf{x}^*) = \emptyset$. This can be formulated as the unsolvability of the system

$$\nabla f(\mathbf{x}^*)^{\mathrm{T}} \mathbf{p} < 0, \tag{25}$$

$$\nabla g_i(\mathbf{x}^*)^{\mathrm{T}} \mathbf{p} \leq 0, \quad i \in \mathcal{I}(\mathbf{x}^*). \tag{26}$$

This can, just as in the Fritz-John conditions, be turned into a solvable system by Farkas Lemma (here we think of the matrix $A$ as consisting the $g_i$-gradients, and $b$ as $-\nabla f(\mathbf{x}^*)$),

$$\sum_{i \in \mathcal{I}} \nabla g_i(\mathbf{x}^*) = -\nabla f(\mathbf{x}^*), \tag{27}$$

$$\mu_i \geq 0, i \in \mathcal{I}(\mathbf{x}^*). \tag{28}$$

This is equivalent to the claimed system by the same tricks we used in the Fritz-John conditions. $\qquad\square$

Note now that the KKT conditions are precisely the Fritz-John conditions, with the added requirement that $\mu_0 = 1$. Also a note on terminology, we call the vector $\boldsymbol{\mu}$ solving the KKT system for some fixed $\mathbf{x} \in S$ a *Lagrange multiplier*. Beware that this terminology will unfortunately conflict with the terminology of Lagrangian dualoty later on!

## 5   Constraint qualifications

Our final task in showing that the KKT conditions is now to find practical, usable conditions under which we can guarantee that Abadies CQ holds. We start with one of the simplest and most useful ones.

**Definition 6** (LICQ). *We say that the linear independence constraint qualification (LICQ) holds at* $\mathbf{x} \in S$ *if the set* $\{\nabla g_i(\mathbf{x}), i \in \mathcal{I}\}$ *is linearly independent.*

**Proposition 1.** *The LICQ implies the validity of the KKT conditions*

*Proof.* Note: the book does a different (in my view harder) proof. What we want to show is that KKT conditions holds at any locally optimal point, assuming the LICQ. But the Fritz-John conditions are *always* valid, so there is a solution to the Fritz-John system at $\mathbf{x}^*$. But this solution must satisfy $\mu_0 \neq 0$, since otherwise we have a nonzero solution to $\sum_{i \in \mathcal{I}} \mu_i \nabla g_i(\mathbf{x}^*) = 0$, which contradicts the linear independence assumption. $\qquad \square$

The other constraint qualifications we will consider in this course comes from comparing the cones of (10). As a first example we consider the Mangasarian-Fromowitz CQ.

**Definition 7** (MFCQ). *The Mangarasarian-Fromowitz constraint qualification holds at $\mathbf{x} \in S$ if $\overset{\circ}{G}(\mathbf{x}) \neq \emptyset$.*

**Proposition 2.** *The MFCQ implies the Abadie CQ.*

*Proof.* The main idea is to show that $\mathrm{cl}\ \overset{\circ}{G}(\mathbf{x}) = G(\mathbf{x})$, since if this holds then by (10) $T_S(\mathbf{x}) = G(\mathbf{x})$ holds. To show this claim, we pick an arbitrary $\mathbf{p}_1 \in G(\mathbf{x})$, and an arbitrary $\mathbf{p}_0 \in \overset{\circ}{G}(\mathbf{x})$ (which exists since $\overset{\circ}{G}(\mathbf{x}) \neq \emptyset$). We then let, for $t \in (0, 1)$, $\mathbf{p}_t := (1 - t)\mathbf{p}_0 + t\mathbf{p}_1$. Then, for any $i \in \mathcal{I}$,

$$\nabla g_i(\mathbf{x})^{\mathrm{T}} \mathbf{p}_t = \underbrace{(1 - t)\nabla g_i(\mathbf{x})^{\mathrm{T}} \mathbf{p}_0}_{<0,\ \mathbf{p}_0 \in \overset{\circ}{G}(\mathbf{x})} + \underbrace{t\nabla g_i(\mathbf{x})^{\mathrm{T}} \mathbf{p}_1}_{\leq 0,\ \mathbf{p}_1 \in G(\mathbf{x})} < 0. \tag{29}$$

This shows that $\mathbf{p}_t \in \overset{\circ}{G}(\mathbf{x})$ for all $t \in [0, 1)$. Since clearly $\mathbf{p}_t \to \mathbf{p}_1$ as $t \to \infty$, we have shown that $\mathbf{p}_1 \in \mathrm{cl}\ \overset{\circ}{G}(\mathbf{x})$. $\qquad \square$

The MFCQ can be used to get other constraint qualifications as well.

**Definition 8.** *Slaters CQ holds for* (6) *if $g_i$ are all convex functions, and there is an interior point, i.e., a point $\mathbf{x}_0$ such that $g_i(\mathbf{x}_0) < 0$ for all $\mathbf{x} \in S$.*

**Proposition 3.** *Slaters CQ implies Abadies CQ.*

*Proof.* For any $\mathbf{x} \in S$, let $i \in \mathcal{I}(\mathbf{x})$. Then we have by convexity

$$0 > g_i(\mathbf{x}_0) \geq \underbrace{g_i(\mathbf{x})}_{=0} + \nabla g_i(\mathbf{x})^{\mathrm{T}}(\mathbf{x}_0 - \mathbf{x}). \tag{30}$$

Rearranging yields $\nabla g_i(\mathbf{x})^{\mathrm{T}}(\mathbf{x}_0 - \mathbf{x}) < 0$ for all $i \in \mathcal{I}(\mathbf{x})$. Hence $\mathbf{x}_0 - \mathbf{x} \in \overset{\circ}{G}(x)$, so MFCQ holds at $\mathbf{x}$, and thus also Abadie. $\qquad \square$

**Definition 9** (Affine constraints CQ). *The affine constraints holds for* (6) *if all the functions $g_i$ are affine, $i = 1, \ldots, m$.*

**Proposition 4.** *The affine constraints CQ imply Abadies CQ.*

*Proof.* If the constraints are affine then for any $\mathbf{p}$ and any $\mathbf{x} \in S$ and any $i \in \mathcal{I}(\mathbf{x})$ we have $g_i(\mathbf{x} + t\mathbf{p}) = g_i(\mathbf{x}) + t\nabla g_i(\mathbf{x})^{\mathrm{T}}\mathbf{p} = t\nabla g_i(\mathbf{p})^{\mathrm{T}}\mathbf{p}$. So if $\mathbf{p} \in G(\mathbf{x})$ we have $g_i(\mathbf{x} + t\mathbf{p}) \leq 0$, and thus $\mathbf{x} + t\mathbf{p}$ is feasible for all small enough $t \geq 0$, i.e., $\mathbf{p} \in R_S(\mathbf{x})$. Hence $G(\mathbf{x}) = R_S(\mathbf{x})$ from which it follows that $T_S(\mathbf{x}) = G(\mathbf{x})$ from (10). □

# 6 Equality constraints

So far we have only talked about problems with inequality constraints, we briefly outline how to apply the above theory for the problem to

$$\min f(\mathbf{x}), \tag{31a}$$
$$\text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m, \tag{31b}$$
$$h_j(\mathbf{x}) = 0, \quad j = 1, \ldots, l. \tag{31c}$$

where are the functions above are $C^1$. The main idea is simply replace the equality constraints $h_j(\mathbf{x}) = 0$, with two inequality constraints, i.e., by $h_j(\mathbf{x}) \leq 0$ and $h_j(\mathbf{x}) \leq 0$, and apply the KKT theory to the problem

$$\min f(\mathbf{x}), \tag{32a}$$
$$\text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m, \tag{32b}$$
$$h_j(\mathbf{x}) \leq 0, \quad j = 1, \ldots, l, \tag{32c}$$
$$-h_j(\mathbf{x}) \leq 0, \quad j = 1, \ldots, l. \tag{32d}$$

The details of what happens to KKT system under this rewriting can be found in the book, but the main observation is just that the equality constraints are *always* active in any feasible solution, and they will enter the KKT system with non-negative multipliers of opposite sign, which we can rewrite simply as a multiplier with any sign restrictions. The KKT system becomes

$$f(\mathbf{x}) + \sum_{i=1}^{m} \mu_i \nabla g_i(\mathbf{x}) + \sum_{j=1}^{l} h_j(\mathbf{x}) = 0, \tag{33a}$$
$$\mu_i g_i(\mathbf{x}) = 0, \quad i = 1, \ldots, m, \tag{33b}$$
$$\mu_i \geq 0, \quad i = 1, \ldots, m \tag{33c}$$

The main difficulty that we don't really address in this course is what happens to the constraint qualifications when we add eequality constraints. In the previous sections our main "useful" CQ was the MFCQ, i.e., that $\overset{\circ}{G}(\mathbf{x}) \neq \emptyset$. But in the presence of equality constraints, this can *never* hold (think about why?). The way out is to introduce (yet another) cone $H(\mathbf{x}) := \{\mathbf{p} \in \mathbb{R}^n \mid$

$\nabla h_j(\mathbf{x})^{\mathrm{T}}\mathbf{p} = 0, \ j = 1, \ldots, l\}$ and try to develop the theory of the previous sections when replacing the cones $\overset{\circ}{G}(\mathbf{x})$ and $G(\mathbf{x})$ by $\overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x})$ and $G(\mathbf{x}) \cap H(\mathbf{x})$ respectively. THe difficulty lies in that we cannot immediately generalize the statement $\overset{\circ}{G}(\mathbf{x}) \subset T_S(\mathbf{x}^*)$ to the statement $\overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x}) \subset T_S(\mathbf{x})$, which make creating useful constraint qualifications somewhat trickier. However, it turns out that we need to require to make the above theory work is that the set of gradients $\{\nabla h_j(\mathbf{x})\}_{j=1}^l$ is linearly independent, but this lies beyond the scope of this course. The constraint qualifications above then have to get modified to include the statement "and the set of gradients $\{\nabla h_j(\mathbf{x})\}$ is linearly independent"[2] . We refer the reader to the book for detailed statements of all the CQs.

# 7 Sufficiency under convexity

All of what we have done above is about developing *neccessary* optimality conditions for the problem (6). However, it is very natural to ask whether the KKT conditions are ever *sufficient* for optimality, that is, can we say that if the KKT system is solvable at $\mathbf{x}^*$, that we can conclude that $\mathbf{x}^*$ is optimal in (6)? In the unconstrained case, we saw that the property that allows such statements is convexity, and it turns out that what we need in the constrained case is also convexity, but in terms of both the objective function and the constraints.

**Theorem 5.** *If, in* (6), *the objective function $f$ and all constraint functions $g_i$, $i = 1, \ldots, m$ are convex, then the KKT conditions are a sufficient optimality condition.*

*Proof.* Suppose that the KKT conditions hold at $\mathbf{x}^*$, with Lagrange multipliers $\mu_i^*$, $i = 1, \ldots, m$, and pick an arbitary feasible $\mathbf{x}$.

We have

$$
\begin{aligned}
f(\mathbf{x}) - f(\mathbf{x}^*) &\geq \nabla f(\mathbf{x}^*)^{\mathrm{T}}(\mathbf{x} - \mathbf{x}^*) \\
&= -\sum_i \mu_i^* \nabla g_i(\mathbf{x})^{\mathrm{T}}(\mathbf{x} - \mathbf{x}^*) \\
&\geq -\sum_i \mu_i^* \left( g_i(\mathbf{x}) - g_i(\mathbf{x}^*) \right) \\
&= -\sum_i \mu_i^* g_i(\mathbf{x}) \geq 0
\end{aligned}
\tag{34}
$$

So that $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all feasible $\mathbf{x}$, i.e., $\mathbf{x}^*$ is optimal. $\qquad\square$

Note that if we apply the above theorem to problems with equality constraints, then we must require that the equality constraints are affine.

---

[2]Except the affine constraints CQ and Abadies CQ. Think about why!

# 8 Lagrangian Duality

We will consider one the most classical versions of a general method for optimization, namely relaxation/duality based methods. The basic premise is to take a "difficult" optimization problem, and replace it with something simpler. The "simpler" here will be what we call the Lagrangian relaxation, but first we are going to state an obvious theorem. We are now again working with the abstract problem

$$f^* = \inf f(\mathbf{x}) \tag{35}$$
$$\text{subject to } \mathbf{x} \in S \tag{36}$$

We define a *relaxation* of the problem above to be a problem of the form

$$f_R^* = \inf f_R(\mathbf{x}), \tag{37}$$
$$\text{subject to } \mathbf{x} \in S_R \tag{38}$$

where the function $f_R(\mathbf{x}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in S$, and where $S_R(\mathbf{x}) \subseteq S$. That is, we have replaced the feasible set with larger one, and the objective with something smaller. The following should should then be obvious (but nevertheless, a proof can be found in the book).

**Theorem 6** (The relaxation theorem). *a) $f_R^* \leq f^*$ b) If the relaxed problem in infeasible, then so is (36) c) If $\mathbf{x}_R^*$ is optimal in the relaxed problem, and $\mathbf{x}_R^* \in S$, then $\mathbf{x}_R^*$ is optimal also in (36)*

## 8.1 Lagrangian relaxation

Now we consider a problem of the form

$$f^* = \inf f(\mathbf{x}) \tag{39a}$$
$$\text{subject to } \mathbf{x} \in X, \tag{39b}$$
$$g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m, \tag{39c}$$

where $f$ and $g_i$ are some given functions[3], and $X \subseteq \mathbb{R}^n$ is some subset. The basic idea of Lagrangian relaxation is to replace *constraints* (i.e., things that acceptable solutions *must* satisfy), with a *price* in the objective function for violation. That is, for any $\boldsymbol{\mu} \in \mathbb{R}^m$ we define the Lagrangian relaxation of (the constraints (39c) of) the problem (39) as the problem

---

[3]note that we do not say anything about smoothness!

$$q(\boldsymbol{\mu}) = \inf f(\mathbf{x}) + \sum_{i=1}^{m} \mu_i g_i(\mathbf{x}) \tag{40a}$$

$$\text{subject to } \mathbf{x} \in X. \tag{40b}$$

It is immediate that whenever $\mu_i \geq 0$ for $i = 1, \ldots, m$, the above is indeed a relaxation of (39). The objective function appearing above is important enough to merit its own name, and we call it the *Lagrange function* of (39), and denote it by $L(\mathbf{x}, \boldsymbol{\mu}) := f(\mathbf{x}) + \sum_i \mu_i g_i(\mathbf{x})$.

Now we have, in a fairly simple way, defined a family of relaxations of (39), parametrized by the "price" vector $\boldsymbol{\mu}$. We immediately have the following, very important, result.

**Theorem 7** (Weak duality). *For any $\boldsymbol{\mu} \geq 0$, and any $\mathbf{x}$ feasible in (39) we have*

$$q(\boldsymbol{\mu}) \leq f(\mathbf{x}) \tag{41}$$

*Proof.* This is really just a rephrasing of the statement that the Lagrangian relaxation is, indeed, a relaxation. $\square$

The reason why this result is so important is that allows to get lower bounds on the optimal value $f^*$ of (39).

**Example 6.** *Consider the problem to*

$$f^* = \min x^2,$$
$$\text{subject to } x \geq 1.$$

*We can relax the constraint $x \geq 1$, to get a Lagrangian dual function (note the rewriting of $x \geq 1$ as $1 - x \leq 0$)*

$$q(\mu) = \min x^2 + \mu(1 - x) = \min \left(x - \frac{\mu}{2}\right)^2 - \frac{\mu^2}{4} + \mu.$$

*For each fixed $\mu \geq 0$, the above is an unconstrained minimization problem of a convex function of $x$, so we can actually compute*

$$q(\mu) = \mu - \frac{\mu^2}{4}.$$

*Evaluating at, say, $\mu = 0$, we get $q(\mu) = 0$, and we can conclude that the optimal value $f^*$ must satisfy $f^* \geq 0$. If we instead evaluate at $\mu = 1$, we would be able to conclude $f^* \geq q(1) = 3/4$.*

Having the weak duality theorem at the back of our heads it makes sense to try to find the *best* lower bound of $f^*$, which motivates the following definition.

**Definition 10** (The (Lagrange) dual problem). *The Lagrange dual problem to* (39) *(with respect to the relaxation of* (39c)*) is the problem*

$$q^* = \sup q(\boldsymbol{\mu}), \tag{42a}$$
$$\text{subject to } \boldsymbol{\mu} \geq 0 \tag{42b}$$

In other words, the Lagrange dual problem is just the problem of finding the *best* "price" for defining as tight a relaxation as possible. With these definitions out of the way, we can also the an immediate consequence of the weak duality theorem: for any pair of primal/dual problems, we have $q^* \leq f^*$.

A note on terminology: from now we will refer to (42) as the dual problem, and to (39) as the primal problem. So, for example, the phrase "$\mathbf{x}^*$ is primally optimal" should be taken to mean that $\mathbf{x}^*$ is optimal in (39).

**Example 7.** *Consider again the problem from the previous example. The dual problem is to*

$$q^* = \sup_{\mu \geq 0} \mu - \frac{\mu^2}{4},$$

*and one can easily verify that the maximum is attained at $\mu = 2$, $q^* = q(2) = 1$, which can also be noted to be the optimal value $f^* = q^* = 1$*

Note that in the above example we have $f^* = q^*$. If this holds, we say the pair of primal and dual problems has *no duality gap*. In general, we define the duality gap to be the difference $f^* - q^*$. We also define

**Definition 11.** *We call $\boldsymbol{\mu}^*$ a Lagrange multiplier vector if*

$$f^* = \inf_{x \in X} L(\mathbf{x}, \boldsymbol{\mu}^*) \tag{43}$$

Note that by the above definition, we cannot have a Lagrange multiplier vector unless $f^* = q^*$. Also note the conflict of terminology, Lagrange multiplier is also used to talk about the vector $\boldsymbol{\mu}$ appearing in the Fritz-John and KKT conditions. ALthough the meanings are realted it should be kept in mind that the word Lagrange multiplier mean slightly different things in the context of the KKT conditions and in the context of Lagrange multipliers!

## 8.2 The dual problem

One might ask oneself why one should bother with Lagrangian duality. Really, all we have done is to take an optimization problem (the primal) and replaced it with another problem (the dual). This only makes sense if the dual problem is in some sense "easier" than the first. However we have the following (fantastic!) theorem.

**Theorem 8.** *The dual function $q(\boldsymbol{\mu})$ is concave, and its effective domain $D_q = \{\boldsymbol{\mu} \mid q(\boldsymbol{\mu}) > -\infty\}$ is convex*

*Proof.* This will essentially be for free, since we have defined $q$ as a infimum of things. Spelling out the details, we have, for any pair $\boldsymbol{\mu}, \boldsymbol{\nu}$, and any $\lambda \in (0, 1)$

$$
\begin{aligned}
q((1-\lambda)\boldsymbol{\mu} + \lambda\nu) &= \inf_{\mathbf{x} \in X} f(\mathbf{x}) + (1-\lambda)\sum_i \mu_i g_i(\mathbf{x}) + \lambda \sum_i \nu_i g_i(\mathbf{x}) \\
&= \inf_{\mathbf{x} \in X} \left\{ (1-\lambda)\left[ f(\mathbf{x}) + \sum_i \mu_i g_i(\mathbf{x}) \right] + \lambda \left[ f(\mathbf{x}) + \sum_i \nu_i g_i(\mathbf{x}) \right] \right\} \\
&\geq (1-\lambda) \inf_{bx \in X} \left[ f(\mathbf{x}) + \sum_i \mu_i g_i(\mathbf{x}) \right] + \lambda \inf_{\mathbf{x} \in X} \left[ f(\mathbf{x}) + \sum_i \nu_i g_i(\mathbf{x}) \right], \\
&= \lambda q(\boldsymbol{\mu}) + (1-\lambda)q(\boldsymbol{\nu}),
\end{aligned}
\tag{44}
$$

by simply noting that the infimum of a sum is greater than the sum of the infima. Note that the above also implies that the effective domain is convex! $\qquad\square$

Lets pause for a minute and think about why the above theorem is fantastic news! It really tells us that we can think of the dual problem as a maximization of a concave function over a convex set! In other words, the dual problem is *always* a convex problem! This means that the dual problem is amenable to attack by convex optimization methods (which are the subject of lecture 10), and that the nice results about convex problems we have had so far in this course apply!

## 8.3   Global optimality conditions

If there is no duality gap, i.e., $f^* = q^*$, then it turns out that we can actually use the Lagrangian relaxation to get a *sufficient* condition for optimality.

**Theorem 9.** *Consider the primal/dual pair of vectors $(\mathbf{x}^*, \mu^*)$. Then $\mathbf{x}^*$ is optimal and $\mu^*$ is a Lagrange multiplier vector if and only if*

$$
\mathbf{x}^* \in argmin\, L(\mathbf{x}, \mu^*), \tag{45a}
$$

$$
\boldsymbol{\mu}^* \geq 0, \tag{45b}
$$

$$
\mathbf{x}^* \in X, \quad g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m, \tag{45c}
$$

$$
\mu_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \ldots, m. \tag{45d}
$$

*Remark:* The conditions are, in order, often called Lagrangian optimality, dual feasibility, primal feasibility and complementary slackness. Note also the similarity to the KKT conditions; the only difference is that we have minimization of the Lagrangian instead of stationarity.

*Proof.* Assuming that the above conditions hold we have $f(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*) = f(\mathbf{x})$ for any feasible $x$, since, in order, complementary slackness allows us to add the term $\sum_i \mu_i g_i(\mathbf{x}^*) = 0$, $\mathbf{x}^*$ minimizes $L(\mathbf{x}, \boldsymbol{\mu})$, and $\sum_i \mu_i^* g_i(\mathbf{x}) \leq 0$ for any pair of primal/dual feasible vectors. This shows that $f^* = f(\mathbf{x}^*) = \min L(\mathbf{x}, \boldsymbol{\mu}^*)$, i.e. that $\mathbf{x}^*$ is primally optimal and $\boldsymbol{\mu}^*$ is a Lagrange multiplier vector.

If we instead assume that $\mathbf{x}^*$ is optimal and $\boldsymbol{\mu}^*$ is a Lagrange multiplier vector, the only non-trivial thing to verify is the complementary slackness. But this follows since by definition $f(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\mu}^*) = f(\mathbf{x}^*) + \sum_i \mu_i^* g_i(\mathbf{x}^*)$. Subtracting $f(\mathbf{x}^*)$ from both sides yields $\sum_i \mu_i^* g_i(\mathbf{x}^*) = 0$. But by primal and dual feasibility each term in this is sum is non-positive, so all the terms must be zero individually. $\square$

We can in fact formulate the above conditions in more compact way, as what is called saddle-point optimality conditions, meaning that the pair $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ simultaneously maximizes $L$ (over $\boldsymbol{\mu}$) and minimizes $L$ (over $\mathbf{x}$)

**Theorem 10.** $\mathbf{x}^*$ *is primally optimal and $\boldsymbol{\mu}^*$ is a Lagrange multiplier if and only if $\mathbf{x}^* \in X$, $\boldsymbol{\mu}^* \geq 0$ and*

$$L(\mathbf{x}^*, \boldsymbol{\mu}) \leq L(\mathbf{x}^*, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*), \quad (\mathbf{x}, \boldsymbol{\mu}) \in X \times \mathbb{R}_+^m \qquad (46)$$

*Proof.* The first inequality is equivalent to requiring that $\boldsymbol{\mu}^*$ maximizes $L(\mathbf{x}^*, \boldsymbol{\mu})$ over $\mathbb{R}_+^m$. Using optimality conditions convex set $\mathbb{R}_+^m$ we have that this is equivalent[4] to $\nabla_{\boldsymbol{\mu}} L(\mathbf{x}, \boldsymbol{\mu}) \in N_{\mathbb{R}_+^m}(\boldsymbol{\mu}^*)$, where $\nabla_{\boldsymbol{\mu}}$ denotes the gradient with respect to the $\boldsymbol{\mu}$-variables of $L$. But since $\nabla_{\mu}(L(\mathbf{x}^*, \mu^*)) = \mathbf{g}(\mathbf{x}^*)$, we get $\mathbf{g}(\mathbf{x}^*) \in N_{\mathbb{R}_+^m}(\boldsymbol{\mu}^*)$, i.e., that $\mu_i g_i(\mathbf{x}^*) = 0$ and $g_i(\mathbf{x}^*) \leq 0$, for all $i = 1, \ldots, m$.

The second inequlity is the statement $\mathbf{x}^* \in \text{argmin}_{x \in X} L(\mathbf{x}, \boldsymbol{\mu}^*)$. Hence the conditions of the theorem are equivalent to the conditions of the previous theorem. $\square$

## 8.4 Strong Lagrangian Duality

A natural question to ask is now under what conditions one can guarantee that $q^* = f^*$, i.e., that the primal and dual optimal values coincide. Not surprisingly, it turns out that what we need to require is convexity. However we also need to assume some regularity, which will here be a variant of the Slater CQ. That is, we now require that $X$ is a convex set, $g_i$ are convex for $i = 1 \ldots, m$ and there is some point $\mathbf{x} \in X$ such that $g_i(\mathbf{x}) \leq 0$, $i = 1, \ldots, m$.

**Theorem 11.** *Assume that the problem (39) satisfies the Slater CQ, and that $f^* \geq -\infty$. Then strong Lagrangian duality holds, and there exists at least one Lagrange multiplier vector.*

*Proof.* Consider the set

$$A := \{(z_1, z_2, \ldots, z_n, w) \in \mathbb{R}^{m+1} \mid \exists \mathbf{x} \in X : f(\mathbf{x}) \leq w, g_i(\mathbf{x}) \leq z_i, i = 1, \ldots, m\} \qquad (47)$$

---

[4]Since $L(\mathbf{x}, \boldsymbol{\mu})$ is convex in the $\boldsymbol{\mu}$-variable

Since all the functions above are convex, it follows that $A$ is convex. Further, since $-\infty < f^* < \infty$, we have that $(0^m, f^*)$ lies on the boundary of $A$, by the very definition of $f^*$. This allows us to find a supporting hyperplane of $A$ passing through the point $(0^m, f^*)$, i.e., a vector $(\boldsymbol{\mu}, \beta)$ such that

$$\beta f^* \leq \beta w + \sum_{i=1}^{m} \mu_i z_i, \quad \forall (z_1, z_2, \ldots, z_m, w) \in A. \tag{48}$$

Since by how $A$ is defined, we may take $w \to \infty$ and remain in $A$, which shows that for the above to hold we must have $\beta \geq 0$. The same argument applied to letting $z_i \to \infty$ shows that $\mu_i \geq 0$ for $i = 1, \ldots, m$. Further, using the point $\bar{\mathbf{x}}$ given by the Slater CQ shows that there are points in $A$ for which $z_i < 0$ for all $i = 1, \ldots, m$. This shows that for (48) to hold we must have $\beta > 0$, and we may take without loss of generality $\beta = 1$. Finally, since $(g_1()\mathbf{x}, g_2(\mathbf{x}), \ldots, g_m(\mathbf{x}), f(\mathbf{x})) \in A$, for any $\mathbf{x} \in X$ we thus conclude that

$$f^* \leq f(\mathbf{x}) + \sum_{i=1}^{m} \mu_i g_i(\mathbf{x}), \quad \forall \mathbf{x} \in X \tag{49}$$

Taking the infimum over $\mathbf{x} \in X$ yields $f^* \leq q(\boldsymbol{\mu}) \leq q^*$. But by the weak duality theorem we always have $q^* \leq f^*$. Thus, there is no duality gap, and the vector $\boldsymbol{\mu}$ is a Lagrange multiplier vector. $\qquad\square$

We finally note what happens if we assume that $f$, and $g_i$ are also $C^1$, $X = \mathbb{R}^m$, and the problem (39) satisfies Slaters CQ and has some optimal solution $\mathbf{x}^*$. The above theorem then gives a Lagrange multiplier vector $\boldsymbol{\mu}^*$. Then the pair $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ are a pair of a primally optimal solution and a Lagrange multiplier vector, so they satsify the system (45). But the Lagrange function $L(\mathbf{x}, \boldsymbol{\mu})$ is convex in $\mathbf{x}$ for any $\boldsymbol{\mu}$, so the condition that $\mathbf{x}^* \in \text{argmin}_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*)$ can be replaced by the neccessary and sufficient condition that $\nabla_{\mathbf{x}}^* L(\mathbf{x}, \boldsymbol{\mu}^*) = 0$. But $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\mu}^*) = \nabla f(\mathbf{x}) + \sum_{i=1}^{m} \mu_i^* \nabla g_i(\mathbf{x}^*)$. Thus in this case the golbal optimality conditions just reduces to the Karush-Kuhn-Tucker conditions!