# FOUNDATIONS OF OPTIMIZATION

## Basics

- *Optimization problems*
  - An optimization problem is
    $$\text{minimise } f(\boldsymbol{x}) \text{ subject to } \boldsymbol{x} \in \mathcal{C}$$
    $f$ is the *objective* (real) $\mathcal{C}$ is the constraint set/feasible set/search space.
  - $\boldsymbol{x}^*$ is an *optimal solution* (*global minimizer*) if and only if
    $$f(\boldsymbol{x}^*) \le f(\boldsymbol{x}) \qquad \forall \boldsymbol{x} \in \mathcal{C}$$
  - Maximizing $f(\boldsymbol{x})$ is equivalent to minimizing $-f(\boldsymbol{x})$.
  - We consider problems in the following form
    $$\begin{aligned}
    \text{minimize} \quad & f(\boldsymbol{x}) \\
    \text{subject to} \quad & h_i(\boldsymbol{x}) = 0 \qquad \forall\ 1 \le i \le m \\
    & g_i(\boldsymbol{x}) \le 0 \qquad \forall\ m \le i \le r \\
    & \boldsymbol{x} \in \mathbb{R}^n
    \end{aligned}$$
  - We consider the following subsets of the problem
    - In *linear programming*, all functions are linear.
    - In convex programming, the $f$ and $g$ are convex, and the $h$ are linear.
  - If $\mathcal{C}$ is the feasible set of a problem, a point $\boldsymbol{x} \in \mathcal{C}$ is a *local minimum* if there exists a neighborhood $N_r(\boldsymbol{x})$ such that $f(\boldsymbol{x}) \le f(\boldsymbol{y})\ \forall \boldsymbol{y} \in \mathcal{C} \cap N_r(\boldsymbol{x})$. It is an *unconstrained local minimum* if $f(\boldsymbol{x}) \le f(\boldsymbol{y})\ \forall \boldsymbol{y} \in N_r(\boldsymbol{x})$. (Strict equivalents exist).

- *Topology*
  - An *open ball* around a point $\boldsymbol{x} \in \mathbb{R}^n$ with radius $r > 0$ is the set $N_r(\boldsymbol{x}) = \left\{ \boldsymbol{y} \in \mathbb{R}^n : \left\| \boldsymbol{x} - \boldsymbol{y} \right\| < r \right\}$, where $\left\| \boldsymbol{x} \right\| = \sqrt{\sum x_i^2}$ .
  - A point $\boldsymbol{x} \in \mathcal{E} \subset \mathbb{R}^n$ is an *interior point* if there exists an open ball such that $N_r(\boldsymbol{x}) \subset \mathcal{E}$. A set $\mathcal{E} \subset \mathbb{R}^n$ is *open* if $\mathcal{E} = \operatorname{int} \mathcal{E}$ .

○ A point $\boldsymbol{x} \in \mathcal{E} \subset \mathbb{R}^n$ is a *closure point* if, for every open ball $N_r(\boldsymbol{x})$, there exists $\boldsymbol{y} \in \mathcal{E}$ with $\boldsymbol{y} \in N_r(\boldsymbol{x})$. A set $\mathcal{E} \subset \mathbb{R}^n$ is *closed* if $\mathcal{E} = \text{cl } \mathcal{E}$.

○ The set of reals is both closed and open.

○ ***Theorems***:

- The union of open sets is open. The intersection of a *finite* number of open sets is open.

- The intersection of closed sets is closed. The union of a *finite* number of closed sets is closed.

- **Analysis**

○ A sequence of vectors $\left\{\boldsymbol{x}_n\right\} \subset \mathbb{R}^n$ converges to a limit $\boldsymbol{x} \in \mathbb{R}^n$ if $\lim_{k\to\infty} \left\| \boldsymbol{x} - \boldsymbol{x}_k \right\| = 0$, and we say that $\boldsymbol{x}_k \to \boldsymbol{x}$.

○ A set $\mathcal{E} \subset \mathbb{R}^n$ is (sequentially) compact if, given a sequence $\left\{\boldsymbol{x}_k\right\} \subset \mathcal{E}$, there is a subsequence $\left\{\boldsymbol{x}_{k_i}\right\}$ converging to an element $\boldsymbol{x} \in \mathcal{E}$.

- ***Theorem*** (Heine-Borel): A set $\mathcal{E} \subset \mathbb{R}^n$ is compact if and only if it is closed and bounded.

- ***Theorem***: A closed subset of a compact set is compact.

- ***Theorem***: Suppose $\left\{\mathcal{E}_n\right\}$ are a sequence of non-empty, compact sets that are nested (ie: $\mathcal{E}_{n+1} \subset \mathcal{E}_n$) – then their intersection is non-empty.

○ A real-valued function $f$ defined on a domain $\mathcal{X} \subset \mathbb{R}^n$ is continuous at the point $x \in \mathcal{X}$ if, for every sequence $\left\{\boldsymbol{x}_k\right\} \subset \mathcal{X}$ with $\boldsymbol{x}_k \to \boldsymbol{x}$, $\lim_{k\to\infty} f(\boldsymbol{x}_k) = f(\boldsymbol{x})$. $f$ is *continuous* if it is continuous at all points in $\mathcal{X}$.

○ A function $f$ is *coercive* over a set $\mathcal{C} \subset \mathbb{R}^n$ if, for every sequence $\left\{\boldsymbol{x}_k\right\} \subset \mathcal{C}$ with $\left\|\boldsymbol{x}_k\right\| \to \infty$, we have $\lim_{k\to\infty} f(\boldsymbol{x}_k) = \infty$.

○ The *inverse image* of the set $\mathcal{A} \subset \mathbb{R}$ is defined by $f^{-1}(\mathcal{A}) = \left\{\boldsymbol{x} \in \mathcal{X} : f(\boldsymbol{x}) \in \mathcal{A}\right\}$.

- ***Theorem***: If $f$ is continuous and $\mathcal{X}$ is open ^closed and $\mathcal{A}$ is open ^closed, then $f^{-1}(\mathcal{A})$ is also open ^closed. This is the standard way to prove that a set is open/closed.

- o **Definition**: A function is *convex* if $f\left(\lambda x_1 + (1-\lambda)x_2\right) \le \lambda f(x_1) + (1-\lambda)f(x_2)$. It is *strictly convex* if the inequality is strict for $x_1 \ne x_2$. We say $f$ is convex over $\mathcal{X} = \operatorname{dom} f$ if it is convex when restricted to $\mathcal{X}$. $f$ is (strictly )concave if $-f$ is (strictly) convex.

- o **Definition**: If $f$ is convex with a convex domain $\mathcal{X}$, we define the *extended-value extension* $\tilde{f} : \mathbb{R}^n \to \mathbb{R} \cup \{\infty\}$ by

$$\tilde{f}(x) = \begin{cases} f(x) & \text{if } x \in \mathcal{X} \\ \infty & \text{otherwise} \end{cases}$$

   and we let

$$\operatorname{dom} \tilde{f} = \left\{ x \in \mathbb{R}^n : \tilde{f}(x) < \infty \right\}$$

   An extended-value function is convex if

   - ▪ Its domain is convex.

   - ▪ The standard convexity property holds.

- o Given $\mathcal{C} \subset \mathbb{R}^n$, the *indicator function* $I_{\mathcal{C}} : \mathbb{R}^n \to \mathbb{R} \cup \{\infty\}$ as

$$I_{\mathcal{C}}(x) = \begin{cases} 0 & \text{if } x \in \mathcal{C} \\ \infty & \text{otherwise} \end{cases}$$

   If $\mathcal{C}$ is a convex set, then $I_{\mathcal{C}}$ is a convex function.

- o **Theorem**: If $f$ is convex over a convex set $\mathcal{C} \subset \mathbb{R}^n$, then every sublevel set $\left\{ x \in \mathcal{C} : f(x) \le \gamma \right\}$ is a convex subset of $\mathbb{R}^n$. The converse is *not* true (eg: $\log x$ on $(0, \infty)$). However, we define...

- o ...**Definition**: A extended real valued function $f : \mathbb{R}^n \to \mathbb{R} \cup \{\infty\}$ is *quasiconvex* if, every one of its sublevel sets (ie: for every $\gamma \in \mathbb{R}$) is convex.

- • *Calculus*

  - o A function $f : \mathcal{X} \to \mathbb{R}$ with $\mathcal{X} \subset \mathbb{R}^n$ is *differentiable* at $x \in \operatorname{int} \mathcal{X}$ if there exists a vector $\nabla f(x) \in \mathbb{R}^n$, known as the *gradient*, such that

$$\lim_{d \to 0} \frac{f(x+d) - f(x) - \nabla f(x) \cdot d}{\|d\|} = 0$$

   And

$$\nabla f(\boldsymbol{x}) = \left[\frac{\partial f(\boldsymbol{x})}{\partial x_1}, \cdots, \frac{\partial f(\boldsymbol{x})}{\partial x_n}\right]^T \in \mathbb{R}^n \qquad \frac{\partial f(\boldsymbol{x})}{\partial x_i} = \lim_{h \to 0} \frac{f(\boldsymbol{x} + h\boldsymbol{e}_i) - f(\boldsymbol{x})}{h}$$

$f$ is *differentiable* over an open set $\mathcal{U} \in \mathcal{X}$ if it is differentiable at every point in the set. If, in addition, the components of the gradient are continuous over $\mathcal{U}$, then $f$ is *continuously differentiable* over $\mathcal{U}$.

o   If, for a point $\boldsymbol{x} \in \operatorname{int} \mathcal{X}$, each component of the gradient is differentiable, we say $f$ is *twice differentiable* at $\boldsymbol{x}$, and we define the *Hessian Matrix* $\nabla^2 f(\boldsymbol{x}) \in \mathbb{R}^{n \times n}$ by

$$\nabla^2 f(\boldsymbol{x}) = \left[\frac{\partial^2 f(\boldsymbol{x})}{\partial x_i \partial x_j}\right]_{ij}$$

If $f$ is twice continuously differentiable in a neighborhood of $\boldsymbol{x}$, then the Hessian is symmetric.

o   Suppose at $f$ is twice continuously differentiable over a neighborhood $N_r(\boldsymbol{x})$, then for all $\boldsymbol{d} \in N_r(\boldsymbol{0})$

$$f(\boldsymbol{x} + \boldsymbol{d}) = f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^T \boldsymbol{d} + \tfrac{1}{2} \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}) \boldsymbol{d} + o\left(\left\|\boldsymbol{d}\right\|^2\right)$$

(Formally, this means that for every $C > 0$, there exists a neighborhood around $\boldsymbol{0}$ such that the estimate of $f(\boldsymbol{x} + \boldsymbol{d})$ differs from the real value by no more than $C \left\|\boldsymbol{d}\right\|^2$.

o   Consider a vector-valued function $\boldsymbol{F} : \mathcal{X} \to \mathbb{R}^m, \mathcal{X} \subset \mathbb{R}^n$ and a point $\boldsymbol{x} \in \operatorname{int} \mathcal{X}$. We define the *gradient* to be the matrix $\nabla \boldsymbol{F}(\boldsymbol{x}) \in \mathbb{R}^{n \times m}$ with

$$\nabla \boldsymbol{F}(\boldsymbol{x}) = \left[\nabla F_1(\boldsymbol{x}), \cdots, \nabla F_m(\boldsymbol{x})\right] \qquad \nabla F(\boldsymbol{x})_{ij} = \frac{\partial F_j(\boldsymbol{x})}{\partial x_i}$$

o   The chain rule states that for interior points, if $\boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{g}(\boldsymbol{f}(\boldsymbol{x}))$, then

$$\nabla \boldsymbol{h}(\boldsymbol{x}) = \nabla \boldsymbol{f}(\boldsymbol{x}) \nabla \boldsymbol{g}(\boldsymbol{f}(\boldsymbol{x}))$$

- *Linear algebra – Kernels and Images*
  o   Consider a matrix $A \in \mathbb{R}^{m \times n}$. Then
    ▪   $\ker A = \left\{\boldsymbol{x} \in \mathbb{R}^n : A\boldsymbol{x} = \boldsymbol{0}\right\}$
    ▪   $\operatorname{im} A = \left\{\boldsymbol{y} \in \mathbb{R}^m : \boldsymbol{y} = A\boldsymbol{x}, \boldsymbol{x} \in \mathbb{R}^n\right\}$
  o   Given a set $\mathcal{S} \in \mathbb{R}^n$, $\mathcal{S}^\perp = \left\{\boldsymbol{x} \in \mathbb{R}^n : \boldsymbol{x} \cdot \boldsymbol{y} = 0 \; \forall \boldsymbol{y} \in \mathcal{S}\right\}$

Daniel Guetta

- o **Lemma**: $\operatorname{im} A = \left[\ker(A^\top)\right]^\perp$. In other words, given $\boldsymbol{z} \in \mathbb{R}^m$,

$$\boldsymbol{z} = A\boldsymbol{x} \text{ for some } \boldsymbol{x} \in \mathbb{R}^n \Leftrightarrow \boldsymbol{z} \cdot \boldsymbol{y} = 0 \ \forall \boldsymbol{y} \text{ with } A^\top \boldsymbol{y} = 0$$

- **Sets, etc...**

  - o **Affine sets**

    - ▪ **Definition**: A set $\mathcal{C} \subset \mathbb{R}^n$ is *affine* if, for all points $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{C}$ and a scalar $\lambda \in \mathbb{R}$, $\lambda \boldsymbol{x}_1 + (1-\lambda)\boldsymbol{x}_2 \in \mathcal{C}$.

      - • **Example**: The empty space, a line and any subspace are affine. Similarly, $\mathcal{C} = \left\{\boldsymbol{x} \in \mathbb{R}^n : A\boldsymbol{x} = \boldsymbol{b}\right\}$ is affine.

    - ▪ **Definition**: Given a set of points $\mathcal{X} \subset \mathbb{R}^n$, the *affine hull* aff $\mathcal{X}$ is the set of points $\lambda_1 \boldsymbol{x}_1 + \cdots + \lambda_k \boldsymbol{x}_k$, where $k \geq 1$, $\{\boldsymbol{x}_i\} \subset \mathcal{X}$ and $\lambda_+ = 1$. The affine hull is affine and is the smallest affine set containing $\mathcal{X}$.

  - o **Convex sets**

    - ▪ **Definition**: The set $\mathcal{C}$ is *convex* if, for all points $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathcal{C}$ and a scalar $\lambda \in (0,1)$, $\lambda \boldsymbol{x}_1 + (1-\lambda)\boldsymbol{x}_2 \in \mathcal{C}$. Clearly, affine sets are also convex.

    - ▪ **Definition**: Given a set of points $\mathcal{X} \subset \mathbb{R}^n$, the *convex hull* conv $\mathcal{X}$ is the set of points $\lambda_1 \boldsymbol{x}_1 + \cdots + \lambda_k \boldsymbol{x}_k$, where $k \geq 1$, $\{\boldsymbol{x}_i\} \subset \mathcal{X}$, $\lambda_i \geq 0$ and $\lambda_+ = 1$.

    - ▪ **Theorem (scalar multiplication)**: if $\mathcal{C} \subset \mathbb{R}^n$ is convex and $\alpha \in \mathbb{R}$, then $\alpha \mathcal{C} = \left\{\alpha \boldsymbol{x} : \boldsymbol{x} \in \mathcal{C}\right\}$ is convex.

    - ▪ **Theorem (vector sum)**: If $\mathcal{C}, \mathcal{D} \subset \mathbb{R}^n$ are convex sets, then the set $\mathcal{C} + \mathcal{D} = \left\{\boldsymbol{x} + \boldsymbol{y} : \boldsymbol{x} \in \mathcal{C}, \boldsymbol{y} \in \mathcal{D}\right\}$ is also convex.

    - ▪ **Theorem (affine transformations)**: If $\mathcal{C} \subset \mathbb{R}^n$ is a convex set, $A \in \mathbb{R}^{m \times n}$ is a matrix and $\boldsymbol{b} \in \mathbb{R}^m$ is a vector, then the set $\left\{A\boldsymbol{x} + \boldsymbol{b} : \boldsymbol{x} \in \mathcal{C}\right\}$ is a convex subset of $\mathbb{R}^m$.

    - ▪ **Theorem**: If $\mathcal{K}$ is an arbitrary collection of convex sets, then the intersection $\bigcap_{\mathcal{C} \in \mathcal{K}} \mathcal{C}$ is also convex.

  - o **Miscellaneous definitions**

- ▪ **Definition**: A set $\mathcal{P}$ is a *polyhedron* if it is of the form $\mathcal{P} = \left\{ x \in \mathbb{R} : A x \leq b \right\}$, for some $A \in \mathbb{R}^{m \times n}$ and a vector $b \in \mathbb{R}^n$. Polyhedra are convex.

- ▪ **Definition**: A set $\mathcal{C} \subset \mathbb{R}^n$ is a *cone* if for all $x \in \mathcal{C}, \lambda \geq 0$, $\lambda x \in \mathcal{C}$. If the cone is convex (ie, if $\lambda_1 x_1 + \lambda_2 x_2 \in \mathcal{C}$ for all $x_1, x_2 \in \mathcal{C}$, $\lambda_1, \lambda_2 \geq 0$), it is a *convex cone*.

  - ● **Example**: The *conic hull* of the points $\mathcal{X} \subset \mathbb{R}^n$ consists of the points $\lambda_1 x_1 + \cdots + \lambda_k x_k$, where $k \geq 1$, $\{x_i\} \subset \mathcal{X}$ and $\boldsymbol{\lambda} \geq \boldsymbol{0}$. It is a convex cone.

  - ● **Example**: Given a norm on $\mathbb{R}^n$, the *norm cone*
    $$\left\{ (x, t) \in \mathbb{R}^{n+1} : \left\| x \right\| \leq t \right\}$$
    is a convex cone in $\mathbb{R}^{n+1}$. When the norm in question is the 2-norm, the cone is called the *second-order cone*.

  - ● **Example**: The set of positive semidefinite matrices
    $$S_+^n = \left\{ X \in \mathbb{R}^{n \times n} : X^\top = X, X \succeq 0 \right\}$$
    is a convex cone in $\mathbb{R}^{n \times n}$ called the *positive semidefinite cone*.

- ○ **Hyperplanes and halfspaces**

  - ▪ **Definition**: A hyperplane is a set of the form $\left\{ x \in \mathbb{R}^n : a \cdot x = b \right\}$, where $a \in \mathbb{R}^n \setminus \{\boldsymbol{0}\}$ is called the *normal vector*. Hyperplanes are affine and therefore convex.

  - ▪ **Definition**: A halfspace is a set of the form $\left\{ x \in \mathbb{R}^n : a \cdot x \leq b \right\}$. Halfspaces are convex but not affine.

- ○ **Definition**: A *norm* is a real valued function $\left\| \cdot \right\|$ on $\mathbb{R}^n$ such that

  - ▪ $\left\| x \right\| = 0 \Leftrightarrow x = \boldsymbol{0}$

  - ▪ For all $x \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}$, $\left\| \lambda x \right\| = \left| \lambda \right| \left\| x \right\|$

  - ▪ For all $x_1, x_2 \in \mathbb{R}^n$, $\left\| x_1 + x_2 \right\| \leq \left\| x_1 \right\| + \left\| x_2 \right\|$

  Examples of norms:

- The L2-norm (Euclidean norm): $\left\|x\right\|_2 = \sqrt{\sum_{i=1}^{n} x_i^2} = \sqrt{x \cdot x}$

- $\left\|x\right\|_\Gamma = \sqrt{x^\top \Gamma x}$ (when $\Gamma \succ 0$ and symmetric)

- The $p$-norm: $\left\|x\right\| = \left(\sum_{i=1}^{n} |x_i|^p\right)^{1/p}$ for $p \geq 1$

- $\left\|x\right\|_\infty = \max\left\{|x_1|, \cdots, |x_n|\right\}$

Given a norm, the (closed) ball with centre $x_0$ and radius $r \geq 0$ is $\left\{x \in \mathbb{R}^n : \left\|x - x_0\right\| \leq r\right\}$, and it is convex.

o **Theorem (projection)**: Let $\mathcal{C} \subset \mathbb{R}^n$ be a *closed* and non-empty convex set, and consider the Euclidean norm. Fix the vector $x \in \mathbb{R}^n$. Consider the problem

$$\begin{aligned}\min \quad & \left\|z - x\right\| \\ \text{s.t.} \quad & z \in \mathcal{C} \subset \mathbb{R}^n\end{aligned}$$

For every $x \in \mathbb{R}^n$, the problem has a unique global minimum $x^*$ called the *projection* of $x$ onto $\mathcal{C}$. A vector $x' \in \mathcal{C}$ is equal to $x^*$ if and only if

$$(x - x') \cdot (z - x') \leq 0 \qquad \forall z \in \mathcal{C}$$

Geometrically, the angle between $x' \to x$ and $x' \to z$ must be larger than $90°$ for all points in the set:



**Proof**: Existence follows from the fact $\left\|z - x\right\|$ is coercive and $\mathcal{C}$ is closed. Uniqueness follows because minimizing $\left\|z - x\right\|$ is equivalent to minimizing $\left\|z - x\right\|^2 = z \cdot z - 2z \cdot x + x \cdot x$, which is strictly convex.

Now, consider that $\nabla f(x^*) = 2(x^* - x)$. By necessary and sufficient conditions for convex optimization problems (derived later), the condition in the theorem must hold.

***Application***: Suppose we want to approximate $f(\boldsymbol{x})$ over a set of points $\left\{\boldsymbol{x}_1, \cdots, \boldsymbol{x}_m\right\}$ using $g(\boldsymbol{x}) = \sum_{\ell=1}^{k} r_\ell \phi_\ell(\boldsymbol{x})$, where the $\phi_i$ are basis functions and $\boldsymbol{r}$ is a vector of weights. One way to do this is to solve the problem

$$\begin{aligned} \min \quad & \sum_{i=1}^{m}\left[f(\boldsymbol{x}_i) - g(\boldsymbol{x}_i)\right]^2 \\ \text{s.t.} \quad & g(\cdot) \text{ is a linear combination of } \left\{\phi_\ell(\cdot)\right\} \end{aligned}$$

Consider the matrix $\Phi_{i,\ell} = \phi_\ell(\boldsymbol{x}_i)$ and the vector $\boldsymbol{y}$, $y_i = f(\boldsymbol{x}_i)$. This problem is equivalent to

$$\begin{aligned} \min \quad & \left\|\boldsymbol{y} - \boldsymbol{z}\right\| \\ \text{s.t.} \quad & \boldsymbol{z} \in \left\{\Phi\boldsymbol{r} : \boldsymbol{r} \in \mathbb{R}^k\right\} \end{aligned}$$

This is a projection problem, and so a unique optimizer exists.

# Existence of solutions

- ***Theorem*** – ***Sufficient Conditions*** (Weierstrass): Consider the problem $\min f(\boldsymbol{x})$ s.t. $x \in \mathcal{C} \subset \mathbb{R}^n$. Then if
  - $\mathcal{C}$ is non-empty
  - $f$ is lower semicontinuous over $\mathcal{C}$

  and one of the following conditions hold:
  1. $\mathcal{C}$ is compact
  2. $\mathcal{C}$ is closed, and $f$ is coercive
  3. There exists a scalar $\gamma$ such that the level set $\mathcal{C}(\gamma) = \{\boldsymbol{x} \in \mathcal{C} : f(\boldsymbol{x}) \le \gamma\}$ is nonempty and compact.

  then the set of optimal minimizing solutions of $f$ is non-empty and compact.

  ***Proof***:
  - $\boxed{1 \to 3}$: define $f^* = \inf_{x \in \mathcal{C}} f(x) \in \mathbb{R} \cup \{-\infty\}$ (this always exists). Then, given $\gamma > f^*$, the level set $\{\boldsymbol{x} \in \mathcal{C} : f(\boldsymbol{x}) \le \gamma\}$ must be non-empty. By the continuity of f, it is also closed. Thus, since $\mathcal{C}$ is compact, so is this set.
  - $\boxed{2 \to 3}$: Define $\mathcal{C}(\gamma) = \{\boldsymbol{x} \in \mathcal{C} : f(\boldsymbol{x}) \le \gamma\}$. Since f is coercive, $\mathcal{C}(\gamma)$ is non-empty and bounded for any $\gamma > f^*$. Furthermore, since the domain of f (ie: $\mathcal{C}$) is

closed, $(-\infty, \gamma]$ is closed and $\mathcal{C}(\gamma) = f^{-1}\left((-\infty, \gamma]\right)$, $\mathcal{C}(\gamma)$ is also closed. Thus, $\mathcal{C}(\gamma)$ is compact.

o $\boxed{3}$: Given a sequence of real numbers $\{\gamma_k\}$ with $\gamma_k \downarrow f^*$, the set of optimal solutions is

$$\mathcal{X}^* = \bigcap_{k=1}^{\infty} \mathcal{C}\left(\gamma_k\right)$$

By a theorem stated above, the intersection of these nested, non-empty compact sets is also non-empty, and compact.

Note that the lower semicontinuity of $f$ was only used in proving that $\mathcal{C}(\gamma)$ is *closed*.

- ***Example***: consider $\min \frac{1}{2}\boldsymbol{x}^\top \Gamma \boldsymbol{x} - \boldsymbol{b}^\top \boldsymbol{x}, \boldsymbol{x} \in \mathbb{R}$. If $\lambda$ is the smallest eigenvalue of $\Gamma$, we have $\frac{1}{2}\boldsymbol{x}^\top \Gamma \boldsymbol{x} - \boldsymbol{b}^\top \boldsymbol{x} \geq \frac{\lambda}{2}\left\|\boldsymbol{x}\right\|^2 - \left\|\boldsymbol{b}\right\|\left\|\boldsymbol{x}\right\|$, which is coercive if $\lambda > 0$. Thus, a solution exists if $\Gamma \succ 0$ (ie: it is positive definite).

## Unconstrained local optimality

- In this section, we consider the problem $\min f(\boldsymbol{x})$ s.t. $x \in \mathcal{C} \subset \mathbb{R}^n$, but we focus on minima that lie in the interior of $\mathcal{C}$. In other words, unconstrained minima.
- *Local optimality*
    - ***Theorem – necessary conditions***: let $\boldsymbol{x}^* \in \text{int}\,\mathcal{C}$ be an unconstrained local minimum. Then
        - If $f$ is continuously differentiable in a neighborhood of $\boldsymbol{x}^*$, then
        $$\nabla f(\boldsymbol{x}^*) = \boldsymbol{0}$$
        - If $f$ is twice continuously differentiable in a neighborhood of $\boldsymbol{x}^*$, then
        $$\nabla^2 f(\boldsymbol{x}^*) \succeq 0 \quad \text{[Positive semidefinite]}$$

        Geometrically, we simply require that the tangent at the said point be horizontal, and also that the tangent *underestimate* the curve in the neighborhood of the stationary point (in other words, the curve should lie above the tangent).

        ***Proof***:
        - $\boxed{\text{First order}}$: fix $\boldsymbol{d} \in \mathbb{R}^n \setminus \{\boldsymbol{0}\}$. By the definition of the gradient,

$$\lim_{\alpha \to 0} \frac{f(\boldsymbol{x}^* + \alpha\boldsymbol{d}) - f(\boldsymbol{x}^*) - \alpha\nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d}}{\alpha \|\boldsymbol{d}\|} = 0$$

$$\lim_{\alpha \to 0} \frac{f(\boldsymbol{x}^* + \alpha\boldsymbol{d}) - f(\boldsymbol{x}^*)}{\alpha} = \nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d}$$

If $\boldsymbol{x}^*$ is a global optimum, the LHS must be positive for small enough $\alpha$. Thus, $\nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d} \geq 0$. Since $\boldsymbol{d}$ is arbitrary, we must have $\nabla f(\boldsymbol{x}^*) = 0$.

- ▪ $\boxed{\text{Second order}}$ : fix $\boldsymbol{d} \in \mathbb{R}^n$. For sufficiently small $\alpha$:

$$f(\boldsymbol{x}^* + \alpha\boldsymbol{d}) - f(\boldsymbol{x}^*) = \alpha\nabla f(\boldsymbol{x}) \cdot \boldsymbol{d} + \tfrac{1}{2}\alpha^2 \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} + o(\alpha^2)$$
$$= \tfrac{1}{2}\alpha^2 \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} + o(\alpha^2)$$

If $\boldsymbol{x}^*$ is a global optimum, the LHS must be positive for small enough $\alpha$, and so

$$\tfrac{1}{2}\alpha^2 \boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} + o(\alpha^2) \geq 0$$

$$\tfrac{1}{2}\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} + \frac{o(\alpha^2)}{\alpha^2} \geq 0$$

Taking limits as $\alpha \to 0$:

$$\boldsymbol{d}^T \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} \geq 0$$

Since $\boldsymbol{d}$ is arbitrary, this leads to our result.

- ○ ***Theorem – sufficient conditions***: Consider a point $\boldsymbol{x}^* \in \operatorname{int}\mathcal{C}$. If $f$ is twice continuously differentiable in a neighborhood of $\boldsymbol{x}^*$, and

$$\nabla f(\boldsymbol{x}^*) = 0 \qquad\qquad \nabla^2 f(\boldsymbol{x}^*) \succ 0$$

Then $\boldsymbol{x}^*$ is a *strict unconstrained local minimum*. The geometric interpretation is as above – the only difference is that we now require a positive definite instead of a positive *semi*definite matrix.

***Proof***: Let $\lambda > 0$ be the smallest eigenvalue of $\nabla^2 f(\boldsymbol{x}^*)$, and let $\boldsymbol{d} \in N_r(\boldsymbol{0}) \setminus \{\boldsymbol{0}\}$

$$f(\boldsymbol{x}^* + \boldsymbol{d}) - f(\boldsymbol{x}^*) = \nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d} + \tfrac{1}{2}\boldsymbol{d}^T \nabla f(\boldsymbol{x}^*)\boldsymbol{d} + o\left(\|\boldsymbol{d}\|^2\right)$$
$$= \tfrac{1}{2}\boldsymbol{d}^T \nabla f(\boldsymbol{x}^*)\boldsymbol{d} + o\left(\|\boldsymbol{d}\|^2\right)$$
$$\geq \tfrac{1}{2}\lambda\|\boldsymbol{d}\|^2 + o\left(\|\boldsymbol{d}\|^2\right)$$
$$= \left(\frac{\lambda}{2} + \frac{o\left(\|\boldsymbol{d}\|^2\right)}{\|\boldsymbol{d}^2\|}\right)\|\boldsymbol{d}\|^2$$

Now, for any $\gamma \in (0, \lambda)$, there exists $\varepsilon \in (0, r]$ such that

$$\frac{\lambda}{2} + \frac{o(|| \, \boldsymbol{d}^2 \, ||)}{|| \, \boldsymbol{d} \, ||^2} \geq \frac{\gamma}{2} \qquad \forall \boldsymbol{d} \text{ with } || \, \boldsymbol{d} \, || < \varepsilon$$

And this means that

$$f(\boldsymbol{x}^* + \boldsymbol{d}) \geq f(\boldsymbol{x}^*) + \frac{\gamma}{2} || \, \boldsymbol{d} \, ||^2 > f(\boldsymbol{x}^*)$$

- *Using the necessary conditions*
    - Verify there is a global minimum (using the existence theorem).
    - Find the set of possible unconstrained local minima using $\nabla f(\boldsymbol{x}) = \boldsymbol{0}$.
    - Compare these points with all points on the boundary $\mathcal{C} \setminus \text{int}\, \mathcal{C}$.
    - ***Example***: Consider $\min_{\boldsymbol{x} \in \mathbb{R}^n} \frac{1}{2} \boldsymbol{x}^\top \Gamma \boldsymbol{x} - \boldsymbol{b}^\top \boldsymbol{x}$ and $\Gamma \succ 0$. By an earlier theorem, global minima must exist. Furthermore, $\mathcal{C} \setminus \text{int}\, \mathcal{C}$ is empty, and so the global minimum must be an unconstrained local minimum. The first order necessary conditions immediately allow us to characterize that point as $\Gamma \boldsymbol{x}^* - \boldsymbol{b} = \boldsymbol{0}$.

- *Sensitivity analysis*
    - Consider the problem $\min f(\boldsymbol{x}, \boldsymbol{a})$ s.t. $x \in \mathbb{R}^n$. We let $\boldsymbol{x}^*$ be a local optimum, and $f^*(\boldsymbol{a}) = f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a})$. The first-order conditions are

    $$\nabla_x f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) = \boldsymbol{0}$$

    Taking the derivative with respect to $\boldsymbol{a}$, we obtain

    $$\nabla \boldsymbol{x}^*(\boldsymbol{a}) \nabla_{xx}^2 f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) + \nabla_{xa}^2 f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) = \boldsymbol{0}$$

    From this expression, we can obtain expressions for the sensitivity of the optimum, and of the optimal value:

    $$\nabla \boldsymbol{x}^*(\boldsymbol{a}) = -\nabla_{xa}^2 f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) \left\{ \nabla_{xx}^2 f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) \right\}^{-1}$$

    $$\nabla f^*(\boldsymbol{a}) = \nabla_a f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) = \nabla \boldsymbol{x}^*(\boldsymbol{a}) \nabla_x f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) + \nabla_a f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a}) = \nabla_a f(\boldsymbol{x}^*(\boldsymbol{a}), \boldsymbol{a})$$

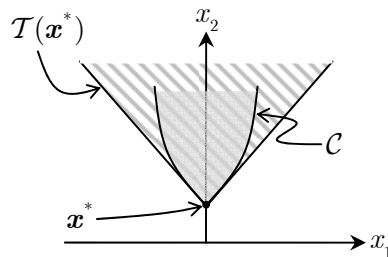    - The implicit function theorem tells us when this exists.

## Constrained local optimality

- Consider the problem $\min f(\boldsymbol{x})$ s.t. $x \in \mathcal{C} \subset \mathbb{R}^n$. We are interested in characterizing local minima that are *not* in $\text{int}\, \mathcal{C}$. We will assume, though, that $f$ is continuously differentiable in a neighborhood of the point considered.

- **Definition**: The set of *descent directions* of $f$ at $x^*$ is the set $\mathcal{D}(x^*) = \left\{ d \in \mathbb{R}^n : \nabla f(x^*) \cdot d < 0 \right\}$.

- **Definition**: The *tangent cone* $\mathcal{T}(x^*)$ of the constraint set $\mathcal{C}$ at $x^*$ is the set of directions $d \in \mathbb{R}^n$ such that either

  - $d = 0$

  - There exists a sequence $\{x_k\} \subset \mathcal{C}, x_k \to x^*$ such that

  $$\frac{x_k - x^*}{\left\| x_k - x^* \right\|} \to \frac{d}{\|d\|}$$

  Geometrically, this is simply the statement that $d$ is tangent to $\mathcal{C}$ if and only if there is some walk we can take in $\mathcal{C}$ that leads us to $x^*$ and that ends up being in direction $d$. For example:

  

- **Theorem – necessary condition**: If $x^*$ is a local minimum, there is no descent direction in the tangent cone:

  $$\mathcal{D}(x^*) \cap \mathcal{T}(x^*) = \varnothing$$

  Geometrically, the tangent cone contains the "directions in which we can move". The set of descent directions contains the "directions in which we can improve our objective". If any direction fulfils both these conditions, then we can clearly improve on our current point.

  **Proof**: Consider $d \in \mathcal{T}(x^*) \setminus \{0\}$ and an appropriate sequence $\{x_k\} \subset \mathcal{C}, x_k \to x^*$. Define

  $$\zeta_k = \frac{x_k - x^*}{\left\| x_k - x^* \right\|} - \frac{d}{\|d\|} \to 0 \qquad\qquad d_k = d + \|d\| \zeta_k \to d$$

  Now, if $\tilde{x}_k$ is a point on the line segment between $x_k$ and $x^*$, the mean value theorem tells us that

  $$f(x_k) = f(x^*) + \nabla f(\tilde{x}_k) \cdot (x_k - x^*)$$

  Note, however, that we can write

$$\left(\boldsymbol{x}_k - \boldsymbol{x}^*\right) = \frac{\left\|\boldsymbol{x}_k - \boldsymbol{x}^*\right\|}{\|\boldsymbol{d}\|}\left(\|\boldsymbol{d}\|\boldsymbol{\zeta}_k + \boldsymbol{d}\right) = \frac{\left\|\boldsymbol{x}_k - \boldsymbol{x}^*\right\|}{\|\boldsymbol{d}\|}\boldsymbol{d}_k$$

And so we can re-write the above as

$$f(\boldsymbol{x}_k) = f(\boldsymbol{x}^*) + \frac{\left\|\boldsymbol{x}_k - \boldsymbol{x}^*\right\|}{\|\boldsymbol{d}\|}\nabla f(\tilde{\boldsymbol{x}}_k)\cdot\boldsymbol{d}_k$$

Now, if $\boldsymbol{d} \in \mathcal{D}(\boldsymbol{x}^*)$ as well, then $\nabla f(\boldsymbol{x}^*)\cdot\boldsymbol{d} < 0$. The strict inequality implies that this is also true in a neighborhood of $\boldsymbol{x}^*$, and so for $k$ large enough, we get $f(\boldsymbol{x}_k) < f(\boldsymbol{x}^*)$. This contradicts the local minimality of $\boldsymbol{x}^*$.

- Unfortunately, $\mathcal{T}$ is hard to characterize algebraically, unless we focus on the particular example where $\mathcal{C}$ is the intersection of equality constraints.

## Equality constrained optimization

- Consider the problem $\min f(\boldsymbol{x})$ s.t. $\boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0}, \boldsymbol{x} \in \mathbb{R}^n$ where $h : \mathbb{R}^n \to \mathbb{R}^m$. We assume the $f$ and $h_i$ are continuously differentiable in a neighborhood of the local minimum.

- In this particular case, we will show we can characterize $\mathcal{T}$ in a simple way. The intuition behind our result is that for any feasible $\boldsymbol{x}$, $\boldsymbol{d} \in \mathbb{R}^n$ and $\alpha > 0$

$$\boldsymbol{h}(\boldsymbol{x} + \alpha\boldsymbol{d}) \approx \boldsymbol{h}(\boldsymbol{x}) + \alpha\nabla\boldsymbol{h}(\boldsymbol{x})^\top\boldsymbol{d} = \alpha\nabla\boldsymbol{h}(\boldsymbol{x})^\top\boldsymbol{d}$$

So intuitively, one might expected that any direction for which $\nabla\boldsymbol{h}(\boldsymbol{x})^\top\boldsymbol{d} = \boldsymbol{0}$ to maintain feasibility. We now formalize this statement…

- **Definition**: the cone of first-order feasible variations at $\boldsymbol{x}^* \in \mathbb{R}^n$ is the set

$$\mathcal{V}(\boldsymbol{x}^*) = \left\{d \in \mathbb{R}^n : \nabla\boldsymbol{h}(\boldsymbol{x}^*)^\top\boldsymbol{d} = \boldsymbol{0}\right\} = \left[\ker\nabla\boldsymbol{h}(\boldsymbol{x}^*)^\top\right]$$

Note that $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*) \Rightarrow -\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$. As such, $\mathcal{V}(\boldsymbol{x}^*)$ is actually a *subspace* of $\mathbb{R}^n$.

- ***Definition***: A point $\boldsymbol{x}^* \in \mathbb{R}^n$ is a *regular point* if it is feasible and the constraint gradients $\nabla h_i(\boldsymbol{x}^*)$ are linearly independent. In other words, $\left|\nabla\boldsymbol{h}(\boldsymbol{x}^*)\right| \neq 0$. If $m > n$, no regular points exist, and if $m = 1$, this reduces to $\nabla h_1(\boldsymbol{x}^*) \neq \boldsymbol{0}$.

- ***Lemma (regularity)***: Let $\boldsymbol{x}^*$ be a regular point. Then $\mathcal{T}(\boldsymbol{x}^*) = \mathcal{V}(\boldsymbol{x}^*)$

  ***Proof***: This theorem is hard. The intuition behind the proof is

  o Consider the curve we would trace if we were sitting at a point $\boldsymbol{x}^*$ and we started walking forward or backwards *while staying on the constraint* (ie: while keeping the constraint satisfied). We'll start by showing that for any direction

$d \in \mathcal{V}(x^*)$, there is such a path that starts by walking forward or backward along the direction $d$.

o   Once we've established this, the result is relatively easy, because the path constitutes a "walk" fully contained in our set $\mathcal{C}$ which eventually ends up being in the direction $d$. It's therefore in $\mathcal{T}$.

And now the painful details! First, let's find the curve in question:

o   Begin by choosing $d \in \mathcal{V}(x^*)$. Given a scalar $t$, consider the curve $x(t) = x^* + td$. This satisfies our requirement that we be moving either side of $x^*$, and that we start by going in direction $d$. However, there's no guarantee we stay on the constraints.

o   Instead, consider the path $x(t) = x^* + td + \nabla h(x^*)u(t)$ for some unknown vector $u(t) \in \mathbb{R}^m$. This seems sensible – we are correcting our path to reflect how $h$ might change. For $x(t)$ to be "valid", we require it to satisfy the $m$ equations

$$\boxed{h\Big(x^* + td + \nabla h(x^*)u(t)\Big) = 0}$$

For $t = 0$, $u(0) = 0$ is clearly a solution.

Now, take the gradient of the boxed equation with respect to $u$ and evaluate it at $(t, u) = 0$. We get

$$\nabla h(x^*)^\top \nabla h(x^*)$$

Since the columns of $\nabla h(x^*)$ are linearly independent, this matrix is invertible.

The two results above allow us to use the implicit function theorem to deduce that a solution $u(t)$ to the boxed equation exists for all $t \in (-\tau, \tau)$, for some $\tau$.

Thus, we have managed to find a curve $x(t)$ that keeps us on the constraints and that is defined over $t \in (-\tau, \tau)$ with $x(0) = x^*$ (this implies that the curve represents moving forward and backward from $x^*$).

o   All we now need to prove is that the initial direction in which we move is $d$. To do that, differentiate the boxed equation above with respect to $t$ and evaluate at $t = 0$. We get

$$\left(\boldsymbol{d}^\top + \dot{\boldsymbol{u}}(0)^\top \nabla \boldsymbol{h}(\boldsymbol{x}^*)^\top\right) \nabla \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0}$$

But since $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$, we also have that $\boldsymbol{d}^\top \nabla \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0}$, and so $\dot{\boldsymbol{u}}(0) = \boldsymbol{0}$, and $\dot{\boldsymbol{x}}(0) = \boldsymbol{d}$, as required.

o  [It will be useful for later to note that if $\boldsymbol{h}(\cdot)$ is twice continuously differentiable, then so is $\boldsymbol{x}(\cdot)$. Though I'm not quite sure how to prove that result].

We now have our elusive curve! Let's now prove the theorem.

o  $\boxed{\mathcal{V}(\boldsymbol{x}^*) \subset \mathcal{T}(\boldsymbol{x}^*)}$: choose $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*) \setminus \{\boldsymbol{0}\}$ and let $\boldsymbol{x}(t)$ be the curve discussed above. Take a sequence $t_k \subset (0, \tau), t_k \to 0$, so that $\boldsymbol{x}(t_k) \neq \boldsymbol{x}^*$. Then, by the mean value theorem, there is some $\tilde{t} \in [0, t_k]$ such that

$$\boldsymbol{x}(t_k) - \boldsymbol{x}(0) = \dot{\boldsymbol{x}}(\tilde{t})(t_k - 0)$$
$$\frac{\boldsymbol{x}(t_k) - \boldsymbol{x}^*}{\left\|\boldsymbol{x}(t_k) - \boldsymbol{x}^*\right\|} = \frac{\dot{\boldsymbol{x}}(\tilde{t})}{\left\|\boldsymbol{x}(t_k) - \boldsymbol{x}^*\right\| / t_k}$$

As $t_k \to 0$ and therefore $\tilde{t} \to 0$, this tends to

$$\to \frac{\dot{\boldsymbol{x}}(0)}{\left\|\dot{\boldsymbol{x}}(0)\right\|} = \frac{\boldsymbol{d}}{\left\|\boldsymbol{d}\right\|}$$

So $\boldsymbol{d} \in \mathcal{T}(\boldsymbol{x}^*)$.

o  $\boxed{\mathcal{T}(\boldsymbol{x}^*) \subset \mathcal{V}(\boldsymbol{x}^*)}$: consider $\boldsymbol{d} \in \mathcal{T}(\boldsymbol{x}^*) \setminus \{\boldsymbol{0}\}$ and an associated sequence $\{\boldsymbol{x}_k\}$ in the feasible set, as defined in the definition of $\mathcal{T}(\boldsymbol{x}^*)$. By the mean value theorem, there is some $\tilde{\boldsymbol{x}} \in [\boldsymbol{x}_k, \boldsymbol{x}^*]$ such that

$$\boldsymbol{h}(\boldsymbol{x}_k) - \boldsymbol{h}(\boldsymbol{x}^*) = \nabla \boldsymbol{h}(\tilde{\boldsymbol{x}})^\top (\boldsymbol{x}_k - \boldsymbol{x}^*)$$

But since $\boldsymbol{x}^*$ and every $\boldsymbol{x}_k$ are in the feasible set, $\boldsymbol{h}(\boldsymbol{x}_k) = \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0}$, and

$$0 = \boldsymbol{h}(\boldsymbol{x}_k) = \boldsymbol{h}(\boldsymbol{x}^*) + \nabla \boldsymbol{h}(\tilde{\boldsymbol{x}}_k)^\top (\boldsymbol{x}_k - \boldsymbol{x}^*) = \nabla \boldsymbol{h}(\tilde{\boldsymbol{x}}_k)^\top (\boldsymbol{x}_k - \boldsymbol{x}^*)$$
$$\nabla \boldsymbol{h}(\tilde{\boldsymbol{x}}_k)^\top \frac{\boldsymbol{x}_k - \boldsymbol{x}^*}{\left\|\boldsymbol{x}_k - \boldsymbol{x}^*\right\|} = 0$$
$$\to \nabla \boldsymbol{h}(\boldsymbol{x}^*)^\top \boldsymbol{d} = \boldsymbol{0}$$

And so $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$.

{*Done! Take a deep breath!*}

•  ***Theorem – necessary condition***: If $\boldsymbol{x}^*$ is a local minimum that is a regular point, then there is no descent direction that is also a first order feasible variation:

$$\nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d} = 0 \qquad \forall \; \boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$$

Daniel Guetta

Or in other words, we require $\nabla f(\boldsymbol{x}^*)$ to be in $\mathcal{V}(\boldsymbol{x}^*)^\perp$:

$$\nabla f(\boldsymbol{x}^*) \in \mathcal{V}(\boldsymbol{x}^*)^\perp = \left[\ker \nabla \boldsymbol{h}(\boldsymbol{x}^*)^\top\right]^\perp = \operatorname{im} \nabla \boldsymbol{h}(\boldsymbol{x}^*)$$

Or in other words, there exists $\boldsymbol{\lambda} \in \mathbb{R}^m$ such that $\nabla f(\boldsymbol{x}^*) = \nabla \boldsymbol{h}(\boldsymbol{x}^*)\,\boldsymbol{\lambda}$.

**Proof:** Since $\boldsymbol{x}^*$ is a local minimum, $\mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{T}(\boldsymbol{x}^*) = \varnothing$, and since $\boldsymbol{x}^*$ is regular, $\mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{V}(\boldsymbol{x}^*) = \varnothing$.

Now, assume $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$ – by what we have such said, $\boldsymbol{d} \notin \mathcal{D}(\boldsymbol{x}^*)$, and so $\nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d} \geq 0$. However, since we also have $-\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$, we must have $\nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d} = 0$.

For the last part of the theorem, note that $\operatorname{im} A = (\ker A^\top)^\perp$, as proved in the introductory section of these notes.

- The last part of the previous theorem is important, because it provides a "simple" way to characterize the tangent cone, and a "recipe" to find optimal points. This can be formalized further using…

- *...Lagrange Multipliers*

  o **Theorem – necessary conditions**: If $\boldsymbol{x}^*$ is a local minimum that is a regular point, then there exists a unique vector $\boldsymbol{\lambda}^* \in \mathbb{R}^m$ called a *Lagrange multiplier* such that

$$\nabla f(\boldsymbol{x}^*) + \boldsymbol{\lambda}^{*\top} \nabla \boldsymbol{h}(\boldsymbol{x}^*) = \nabla f(\boldsymbol{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla h_i(\boldsymbol{x}^*) = \boldsymbol{0}$$

  In addition, if $f$ and $\boldsymbol{h}$ are twice continuously differentiable

$$\boldsymbol{d}^\top \left(\nabla^2 f(\boldsymbol{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla^2 h_i(\boldsymbol{x}^*)\right) \boldsymbol{d} \geq 0 \qquad \forall \boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$$

  There is an interesting geometrical interpretation of the first-order condition. It effectively states that $\nabla f(\boldsymbol{x}^*)$ [the direction in which we might increase our objective] must be a linear combination of the $\nabla h_i(\boldsymbol{x}^*)$ [the *perpendicular* to the constraints $h_i(\boldsymbol{x}^*) = 0$]. Since we cannot move along any of those perpendiculars without leaving the constraints, we clearly cannot move along $\nabla f(\boldsymbol{x}^*)$. Here is an example, in which $\nabla f(\boldsymbol{x})$ is constant:

**Proof:** The existence of $\boldsymbol{\lambda}^*$ is simply a restatement of the previous theorem. The uniqueness of $\boldsymbol{\lambda}^*$ follows from the fact that the columns of $\nabla h(\boldsymbol{x}^*)$ are linearly independent.

For the second-order condition, consider a $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$, and use the first part of the regularity lemma to define a path $\boldsymbol{x}(t)$ either side of $\boldsymbol{x}^*$, which always stays on the constraints and such that $\dot{\boldsymbol{x}}(0) = \boldsymbol{d}$. Now, define $g(t) = f(\boldsymbol{x}(t))$ and take a double derivative

$$\ddot{g}(t) = \dot{\boldsymbol{x}}(t)^\top \nabla^2 f(\boldsymbol{x}(t))\dot{\boldsymbol{x}}(t) + \ddot{\boldsymbol{x}}(t)^\top \nabla f(\boldsymbol{x}(t))$$

Since all points $\boldsymbol{x}(t)$ satisfy the constraints of the problem, and $\boldsymbol{x}^*$ is a local minimum, $t = 0$ must be an unconstrained *local* minimum of $g(t)$. Thus

$$\ddot{g}(0) = \boldsymbol{d}^\top \nabla^2 f(\boldsymbol{x}^*)\boldsymbol{d} + \ddot{\boldsymbol{x}}(0)^\top \nabla f(\boldsymbol{x}^*) \geq 0$$

Finally, consider $\ell(t) = \boldsymbol{\lambda}^{*\top} h(\boldsymbol{x}(t)) = 0$ and differentiate it twice, to get

$$\ddot{\ell}(0) = \boldsymbol{d}^\top \left( \sum_{i=1}^m \lambda_i \nabla^2 h_i(\boldsymbol{x}^*) \right) \boldsymbol{d} + \ddot{\boldsymbol{x}}(0)^\top \nabla h(\boldsymbol{x}^*)\boldsymbol{\lambda}^* = 0$$

Finally, add the last two equations, and apply the first order condition.

o   We define the *Lagrangian* as

$$\mathcal{L}\left(\boldsymbol{x}, \boldsymbol{\lambda}\right) = f(\boldsymbol{x}) + \boldsymbol{\lambda} \cdot h(\boldsymbol{x})$$

The first and second order conditions then reduce to

$$\nabla_{\boldsymbol{x}} \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*) = 0$$
$$\boldsymbol{d}^\top \nabla^2_{\boldsymbol{xx}} \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*)\boldsymbol{d} \geq 0 \quad \forall d \in \mathcal{V}(\boldsymbol{x}^*)$$

And the feasibility condition is given by

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

This formulation of the first order conditions also has an interesting interpretation. We effectively allow the constraint $h(\boldsymbol{x}) = \boldsymbol{0}$ to be broken, *but* we associate a cost $\boldsymbol{\lambda}$ with breaking it. We then apply unconstrained first-order conditions to the resulting problem.

o   As such, we have the following simple recipe to solve an equality-constrained problem (assuming $f$ and $h$ are continuously differentiable on $\mathbb{R}^n$):

   ▪   Check that global minima exist

   ▪   Find the set of $(\boldsymbol{x}^*, \boldsymbol{\lambda}^*)$ satisfying $\nabla_{\boldsymbol{x}} \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*) = 0$ and $\nabla_{\boldsymbol{\lambda}} \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*) = 0$.

   ▪   Find the set of non-regular points

   ▪   The global minima are amongst these points.

A few examples:

   ▪   $\min f(\boldsymbol{x}) = x_1 + x_2$ s.t. $h(\boldsymbol{x}) = x_1^2 + x_2^2 - 2 = 0, \boldsymbol{x} \in \mathbb{R}^2$. In this case, note that $\nabla h = (2x_1, 2x_2)$; thus, provided the constraints are met, all points will be regular. In this case, $\mathcal{L}(\boldsymbol{x}, \lambda) = x_1 + x_2 + \lambda(x_1^2 + x_2^2 - 2)$, and since all points are regular, we simply need to consider all points satisfying the first-order conditions and choose the smallest one.

   ▪   Consider the problem

$$\begin{aligned} \min \quad & f(\boldsymbol{x}) = x_1 + x_2 \\ \text{s.t.} \quad & h_1(\boldsymbol{x}) = (x_1 - 1)^2 + x_2^2 - 1 = 0 \\ & h_2(\boldsymbol{x}) = (x_1 - 2)^2 + x_2^2 - 4 = 0 \\ & \boldsymbol{x} \in \mathbb{R}^2 \end{aligned}$$

The only feasible point is $(0, 0)$ – and so it is the global minimum. However, it is not a regular point. And indeed, we find that there is no solution to the first-order necessary conditions.

o   ***Constraint qualification***: In summary, what we have proved so far is

$$\boldsymbol{x}^* \text{ local minimum}$$
$$\Rightarrow \mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{T}(\boldsymbol{x}^*) = \varnothing$$
$$\Rightarrow \mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{V}(\boldsymbol{x}^*) = \varnothing$$
$$\Rightarrow \text{Lagrange multipliers exist}$$

In proving this sequence of statements, we made use of the fact that $\boldsymbol{x}^*$ was a regular point. However, it is possible to prove the existence of Lagrange

multipliers under weaker assumptions called *constraint qualifications*. If the constraints are linear, for example, Lagrange multipliers are guaranteed to exist. The weakest form of constraint qualification is *quasiregularity*, which requires that $\mathcal{V}(\boldsymbol{x}^*) = \mathcal{T}(\boldsymbol{x}^*)$.

o   ***Theorem  –  Sufficient  Conditions***: Assume that $f$ and $\boldsymbol{h}$ are both twice continuously differentiable, and that $\boldsymbol{x}^* \in \mathbb{R}^n$ and $\boldsymbol{\lambda}^* \in \mathbb{R}^m$ satisfy

$$\nabla_x L(\boldsymbol{x}^*, \boldsymbol{\lambda}^*) = \boldsymbol{0}$$
$$\nabla_\lambda L(\boldsymbol{x}^*, \boldsymbol{\lambda}^*) = \boldsymbol{0}$$
$$\boldsymbol{d}^\top \nabla^2_{xx} L(\boldsymbol{x}^*, \boldsymbol{\lambda}^*) \boldsymbol{d} > 0 \quad \forall \boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*) \setminus \{\boldsymbol{0}\}$$

Then $\boldsymbol{x}^*$ is a strict local minimum.

***Proof***: The second condition above implies that $\boldsymbol{x}^*$ is clearly feasible. Suppose it is not a strict local minimum; then there exists a sequence $\{\boldsymbol{x}_k\} \subset \mathbb{R}^n$ such that $\boldsymbol{x}_k \neq \boldsymbol{x}^*$ and $\boldsymbol{x}_k \to \boldsymbol{x}^*$ which lies entirely in the feasible region of the problem [ie: $\boldsymbol{h}(\boldsymbol{x}_k) = \boldsymbol{0}$] and $f(\boldsymbol{x}_k) \leq f(\boldsymbol{x}^*)$. We define, for some $\boldsymbol{d}$

$$\boldsymbol{d}_k = \frac{\boldsymbol{x}_k - \boldsymbol{x}^*}{\left\| \boldsymbol{x}_k - \boldsymbol{x}^* \right\|} \to \boldsymbol{d} \qquad\qquad \delta_k = \left\| \boldsymbol{x}_k - \boldsymbol{x}^* \right\| \to 0$$

Now, by the mean value theorem, there exists $\tilde{\boldsymbol{x}} \in [\boldsymbol{x}^*, \boldsymbol{x}_k]$ with

$$\boldsymbol{h}(\boldsymbol{x}_k) - \boldsymbol{h}(\boldsymbol{x}^*) = \nabla \boldsymbol{h}(\tilde{\boldsymbol{x}}_k)^\top (\boldsymbol{x}_k - \boldsymbol{x}^*) = \nabla \boldsymbol{h}(\tilde{\boldsymbol{x}}_k)^\top (\delta_k \boldsymbol{d}_k)$$

But since $\boldsymbol{x}^*$ and $\boldsymbol{x}_k$ are feasible, $\boldsymbol{h}(\boldsymbol{x}_k) = \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0}$, so.

$$\nabla \boldsymbol{h}(\tilde{\boldsymbol{x}}_k)^\top \boldsymbol{d}_k = 0$$

Taking the limit as $k \to \infty$, we get $\nabla \boldsymbol{h}(\boldsymbol{x}^*)^\top \boldsymbol{d} = 0$, and so $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)$.

Now, we know that

$$\boldsymbol{h}(\boldsymbol{x}_k) = \boldsymbol{0} \qquad\qquad f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) \leq 0$$

Using a second order Taylor expansion (with remainder) with some set of $\hat{\boldsymbol{x}}^i \in [\boldsymbol{x}_k, \boldsymbol{x}^*]$, we can re-write these as

$$h_i(\boldsymbol{x}_k) = h_i(\boldsymbol{x}^*) + \delta_k \nabla h_i(\boldsymbol{x}^*) \cdot \boldsymbol{d}_k + \tfrac{1}{2} \delta_k^2 \boldsymbol{d}_k^\top \nabla^2 h_i(\hat{\boldsymbol{x}}^i) \boldsymbol{d}_k = 0$$

$$f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) = \boxed{\delta_k \nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d}_k + \tfrac{1}{2} \delta_k^2 \boldsymbol{d}_k^\top \nabla^2 f(\hat{\boldsymbol{x}}^0) \boldsymbol{d}_k \leq 0}$$

We can modify the first set of equations slightly by remembering that $\boldsymbol{h}(\boldsymbol{x}^*) = 0$, and multiplying both sides of the equation by $\lambda_i^*$. This gives

$$h_i(\boldsymbol{x}_k) = \boxed{\delta_k \lambda_i^* \nabla h_i(\boldsymbol{x}^*) \cdot \boldsymbol{d}_k + \tfrac{1}{2}\delta_k^2 \boldsymbol{d}_k^\top \lambda_i^* \nabla^2 h_i(\hat{\boldsymbol{x}}^i)\boldsymbol{d}_k = 0}$$

Adding these $m+1$ equations, we get

$$\delta_k\left(\nabla f(\boldsymbol{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla h_i(\boldsymbol{x}^*)\right)\cdot \boldsymbol{d}_k + \tfrac{1}{2}\delta_k^2 \boldsymbol{d}_k^\top \left(\nabla^2 f(\hat{\boldsymbol{x}}^0) + \sum_{i=1}^m \lambda_i^* \nabla^2 h_i(\hat{\boldsymbol{x}}^i)\right)\boldsymbol{d}_k \le 0$$

$$\delta_k \nabla_{\boldsymbol{x}} L(\boldsymbol{x}^*,\boldsymbol{\lambda}^*)\cdot \boldsymbol{d}_k + \tfrac{1}{2}\delta_k^2 \boldsymbol{d}_k^\top \left(\nabla^2 f(\hat{\boldsymbol{x}}^0) + \sum_{i=1}^m \lambda_i^* \nabla^2 h_i(\hat{\boldsymbol{x}}^i)\right)\boldsymbol{d}_k \le 0$$

Noting that, by the first order conditions, $\nabla_{\boldsymbol{x}} L(\boldsymbol{x}^*,\boldsymbol{\lambda}^*)$ and then dividing by $\tfrac{1}{2}\delta_k^2$ and taking the limit as $k \to \infty$, this becomes

$$\boldsymbol{d}^\top\left(\nabla^2 f(\boldsymbol{x}^*) + \sum_{i=1}^m \lambda_i^* \nabla^2 h_i(\boldsymbol{x}^*)\right)\boldsymbol{d} \le 0$$

But since $\boldsymbol{d} \in \mathcal{V}(\boldsymbol{x}^*)\setminus\{\boldsymbol{0}\}$, this violates our assumed second order condition.

o We now consider an application of these conditions. Consider the program

$$\min_{\boldsymbol{x}\in\mathbb{R}^n} \sigma^2 = \boldsymbol{x}^\top \Gamma \boldsymbol{x} \text{ s.t. } \boldsymbol{1}^\top \boldsymbol{x} = 1, \boldsymbol{\mu}^\top \boldsymbol{x} = \bar{\mu}$$

which might represent minimizing the variance in a portfolio while keeping total sales equal to 1 unit, and keeping the expected return equal to a certain value $\bar{\mu}$. The first-order conditions give

$$2\Gamma \boldsymbol{x}^* + \lambda_1^* \boldsymbol{1} + \lambda_2^* \boldsymbol{\mu} = \boldsymbol{0} \qquad \boldsymbol{1}^\top \boldsymbol{x}^* = 1, \boldsymbol{\mu}^\top \boldsymbol{x}^* = \bar{\mu}$$

From the first equation, we obtain

$$\boldsymbol{x}^* = -\tfrac{1}{2}\Gamma^{-1}\boldsymbol{1}\lambda_1^* - \tfrac{1}{2}\Gamma^{-1}\boldsymbol{\mu}\lambda_2^*$$
$$\boldsymbol{1}^\top \boldsymbol{x}^* = -\tfrac{1}{2}\boldsymbol{1}^\top \Gamma^{-1}\boldsymbol{1}\lambda_1^* - \tfrac{1}{2}\boldsymbol{1}^\top \Gamma^{-1}\boldsymbol{\mu}\lambda_2^* = 1$$
$$\boldsymbol{\mu}^\top \boldsymbol{x}^* = -\tfrac{1}{2}\boldsymbol{\mu}^\top \Gamma^{-1}\boldsymbol{1}\lambda_1^* - \tfrac{1}{2}\boldsymbol{\mu}^\top \Gamma^{-1}\boldsymbol{\mu}\lambda_2^* = \bar{\mu}$$

The last two equations are a system of equations for $(\lambda_1^*, \lambda_2^*)$:

$$-\frac{1}{2}\begin{pmatrix} \boldsymbol{1}^\top \Gamma^{-1}\boldsymbol{1} & \boldsymbol{1}^\top \Gamma^{-1}\boldsymbol{\mu} \\ \boldsymbol{\mu}^\top \Gamma^{-1}\boldsymbol{1} & \boldsymbol{\mu}^\top \Gamma^{-1}\boldsymbol{\mu} \end{pmatrix}\begin{pmatrix} \lambda_1^* \\ \lambda_2^* \end{pmatrix} = \begin{pmatrix} 1 \\ \bar{\mu} \end{pmatrix}$$

this system is nonsingular provided that $\Gamma \succ 0$ and $\boldsymbol{1}$ and $\boldsymbol{\mu}$ are linearly independent. We then get

$$\begin{pmatrix} \lambda_1^* \\ \lambda_2^* \end{pmatrix} = \begin{pmatrix} \eta_1 + \zeta_1\bar{\mu} \\ \eta_2 + \zeta_2\bar{\mu} \end{pmatrix}$$
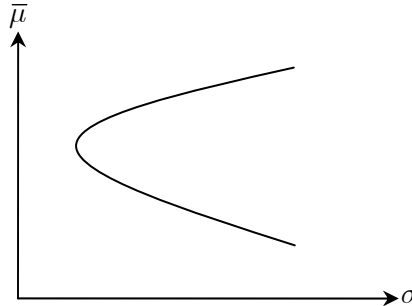
Where the constants depend on $\Gamma$ and $\boldsymbol{\mu}$. Now, using the first equation in the FOCs, we obtain, for some vectors $\boldsymbol{v}$ and $\boldsymbol{w}$

$$\boldsymbol{x}^* = \bar{\mu}\boldsymbol{v} + \boldsymbol{w}$$

which finally gives, for some constants $\alpha, \beta, \gamma$ depending on $\Gamma$ and $\boldsymbol{\mu}$

$$\sigma^2 = \left(\overline{\mu}\boldsymbol{v} + \boldsymbol{w}\right)^\top \Gamma\left(\overline{\mu}\boldsymbol{v} + \boldsymbol{w}\right) = \left(\alpha\overline{\mu} + \beta\right)^2 + \gamma$$

This implies that of all the possible asset combinations, those that provide the maximum expected returns for a particular variance lie on a parabola that looks like this:



The upper part of this parabola is called the *efficient frontier*.

- *Sensitivity analysis*
  - We now look at what happens when we slightly vary our constraint.
  - ***Theorem***: Consider a family of equality-constrained problems $\min_{\boldsymbol{x}\in\mathbb{R}^n} f(\boldsymbol{x})$ s.t. $\boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{u}$. Suppose there exists a local pair $(\boldsymbol{x}^*, \boldsymbol{\lambda}^*)$ satisfying the second-order sufficient conditions when $\boldsymbol{u} = \boldsymbol{0}$, and $\boldsymbol{x}^*$ is regular. Then there exists a neighborhood $N_\varepsilon(0) \subset \mathbb{R}^m$ of $\boldsymbol{u} = \boldsymbol{0}$ and a function $\boldsymbol{x}^*(\cdot)$ defined on that neighborhood such that

      - $\boldsymbol{x}^*(\boldsymbol{0}) = \boldsymbol{x}^*$, and for each $\boldsymbol{u} \in N_\varepsilon(0)$, $\boldsymbol{x}^*(\boldsymbol{u})$ is a strict local minimum.
      - $\boldsymbol{x}^*(\boldsymbol{u})$ is continuously differentiable.
      - If $p(\boldsymbol{u}) = f(\boldsymbol{x}^*(\boldsymbol{u}))$, then $\nabla p(\boldsymbol{0}) = -\boldsymbol{\lambda}^*$.

    The last part of the theorem is the interesting one – it tells us that changing $\boldsymbol{u}$ by an infinitesimal amount around our solution point will lead in a change of $-\boldsymbol{\lambda}^*$ in the objective.

    ***Proof***: Consider, for $\boldsymbol{u} \in \mathbb{R}^m$, the following system of equations in $(\boldsymbol{x}, \boldsymbol{\lambda})$, representing the first-order sufficient conditions:

    $$\nabla f(\boldsymbol{x}) + \nabla \boldsymbol{h}(\boldsymbol{x})\boldsymbol{\lambda} = \boldsymbol{0}$$
    $$\boldsymbol{h}(\boldsymbol{x}) - \boldsymbol{u} = \boldsymbol{0}$$

    Let's carefully calculate the gradient of this system

$$\begin{bmatrix} \nabla_x \{\nabla f(x) + \nabla h(x)\lambda\} & \nabla_x \{h(x) - u\} \\ \nabla_\lambda \{\nabla f(x) + \nabla h(x)\lambda\} & \nabla_\lambda \{h(x) - u\} \end{bmatrix} = \begin{bmatrix} \nabla_{xx} L(x, \lambda) & \nabla_x h(x) \\ \nabla h(x)^\top & 0 \end{bmatrix}$$

Clearly, the system above has a solution at $(x^*, \lambda^*)$, and at that point, the gradient is given by

$$J = \begin{bmatrix} \nabla_{xx} L(x^*, \lambda^*) & \nabla_x h(x^*) \\ \nabla h(x^*)^\top & 0 \end{bmatrix}$$

Now, if this matrix were singular, there would exist some nonzero vector $x = (y^\top, z^\top)^\top$ such that $Jx = 0$. This implies that

$$\nabla_{xx}^2 L(x^*, \lambda^*)y + \nabla_x h(x^*) = 0$$
$$\nabla h(x^*)^\top y = 0$$

The second equation implies that $y \in \mathcal{V}(x^*)$. Multiplying the first equation by $y^\top$ and using the second equation implies that we would then have to have $y^\top \nabla_{xx}^2 L(x^*, \lambda^*)y = 0$ for some $y \in \mathcal{V}(x^*)$; this violates the second-order sufficient conditions. Thus, our matrix is nonsingular.

Using the implicit function theorem, this implies that we can define $\left(x^*(u), \lambda^*(u)\right)$ satisfying first-order conditions for all $u$ in some $N_\varepsilon(0)$. Second-order conditions follow, for $u$ sufficiently close to $0$, by continuity assumptions.

For the last part of our theorem, consider the FOCs in terms of this $\left(x^*(u), \lambda^*(u)\right)$, and multiply them by $\nabla x^*(u)$:

$$\nabla x^*(u)\left[\nabla f(x^*(u))\right] + \nabla x^*(u)\left[\nabla h(x^*(u))\right]\lambda^*(u) = 0$$
$$\nabla_u f(x^*(u) + \nabla x^*(u)\left[\nabla h(x^*(u))\right]\lambda^*(u) = 0$$

Note also that differentiating $h(x^*(u)) = u$ with respect to $u$ yields $\nabla_u x^*(u)\nabla_x h(x^*(u)) = I$. Using that result, we can re-write the previous equation as

$$\nabla_u f(x^*(u) + I\lambda^*(u) = 0$$
$$\nabla_u f(x^*(u) = -\lambda^*(u)$$

As required.

## Inequality constrained optimization

- Consider the program

$$\min_{x \in \mathbb{R}^n} f(x) \text{ s.t. } h(x) = 0, g(x) \le 0$$

Where $f : \mathbb{R}^n \to \mathbb{R}$, $h : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^n \to \mathbb{R}^r$. Note also that when we say $x \le 0$, we mean that every component of $x$ is less than or equal to 0.

We also assume that $f$, $h$ and $g$ are continuously differentiable on $\mathbb{R}^n$, though the necessary and sufficient conditions also hold if these are defined and continuously differentiable only in a neighborhood of the local minimum.

- ***Definition***: Given a feasible point $x^* \in \mathbb{R}^n$, the set of active inequality constraints $\mathcal{A}(x^*)$ is defined as

$$\mathcal{A}(x^*) = \left\{ j : g_j(x^*) = 0 \right\} \subseteq \left\{ 1, \cdots, r \right\}$$

- ***Lemma***: Let $x^*$ be a local minimum for the inequality constrained program above (IneqCP). Then it is also a local minimum for the following equality constrained program (EqCP):

$$\min_{x \in \mathbb{R}^n} f(x) \text{ s.t } h(x) = 0, g_j(x) = 0 \ \forall j \in \mathcal{A}(x^*)$$

***Proof***: Suppose $x^*$ is not a local optimum for EqCP. Then there exists a sequence of points $\{x_k\}$ feasible for EqCP such that $x_k \to x^*$ and $f(x_k) < f(x^*)$.

Since $g$ is continuous, we have $g(x_k) \to g(x^*)$. In particular, if $j \notin \mathcal{A}(x^*)$, $g_j(x_k) \to g_j(x^*) < 0$. Thus, for sufficiently large $k$, $x_k$ is feasible for IneqCP. This contradicts the local optimality of $x^*$ for IneqCP.

- ***Definition***: Consider the inequality constrained program above. A point $x^* \in \mathbb{R}^n$ is a *regular point* if it is feasible and if the set of constraint gradients

$$\left\{ \nabla h_i(x^*) : 1 \le i \le m \right\} \cup \left\{ \nabla g_j(x^*) : j \in \mathcal{A}(x^*) \right\}$$

are linearly independent.

- ***Definition***: The cone $\mathcal{V}^{\text{EQ}}(x^*)$ at the point $x^* \in \mathbb{R}^n$ is the set of vectors $d \in \mathbb{R}^n$ such that

$$\begin{aligned} d \cdot \nabla h_i(x^*) &= 0 & \forall 1 \le i \le m \\ d \cdot \nabla g_j(x^*) &= 0 & \forall j \in \mathcal{A}(x^*) \end{aligned}$$

Intuitively, it is the set of points in which we can move while keeping the equality constraints and the active inequality constraints satisfied.

- ***Theorem (Karush-Kuhn-Tucker Necessary Conditions)***: If $\boldsymbol{x}^*$ is a local minimum that is a regular point, then there exists unique Lagrange multiplier vectors $\boldsymbol{\lambda}^* \in \mathbb{R}^m$ and $\boldsymbol{\mu}^* \in \mathbb{R}^r$ such that

$$\nabla f(\boldsymbol{x}^*) + \sum_{i=1}^{m} \lambda_i^* \nabla h_i(\boldsymbol{x}^*) + \sum_{j=1}^{r} \mu_i^* \nabla g_j(\boldsymbol{x}^*) = \boldsymbol{0}$$
$$\boldsymbol{\mu}_j \geq \boldsymbol{0}$$
$$\mu_j^* = 0 \qquad \forall j \notin \mathcal{A}(\boldsymbol{x}^*)$$

In addition, if *f*, $\boldsymbol{h}$ and $\boldsymbol{g}$ are twice continuously differentiable, then

$$\boldsymbol{d}^\top \left( \nabla^2 f(\boldsymbol{x}^*) + \sum_{i=1}^{m} \lambda_i^* \nabla^2 h_i(\boldsymbol{x}^*) + \sum_{j=1}^{r} \mu_j^* \nabla^2 g_j(\boldsymbol{x}^*) \right) \boldsymbol{d} \geq 0 \qquad \forall \boldsymbol{d} \in \mathcal{V}^{\text{EQ}}(\boldsymbol{x}^*)$$

There is an interesting geometrical explanation of the first-order conditions. First, consider that

- o Because the *h* are equalities, we cannot move along $\nabla h_i(\boldsymbol{x})$ [the perpendiculars to the equality constraints] in *any* direction (either forwards or backwards).

- o Because the active constraints are $\leq$ inequalities, we can move along $\nabla g_i(\boldsymbol{x}), i \in \mathcal{A}(\boldsymbol{x})$ [the perpendiculars to the active inequality constraints] only in the *negative* direction. Moving along $+\nabla g_i(\boldsymbol{x})$ would leave the feasible region.

- o Inactive constraints are not tight, and so we can move both forwards and backwards along $\nabla g_i(\boldsymbol{x}), i \notin \mathcal{A}(\boldsymbol{x})$.

All the first-order condition states is that there is no way to move in such a way that we remain in the feasible region and decrease $f(\boldsymbol{x}^*)$. Mathematically, the conditions states that $-\nabla f(\boldsymbol{x}^*)$ [the direction which would decrease $f(\boldsymbol{x}^*)$] must be composed entirely of

- o A linear combination of the perpendiculars to the equality constraints (along which we cannot move).

- o A *positive* linear combination of the perpendiculars to the active inequality constraints (along which we cannot move).

- o [No perpendiculars to the inactive inequality constraints, because we *could* move along those].

***Proof.*** Everything follows from applying Lagrange multipliers to the equality constrained program, except for the fact that $\mu_j^* \geq 0$ for all $j \in \mathcal{A}(\boldsymbol{x}^*)$.

To see why this is the case, assume that this does *not* hold for some $j \in \mathcal{A}(\boldsymbol{x}^*)$. Then let $\mathcal{C}_j \subset \mathbb{R}^n$ be the set of points feasible for all *other* active constraints, and let $\mathcal{V}_j^{\text{EQ}}(\boldsymbol{x}^*)$ be the corresponding cone of first-order feasible directions:

$$\mathcal{C}_j = \Big\{ \boldsymbol{x} : \boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0}, g_\ell(\boldsymbol{x}) = 0 \ \forall \ell \in \mathcal{A}(\boldsymbol{x}^*) \setminus \{j\} \Big\}$$
$$\mathcal{V}_j^{\text{EQ}}(\boldsymbol{x}^*) = \Big\{ \boldsymbol{d} : \boldsymbol{d} \cdot \nabla \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0}, \boldsymbol{d} \cdot \nabla g_\ell(\boldsymbol{x}^*) = 0 \ \forall \ell \in \mathcal{A}(\boldsymbol{x}^*) \setminus \{j\} \Big\}$$

By regularity, there must exist a $\boldsymbol{d} \in \mathcal{V}_j^{\text{EQ}}(\boldsymbol{x}^*)$ with $\boldsymbol{d} \cdot \nabla g_j(\boldsymbol{x}^*) < 0$ (if it were the case that $\boldsymbol{d} \cdot \nabla g_j(\boldsymbol{x}^*) = 0 \ \forall \boldsymbol{d} \in \mathcal{V}_j^{\text{EQ}}(\boldsymbol{x}^*)$, then the point would not be regular, and if $\boldsymbol{d}' \cdot \nabla g_j(\boldsymbol{x}^*) > 0$ for some $\boldsymbol{d}' \in \mathcal{V}_j^{\text{EQ}}(\boldsymbol{x}^*)$, then just take $\boldsymbol{d} = -\boldsymbol{d}' \in \mathcal{V}_j^{\text{EQ}}(\boldsymbol{x}^*)$).

Then, by the regularity lemma, there exists a curve $\boldsymbol{x}(t) \in \mathcal{C}_j$ with $\boldsymbol{x}(0) = \boldsymbol{x}^*$, $\dot{\boldsymbol{x}}(0) = \boldsymbol{d}$ and such that $\boldsymbol{x}(t)$ is feasible for small $t$. Now, if we let $\ell(t) = f[\boldsymbol{x}(t)]$, we find that

$$\dot{\ell}(t) = \dot{\boldsymbol{x}}(t)^\top \nabla f[\boldsymbol{x}(t)]$$

Evaluating at $t = 0$

$$\begin{aligned}
\dot{\ell}(0) &= \boldsymbol{d}^\top \nabla f(\boldsymbol{x}^*) \\
&= -\boldsymbol{d}^\top \Big\{ \sum_{i=1}^m \lambda_i^* h_i(\boldsymbol{x}^*) + \sum_{j=1}^m \mu_j^* g_j(\boldsymbol{x}^*) \Big\} \\
&= -\boldsymbol{d}^\top \nabla g_j(\boldsymbol{x}^*) \mu_j^* \\
&< 0
\end{aligned}$$

This contradicts the local optimality of $\boldsymbol{x}^*$.

- Practically, we typically define a Lagrangian

$$\mathcal{L}\big(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}\big) = f(\boldsymbol{x}) + \boldsymbol{\lambda} \cdot \boldsymbol{h}(\boldsymbol{x}) + \boldsymbol{\mu} \cdot \boldsymbol{g}(\boldsymbol{x})$$

The first-order conditions then reduce to

$$\nabla_x \mathcal{L}\big(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\big) = \boldsymbol{0} \quad \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0} \quad \boldsymbol{g}(\boldsymbol{x}^*) \leq \boldsymbol{0}$$
$$\boldsymbol{\mu}^* \geq \boldsymbol{0}$$
$$\mu_j^* g_j(\boldsymbol{x}^*) = 0 \qquad \forall 1 \leq j \leq r$$

For any given problem, then, we

  o  Prove global minima exist

  o  Find the set of $(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ satisfying these FOC

o   Find the set of non-regular points

o   Choose the point with the lowest objective.

For example, consider the problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^3} \tfrac{1}{2}(x_1^2 + x_2^2 + x_3^2) \text{ s.t. } x_1 + x_2 + x_3 \le -3$$

The objective and constraints are continuously differentiable, and minima exist (by coerciveness). The Lagrangian is

$$\mathcal{L}\left(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}\right) = \tfrac{1}{2}(x_1^2 + x_2^2 + x_3^2) + \mu(x_1 + x_2 + x_3 + 3)$$

The first-order conditions are

$$\nabla_x \mathcal{L}\left(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\right) = \boldsymbol{0} \quad \Rightarrow \quad x_1^* + \mu^* = x_2^* + \mu^* = x_3^* + \mu^* = 0$$
$$\boldsymbol{g}(\boldsymbol{x}^*) \le \boldsymbol{0} \quad \Rightarrow \quad x_1^* + x_2^* + x_3^* \le -3$$
$$\mu_j^* g_j(\boldsymbol{x}^*) = 0 \quad \Rightarrow \quad \mu^*(x_1^* + x_2^* + x_3^* + 3) = 0$$

The solution is $\boldsymbol{x}^* = (-1, -1, -1)$ and $\mu^* = 1$ which satisfies $\mu \ge 0$. Furthermore, all points are regular, so this is the global minimum.

- ***Theorem (KKT Sufficient Conditions)***: Assume that $f$, $\boldsymbol{h}$ and $\boldsymbol{g}$ are twice continuously differentiable, and that $\boldsymbol{x}^* \in \mathbb{R}^n, \boldsymbol{\lambda}^* \in \mathbb{R}^m, \boldsymbol{\mu}^* \in \mathbb{R}^r$ satisfy

$$\nabla_x \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \boldsymbol{0} \quad \boldsymbol{h}(\boldsymbol{x}^*) = \boldsymbol{0} \quad \boldsymbol{g}(\boldsymbol{x}^*) \le \boldsymbol{0}$$
$$\boldsymbol{\mu}^* \ge \boldsymbol{0} \quad \mu_j^* = 0 \; \forall j \notin \mathcal{A}(\boldsymbol{x}^*)$$
$$\boldsymbol{d}^\top \nabla_{xx}^2 \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \boldsymbol{d} > 0 \quad \forall \boldsymbol{d} \in \mathcal{V}^{\mathrm{EQ}}(\boldsymbol{x}^*) \setminus \{\boldsymbol{0}\}$$

Assume also that

$$\mu_j^* > 0 \qquad \forall j \in \mathcal{A}(\boldsymbol{x}^*)$$

Then $\boldsymbol{x}^*$ is a strict local minimum.

***Proof***: We follow the equality case. Suppose $\boldsymbol{x}^*$ is not a strict local minimum. Then the exists $\{\boldsymbol{x}_k\} \subset \mathbb{R}^n, \boldsymbol{h}(\boldsymbol{x}_k) = \boldsymbol{0}, \boldsymbol{g}(\boldsymbol{x}_k) \le \boldsymbol{0}, \boldsymbol{x}_k \ne \boldsymbol{x}^*, \boldsymbol{x}_k \to \boldsymbol{x}^*$ with $f(\boldsymbol{x}_k) \le f(\boldsymbol{x}^*)$. We define

$$\boldsymbol{d}_k = \frac{\boldsymbol{x}_k - \boldsymbol{x}^*}{\left\|\boldsymbol{x}_k - \boldsymbol{x}^*\right\|} \qquad\qquad \delta_k = \left\|\boldsymbol{x}_k - \boldsymbol{x}^*\right\|$$

Without loss of generality, assume $\boldsymbol{d}_k \to \boldsymbol{d}$. Using the same mean-value-theorem argument as in the sufficient conditions proof for Lagrange multipliers, we find that $\nabla \boldsymbol{h}(\boldsymbol{x}^*)^\top \boldsymbol{d} = \boldsymbol{0}$.

Now, consider the $\boldsymbol{g}$. If $j \in \mathcal{A}(\boldsymbol{x}^*)$, then $g_j(\boldsymbol{x}^*) = 0$, and so since $\boldsymbol{x}_k$ is feasible, $g_j(\boldsymbol{x}_k) - g_j(\boldsymbol{x}^*) \le 0$. Furthermore, by the mean-value-theorem,

$$g_j(\boldsymbol{x}_k) - g_j(\boldsymbol{x}^*) = \nabla g_j(\tilde{\boldsymbol{x}}_k)^\top (\boldsymbol{x}_k - \boldsymbol{x}^*) \to \nabla g_j(\boldsymbol{x}^*)^\top \boldsymbol{d}$$

for some $\tilde{\boldsymbol{x}}_k \in [\boldsymbol{x}_k, \boldsymbol{x}^*]$. Together, these imply that

$$\nabla g_j(\boldsymbol{x}^*)^\top \boldsymbol{d} \le 0$$

If the inequality is tight for all $j \in \mathcal{A}(\boldsymbol{x}^*)$, then $\boldsymbol{d} \in \mathcal{V}^{\mathrm{EQ}}(\boldsymbol{x}^*)$, and we proceed as in the equality case.

If, on the other hand, $\nabla g_j(\boldsymbol{x}^*)^\top \boldsymbol{d} < 0$ for some $j \in \mathcal{A}(\boldsymbol{x}^*)$, then

$$\boldsymbol{d}^\top \nabla f(\boldsymbol{x}^*) = -\boldsymbol{d}^\top \left[ \nabla \boldsymbol{h}(\boldsymbol{x}^*) \cdot \boldsymbol{\lambda}^* + \nabla \boldsymbol{g}(\boldsymbol{x}^*) \cdot \boldsymbol{\mu}^* \right] > 0$$

[The inequality holds because $\nabla g_j(\boldsymbol{x}^*)^\top \boldsymbol{d} < 0$ whereas $\boldsymbol{d}$ dotted with every other gradient of $\boldsymbol{h}$ and $\boldsymbol{g}$ is 0, as we found above]. However, by the definition of our sequence, $f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) \le 0$, and by the mean value theorem,

$$f(\boldsymbol{x}_k) - f(\boldsymbol{x}^*) = \nabla f(\tilde{\boldsymbol{x}}_k)^\top (\boldsymbol{x}_k - \boldsymbol{x}^*) \to \nabla f(\boldsymbol{x}^*)^\top \boldsymbol{d}$$

For some $\tilde{\boldsymbol{x}}_k \in [\boldsymbol{x}_k, \boldsymbol{x}^*]$. This implies that $\nabla f(\boldsymbol{x}^*)^\top \boldsymbol{d} \le 0$. This, together with the statement $\boldsymbol{d}^\top \nabla f(\boldsymbol{x}^*) > 0$, is a contradiction.

- Practical point of note: the problem above was a **min** problem with $\le$ inequalities, and it resulted in a *positive* Lagrange multiplier. If either of these facts changed (eg: max and/or $\ge$), the sign of the Lagrange multiplier required would be flipped. For consistency, it is then customary to change the way we define the Lagrangian to ensure the Lagrange multiplier remains positive; specifically, the term $+\mu g_i(\boldsymbol{x})$ changes to $-\mu g_i(\boldsymbol{x})$.

- *Sensitivity analysis*

  o Consider the optimization problem

  $$\min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) \text{ s.t. } \boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{u}, \boldsymbol{g}(\boldsymbol{x}) \le \boldsymbol{v}$$

  o *Theorem*: Suppose there exists a triple $(\boldsymbol{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ satisfying the second order sufficient conditions when $(\boldsymbol{u}, \boldsymbol{v}) = (\boldsymbol{0}, \boldsymbol{0})$, with $\boldsymbol{x}^*$ regular. Then there exists a neighborhood $N$ of $(\boldsymbol{u}, \boldsymbol{v}) = (\boldsymbol{0}, \boldsymbol{0})$ and a function $\boldsymbol{x}^*(\cdot, \cdot)$ defined on $N$ such that

    ▪ $\boldsymbol{x}^*(\boldsymbol{0}, \boldsymbol{0}) = \boldsymbol{x}^*$, and for each $(\boldsymbol{u}, \boldsymbol{v}) \in N$, $\boldsymbol{x}^*(\boldsymbol{u}, \boldsymbol{v})$ is a strict local minimum.

- $\boldsymbol{x}^*(\cdot,\cdot)$ is continuously differentiable.

- If $p(\boldsymbol{u},\boldsymbol{v}) = f[\boldsymbol{x}^*(\boldsymbol{u},\boldsymbol{v})]$

$$\nabla_{\boldsymbol{u}} p(\boldsymbol{0},\boldsymbol{0}) = -\boldsymbol{\lambda}^* \qquad \nabla_{\boldsymbol{v}} p(\boldsymbol{0},\boldsymbol{0}) = -\boldsymbol{\mu}^*$$

- *Existence of Lagrange Multipliers*

  o **Theorem (Farkas' Lemma)**: Consider a matrix $A \in \mathbb{R}^{m \times n}$. Then a vector $\boldsymbol{z} \in \mathbb{R}^m$ satisfies

  $$\boldsymbol{z} \cdot \boldsymbol{y} \leq 0 \text{ for all } \boldsymbol{y} \in \mathbb{R}^m \text{ such that } A^\top \boldsymbol{y} = \boldsymbol{0}$$

  if and only if

  $$\boldsymbol{z} = A\boldsymbol{x} \text{ for some } \boldsymbol{x} \in \mathbb{R}^n \text{ with } \boldsymbol{x} \geq \boldsymbol{0}$$

  o **Theorem**: Consider the standard inequality optimization program above. Let $\boldsymbol{x}^*$ be a local minimum. Then there exists Lagrange multipliers $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ satisfying the FOCs

  $$\nabla f(\boldsymbol{x}^*) + \nabla h(\boldsymbol{x}^*)\boldsymbol{\lambda}^* + \nabla g(\boldsymbol{x}^*)\boldsymbol{\mu}^* = \boldsymbol{0}$$
  $$\boldsymbol{\mu}^* \geq \boldsymbol{0} \qquad \mu_j^* g_j(\boldsymbol{x}^*) = 0 \quad \forall 1 \leq j \leq r$$

  If and only if

  $$\mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{V}^{\mathrm{EQ}}(\boldsymbol{x}^*) = \varnothing$$

  **Proof**: Suppose there are no equality constraints. Then $\mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{V}^{\mathrm{EQ}}(\boldsymbol{x}^*) = \varnothing$ is equivalent to

  $$\boldsymbol{d} \notin \mathcal{D}(\boldsymbol{x}^*) \text{ for all } \boldsymbol{d} \text{ such that } \boldsymbol{d} \in \mathcal{V}^{\mathrm{EQ}}(\boldsymbol{x}^*)$$

  {$\boldsymbol{d}$ must be "uphill} for every $\boldsymbol{d}$ {that keeps us in the feasible set}

  $$\nabla f(\boldsymbol{x}^*) \cdot \boldsymbol{d} \geq 0 \text{ for all } \boldsymbol{d} \text{ such that } \nabla g_j(\boldsymbol{x}^*) \cdot \boldsymbol{d} \leq 0 \ \forall j \in \mathcal{A}(\boldsymbol{x}^*)$$

  By Farkas' Lemma, this is equivalent to

  $$\nabla f(\boldsymbol{x}^*) + \nabla g(\boldsymbol{x}^*)\boldsymbol{\mu} = \boldsymbol{0}$$

  For some $\boldsymbol{\mu}$ with $\boldsymbol{\mu} \geq \boldsymbol{0}$ and $\mu_j = 0$ if $j \notin \mathcal{A}(\boldsymbol{x}^*)$.

  If there are equality constraints $\boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0}$, we can replace them with inequality constraints $\boldsymbol{h}(\boldsymbol{x}) \leq \boldsymbol{0}$ and $-\boldsymbol{h}(\boldsymbol{x}) \leq \boldsymbol{0}$.

  o **Constraint qualification**: In summary, what we have proved so far is

$\boldsymbol{x}^*$ local minimum

$$\Rightarrow \mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{T}(\boldsymbol{x}^*) = \varnothing$$
$$\Rightarrow \mathcal{D}(\boldsymbol{x}^*) \cap \mathcal{V}(\boldsymbol{x}^*) = \varnothing$$
$$\Rightarrow \text{Lagrange multipliers exist}$$

In proving this sequence of statements, we made use of the fact that $\boldsymbol{x}^*$ was a regular point. However, it is possible to prove the existence of Lagrange multipliers under weaker assumptions called *constraint qualifications*. If the constraints are linear, for example, Lagrange multipliers are guaranteed to exist. The weakest form of constraint qualification is *quasiregularity*, which requires that $\mathcal{V}(\boldsymbol{x}^*) = \mathcal{T}(\boldsymbol{x}^*)$.

o **Theorem (Linear Constraint Qualification)**: Consider the optimization program

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) \text{ s.t. } A\boldsymbol{x} \leq \boldsymbol{b}$$

Suppose that $\boldsymbol{x}^*$ is a local minimum. Then $\boldsymbol{x}^*$ is quasi-regular and Lagrange multipliers exist. This trivially also applies to linear equality constraints, and can be extended to linear equality constraints and concave inequality constraints.

o **Theorem (General Sufficiency Condition)**: Consider the optimization program

$$\min_{\boldsymbol{x} \in \Omega} f(\boldsymbol{x}) \text{ s.t. } \boldsymbol{g}(\boldsymbol{x}) \leq \boldsymbol{0}$$

Let $\boldsymbol{x}^* \in \mathbb{R}^n$ be a feasible point, and $\boldsymbol{\mu}^* \in \mathbb{R}^r$ be a vector such that

$$\boldsymbol{\mu}^* \geq \boldsymbol{0}$$
$$\mu_j^* = 0 \qquad \forall j \notin \mathcal{A}(\boldsymbol{x}^*)$$
$$\boldsymbol{x}^* \in \operatorname{argmin}_{\boldsymbol{x} \in \Omega} L(\boldsymbol{x}, \boldsymbol{\mu}^*)$$

Then $\boldsymbol{x}^*$ is a global minimum. Note that no differentiability or continuity assumptions were made.

**Proof**:

$$f(\boldsymbol{x}^*) = f(\boldsymbol{x}^*) + \boldsymbol{\mu}^* \cdot \boldsymbol{g}(\boldsymbol{x}^*)$$
$$= \min_{\boldsymbol{x} \in \Omega} \left\{ f(\boldsymbol{x}) + \boldsymbol{\mu}^* \boldsymbol{g}(\boldsymbol{x}) \right\}$$
$$\leq \min_{\boldsymbol{x} \in \Omega, \boldsymbol{g}(\boldsymbol{x}) \leq 0} \left\{ f(\boldsymbol{x}) + \boldsymbol{\mu}^* \boldsymbol{g}(\boldsymbol{x}) \right\}$$
$$\leq \min_{\boldsymbol{x} \in \Omega, \boldsymbol{g}(\boldsymbol{x}) \leq 0} f(\boldsymbol{x})$$

## Optimization over Convex Sets

- Consider the problem $\min f(\boldsymbol{x})$ subject to $\boldsymbol{x} \in \mathcal{C} \subset \mathbb{R}^n$ where $\mathcal{C}$ is a convex set.

- ***Theorem (necessary conditions)***: If $\boldsymbol{x}^*$ is a local minimum and $f$ is continuously differentiable in a neighborhood of $\boldsymbol{x}^*$, then

$$\nabla f(\boldsymbol{x}^*) \cdot (\boldsymbol{x} - \boldsymbol{x}^*) \geq 0 \qquad \forall \boldsymbol{x} \in \mathcal{C}$$

Geometrically, at a local minimum, the gradient of the objective function must make an acute angle with any improving direction (or else, we could improve along that direction), and since the set is convex, $\boldsymbol{x} - \boldsymbol{x}^*$ must be an improving direction for any $\boldsymbol{x} \in \mathcal{C}$, since the line between $\boldsymbol{x}$ and $\boldsymbol{x}^*$ is in $\mathcal{C}$:



***Proof***: Suppose that there exists an $\boldsymbol{x} \in \mathcal{C}$ with $\nabla f(\boldsymbol{x}^*) \cdot (\boldsymbol{x} - \boldsymbol{x}^*) < 0$. Consider applying the mean value theorem to the function $g(\varepsilon) = f\left[\boldsymbol{x}^* + \varepsilon(\boldsymbol{x} - \boldsymbol{x}^*)\right]$:

$$g(\varepsilon) - g(0) = g'(c)\left[\varepsilon - 0\right]$$
$$f\left[\boldsymbol{x}^* + \varepsilon(\boldsymbol{x} - \boldsymbol{x}^*)\right] = f(\boldsymbol{x}^*) + \varepsilon \nabla f\left(\boldsymbol{x}^* + \varepsilon s(\boldsymbol{x} - \boldsymbol{x}^*)\right) \cdot (\boldsymbol{x} - \boldsymbol{x}^*)$$

Where $s \in (0,1)$. Since $\nabla f$ is continuous, and $\nabla f(\boldsymbol{x}^*) \cdot (\boldsymbol{x} - \boldsymbol{x}^*) < 0$, then fall all sufficiently small $\varepsilon$, this will also be true of $\nabla f\left(\boldsymbol{x}^* + \varepsilon s(\boldsymbol{x} - \boldsymbol{x}^*)\right) \cdot (\boldsymbol{x} - \boldsymbol{x}^*) < 0$. This, however, implies that $f(\boldsymbol{x}^*) > f\left[\boldsymbol{x}^* + \varepsilon(\boldsymbol{x} - \boldsymbol{x}^*)\right]$. Since $\mathcal{C}$ is convex, this point is in $\mathcal{C}$, and so this contradicts the optimality of $\boldsymbol{x}^*$.

- ***Theorem (optimality for convex optimization)***: Suppose $f : \mathbb{R}^n \to \mathbb{R}$ is convex over $\mathcal{C}$. Then

    1. Any local minimum of $f$ is also a global minimum.

    2. If $f$ is strictly convex, then there exists at most one global minimum.

    ***Proof***: Consider:

    1. Suppose $\boldsymbol{x}^*$ is a local minimum, and there exists some $\boldsymbol{x} \neq \boldsymbol{x}^*$ with $f(\boldsymbol{x}) < f(\boldsymbol{x}^*)$. Then for any $\lambda \in [0,1)$, then

$$f\left[\lambda \boldsymbol{x}^* + (1-\lambda)\boldsymbol{x}\right] \le \lambda f(\boldsymbol{x}^*) + (1-\lambda)f(\boldsymbol{x}) < f(\boldsymbol{x}^*)$$

Now:

- Since $\mathcal{C}$ is convex, $\lambda \boldsymbol{x}^* + (1-\lambda)\boldsymbol{x} \in \mathcal{C}$

- For every small $r$, there will be a $\lambda$ such that $\lambda \boldsymbol{x}^* + (1-\lambda)\boldsymbol{x} \in B_r(\boldsymbol{x}^*)$
  . The fact the value of $f$ at that point is more than at $\boldsymbol{x}^*$ contradicts our local optimality assumption.

2. Suppose $\boldsymbol{x}_0 \ne \boldsymbol{x}_1$ are two global minima. Then

$$f\left(\tfrac{1}{2}(\boldsymbol{x}_0 + \boldsymbol{x}_1)\right) < \tfrac{1}{2}\left(f(\boldsymbol{x}_0) + f(\boldsymbol{x}_1)\right) = f(\boldsymbol{x}_0) = f(\boldsymbol{x}_1)$$

   This means there is a point at which $f$ is strictly less than at $\boldsymbol{x}_0$ and $\boldsymbol{x}_1$, which contradicts our local optimality assumption.

- ***Theorem (necessary & sufficient conditions)***: If $f : \mathbb{R}^n \to \mathbb{R}$ is convex over $\mathcal{C}$ and differentiable and $\boldsymbol{x}^* \in \mathcal{C}$ is a feasible point, then $\boldsymbol{x}^*$ is globally optimal if and only if

$$\nabla f(\boldsymbol{x}^*) \cdot (\boldsymbol{x} - \boldsymbol{x}^*) \ge 0 \qquad \forall \boldsymbol{x} \in \mathcal{C}$$

   ***Proof***: Necessity follows from the previous theorem. For sufficiency, note that

$$f(\boldsymbol{x}) \ge f(\boldsymbol{x}^*) + \nabla f(\boldsymbol{x}^*) \cdot (\boldsymbol{x} - \boldsymbol{x}^*) \ge f(\boldsymbol{x}^*) \qquad \forall \boldsymbol{x} \in \mathcal{C}$$

## Duality

- ***Supporting & Separating hyperplanes***
  - ***Definition***: The hyperplane $\left\{\boldsymbol{x} \in \mathbb{R}^n : \boldsymbol{\mu} \cdot \boldsymbol{x} = b\right\}$ with $\boldsymbol{\mu} \in \mathbb{R}^n \setminus \{0\}$ *supports* the convex set $\mathcal{C} \subset \mathbb{R}^n$ at the point $\bar{\boldsymbol{x}}$ if

$$\boldsymbol{\mu} \cdot \boldsymbol{x} \ge \boldsymbol{\mu} \cdot \bar{\boldsymbol{x}} = b \qquad \forall \boldsymbol{x} \in \mathcal{C}$$

    or equivalently

$$\inf_{\boldsymbol{x} \in \mathcal{C}} \boldsymbol{\mu} \cdot \boldsymbol{x} \ge \boldsymbol{\mu} \cdot \bar{\boldsymbol{x}} = b$$

    This statement requires that (a) $\bar{\boldsymbol{x}}$ lies on the hyperplane, and (b) $\mathcal{C}$ lies entirely on one side of the hyperplane.

o **Theorem (supporting hyperplane)**: Let $\mathcal{C} \subset \mathbb{R}^n$ be a convex set and $\bar{x} \in \mathbb{R}^n$ be a point that is not in the interior of $\mathcal{C}$. Then there exists a supporting hyperplane at $\bar{x}$.

**Proof**: The intuition behind this proof will be to consider a sequence $\{x_k\}$ of points *outside* the closure of $\mathcal{C}$ converging to $\bar{x}$. This is possible, since $\bar{x} \notin \text{int}\,\mathcal{C}$. We consider the projection of each of these $x_k$ onto $\bar{\mathcal{C}}$ (denoted $\hat{x}_k$), and we consider the hyperplane *perpendicular* to the line from $x_k$ to $\hat{x}_k$:



We then show that

- The set $\mathcal{C}$ lies on one side of each of these hyperplanes.

- These hyperplanes tend to the supporting hyperplane of interest. In other words, $\bar{x}$ lies on the limiting supporting hyperplane.

Now, to the formal stuff. Define $\bar{\mathcal{C}} = \text{cl}\,\mathcal{C}$, and note that $\bar{\mathcal{C}}$ is convex. Let $\{x_k\}$ be a sequence of points such that $x_k \notin \bar{\mathcal{C}}, x_k \to \bar{x}$, and for each $x_k$, let $\hat{x}_k$ be its projection onto $\bar{\mathcal{C}}$. Then, by the projection theorem,

$$\left(\hat{x}_k - x_k\right) \cdot \left(x - \hat{x}_k\right) \geq 0 \qquad \forall x \in \bar{\mathcal{C}}$$

And this means that for all $k$ and $x \in \bar{\mathcal{C}}$,

$$\left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \boldsymbol{x} \geq \left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \hat{\boldsymbol{x}}_k$$
$$= \left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) + \left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \boldsymbol{x}_k$$
$$\geq \left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \boldsymbol{x}_k$$

More succinctly

$$\left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \boldsymbol{x} \geq \left(\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right) \cdot \boldsymbol{x}_k = \text{constant} \qquad \forall \boldsymbol{x} \in \overline{\mathcal{C}}$$

In other words, $\overline{\mathcal{C}}$ lies on one side of each of those hyperplanes.

But now, set

$$\boldsymbol{\mu}_k = \frac{\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k}{\left\|\hat{\boldsymbol{x}}_k - \boldsymbol{x}_k\right\|}$$

Then the equation above can be written as

$$\boldsymbol{\mu}_k \cdot \boldsymbol{x} \geq \boldsymbol{\mu}_k \cdot \boldsymbol{x}_k \qquad \forall k, x \in \mathcal{C}$$

Since $\left\|\boldsymbol{\mu}_k\right\| = 1$, the sequence $\{\boldsymbol{\mu}_k\}$ is bounded has a non-zero subsequential limit

$\boldsymbol{\mu}$. Letting $k \to \infty$, we get $\boldsymbol{\mu}_k \to \boldsymbol{\mu}$ and $\boldsymbol{x}_k \to \overline{\boldsymbol{x}}$, so

$$\boldsymbol{\mu} \cdot \boldsymbol{x} \geq \boldsymbol{\mu} \cdot \overline{\boldsymbol{x}} \qquad \forall k, \boldsymbol{x} \in \overline{\mathcal{C}}$$

As required.

o **Theorem (separating hyperplane)**: Let $\mathcal{C}_1, \mathcal{C}_2 \in \mathbb{R}^n$ be two disjoint non-empty convex sets. Then there exists a hyperplane that separates them; ie: a vector $\boldsymbol{\mu} \in \mathbb{R}^n, \boldsymbol{\mu} \neq 0$ and a scalar $b \in \mathbb{R}$ with

$$\boldsymbol{\mu} \cdot \boldsymbol{x}_1 \leq b \leq \boldsymbol{\mu} \cdot \boldsymbol{x}_2 \qquad \forall \boldsymbol{x}_1 \in \mathcal{C}_1, \boldsymbol{x}_2 \in \mathcal{C}_2$$



**Proof**: Consider the convex set

$$\mathcal{D} = \mathcal{C}_1 - \mathcal{C}_2 = \left\{x_1 - x_2 : x_1 \in \mathcal{C}_1, x_2 \in \mathcal{C}_2\right\}$$

Since the two sets are disjoint, $\boldsymbol{0} \notin \mathcal{D}$. Thus, by the supporting hyperplane theorem, there exists a vector $\boldsymbol{\mu} \neq \boldsymbol{0}$ with

$$\boldsymbol{0} \leq \boldsymbol{\mu}^\top \left( \boldsymbol{x}_1 - \boldsymbol{x}_2 \right) \qquad \forall \boldsymbol{x}_1 \in \mathcal{C}_1, \boldsymbol{x}_2 \in \mathcal{C}_2$$

Setting $b = \sup_{\boldsymbol{x}_2 \in \mathcal{C}_2} \boldsymbol{\mu}^\top \boldsymbol{x}_2$, we obtain the desired result.

o **Theorem (strictly separating hyperplane)**: Let $\mathcal{C} \subset \mathbb{R}^n$ be a *closed* convex set and $\bar{\boldsymbol{x}} \notin \mathcal{C}$ a point. Then there exists a hyperplane that *strictly* separates the point and the set. In other words, $\exists \boldsymbol{\mu} \in \mathbb{R}^n \setminus \{0\}, b \in \mathbb{R}$ such that

$$\boldsymbol{\mu} \cdot \bar{\boldsymbol{x}} < b < \inf_{\boldsymbol{x} \in \mathcal{C}} \boldsymbol{\mu} \cdot \boldsymbol{x}$$



**Proof**: Define $r = \min_{\boldsymbol{x} \in \mathcal{C}} \left\| \boldsymbol{x} - \bar{\boldsymbol{x}} \right\|$. This will be $> 0$ because $\mathcal{C}$ is closed, and $\bar{\boldsymbol{x}} \notin \mathcal{C}$. Now, let $\bar{\mathcal{C}} = \left\{ \boldsymbol{x} \in \mathbb{R}^n : \left\| \boldsymbol{x} - \bar{\boldsymbol{x}} \right\| \leq r / 2 \right\}$. Clearly, $\mathcal{C}$ and $\bar{\mathcal{C}}$ are disjoint, so we can apply the separating hyperplane theorem. Diagramatically:



o **Corollary**: If $\mathcal{C} \subsetneq \mathbb{R}^n$ is a closed convex set, then it is the intersection of all closed halfspaces that contain it.

**Proof**: Let $\mathcal{H}$ be the collection of all closed halfspaces containing $\mathcal{C}$, and let $\bar{H} = \bigcap_{H \in \mathcal{H}} H$.

Since $\mathcal{C} \neq \mathbb{R}^n$, the strictly separating hyperplane theorem implies that $\mathcal{H}$ is non-empty (since there is a point $\boldsymbol{x} \in \mathbb{R}^n$ and $\boldsymbol{x} \notin \mathcal{C}$) and clearly, $\mathcal{C} \subset \bar{H}$.

Now, suppose there exists an $\boldsymbol{x} \in \bar{H}$ with $\boldsymbol{x} \neq \mathcal{C}$, then

- By the strictly separating hyperplane theorem, there exists a halfplane that strictly separates $\boldsymbol{x}$ and $\mathcal{C}$, with $\boldsymbol{x}$ on the opposite side of the hyperplane as $\mathcal{C}$.

- Thus, $\boldsymbol{x} \notin \bar{H}$. This is a contradiction.

- *Farkas' Lemma*

  - *Lemma (Farkas)*: Consider a matrix $A \in \mathbb{R}^{m \times n}$. Given a vector $\boldsymbol{z} \in \mathbb{R}^m$, the following statements are equivalent
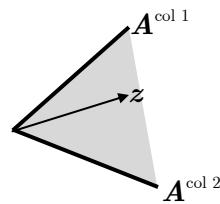
    1.  $\boldsymbol{z} \cdot \boldsymbol{y} \le 0$ for all $\boldsymbol{y} \in \mathbb{R}^m$ with $A^\top \boldsymbol{y} \le \boldsymbol{0}$.

    2.  $\boldsymbol{z} = A\boldsymbol{x}$ for some $\boldsymbol{x} \in \mathbb{R}^n$ with $\boldsymbol{x} \ge \boldsymbol{0}$.

    Geometrically, the two parts of this theorem are as follows:

    1.  $\boldsymbol{z}$ makes an obtuse angle with all vectors $\boldsymbol{y}$ that make an obtuse angle with *every* column of $A$.

    

    2.  $\boldsymbol{z}$ lies in the cone formed by the columns of $A$.

    

    *Proof*: Consider:

    - $\boxed{2 \to 1}$ Dot both sides of the equation with $\boldsymbol{y}$ to get

      $$\boldsymbol{z} \cdot \boldsymbol{y} = (A\boldsymbol{x}) \cdot \boldsymbol{y} = \boldsymbol{x} A^\top \boldsymbol{y} \le 0$$

      Where the last equality follows because $\boldsymbol{x}$ is positive and $A^\top \boldsymbol{y} \le 0$, by assumption.

    - $\boxed{1 \to 2}$ Suppose $\boldsymbol{z}$ satisfies (1) [it makes an obtuse angle to $\boldsymbol{y}$, which makes an obtuse angle to all the columns of $A$], but that there is no $\boldsymbol{x} \ge \boldsymbol{0}$ with $\boldsymbol{z} = A\boldsymbol{x}$ [it lies outside the cone defined by the columns of $A$].

Then define $\mathcal{C} = \left\{A\boldsymbol{x} : \boldsymbol{x} \geq 0\right\}$ to be the closed, convex cone formed by the columns of $A$. Clearly, $\boldsymbol{z} \neq \mathcal{C}$.

By the strictly separating hyperplane theorem, there exists $\boldsymbol{y} \in \mathbb{R}^m \setminus \{0\}$ such that $\boldsymbol{y} \cdot \boldsymbol{z} > \boldsymbol{y} \cdot \boldsymbol{c} \; \forall \boldsymbol{c} \in \mathcal{C}$:



We will show that this $\boldsymbol{y}$ forms an obtuse angle with all the columns of $A$, but an acute angle with $\boldsymbol{z}$ (thus contradicting [1]).

- **_Obtuse angle with the columns of $A$_**: Note that $\lambda \boldsymbol{A}^{\text{col } i} \in \mathcal{C}$, and so $\boldsymbol{y} \cdot \boldsymbol{z} > \lambda \boldsymbol{y} \cdot \boldsymbol{A}^{\text{col } i}$. This means that regardless how far we go along any of the edges of the cone, the point we reach must still be on a different side of the hyperplane to $\boldsymbol{z}$. This implies that the hyperplane must be "titled away" from every edge of the cone (the columns of $A$), and that the normal _must_ make an obtuse angle to the said column. More formally,

$$\boldsymbol{y} \cdot \boldsymbol{A}^{\text{col } i} < \frac{\boldsymbol{y} \cdot \boldsymbol{z}}{\lambda} \to_{\lambda \to \infty} 0$$

- **_Acute angle with $\boldsymbol{z}$_**: Note that since $\boldsymbol{0} \in \mathcal{C}$, $\boldsymbol{y} \cdot \boldsymbol{z} > 0$ and so $\boldsymbol{y}$ must form an acute angle with $\boldsymbol{z}$.

We have therefore reached a contradiction.

o We quickly digress to cover an interesting application of Farkas' Lemma: arbitrage-free pricing.

  ▪ Consider a market with $n$ assets. A portfolio is a vector $\boldsymbol{x} \in \mathbb{R}^n$ in which $x_i$ describes the amount of asset $i$ purchased.

- The current price of asset $i$ is $v_i$, and so the current price of portfolio $\boldsymbol{x}$ is $\boldsymbol{v} \cdot \boldsymbol{x}$.

- The future is modeled by $m$ different scenarios, described by a matrix $R \in \mathbb{R}^{m \times n}$. In scenario $m$, asset $n$ yields $R_{mn}$. The future payoffs are therefore $R\boldsymbol{x}$.

- **Definition**: An *arbitrage opportunity* is a portfolio $\boldsymbol{x}$ such that

$$\boldsymbol{v} \cdot \boldsymbol{x} < 0 \qquad R\boldsymbol{x} \geq \boldsymbol{0}$$

  In other words, the portfolio costs a negative amount to buy right now and will always yield nonnegative returns. We say a market is *consistent* if and only if it has no arbitrage opportunity.

- By Farkas' Lemma, a market is consistent if and only if there exists a vector $\boldsymbol{q} \geq \boldsymbol{0}$ with $\boldsymbol{v} = R^\top \boldsymbol{q}$. Suppose there exists such a vector, and define

$$r = \frac{1}{\boldsymbol{1} \cdot \boldsymbol{q}} - 1 \qquad\qquad \boldsymbol{\pi} = (1 + r)\boldsymbol{q}$$

  $\boldsymbol{\pi}$ is then a probability distribution, because $\boldsymbol{\pi} \geq \boldsymbol{0}$ and $\boldsymbol{1} \cdot \boldsymbol{\pi} = 1$. We also have that

$$v_i = \frac{1}{1+r}\left(\boldsymbol{\pi} \cdot \boldsymbol{R}^{\text{col } i}\right)$$

  In other words, the prices today are the discounted value of the future prices under the distribution $\boldsymbol{\pi}$. We call $\boldsymbol{\pi}$ the *risk-neutral distribution*.

- ***The primal and the dual***
  - **Definition**: Consider the *primal problem* $\min f(\boldsymbol{x})$ s.t. $\boldsymbol{g}(\boldsymbol{x}) \leq \boldsymbol{0}_r, \boldsymbol{x} \in \Omega \subset \mathbb{R}^n$, and define

$$f^* = \inf_{\boldsymbol{x} \in \Omega, \boldsymbol{g}(\boldsymbol{x}) \leq \boldsymbol{0}} f(\boldsymbol{x})$$

    **Assumption**: The feasible set is non-empty and the optimal cost is bounded below. In other words, $-\infty < f^* < \infty$.

    **Definition (*Lagrangian function*)**: For $\boldsymbol{x} \in \Omega, \boldsymbol{\mu} \in \mathbb{R}^r$

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu}) = f(\boldsymbol{x}) + \boldsymbol{\mu} \cdot \boldsymbol{g}(\boldsymbol{x})$$

    We say $\boldsymbol{\mu}^*$ is a *geometric multiplier*

- $\mu^* \geq 0$
- $f^* = \inf_{x \in \Omega} L(x, \mu^*)$

o **Definition**: The set $\mathcal{S} \subset \mathbb{R}^{r+1}$ of *constraint-cost pairs* is defined by

$$\mathcal{S} = \left\{ \big(g(x), f(x)\big) : x \in \Omega \right\}$$

Given $(\mu, \mu_0) \in \mathbb{R}^{r+1} \setminus \{0\}$, the *hyperplane passing through* $(\bar{z}, \bar{w})$ is

$$\left\{ (z, w) : \mu \cdot z + \mu_0 w = \mu \cdot \bar{z} + \mu_0 \bar{w} \right\}$$

The positive/negative halfspace has the $=$ replaced by $\geq/\leq$. The hyperplane is *non-vertical* if $\mu_0 \neq 0$. In such a case, it can be normalized such that $\mu_0 = 1$.

o **Lemma (visualization)**:

- The hyperplane with normal $(\mu, 1)$ that passes through $\big(g(x), f(x)\big)$ intercepts the vertical axis $\left\{ (0, w) \in \mathbb{R}^{r+1} : x \in \mathbb{R} \right\}$ at the level $L(x, \mu)$.



In other words, the simple interpretation of $L(x, \mu)$ is the point at which a given hyperplane crosses the vertical axis.

- Amongst all hyperplanes with normal $(\mu, 1)$ that contain $\mathcal{S}$ in the positive halfspace, the highest interception of the vertical axis is attained at $\inf_{x \in \Omega} L(x, \mu)$.

$$\left(0, \inf_{x \in \Omega} L(x, \bar{\mu})\right)$$

- $\mu^*$ is a geometric multiplier if and only if $\mu^* \geq 0$ and, among all hyperplanes with normal $(\mu^*, 1)$ that contain $\mathcal{S}$ in the positive halfspace, the highest interception of the vertical axis is attained at $f^*$ (in other words $\inf_{x \in \Omega} L(x, \mu^*) = f^*$).



$$\left(0, f^*\right)$$

Here is an example of a problem in which a geometric multiplier does *not* exist:



$$(0, f^*)$$

In this case, there is no hyperplane that contains the whole of $\mathcal{S}$ in one of its halfspaces that also passes through the point $(0, f^*)$. Thus, there is no $\mu$ for which $\inf_{x \in \Omega} L(x, \mu) = f^*$ because, as we saw in part (2),

$\inf_{x \in \Omega} L(x, \mu)$ considers the highest intercept for planes *that contain the whole of $\mathcal{S}$ in one of their halfpsaces.*

**Proof:**

1. The hyperplane is the set $(z, w)$ satisfying

$$\mu \cdot z + w = \mu \cdot g(x) + f(x)$$

Clearly, if $z = 0$, we must have $w = L(x, \mu)$.

2. The hyperplane with normal $(\mu, 1)$ that intercepts the axis at level $c$ is the set $(z, w)$ with

$$\mu \cdot z + w = c$$

If $\mathcal{S}$ lies in the positive halfspace, then

$$L(x, \mu) = \mu \cdot g(x) + f(x) \geq c \qquad \forall x \in \Omega$$

Thus, the maximum intercept is $\inf_{x \in \Omega} L(x, \mu)$.

3. We need to show that $f^* = \inf_{x \in \Omega} L(x, \mu^*)$. It is obvious from part (2) that this is true if and only if among all hyperplanes with normal $(\mu^*, 1)$, the highest interception with the vertical axis is at $f^*$.

o **Theorem:** Let $\mu^*$ be a geometric multiplier. Then $x^*$ is a global minimum if and only if $x^*$ is feasible and

$$x^* \in \operatorname{argmin}_{x \in \Omega} L(x, \mu^*) \qquad\qquad \mu_j^* g_j(x^*) = 0 \quad \forall 1 \leq j \leq r$$

Geometrically, the first statement is that $x^*$ is, indeed, the value of $x$ that minimizes $L$ at that value of $\mu$ (and recall that $f^* = \inf_{x \in \Omega} L(x, \mu^*)$), and the second is that *either* $x^*$ is on the *boundary* of the feasible set (in which case we can "improve no further") or that the geometric multiplier is horizontal (in which case the minimum is attained on the interior of the feasible set).

***Proof***: Assume $x^*$ is a global minimum. Then

$$f^* = f(x^*) \geq f(x^*) + \mu^* \cdot g(x^*) = L(x^*, \mu^*) \geq \inf_{x \in \Omega} L(x, \mu^*) = f^*$$

Since the two ends are equal, we have equality the whole way through, and

$$\mathcal{L}(x^*, \mu^*) = \inf_{x \in \Omega} \mathcal{L}(x, \mu^*)$$

Similarly, we have $f(x^*) \geq f(x^*) + \mu^* \cdot g(x^*)$. Thus, since $\mu^* \geq 0$, we must have

$$\mu_j^* g_j(x^*) = 0$$

Conversely, if $\mathcal{L}(x^*, \mu^*) = \inf_{x \in \Omega} \mathcal{L}(x, \mu^*)$ and $\mu^* \cdot g(x^*)$:

$$f(x^*) = f(x^*) + \mu^* \cdot g(x^*) = L(x^*, \mu^*) = \min_{x \in \Omega} L(x, \mu^*) = f^*$$

o We define the *dual function* $q(\mu) = \inf_{x \in \Omega} \mathcal{L}(x, \mu)$; it is concave over its domain, which is convex [because $q$ is a point-wise minimum of concave – actually linear – functions]. Geometrically, $q$ simply does the following:

- Take a hyperplane with normal $(\mu, 1)$

- Push it up as far as possible until it hits the set $\mathcal{S}$.

- Find the value at which it intercepts the $f(x)$ axis. By the second part of the visualization lemma, this is $q(\mu)$.
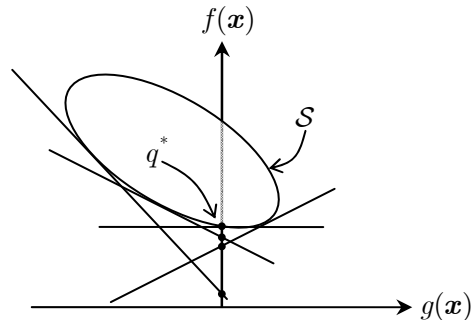
Diagrammatically:



o We then define the *dual problem* as

$$\max q(\mu) \text{ s.t. } \mu \geq 0$$

And the dual optimal value is given by $q^* = \max_{\mu \geq 0} q(\mu)$. [It is possible, however, that $q(\mu) = -\infty$ for all $\mu \geq 0$, in which case the dual problem is infeasible and $q^* = -\infty$].

Geometrically, the dual problem takes all of the snuggly fitting hyperplanes defined by different values of $\boldsymbol{\mu}$ and finds the one with the highest intercept:



[Parenthetically, recall that we proved above that a convex set can be described as the intersection of half-spaces containing it. If this is true, then the dual problem can – in some sense – be seen as "exploring" this set via its halfspaces rather than directly].

o We now consider a number of examples to illustrate the concepts above

- Consider the program

$$\min_{\boldsymbol{x} \in \mathbb{R}_+^2} f(\boldsymbol{x}) = x_1 - x_2 \text{ s.t. } g(\boldsymbol{x}) = x_1 + x_2 - 1 \le 0$$

It turns out the solution to this problem is $f^* = -1, \boldsymbol{x}^* = (0,1)$.

Now, the Lagrangian is given by

$$\mathcal{L}(\boldsymbol{x}, \mu) = x_1 - x_2 + \mu(x_1 + x_2 - 1)$$

with $\mu \ge 0$. The dual function is

$$q(\mu) = \min_{\boldsymbol{x} \in \mathbb{R}_+^2} \mathcal{L}(\boldsymbol{x}, \mu)$$
$$= \min_{\boldsymbol{x} \in \mathbb{R}_+^2} (\mu + 1)x_1 + (\mu - 1)x_2 - \mu$$

Clearly, this is only bounded if both the coefficients of $x_1$ and $x_2$ are positive. This means we need $\mu \ge 1$. Thus

$$q(\mu) = \begin{cases} -\infty & \mu < 1 \\ -\mu & \mu \ge 1 \end{cases}$$

Maximising this function, we clearly find $\mu^* = 1 \Rightarrow q^* = -1 = f^*$. Thus, there is no duality gap.

- Consider the program

$$\min_{x \in \left\{(x_1, x_2) \in \mathbb{R}^2 : x_2 \geq 0\right\}} f(\boldsymbol{x}) = |x_1| + x_2 \text{ s.t. } g(\boldsymbol{x}) = x_1 \leq 0$$

It turns out the solution to this problem is $f^* = 0, \boldsymbol{x}^* = (0,0)$.

Now, the Lagrangian is given by $\mathcal{L}(\boldsymbol{x}^*, \mu) = |x_1| + x_2 + \mu x_1$, with $\mu \geq 0$. If $|\mu| \leq 1$, the first term will dominate over the last term, and $x_1 = 0$ will be the bounded solution. If, on the other hand, $\mu$ falls outside this range, the last term dominate, and the Lagrangian can be shrunk ad infinitum, Thus

$$q(\mu) = \begin{cases} -\infty & |\mu| > 1 \\ 0 & |\mu| \leq 1 \end{cases}$$

The optimal solution is clearly $\mu \in [0,1]$ (since $\mu \geq 0$), with $q^* = 0 = f^*$. So once again, no duality gap. Notice that primal degeneracy leads to dual non-uniqueness; see linear programming notes for more on this.

- Consider the program

$$\min_{x \in \mathbb{R}} f(x) = x \text{ s.t. } g(x) = x^2 \leq 0$$

It turns out the solution to this problem is $f^* = 0, x = 0$.

The Lagrangian is $\mathcal{L}(x, \mu) = x + \mu x^2$ with $\mu \geq 0$ and

$$q(\mu) = \begin{cases} -\frac{1}{4\mu} & \mu > 0 \\ -\infty & \mu \leq 0 \end{cases}$$

In this case, $q^* = 0$, so there is no duality gap, but there is also no optimal dual solution; it is attained as $\mu \to \infty$.

- Consider the program

$$\min_{x \in \{0,1\}} f(x) = -x \text{ s.t. } g(x) = x - \tfrac{1}{2} \leq 0$$

By inspection, the solution is $f^* = 0, x^* = 0$.

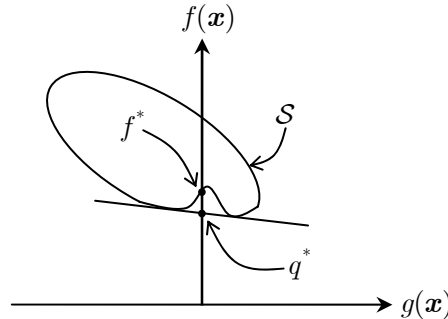The Lagrangian is $\mathcal{L}(x, \mu) = (\mu - 1)x - \tfrac{1}{2}\mu$ with $\mu \geq 0$ and

$$q(\mu) = \begin{cases} -\frac{1}{2}\mu & \mu \geq 1 \\ \frac{1}{2}\mu - 1 & \mu < 1 \end{cases} = \min\left\{-\tfrac{\mu}{2}, \tfrac{\mu}{2} - 1\right\}$$

And so the solution of the dual problem is $q^* = -\tfrac{1}{2}$. There is a duality gap.

- *Weak duality*

  o **Theorem (Weak Duality)**: $q^* \leq f^*$

    Geometrically, this is an obvious statement. Whatever the shape of the set $\mathcal{S}$, the "snug hyperplane" with the highest cutting point will be lower than, or at the same place as $f^*$:

    

    ***Proof***: Consider feasible $\boldsymbol{x}$ and $\boldsymbol{\mu}$; in other words, $\boldsymbol{\mu} \geq \boldsymbol{0}$, $\boldsymbol{x} \in \Omega$ and $\boldsymbol{g}(\boldsymbol{x}) \leq \boldsymbol{0}$.

    $$q(\boldsymbol{\mu}) = \inf\nolimits_{z \in \Omega} L(\boldsymbol{x}, \boldsymbol{\mu}) \leq f(\boldsymbol{x}) + \boldsymbol{\mu}^\top \boldsymbol{g}(\boldsymbol{x}) \leq f(\boldsymbol{x})$$

    Thus

    $$q^* = \sup\nolimits_{\boldsymbol{\mu} \geq \boldsymbol{0}} q(\boldsymbol{\mu}) \leq \inf\nolimits_{x \in \Omega, g(x) \leq 0} f(\boldsymbol{x}) = f^*$$

  o **Definition (Duality gap)**: If $q^* = f^*$, we say there is no *duality gap*. If $q^* < f^*$ then there is a duality gap.

  o **Theorem**: If there is no duality gap, the set of geometric multipliers is equal to the set of dual optimal solutions. If there is a duality gap, the set of geometric multipliers is empty.

    ***Proof***: By definition, $\boldsymbol{\mu}^* \geq \boldsymbol{0}$ is a geometric multiplier if and only if $f^* = \inf\nolimits_{\boldsymbol{x} \in \Omega} \mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu}^*) = q(\boldsymbol{\mu}^*) \leq q^*$. By weak duality, this holds if and only if there is no duality gap.

  o Note also that if the primal problem is unbounded, $f^* = -\infty$. By weak duality, this means that $q(\boldsymbol{\mu}) = -\infty \ \forall \boldsymbol{\mu} \geq \boldsymbol{0}$, and so the dual problem is infeasible. Similarly, if the dual problem is feasible, the primal problem is bounded. However, if the primal is infeasible, we can say nothing about the dual.

- *Primal & Dual Optimality Conditions*

o **Theorem (optimality conditions)**: $(x^*, \mu^*)$ is an optimal solution/geometric multiplier pair to the problem $\min_{x \in \Omega} f(x)$ s.t $g(x) \leq 0$ if and only if

1. <u>PRIMAL FEASIBILITY</u>: $x \in \Omega$ and $g(x) \leq 0$

2. <u>DUAL FEASIBILITY</u>: $\mu^* \geq 0$

3. <u>LAGRANGIAN OPTIMALITY</u>: $x^* \in \text{argmin}_{x \in \Omega} \mathcal{L}(x, \mu^*)$

4. <u>COMPLEMENTARY SLACKNESS</u>: $\mu_j^* g_j^*(x^*) = 0 \quad \forall 1 \leq j \leq r$

Note that this theorem is only useful if there is no duality gap.

**Proof**: Clearly, if $(x^*, \mu^*)$ is an optimal solution/geometric multiplier pair, 1 and 2 must hold. 3 and 4 hold by earlier theorems.

Conversely, if 1-4 hold, then

$$f^* \leq f(x^*) \stackrel{\text{by (4)}}{=} \mathcal{L}(x^*, \mu^*) \stackrel{\text{by (3)}}{=} \min_{x \in \Omega} \mathcal{L}(x, \mu^*) = q(\mu^*) \leq q^*$$

By weak duality, the opposite inequality must hold. Therefore $f^* = q^*$, $f^*$ is primal optimal and $q^*$ is dual optimal.

o **Theorem (Saddle point)**: $(x^*, \mu^*)$ is an optimal solution/geometric multiplier pair of the problem above if and only if $x^* \in \Omega, \mu^* \geq 0$ and $(x^*, \mu^*)$ is a *saddle point* of the Lagrangian in the sense that

$$\mathcal{L}(x^*, \mu) \leq \mathcal{L}(x^*, \mu^*) \leq \mathcal{L}(x, \mu^*) \qquad\qquad \forall x \in \Omega, \mu \geq 0$$

Again, this theorem is only useful if there is no duality gap.

**Proof**: Two parts:

▪ $\boxed{\text{Optimal} \Rightarrow \text{Saddle}}$ If $(x^*, \mu^*)$ is an optimal solution/geometric multiplier pair, then from optimality condition (3), we have

$$\mathcal{L}(x^*, \mu^*) \leq \mathcal{L}(x, \mu^*) \quad \forall x \in \Omega$$

Furthermore, since $g(x) \leq 0$, for any $\mu \geq 0$ we must have $\mathcal{L}(x^*, \mu) \leq f(x) = \mathcal{L}(x^*, \mu^*)$.

▪ $\boxed{\text{Saddle} \Rightarrow \text{Optimal}}$ If $x^* \in \Omega, \mu^* \geq 0$ and $(x^*, \mu^*)$ is a saddle point, then

$$\sup_{\mu \geq 0} \mathcal{L}(x^*, \mu) = \sup_{\mu \geq 0} \left\{ f(x^* + \mu^\top g(x^*)) \right\} = \begin{cases} f(x^*) & \text{if } g(x^*) \leq 0 \\ +\infty & \text{otherwise} \end{cases}$$

We show that each optimality condition holds:

- _Condition 1_: By the definition of a saddle point, $\mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\mu}) \leq \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\mu}^*)$. Therefore, the first of the two options for $\sup_{\boldsymbol{\mu} \geq \boldsymbol{0}} \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\mu})$ options must be correct, and $\boldsymbol{g}(\boldsymbol{x}^*) \leq \boldsymbol{0}$.

- _Condition 2_: The statement of the theorem says that $\boldsymbol{\mu}^* \geq \boldsymbol{0}$.

- _Condition 3_: Obvious, by the definition of a saddle point.

- _Condition 4_: We have established that the first option for $\sup_{\boldsymbol{\mu} \geq \boldsymbol{0}} \mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\mu})$ above is correct. Thus, $\mathcal{L}(\boldsymbol{x}^*, \boldsymbol{\mu}^*) = f(\boldsymbol{x}^*)$ and this means that $\boldsymbol{\mu}^{*\top} \boldsymbol{g}(\boldsymbol{x}^*) = 0$. But since $\boldsymbol{\mu}^* \geq \boldsymbol{0}$ and $\boldsymbol{g}(\boldsymbol{x}^*) \leq \boldsymbol{0}$, we get $\mu_j^* g_j(\boldsymbol{x}^*) = 0$.

- **_Extension to equality constraints_**

  o The theory above can be extended to equality constraints by simply introducing a pair of inequality constraints for each equality constraints. Thus, $\boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{0}$ becomes $\boldsymbol{h}(\boldsymbol{x}) \geq \boldsymbol{0}$ and $\boldsymbol{h}(\boldsymbol{x}) \leq \boldsymbol{0}$.

  o It turns out this is equivalent to simply eliminating the non-negativity constraints on the Lagrange multipliers for equality constraints. See the discussion of the geometrical interpretation of the KTT conditions for some intuition behind this result.

- **_Strong Duality_**

  o **_Theorem (Linear Constraints)_**: Consider the following problem:
  $$(P) : \min_{\boldsymbol{x}} f(\boldsymbol{x}) \text{ s.t. } A\boldsymbol{x} \leq \boldsymbol{b}, \boldsymbol{x} \in \mathbb{R}^n$$
  The Lagrangian is $\mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu}) = f(\boldsymbol{x}) + \boldsymbol{\mu}^\top (A\boldsymbol{x} - \boldsymbol{b})$ and $q(\boldsymbol{\mu}) = \inf_{\boldsymbol{x} \in \mathbb{R}^n} \mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu})$ and the dual is
  $$\max_{\boldsymbol{\mu}} q(\boldsymbol{\mu}) \text{ s.t. } \boldsymbol{\mu} \geq 0, \boldsymbol{\mu} \in \mathbb{R}^m$$
  If $f$ is _convex_ over $\mathbb{R}^n$ and _continuously differentiable_, then if the primal has an optimal solution, then there is no duality gap, and at least one geometric multiplier exists.

  **_Proof_**: Let $\boldsymbol{x}^*$ be an optimal solution for the primal. Then, by the KKT conditions, there exists $\boldsymbol{\mu}^* \in \mathbb{R}^m$ such that
  $$\boldsymbol{\mu}^* \geq \boldsymbol{0} \qquad \boldsymbol{\mu}^{*\top}(A\boldsymbol{x} - \boldsymbol{b}) = 0 \qquad \nabla f(\boldsymbol{x}^*) + A\boldsymbol{\mu}^* = \boldsymbol{0}$$

(We do not need to check for regularity, since the constraints are linear).

Now, since $\mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu})$ is convex and $\nabla f(\boldsymbol{x}^*) + A\boldsymbol{\mu}^* = \boldsymbol{0}$, we have that

$$\boldsymbol{x}^* \in \operatorname{argmin}_{\boldsymbol{x} \in \mathbb{R}^n} \mathcal{L}(\boldsymbol{x}, \boldsymbol{\mu}^*)$$

As such

$$f(\boldsymbol{x}^*) = \min_{\boldsymbol{x} \in \mathbb{R}^n} \left\{ f(\boldsymbol{x}) + \boldsymbol{\mu}^{*\top}(A\boldsymbol{x} - \boldsymbol{b}) \right\} = q(\boldsymbol{\mu}^*)$$

By weak duality, however, we have that

$$q(\boldsymbol{\mu}^*) \leq q^* \leq f^* \leq f(\boldsymbol{x}^*)$$

However, by the previous statement, $f(\boldsymbol{x}^*) = q(\boldsymbol{\mu}^*)$, equality holds throughout, and so $f^* = q^*$.

- o **_Extension to equality constraints_**: The above trivially extends to linear equality constraints. More generally, consider the problem

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) \text{ s.t. } A\boldsymbol{x} = \boldsymbol{b}, g(\boldsymbol{x}) \leq \boldsymbol{0}$$

With Lagrangian

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\boldsymbol{x}) + \boldsymbol{\lambda}^\top(A\boldsymbol{x} - \boldsymbol{b}) + \boldsymbol{\mu}^\top g(\boldsymbol{x})$$

(With $\boldsymbol{\mu} \geq \boldsymbol{0}$ and $\boldsymbol{\lambda}$ unrestricted).

Then, provided that

- ▪ There exists an optimal solution $\boldsymbol{x}^*$
- ▪ $f$ and $\boldsymbol{g}$ are continuously differentiable
- ▪ There exists multipliers $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ satisfying the KKT conditions (ie: some sort of regularity).
- ▪ $f$ and $\boldsymbol{g}$ are convex over $\mathbb{R}^n$

Then there is no duality gap and geometric multipliers exists.

(Note that we are effectively requiring inequality constraints to be convex and equality constraints the be linear. One way to look at this requirement is as a requirement that $\mathcal{L}$ be convex. Indeed, since $\boldsymbol{\mu}$ is positive, $\boldsymbol{\mu}^\top g(\boldsymbol{x})$ is convex for all $\boldsymbol{g}$. Since $\boldsymbol{\lambda}$ can take any value, however, it needs to multiply a linear function to retain convexity).

- o **_Theorem (Slater's Condition)_**: Consider the problem

$$\min_x f(\boldsymbol{x}) \text{ s.t } \boldsymbol{g}(\boldsymbol{x}) \le \boldsymbol{0}, \boldsymbol{x} \in \Omega$$

Then, suppose that

- The problem is bounded, in other words
$$-\infty < f^* = \inf_{\boldsymbol{x} \in \Omega, \boldsymbol{g}(\boldsymbol{x}) \le \boldsymbol{0}} f(\boldsymbol{x})$$

- The set $\Omega$ is convex, and $f$ and $\boldsymbol{g}$ are convex over $\Omega$.

- There exists a vector $\bar{\boldsymbol{x}}$ with $\boldsymbol{g}(\bar{\boldsymbol{x}}) < \boldsymbol{0}$

Then there is no duality gap, and there exists at least one geometric multiplier.

**Proof**: Define

$$\mathcal{A} = \left\{ (\boldsymbol{z}, w) \in \mathbb{R}^{r+1} : \exists \boldsymbol{x} \in \Omega \text{ with } \boldsymbol{g}(\boldsymbol{x}) \le \boldsymbol{z}, f(\boldsymbol{x}) \le w \right\}$$

By the convexity of $\boldsymbol{g}$, $f$ and $\Omega$, this set is convex.

Note also that $\left( 0, f^* \right)$ is not in the interior of $\mathcal{A}$. Otherwise, for some $\varepsilon > 0$, $\left( 0, f^* - \varepsilon \right) \in \mathcal{A}$, which contradicts the definition of $f^*$.

By the supporting hyperplane theorem, there exists a normal vector $(\boldsymbol{\mu}, \beta) \ne (\boldsymbol{0}, 0)$ such that

$$\beta f^* \le \beta w + \boldsymbol{\mu} \cdot \boldsymbol{z} \qquad \forall \ (\boldsymbol{z}, w) \in \mathcal{A}$$

Now, consider; if $(\boldsymbol{z}, w) \in \mathcal{A}$, then $(\boldsymbol{z}, w + \gamma) \in \mathcal{A}$, for all $\gamma \ge 0$. Thus, $\beta \ge 0$. Similarly, $\boldsymbol{\mu} \ge \boldsymbol{0}$.

Consider, however, that if $\beta = 0$, the equation above becomes

$$0 \le \boldsymbol{\mu} \cdot \boldsymbol{z} \qquad \forall \ (\boldsymbol{z}, w) \in \mathcal{A}$$

But since $\left( \boldsymbol{g}(\bar{\boldsymbol{x}}), f(\bar{\boldsymbol{x}}) \right) \in \mathcal{A}$, this would mean that

$$0 \le \boldsymbol{\mu} \cdot \boldsymbol{g}(\bar{\boldsymbol{x}}) \qquad \forall \ (\boldsymbol{z}, w) \in \mathcal{A}$$

Since, by definition, $\boldsymbol{g}(\bar{\boldsymbol{x}}) < \boldsymbol{0}$, this means that $\boldsymbol{\mu} = \boldsymbol{0}$. This is a contradiction, which means that we must have $\beta > 0$.

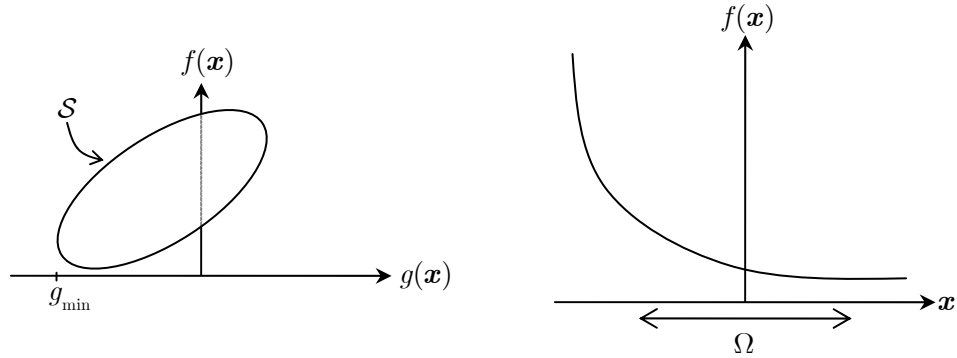We can therefore divide by $\beta$ and normalize, so that $\beta = 1$; we then get

$$f^* \le w + \boldsymbol{\mu} \cdot \boldsymbol{z} \qquad \forall \ (\boldsymbol{z}, w) \in \mathcal{A}$$
$$\Rightarrow f^* \le f(\boldsymbol{x}) + \boldsymbol{\mu} \cdot \boldsymbol{g}(\boldsymbol{x}) \qquad \forall \boldsymbol{x} \in \Omega$$

Minimizing over $\boldsymbol{x} \in \Omega$, we get

$$f^* \leq \inf_{x \in \Omega} \left[ f(x) + \mu \cdot g(x) \right] = q(\mu) \leq q^*$$

Thus, by weak duality, $f^* = q^*$ and $\mu$ is a geometric multiplier.

Diagrammatically, if the set $\mathcal{S} = \left\{ \left( g(x), f(x) \right) : x \in \Omega \right\}$ and the function $f(x)$ look like this:



Then the set $\mathcal{A} = \left\{ (z, w) \in \mathbb{R}^{r+1} : \exists x \in \Omega \text{ with } g(x) \leq z, f(x) \leq w \right\}$ looks like
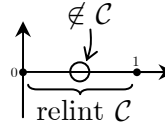


The line must be sloping as above (ie: $\mu$ and $\beta$ must be positive). Now, consider – if there is a point $\bar{x}$ such that $g(\bar{x}) < 0$, it means $g_{\min} < 0$. Therefore, we can't have $\beta = 0$ (ie: the line can't be vertical). Now, the Lagrangian is basically the projection of any point in this set onto the vector $(\mu, 1)$ – which is clearly minimized at $f^*$.

o   Note that the requirement that there be a $g(\bar{x}) < 0$ is crucial. In particular, this condition is *not* satisfied for problems involving equality constraints. These are harder to deal with and require a new definition

o   ***Definition (Relative interior)***: Suppose $\mathcal{C} \subset \mathbb{R}^n$ is a convex set. The *relative interior* of $\mathcal{C}$ is the set relint $\mathcal{C}$ of all $\boldsymbol{x} \in \mathbb{R}^n$ for which there exists an $\varepsilon > 0$ such that if $\boldsymbol{z} \in \text{aff } \mathcal{C}$ with $\left\| \boldsymbol{z} - \boldsymbol{x} \right\| < \varepsilon$, then $\boldsymbol{z} \in \mathcal{C}$.

For example, consider the set $\mathcal{C} = [0,1] \in \mathbb{R}^2$. This set has no interior, because there is no $\mathbb{R}^2$ ball lying totally in the set. However, it has a relative interior:



o   ***Theorem (Slater's Conditions with Mixed Constraints)***: Consider the program ($E$ is a matrix)

$$\min_{\boldsymbol{x} \in \Omega} f(\boldsymbol{x}) \text{ s.t. } \boldsymbol{g}(\boldsymbol{x}) \leq \boldsymbol{0}, A\boldsymbol{x} \leq \boldsymbol{b}, E\boldsymbol{x} = \boldsymbol{d}$$

Suppose that the optimal value $f^*$ is finite and that

- $\Omega$ is the intersection of a convex set $\mathcal{C}$ and a polyhedron

- The functions $f$ and $\boldsymbol{g}$ are convex over $\Omega$

- There is a feasible vector $\bar{\boldsymbol{x}}$ with $\boldsymbol{g}(\bar{\boldsymbol{x}}) < \boldsymbol{0}$

- There is a vector $\boldsymbol{x}$ with $A\boldsymbol{x} \leq \boldsymbol{b}$, $E\boldsymbol{x} = \boldsymbol{d}$, $\boldsymbol{x} \in \text{relint } \mathcal{C}$ and $\boldsymbol{x} \in \Omega$

Then there is no duality gap, and there exists at least one Lagrange multiplier.