# 1. Introduction

## Optimization problem

Let $W \subset \mathbb{R}^n$ be a nonempty set and $f\colon W \to \mathbb{R}$ a function. We consider the problem of finding minima of $f$ in $W$, taking in particular

- $W = \mathbb{R}^n$ (unconstrained optimization),

- $W = \{\, x \in \mathbb{R}^n \colon g_1(x) = 0, \ldots, g_m(x) = 0 \,\}$, where $g_1, \ldots, g_m$ are functions $\mathbb{R}^n \to \mathbb{R}^n$ (equality constraints),

- $W = \{\, x \in \mathbb{R}^n \colon g_1(x) \leqslant 0, \ldots, g_m(x) \leqslant 0 \,\}$, where $g_1, \ldots, g_m$ are functions $\mathbb{R}^n \to \mathbb{R}^n$ (inequality constraints).

The set $W$ is called a <u>feasible set/region</u>.

<u>Definition 1</u> *A point $x_0 \in W$ is called a <u>global minimum</u> of $f$ in $W$ if*

$$f(x) \geqslant f(x_0) \quad \text{for all } x \in W.$$

<u>Definition 2</u> *A point $x_0 \in W$ is called a <u>local minimum</u> of $f$ in $W$ if there exists $\varepsilon > 0$ such that*

$$f(x) \geqslant f(x_0) \quad \text{for all } x \in W \cap B(x_0, \varepsilon),$$

*where $B(x_0, \varepsilon)$ is the ball whose centre is $x_0$ and the radius is $\varepsilon$.*

Any global minimum is a local minimum. A minimum is called <u>strict</u> if in the definitions above there is $f(x) > f(x_0)$ for $x \neq x_0$. In a similar way we define global and local maxima. A point $x_0$ is a (global or local) <u>extremum</u> if it is a minimum or a maximum.

Minima need not exist, if no point $x_0$ fulfills the definitions. A global minimum does not exist if $\inf_{x \in W} f(x) = -\infty$ or $\inf_{x \in W} f(x) = c$ and $f(x) > c$ for all $x \in W$.

<u>Example.</u> Let $f(x) = x \cos x$. If $W = \mathbb{R}$ then $\inf_{x \in W} f(x) = -\infty$, and there is no global minimum and an infinite set of local minima. If $W = [a, b]$, where $a, b \in \mathbb{R}$,

then a global minimum exists. If $W = (a, b)$ then minima either exist or not, depending on the choice of $a, b$. In general, a continuous function is not guaranteed to have extrema if the feasible set is not compact, e.g. if it is open.

## Existence of minima of a continuous function

<u>Theorem 1</u> *If the set $W \in \mathbb{R}^n$ is compact and $f\colon W \to \mathbb{R}$ is a continuous function, then $f$ reaches its infimum and supremum in $W$, i.e., there exist $x_0, y_0 \in W$ such that*

$$f(x_0) \leqslant f(x) \leqslant f(y_0) \quad \text{for all } x \in W.$$

<u>Definition 3</u> *A function $f\colon W \to \mathbb{R}$ is called <u>coercive</u> if $f(x) \to \infty$ for $\|x\| \to \infty$. Equivalently,*

$$\forall_{r>0} \exists_{s>0} \forall_{x \in W} \; \|x\| > s \;\Rightarrow\; f(x) > r.$$

If $W$ is a bounded set, then any function $f\colon W \to \mathbb{R}$ is coercive.

<u>Theorem 2</u> *If $W \subset \mathbb{R}^n$ is a closed set and $f\colon W \to \mathbb{R}$ is continuous and coercive, then there exists a minimum $x_0$ of $f$ in $W$.*

<u>Proof.</u> For a point $y \in W$ we define the set $U_y = \{\, x \in W \colon f(x) \leqslant f(y) \,\}$. The set $U_y$ is nonempty and closed, as the function $f$ is continuous and the inequality in the definition of $U_y$ is nonsharp and $W$ is closed. This set is also bounded: for $r = f(y)$, from the coercivity of $f$ there exists $s > 0$ such that if $\|x\| > s$, then $f(x) > r = f(y)$; hence, $x \notin U_y$ and $U_y \subset B(0, s)$. It follows that $U_y$ is a closed and bounded set, i.e., it is compact. Therefore there exists a global minimum $x_0$ of $f$ in $U_y$. Due to $f(x) > f(y) \geqslant f(x_0)$ for $x \notin U_y$, $x_0$ is also a global minimum of $f$ in $W$. $\square$

<u>Theorem 3</u> *Let $W \subset \mathbb{R}^n$ be nonempty and let $f\colon W \to \mathbb{R}$ be a continuous function. If there exists $y \in W$ such that for any sequence $(x_n)_n \subset W$ such that*

$$x_n \to \operatorname{cl} W \setminus W \quad \text{or} \quad \|x_n\| \to \infty$$

*there is $\liminf_{n \to \infty} f(x_n) > f(y)$, then there exists a minimum $x_0$ of the function $f$.*

Proof. The set $U_y$ is defined as before. To show that it is closed, we take any sequence $(x_n)_n \subset U_y$ which converges to $\overline{x}$. It suffices to show that $\overline{x} \in U_y$. From $x_n \in U_y$ we have $f(x_n) \leqslant f(y)$ and if $\overline{x} \notin W$, then we have an inconsistency with the assumption. Hence, $\overline{x} \in W$. As the function $f$ is continuous in $W$, there is $f(\overline{x}) \leqslant f(y)$, hence $\overline{x} \in U_y$. The set $U_y$ is also bounded, which follows from the assumed implication $\|x_n\| \to \infty \Rightarrow \liminf_{n \to \infty} f(x_n) > f(y)$. The proof is completed just like the proof of the previous theorem. $\square$

## Local minima of functions of one variable

Let $W \subset \mathbb{R}$ be an open set.

Theorem 4 (*necessary condition of the 1st order*) *If* $x_0 \in W$ *is a local minimum or maximum of* $f$ *and* $f'(x_0)$ *exists, then* $f'(x_0) = 0$.

Proof. Let $x_0$ be a local minimum. For sufficiently small $h > 0$ there is $f(x_0 - h) \geqslant f(x_0) \leqslant f(x_0 + h)$ and then

$$\frac{f(x_0 - h) - f(x_0)}{-h} \leqslant 0 \quad \Rightarrow \quad \lim_{h \to 0} \frac{f(x_0 - h) - f(x_0)}{-h} \leqslant 0 \quad \Rightarrow \quad f'(x_0) \leqslant 0,$$

$$\frac{f(x_0 + h) - f(x_0)}{+h} \geqslant 0 \quad \Rightarrow \quad \lim_{h \to 0} \frac{f(x_0 + h) - f(x_0)}{+h} \geqslant 0 \quad \Rightarrow \quad f'(x_0) \geqslant 0,$$

hence, $f'(x_0) = 0$. $\square$

Theorem 5 (*necessary condition of the 2nd order*) *If* $f : W \to \mathbb{R}$ *is of class* $C^2(W)$ *and* $x_0$ *is a local minimum, then* $f''(x_0) \geqslant 0$.

If the set $W$ is not open, then we cannot use the above theorems for $x_0 \in \partial W$. But the theorem below applies also in this case.

Theorem 6 (*sufficient condition of the 2nd order*) *If* $f : W \to \mathbb{R}$ *is of class* $C^2(W)$ *and* $f'(x_0) = 0$, $f''(x_0) > 0$ *at a point* $x_0 \in W$, *then* $f$ *has a strict local minimum at* $x_0$.

Theorem 7 *If* $W \subset \mathbb{R}$ *is open,* $f \in C^k(W)$ *and* $f'(x_0) = f''(x_0) = \cdots = f^{(k-1)}(x_0) = 0$, $f^{(k)}(x_0) \neq 0$ *for* $x_0 \in W$, *then if* $k$ *is odd, there is no extremum of* $f$ *at* $x_0$, *and if* $k$ *is even, then there is a local minimum if* $f^{(k)}(x_0) > 0$ *and a local maximum if* $f^{(k)}(x_0) < 0$.

## Taylor's formulae

Theorem 8 (*Rolle's theorem*) *If a function* $f : [a, b] \to \mathbb{R}$ *is continuous in* $[a, b]$, *differentiable in* $(a, b)$ *and* $f(a) = f(b)$, *then there exists a point* $x_0 \in (a, b)$ *such that* $f'(x_0) = 0$.

Proof. If $f$ is constant, then the claim is obvious. Otherwise there exists an extremum $x_0$ of $f$ in $[a, b]$ other than $a$ and $b$: there is $f(x_0) = \sup_{x \in [a,b]} f(x) > f(a)$ or $f(x_0) = \inf_{x \in [a,b]} f(x) < f(a)$. Let $x_0$ be a maximum. Then $f(x) \leqslant f(x_0)$ for all $x \in [a, b]$ and

$$\frac{f(x) - f(x_0)}{x - x_0} \geqslant 0 \text{ if } x < x_0, \quad \frac{f(x) - f(x_0)}{x - x_0} \leqslant 0 \text{ if } x > x_0.$$

Hence,

$$f'(x_0) = \lim_{x \nearrow x_0} \underbrace{\frac{f(x) - f(x_0)}{x - x_0}}_{\geqslant 0} = \lim_{x \searrow x_0} \underbrace{\frac{f(x) - f(x_0)}{x - x_0}}_{\leqslant 0},$$

therefore, $f'(x_0) = 0$. If $x_0$ is a minimum, the proof is similar. $\square$

Theorem 9 (*mean value theorem*) *If a function* $f : [a, b] \to \mathbb{R}$ *is continuous in* $[a, b]$ *and differentiable in* $(a, b)$, *then there exists a point* $x_0 \in (a, b)$ *such that*

$$f(b) - f(a) = f'(x_0)(b - a).$$

Proof. Let $g(x) \stackrel{\text{def}}{=} \big(f(b) - f(a)\big)x - (b - a)f(x)$. The function $g$ is continuous in $[a, b]$ and differentiable in $(a, b)$, moreover,

$$g(a) = f(b)a - f(a)b = g(b).$$

By Rolle's theorem, there exists $x_0 \in (a, b)$ such that $g'(x_0) = 0$. Hence,

$$0 = g'(x_0) = f(b) - f(a) - (b - a)f'(x_0).$$

The proof is completed by rearranging this formula. $\square$

<u>Theorem 10</u> *(Taylor's formula with the remainder in Peano form) Let* f: $[a, b] \to \mathbb{R}$ *be a function differentiable in* $[a, b]$ *and twice differentiable at some point* $x_0 \in (a, b)$. *Then for all* $x \in [a, b]$ *there is*

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2}(x - x_0)^2 + o\big((x - x_0)^2\big).$$

<u>Proof.</u> Without loss of generality we assume $x_0 = 0$. Let

$$R(x) \stackrel{\text{def}}{=} f(x) - f(0) - f'(0)x - \frac{f''(0)}{2}x^2.$$

We need to show that $R(x) = o(x^2)$. From the continuity of $f'$ we obtain

$$f(x) - f(0) = \int_0^x f'(y)\, dy.$$

The function $f'$ is differentiable at $0$. Hence, $f'(y) = f'(0) + f''(0)y + r(y)$, where $r(y) = o(y)$. This means that

$$\lim_{y=0} \frac{r(y)}{y} = 0,$$

i.e., for any $\varepsilon > 0$ there exists $\delta > 0$ such that $|y| < \delta \Rightarrow |r(y)| < \varepsilon |y|$.

Now we fix an $\varepsilon > 0$ and the related $\delta > 0$. For $|x| < \delta$ we integrate $f'(y)$:

$$f(x) - f(0) = \int_0^x \big(f'(0) + f''(0)y + r(y)\big)\, dy = f'(0)x + \frac{f''(0)}{2}x^2 + \int_0^x r(y)\, dy.$$

Hence, $R(x) = \int_0^x r(y)\, dy$. Using the estimate $|r(y)| < \varepsilon |y|$ for $|y| < \delta$, we obtain

$$|R(x)| \leqslant \int_0^x |r(y)|\, dy < \int_0^x \varepsilon |y|\, dy = \frac{\varepsilon x^2}{2}.$$

Hence,

$$\left| \frac{R(x)}{x^2} \right| < \frac{\varepsilon}{2}.$$

As $\varepsilon > 0$ may be arbitrary, $\lim_{x \to 0} \left| \frac{R(x)}{x^2} \right| = 0$, i.e., $R(x) = o(x^2)$. □

Just a little more effort is needed to prove the formula with more terms, applicable for functions having derivatives up to the order $k - 1$ in $(a, b)$ and the k-th order derivative at $x_0$:

$$f(x) = f(x_0) + \sum_{i=1}^k \frac{f^{(i)}(x_0)}{i!}(x - x_0)^i + o\big((x - x_0)^k\big).$$

<u>Theorem 11</u> *(Taylor's formula with the remainder in Lagrange form) Let* f: $[a, b] \to \mathbb{R}$ *be a function of class* $C^{k-1}[a, b]$ *and* k *times differentiable in* $(a, b)$. *For* $x_0 \in (a, b)$ *and* $x \in [a, b]$ *there is*

$$f(x) = f(x_0) + \sum_{i=1}^{k-1} \frac{f^{(i)}(x_0)}{i!}(x - x_0)^i + \frac{f^{(k)}(\overline{x})}{k!}(x - x_0)^k,$$

*where* $\overline{x}$ *is a point between* $x_0$ *and* $x$.

<u>Proof.</u> The function $h(x) \stackrel{\text{def}}{=} f(x_0) + \sum_{i=1}^{k-1} \frac{f^{(i)}(x_0)}{i!}(x - x_0)^i$ is a polynomial of degree less than k. For $x \neq x_0$ let $g_x(y) \stackrel{\text{def}}{=} f(y) - h(y) - z_x(y - x_0)^k$, where $z_x = \frac{f(x) - h(x)}{(x - x_0)^k}$. It is easy to verify that $g_x(x_0) = g_x'(x_0) = \cdots = g_x^{(k-1)}(x_0) = g_x(x) = 0$. By Rolle's theorem, the derivative $g_x'$ is equal to $0$ at some point $x_1$ between $x_0$ and $x$; note that the point $x_0$ is a zero of multiplicity $k - 1$ of $g_x'$. Using the induction and Rolle's theorem in the similar way, we show the existence of the sequence of points, $x_2, \ldots, x_k$ such that $g_x^{(i)}(x_i) = 0$ and each point $x_i$ is between $x_0$ and $x_{i-1}$.

The point $\overline{x} = x_k$ is a zero of $g_x^{(k)}$ located between $x_0$ and $x$, i.e.,

$$0 = g_x^{(k)}(x_k) = f^{(k)}(x_k) - z_x k!.$$

Hence, $z_x = \frac{f^{(k)}(x_k)}{k!}$. By substititing this expression and $y = x$ to the definition of $g_x$, due to $g_x(x) = 0$, we obtain the needed formula. □

## Global extrema

<u>Theorem 12</u> *Let* $I \subset \mathbb{R}$ *be an interval, open or closed at one or both ends, or even unbounded. Let* f: $I \to \mathbb{R}$ *be of class* $C^1(I)$ *and* $C^2(\text{int } I)$. *Let* $x_0 \in I$ *and* $f'(x_0) = 0$. *If* $f''(x) \geqslant 0$ *for all* $x \in I$, *then* $x_0$ *is a global minimum of* f. *If* $f''(x) \leqslant 0$ *for all* $x \in I$, *then* $x_0$ *is a global maximum of* f. *If in addition* $f''(x_0) > 0$ *or respectively* $f''(x_0) < 0$, *then* $x_0$ *is a unique (strict) global minimum or maximum.*

<u>Proof.</u> By the Taylor's formula we have

$$f(x) = f(x_0) + \frac{1}{2}f''(\overline{x})(x - x_0)^2,$$

where $\overline{x}$ is a point between $x_0$ and $x$. Hence, the last term of the formula above determines the inequality between $f(x)$ and $f(x_0)$.

Assume that $f''(x) \geqslant 0$ for all $x \in I$ and $f''(x_0) > 0$. By $f'(x_0) = 0$ we obtain

$$f'(x) = f'(x) - f'(x_0) = \int_{x_0}^{x} f''(y) \, dy \geqslant 0$$

for $x > x_0$. Similarly we show that $f'(x) \leqslant 0$ for $x < x_0$. As $f''(x_0) > 0$ and $f''$ is continuous, it follows that $f''$ is positive in a neighbourhood of $x_0$. Hence, the integrals are positive, which implies $f'(x) > 0$ for $x > x_0$ and $f'(x) < 0$ for $x < x_0$. Thus $f$ is decreasing for $x < x_0$ and increasing for $x > x_0$ and, therefore, $x_0$ is a strict minimum. The proof for the case of maximum is similar. $\square$

# 2. Extrema of functions of two or more variables

Let $f: W \to \mathbb{R}$, where $W \subset \mathbb{R}^n$ is an open set. Points of $\mathbb{R}^n$, $\mathbf{x} = (x_1, \dots, x_n)$, are identified with column matrices, $[x_1, \dots, x_n]^\mathsf{T}$, but it is convenient to write $f(\mathbf{x})$ and $f(x_1, \dots, x_n)$, which denotes the same thing. We use the Euclidean norm, $\|\mathbf{x}\| = \sqrt{\mathbf{x}^\mathsf{T}\mathbf{x}} = \sqrt{x_1^2 + \cdots + x_n^2}$.

The <u>gradient</u> of $f$ is a row matrix, $Df(\mathbf{x}) = \left[ \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right]$.

The <u>Hessian</u> of $f$ at $\mathbf{x} \in W$ is the $n \times n$ matrix,

$$D^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

<u>Definition 4</u> *The function $f$ is <u>differentiable</u> at $\mathbf{x}_0 \in W$ if there exists a vector $\boldsymbol{\alpha} \in \mathbb{R}^n$ such that*

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \boldsymbol{\alpha}^\mathsf{T}(\mathbf{x} - \mathbf{x}_0) + o(\|\mathbf{x} - \mathbf{x}_0\|), \quad \mathbf{x} \in W.$$

*The function $f$ is <u>twice differentiable</u> at $\mathbf{x}_0 \in W$ if in addition there exists a matrix $H \in \mathbb{R}^{n \times n}$ such that*

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \boldsymbol{\alpha}^\mathsf{T}(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\mathsf{T} H (\mathbf{x} - \mathbf{x}_0) + o(\|\mathbf{x} - \mathbf{x}_0\|^2), \quad \mathbf{x} \in W.$$

If a function is twice differentiable, then there exists a symmetric matrix $H$ mentioned in the definition above; if a nonsymmetric matrix $H$ satisfies the formula in this definition, so does the symmetric matrix $\frac{1}{2}(H + H^\mathsf{T})$.

<u>Theorem 13</u> *I) If a function $f$ is differentiable at $\mathbf{x}_0$, then the gradient $Df(\mathbf{x}_0)$ exists and is equal to $\boldsymbol{\alpha}^\mathsf{T}$. Conversely, if $Df(\mathbf{x})$ exists in a neighbourhood of $\mathbf{x}_0$ and is continuous at $\mathbf{x}_0$, then $f$ is differentiable at $\mathbf{x}_0$.*

*II) if the Hessian $D^2 f(\mathbf{x})$ exists in a neighbourhood of $\mathbf{x}_0$ and is continuous at $\mathbf{x}_0$, then $f$ is twice differentiable at $\mathbf{x}_0$; the Hessian is then a symmetric matrix, $H = D^2 f(\mathbf{x}_0)$.*

<u>Remark.</u> If the function $f$ is differentiable at a point $\mathbf{x}$, the (real) value of the product of matrices $Df(\mathbf{x})\mathbf{v}$, where $\mathbf{v} \in \mathbb{R}^n$, is the directional derivative of the function $f$ in the direction of the vector $\mathbf{v}$ at $\mathbf{x}$. If the function $f$ is twice differentiable, then $\mathbf{v}^\mathsf{T} D^2 f(\mathbf{x}) \mathbf{v}$ is equal to the second order directional derivative of $f$ in the direction of $\mathbf{v}$.

<u>Remark.</u> To use second order derivatives *in practice* we need to assume the continuity of the Hessian.

<u>Remark.</u> A function $f$ whose domain is an *open* set $W \subset \mathbb{R}^n$ is said to be of class $C^1$ ($C^2$) in $W$ if it is continuous in $W$ together with its first (and second) order derivatives. If the set $W$ is not open, the function is said to be of class $C^1$ ($C^2$) if there exists an extension $\tilde{f}$ of class $C^1$ ($C^2$) of the function $f$ to an open set $\tilde{W}$ such that $W \subset \tilde{W}$. Then we can consider the derivatives of $f$ at the boundary points of $W$; if $W \subset \mathrm{cl}(\mathrm{int}\, W)$, then (due to their continuity) the derivatives are uniquely determined by the values of $f$ in $W$.

<u>Lemma 1</u> *Let $W \subset \mathbb{R}^n$ be an open set. If a function $f: W \to \mathbb{R}$ is of class $C^2$ and a line segment $\overline{\mathbf{x}_0 \mathbf{x}}$ is contained in $W$, then*

$$f(\mathbf{x}) = f(\mathbf{x}_0) + Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^\mathsf{T} D^2 f(\overline{\mathbf{x}})(\mathbf{x} - \mathbf{x}_0),$$

*where $\overline{\mathbf{x}}$ is an interior point the line segment $\overline{\mathbf{x}_0 \mathbf{x}}$.*

<u>Proof.</u> Apply the Taylor's formula to the function $g(t) = f\big(\mathbf{x}_0 + t(\mathbf{x} - \mathbf{x}_0)\big)$, $t \in [0, 1]$. $\square$
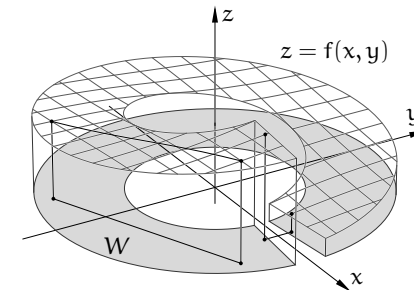


Figure 1: Applicability of the Taylor's formula for a multivariate function

Definition 5 *The set $W \subset \mathbb{R}^n$ is* <u>*convex*</u> *if*

$$\lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in W$$

*for all $\mathbf{x}, \mathbf{y} \in W$, $\lambda \in [0, 1]$.*

Corollary 1 *If $W \subset \mathbb{R}^n$ is open and convex and $f\colon W \to \mathbb{R}$ is of class $C^2(W)$ then for all $\mathbf{x}_0, \mathbf{x} \in W$ there is*

$$f(\mathbf{x}) = f(\mathbf{x}_0) + Df(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0)^T D^2 f(\overline{\mathbf{x}})(\mathbf{x} - \mathbf{x}_0),$$

*where $\overline{\mathbf{x}} = (1 - \lambda)\mathbf{x}_0 + \lambda \mathbf{x}$ for some $\lambda \in (0, 1)$.*

<u>Proof.</u> As the set $W$ is convex, if $\mathbf{x}_0, \mathbf{x} \in W$, then $\overline{\mathbf{x}_0 \mathbf{x}} \subset W$; the claim follows from the lemma. $\square$

## The necessary first order condition

We consider a set $W \subset \mathbb{R}^n$ with a nonempty interior.

Theorem 14 *(<u>necessary 1st order condition</u>) If a function $f\colon W \to \mathbb{R}$ is differentiable at a point $\mathbf{x}_0 \in \operatorname{int} W$ and $\mathbf{x}_0$ is a local extremum of $f$, then $Df(\mathbf{x}_0) = \mathbf{0}^T$.*

<u>Proof.</u> From $\mathbf{x}_0 \in \operatorname{int} W$ it follows that the function $g_i(t) = f(\mathbf{x}_0 + t\mathbf{e}_i)$ (where $\mathbf{e}_i = [0, \dots, 0, \underset{\underset{i}{\uparrow}}{1}, 0, \dots, 0]^T$) is well defined. It has the local extremum at $0$. By the necessary first order condition for functions of one variable there must be $g_i'(0) = 0$, which implies $\frac{\partial f}{\partial x_i} = 0$. As this holds for all $i = 1, \dots, n$, the gradient of $f$ is the zero matrix $1 \times n$. $\square$

Definition 6 *A point $\mathbf{x}_0 \in \operatorname{int} W$ is called a* <u>*critical point*</u> *of the function $f\colon W \to \mathbb{R}$ if $f$ is differentiable at $\mathbf{x}_0$ and $Df(\mathbf{x}_0) = \mathbf{0}^T$.*

## Positive- and negative-definite matrices

Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix, $A = [a_{ij}]$, $a_{ij} = a_{ji}$. It defines a quadratic form

$$F(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} = \sum_{i=1}^{n} \sum_{j=1}^{n} a_{ij} x_i x_j.$$

Definition 7 *The matrix $A$ or the quadratic form $F$ is*

- <u>*positive definite*</u> *if $F(\mathbf{x}) > 0$ for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ (we write $A > 0$),*

- <u>*nonnegative definite*</u> *if $F(\mathbf{x}) \geqslant 0$ for all $\mathbf{x} \in \mathbb{R}^n$ (we write $A \geqslant 0$),*

- <u>*negative definite*</u> *if $F(\mathbf{x}) < 0$ for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ (we write $A < 0$),*

- <u>*nonpositive definite*</u> *if $F(\mathbf{x}) \leqslant 0$ for all $\mathbf{x} \in \mathbb{R}^n$ (we write $A \leqslant 0$),*

- <u>*indefinite*</u> *if there exist vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ such that $F(\mathbf{x}) > 0$, $F(\mathbf{y}) < 0$.*

At the first glance we can notice that if *not all* diagonal elements $a_{ii}$ are positive (nonnegative) then the matrix $A$ is not positive-definite (nonnegative-definite) and if *not all* diagonal coefficients are negative (nonpositive) then the matrix is not negative-definite (nonpositive-definite). The basic characteristic of positive-definite matrices is given by

Theorem 15 *(<u>Sylvester's criterion</u>) Let $A_i$ be a matrix obtained from $A$ by rejecting its last $n - i$ rows and columns (in particular $A_1 = [a_{11}]$, $A_n = A$).*

*I) The matrix $A$ is positive-definite if and only if $\det A_i > 0$ for $i = 1, \dots, n$,*

*II) The matrix $A$ is nonnegative-definite if and only if $\det A_i \geqslant 0$ for $i = 1, \dots, n$.*

A matrix $A$ is negative-definite (nonpositive-definite) if $-A$ is positive-definite (nonnegative-definite). Another characteristic is related to the algebraic eigenproblem. From the linear algebra we know that all eigenvalues of a real symmetric matrix are real numbers; for any such a matrix there exists an orthogonal basis of $\mathbb{R}^n$ whose elements are eigenvectors of this matrix. A positive-,

nonnegative-, negative- or nonpositive-definite symmetric matrix has respectively all eigenvalues positive, nonnegative, negative or nonpositive.

## Conditions of the second order

<u>Theorem 16</u> *(necessary 2nd order condition)* *If* f *is a function of class* $C^2$ *in an open set* $W \subset \mathbb{R}^n$ *and* $x_0 \in W$ *is a local minimum, then the matrix* $D^2 f(x_0)$ *is nonnegative-definite. If* $x_0$ *is a local maximum, then* $D^2 f(x_0)$ *is nonpositive-definite.*

<u>Proof.</u> Let $x_0$ be a local minimum. Let $h \in \mathbb{R}^n \setminus \{0\}$ and $g(t) = f(x_0 + th)$, where $t \in \mathbb{R}$ is chosen so as to obtain $x_0 + th \in W$. The function $g$ has a local minimum at $0$. As $f$ is of class $C^2$, so is $g$. By the second order necessary condition for the univariate case, $g''(0) \geqslant 0$. The second order derivative of the composite function $g$ is

$$g''(0) = h^\mathsf{T} D^2 f(x_0) h.$$

As the vector $h$ may be arbitrary, the matrix $D^2 f(x_0)$ is nonnegative-definite. $\square$

<u>Theorem 17</u> *(sufficient 2nd order condition)* *If* f *is a function of class* $C^2$ *in an open set* $W \subset \mathbb{R}^n$, $Df(x_0) = 0^\mathsf{T}$ *and the matrix* $D^2 f(x_0)$ *is positive-definite (negative-definite), then* $x_0$ *is a local minimum (maximum) of* f.

<u>Proof.</u> Assume that $D^2 f(x_0) > 0$. Let $\alpha \colon W \to \mathbb{R}$ be the function defined by

$$\alpha(x) = \inf_{\|h\|=1} h^\mathsf{T} D^2 f(x) h.$$

The function value $\alpha(x)$ is the minimal eigenvalue of the matrix $D^2 f(x)$; the infimum is the minimum taken at the vector $h$ which is a unit eigenvector corresponding to the minimal eigenvalue of the Hessian. Due to the continuity of the Hessian of $f$, the function $\alpha$ is continuous. Hence, there exists a ball $B(x_0, \varepsilon)$, $\varepsilon > 0$, such that $\alpha(x) > 0$ for all $x \in B(x_0, \varepsilon)$.

For a fixed $x \in B(x_0, \varepsilon)$, due to the Taylor's formula we have

$$f(x) = f(x_0) + Df(x_0)(x - x_0) + \frac{1}{2}(x - x_0)^\mathsf{T} D^2 f(\overline{x})(x - x_0),$$

where $\overline{x}$ is a point of the line segment $\overline{x_0 x} \subset B(x_0, \varepsilon)$. The gradient of $f$ vanishes at $x_0$ and

$$(x - x_0)^\mathsf{T} D^2 f(\overline{x})(x - x_0) = \|x - x_0\|^2 \frac{(x - x_0)^\mathsf{T}}{\|x - x_0\|} D^2 f(\overline{x}) \frac{(x - x_0)}{\|x - x_0\|} \geqslant \|x - x_0\|^2 \alpha(\overline{x}).$$

Hence,

$$f(x) - f(x_0) \geqslant \frac{1}{2}\|x - x_0\|^2 \alpha(\overline{x}) > 0.$$

It follows that $x_0$ is a strict local minimum. The proof for a maximum is similar. $\square$

## Global extrema

Let $W$ be a convex set and $f \colon W \to \mathbb{R}$ a function of class $C^1(W)$ and $C^2(\operatorname{int} W)$.

<u>Theorem 18</u> *If* $x_0 \in \operatorname{int} W$ *is a critical point of* f, *then*

*I) If* $D^2 f(x) \geqslant 0$ *for all* $x \in \operatorname{int} W$, *then* $x_0$ *is a global minimum,*

*II) If* $D^2 f(x) \leqslant 0$ *for all* $x \in \operatorname{int} W$, *then* $x_0$ *is a global maximum.*

*If in addition* $D^2 f(x_0) > 0$ *or* $D^2 f(x_0) < 0$ *respectively, then* $x_0$ *is a strict minimum or maximum.*

<u>Proof.</u> If $x \in W$, then by convexity of $W$ the entire line segment $\overline{x_0 x}$ is contained in $W$. By the Taylor's formula,

$$f(x) = f(x_0) + \frac{1}{2}(x - x_0)^\mathsf{T} D^2 f(\overline{x})(x - x_0),$$

for a point $\overline{x} \in \overline{x_0 x}$. From the inequality $D^2 f(\overline{x}) \geqslant 0$ (or $D^2 f(\overline{x}) \leqslant 0$) it follows that the last term above is nonnegative (or nonpositive), which proves that $x_0$ is a minimum (or a maximum).

If in addition to (I) we have $D^2 f(x_0) > 0$, then we can consider the function $g(t) = f(x_0 + t(x - x_0))$, $t \in [0, 1]$. Due to the convexity of $W$, $x_0 + t(x - x_0) \in W$, so the function $g$ is well defined. From the assumptions it follows that $g'(0) = 0$, $g''(0) > 0$ and $g''(t) \geqslant 0$. Therefore $g$ has a strict global minimum at $0$, i.e. $f(x) > f(x_0)$. As the choice of $x \in W$ is arbitrary, $x_0$ is a strict global minimum of $f$.

The proof for the case of $D^2 f(x_0) < 0$ holding in addition to (II) is similar. $\square$

# 3. Convex sets and functions

<u>Lemma 2</u> *The set $W \subset \mathbb{R}^n$ is convex if and only if for all $m \geqslant 2$ and for all points $x_1, \ldots, x_m \in W$ and numbers $a_1, \ldots, a_m \geqslant 0$, $a_1 + \cdots + a_m = 1$, there is*

$$a_1 x_1 + \cdots + a_m x_m \in W.$$

<u>Lemma 3</u> *Let $W \subset \mathbb{R}^n$ be a convex set with a nonempty interior. Then*

*I) For any $x \in W$ and $x_0 \in \operatorname{int} W$ the line segment $\overline{x_0 x} \setminus \{x\}$ is contained in the interior of $W$:*

$$\lambda x_0 + (1 - \lambda) x \in W \quad \text{for all } \lambda \in (0, 1].$$

*II) $W \subset \operatorname{cl}(\operatorname{int} W)$*

<u>Proof.</u> Let the points $x_0$ and $x$ satisfy the assumptions. As $\operatorname{int} W$ is open, there exists a ball $B(x_0, \varepsilon) \subset \operatorname{int} W$. The union of all line segments, whose one end point is $x$ and the other end point is in this ball, is a "cone" with the vertex $x$ and the base $B(x_0, \varepsilon)$. This cone is a subset of $W$ and its interior contains the line segment $\overline{x_0 x} \setminus \{x\}$. This completes the proof of (I). (II) follows immediately. $\square$
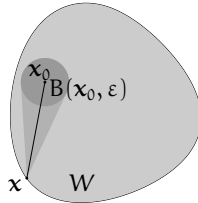


Figure 2: Illustration of Lemma 2

<u>Theorem 19</u> *(weak separation theorem) Let $U, V \subset \mathbb{R}^n$ be nonempty convex sets such that $U \cap V = \emptyset$. There exists a hyperplane separating the sets $U$ and $V$, i.e., there exists a nonzero vector $a \in \mathbb{R}^n$ such that*

$$a^\mathsf{T} x \leqslant a^\mathsf{T} y \quad \text{for all } x \in U, \, y \in V.$$

As the mapping $x \to a^\mathsf{T} x$ is continuous in $\mathbb{R}^n$, from the above we obtain

<u>Corollary 2</u> *Let $U, V \subset \mathbb{R}^n$ be nonempty convex sets such that $\operatorname{int} U \neq \emptyset$ and $(\operatorname{int} U) \cap V = \emptyset$. Then there exists a hyperplane separating the sets $U$ and $V$.*

<u>Theorem 20</u> *(strong separation theorem) Let $U, V \subset \mathbb{R}^n$ be nonempty closed convex sets, let $U$ be compact and let $U \cap V = \emptyset$. Then there exists a hyperplane strictly separating the sets $U$ and $V$, i.e., there exists a nonzero vector $a \in \mathbb{R}^n$ such that*

$$\sup_{x \in U} a^\mathsf{T} x < \inf_{y \in V} a^\mathsf{T} y.$$

There may be more than one hyperplane described by these theorems; one of them is the set defined as follows:

$$\{x \in \mathbb{R}^n : a^\mathsf{T} x = \alpha\}, \quad \alpha = \sup_{x \in U} a^\mathsf{T} x.$$

<u>Proof of the strong separation theorem.</u> Let $d \colon U \times V \to \mathbb{R}$ be a function given by the formula $d(x, y) = \|x - y\|$. As the set $U$ is bounded, the function $d$ is coercive; it may tend to infinity only by taking an appropriate sequence of points $y \in V$. As the function $d$ is continuous and coercive and its domain $U \times V$ is closed, it takes a minimum at a point $(x_0, y_0) \in U \times V$. As $U \cap V = \emptyset$, there is $a = y_0 - x_0 \neq 0$. Below we demonstrate that it is a vector satisfying the claim.

First we show that $a^\mathsf{T} y \geqslant a^\mathsf{T} y_0$ for all $y \in V$. Let

$$g(t) \stackrel{\text{def}}{=} \left( d\big(x_0, y_0 + t(y - y_0)\big) \right)^2, \quad t \in \mathbb{R}.$$

There is

$$g(t) = \|y_0 - x_0\|^2 + 2t(y_0 - x_0)^\mathsf{T}(y - y_0) + t^2 (y - y_0)^\mathsf{T}(y - y_0).$$

This function is differentiable for all $t \in \mathbb{R}$ and, as the set $V$ is convex, $g(0) \leqslant g(t)$ for $t \in [0, 1]$. Hence, $g'(0) \geqslant 0$, i.e.,

$$(y_0 - x_0)^\mathsf{T}(y - y_0) = a^\mathsf{T}(y - y_0) \geqslant 0.$$

In a similar way we can show that $a^\mathsf{T} x \leqslant a^\mathsf{T} x_0$ for all $x \in U$. $\square$

<u>Proof of the weak separation theorem.</u> Consider the set $C = V - U = \{y - x \colon x \in U, \, y \in V\}$. This set is convex and $0 \notin C$. It suffices to find a nonzero vector $a \in \mathbb{R}^n$ such that $a^\mathsf{T} x \geqslant 0$ for all $x \in C$.

Let $A_{\mathbf{x}} \overset{\text{def}}{=} \{\, \mathbf{a} \in \mathbb{R}^n \colon \|\mathbf{a}\| = 1,\ \mathbf{a}^\mathsf{T}\mathbf{x} \geqslant 0 \,\}$. We are going to show that $\bigcap_{\mathbf{x} \in C} A_{\mathbf{x}} \neq \emptyset$. Suppose that $\bigcap_{\mathbf{x} \in C} A_{\mathbf{x}} = \emptyset$. Let $B_{\mathbf{x}} = S \setminus A_{\mathbf{x}}$, where $S$ is the unit sphere in $\mathbb{R}^n$. The sets $B_{\mathbf{x}}$ are open subsets of $S$. If the intersection of all sets $A_{\mathbf{x}}$, where $\mathbf{x} \in C$, is empty, then the family $\{\, B_{\mathbf{x}} \colon \mathbf{x} \in C \,\}$ is an open coverage of $S$, which is a compact set. Hence, there exists a finite coverage $\{\, B_{\mathbf{x}_1}, \ldots, B_{\mathbf{x}_k} \colon \mathbf{x}_1, \ldots, \mathbf{x}_k \in C \,\}$ of $S$. Let

$$\hat{C} \overset{\text{def}}{=} \operatorname{conv}\{\mathbf{x}_1, \ldots, \mathbf{x}_k\} = \Big\{ \sum_{i=1}^{k} \lambda_i \mathbf{x}_i \colon \lambda_1, \ldots, \lambda_k \geqslant 0,\ \sum_{i=1}^{k} \lambda_i = 1 \Big\}.$$

The set $\hat{C}$ is convex and closed and it is a subset of $C$. Hence, $\mathbf{0} \notin \hat{C}$. By the strong separation theorem used to the sets $\{\mathbf{0}\}$ and $\hat{C}$, there exists a nonzero vector $\mathbf{a}$ such that

$$\mathbf{a}^\mathsf{T}\mathbf{x} > 0 \quad \text{for all } \mathbf{x} \in \hat{C}.$$

In particular, $\mathbf{a}^\mathsf{T}\mathbf{x}_i > 0$ i.e., $\frac{\mathbf{a}}{\|\mathbf{a}\|} \in A_{\mathbf{x}_i}$ for $i = 1, \ldots, k$, which contradicts the supposition that $\bigcap_{i=1}^{k} A_{\mathbf{x}_i} = \emptyset$. $\square$

## Convex functions

<u>Definition 8</u> *A function $f\colon W \to \mathbb{R}$, where $W \subset \mathbb{R}^n$ is convex, is called*

- *<u>convex</u>, if for all $\mathbf{x}, \mathbf{y} \in W$ and $\lambda \in (0,1)$ there is*

$$f\big(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}\big) \leqslant \lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{y}),$$

- *<u>strictly convex</u>, if for all $\mathbf{x}, \mathbf{y} \in W$ and $\lambda \in (0,1)$ there is*

$$f\big(\lambda\mathbf{x} + (1-\lambda)\mathbf{y}\big) < \lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{y}),$$

*A function $f$ is (strictly) <u>concave</u> if $-f$ is (strictly) convex.*

<u>Theorem 21</u> *If a function $f\colon W \to \mathbb{R}$, where $W \subset \mathbb{R}^n$ is convex, is Lebesgue-measurable and such that*

$$f\Big(\frac{\mathbf{x}+\mathbf{y}}{2}\Big) \leqslant \frac{f(\mathbf{x}) + f(\mathbf{y})}{2} \quad \text{for all } \mathbf{x}, \mathbf{y} \in W,$$

*then $f$ is a convex function.*

We shall prove a simpler theorem:

<u>Theorem 22</u> *If a function $f\colon W \to \mathbb{R}$, where $W \subset \mathbb{R}^n$ is convex, is continuous and such that*

$$f\Big(\frac{\mathbf{x}+\mathbf{y}}{2}\Big) \leqslant \frac{f(\mathbf{x}) + f(\mathbf{y})}{2} \quad \text{for all } \mathbf{x}, \mathbf{y} \in W,$$

*then $f$ is a convex function.*

<u>Proof.</u> Using induction with respect to $k$, we show that the inequality of the definition of convex functions holds for all $\lambda = \frac{p}{2^k}$, where $p = 0, 1, \ldots, 2^k$. If $k = 1$, then this inequality is satisfied by assumption. Suppose that the inequality is satisfied for some $k \geqslant 1$. Let $p, q \in \mathbb{Z}$, $p, q \geqslant 0$ and $p + q = 2^{k+1}$. Suppose that $p \leqslant q$. Then $p \leqslant 2^k \leqslant q$ and we can write

$$z = \frac{p}{2^{k+1}}\mathbf{x} + \frac{q}{2^{k+1}}\mathbf{y} = \frac{1}{2}\Big(\frac{p}{2^k}\mathbf{x} + \frac{q-2^k}{2^k}\mathbf{y} + \mathbf{y}\Big).$$

Then,

$$f(z) \leqslant \frac{1}{2}f\Big(\frac{p}{2^k}\mathbf{x} + \frac{q-2^k}{2^k}\mathbf{y}\Big) + \frac{1}{2}f(\mathbf{y})$$

$$\leqslant \frac{1}{2}\frac{p}{2^k}f(\mathbf{x}) + \frac{1}{2}\frac{q-2^k}{2^k}f(\mathbf{y}) + \frac{1}{2}f(\mathbf{y}) = \frac{p}{2^{k+1}}f(\mathbf{x}) + \frac{q}{2^{k+1}}f(\mathbf{y}).$$

The first inequality follows from the assumption of the theorem and the second one from the inductive assumption. If $p > q$, then it suffices to exchange $\mathbf{x}$ and $\mathbf{y}$.

The set of numbers $\frac{p}{2^k}$, $k = 1, 2, \ldots$ and $p = 0, \ldots, 2^k$, is dense in the interval $[0,1]$. By the continuity of $f$ we obtain the desired inequality for any $\lambda \in (0,1)$. $\square$

## Properties of convex functions

Below we assume that $W \subset \mathbb{R}^n$ is convex.

<u>Definition 9</u> *The <u>epigraph of a function</u> $f\colon W \to \mathbb{R}$ is the set*

$$\operatorname{epi}(f) = \{\, (\mathbf{x}, z) \in W \times \mathbb{R} \colon z \geqslant f(\mathbf{x}) \,\}.$$

<u>Definition 10</u> *The <u>sublevel set</u> or the <u>trench</u> of a function $f\colon W \to \mathbb{R}$ is the set*

$$W_\alpha(f) = \{\, \mathbf{x} \in W \colon f(\mathbf{x}) \leqslant \alpha \,\}, \quad \alpha \in \mathbb{R}.$$

Theorem 23 *(epigraph theorem) A function* $f$ *is convex if and only if its epigraph is a convex set.*

Theorem 24 *If a function* $f$ *is convex, then its sublevel sets* $W_\alpha(f)$ *are convex for all* $\alpha \in \mathbb{R}$.

Remark. There exist nonconvex functions whose all sublevel sets are convex.

Theorem 25 *If a function* $f$ *is convex, then it is also continuous in* $\operatorname{int} W$.

Theorem 26 *(supporting hyperplane theorem) If* $f$ *is a convex function, then at each point* $\overline{x} \in \operatorname{int} W$ *there exists a supporting hyperplane, i.e., there exists* $\xi \in \mathbb{R}^n$ *such that*

$$f(x) \geqslant f(\overline{x}) + \xi^\mathsf{T}(x - \overline{x}) \quad \text{for all } x \in W.$$

*Moreover, if* $f$ *is strictly convex, then*

$$f(x) > f(\overline{x}) + \xi^\mathsf{T}(x - \overline{x}) \quad \text{for all } x \in W \setminus \{\overline{x}\}.$$

*If* $f$ *is differentiable at* $\overline{x}$, *then in both cases we can take* $\xi = Df(\overline{x})^\mathsf{T}$.

Proof. The set $\operatorname{epi}(f)$ is convex. We apply the weak separation theorem to the sets $U = \operatorname{int} \operatorname{epi}(f)$ and $V = \{(\overline{x}, f(\overline{x}))\}$. There exists a nonzero vector $a = \{(\xi, \alpha)\} \in \mathbb{R}^{n+1}$ such that

$$\xi^\mathsf{T} x + \alpha y \leqslant \xi^\mathsf{T} \overline{x} + \alpha f(\overline{x}) \quad \text{for all } (x, y) \in \operatorname{epi}(f).$$

The inequality above holds for all $y \geqslant f(x)$. Hence, $\alpha \leqslant 0$. It turns out that $\alpha \neq 0$. To prove it, suppose that $\alpha = 0$. Then, for all $x \in W$ there is $\xi^\mathsf{T}(x - \overline{x}) \leqslant 0$. As $\overline{x} \in \operatorname{int} W$, we know that there exists an $\varepsilon > 0$ such that $\overline{x} + \varepsilon \xi \in W$. Let $x = \overline{x} + \varepsilon \xi$. Then $0 \geqslant \xi^\mathsf{T}(x - \overline{x}) = \varepsilon \xi^\mathsf{T} \xi = \varepsilon \|\xi\|^2$; hence, $\xi = 0$. This contradicts the possibility $a \neq 0$, and thus $\alpha < 0$.

We can rescale the vector $a$ to obtain $\alpha = -1$. With that, for all $x \in W$ we obtain

$$\xi^\mathsf{T} x - f(x) \leqslant \xi^\mathsf{T} \overline{x} - f(\overline{x}),$$

which may be rewritten as

$$f(x) \geqslant f(\overline{x}) + \xi^\mathsf{T}(x - \overline{x}),$$

which completes the proof of the first claim.

Suppose that $f$ is strictly convex. Let $\overline{x} \in \operatorname{int} W$. By the first claim, there is $f(x) \geqslant f(\overline{x}) + \xi^\mathsf{T}(x - \overline{x})$ for all $x \in W$. Suppose that there exists $x \in W \setminus \{\overline{x}\}$ such that $f(x) = f(\overline{x}) + \xi^\mathsf{T}(x - \overline{x})$. By the strict convexity of $f$ we obtain

$$f\left(\frac{\overline{x}+x}{2}\right) < \frac{1}{2}\big(f(x) + f(\overline{x})\big) = f(\overline{x}) + \frac{1}{2}\big(f(x) - f(\overline{x})\big) = f(\overline{x}) + \frac{1}{2}\xi^\mathsf{T}(x - \overline{x}).$$

On the other hand, by the existence of the supporting hyperplane, we obtain

$$f\left(\frac{\overline{x}+x}{2}\right) \geqslant f(\overline{x}) + \xi^\mathsf{T}\left(\frac{x+\overline{x}}{2} - x\right) = f(\overline{x}) + \xi^\mathsf{T}\frac{x-\overline{x}}{2}.$$

The two inequalities are inconsistent. Hence, if $f$ is a strictly convex function, there must be $f(x) > f(\overline{x}) + \xi^\mathsf{T}(x - \overline{x})$ and the second claim is proved.

Suppose that $f$ is differentiable at $\overline{x}$. For $x \in W \setminus \{\overline{x}\}$ and $\lambda \in (0, 1)$, by convexity of $f$ we obtain

$$f(x) - f(\overline{x}) = \frac{(1-\lambda)f(\overline{x}) + \lambda f(x) - f(\overline{x})}{\lambda}$$
$$\geqslant \frac{f\big((1-\lambda)\overline{x} + \lambda x\big) - f(\overline{x})}{\lambda} = \frac{f\big(\overline{x} + \lambda(x - \overline{x})\big) - f(\overline{x})}{\lambda}.$$

With this estimation of the divided difference we go to the limit

$$f(x) - f(\overline{x}) \geqslant \lim_{\lambda \searrow 0} \frac{f\big((1-\lambda)\overline{x} + \lambda x\big) - f(\overline{x})}{\lambda} = Df(\overline{x})(x - \overline{x}).$$

The limit exists and is equal to $Df(\overline{x})(x - \overline{x})$ due to the differentiability of $f$. If $f$ is strictly convex, then we can repeat the proof of the second claim with $\xi^\mathsf{T}$ replaced by $Df(\overline{x})$. Then we get the sharp inequality $f(x) - f(\overline{x}) > Df(\overline{x})(x - \overline{x})$ for $x \neq \overline{x}$. $\square$

Corollary 3 *If a function* $f$ *is convex and differentiable at* $\overline{x} \in \operatorname{int} W$, *then* $\overline{x}$ *is a global minimum of* $f$ *if and only if* $Df(\overline{x}) = 0^\mathsf{T}$.

Proof. The gradient of a differentiable function at a minimal point must be equal to $0^\mathsf{T}$; hence, $Df(\overline{x}) = 0^\mathsf{T}$ is a necessary condition. Suppose that it is satisfied. Then, for any $x \in W$ we have

$$f(x) \geqslant f(\overline{x}) + Df(\overline{x})(x - \overline{x}) = f(\overline{x}),$$

which proves that $\overline{\mathbf{x}}$ is a global minimum. $\square$

## Properties of convex functions

<u>Theorem 27</u> *Let $W \subset \mathbb{R}^n$ be a convex set with a nonempty interior. If at each point $\overline{\mathbf{x}} \in \operatorname{int} W$ there exists a vector $\boldsymbol{\xi} \in \mathbb{R}^n$ such that*

$$f(\mathbf{x}) \geqslant f(\overline{\mathbf{x}}) + \boldsymbol{\xi}^\mathsf{T}(\mathbf{x} - \overline{\mathbf{x}}) \quad \text{for all } \mathbf{x} \in W,$$

*then the function $f$ is convex. If the inequality is sharp for $\mathbf{x} \neq \overline{\mathbf{x}}$, then $f$ is strictly convex.*

<u>Proof.</u> Let $\mathbf{x} \in \operatorname{int} W$, $\mathbf{y} \in W$ and $\lambda \in (0,1)$. Denote $\mathbf{x}_\lambda = \lambda \mathbf{x} + (1-\lambda)\mathbf{y}$. We are going to prove that $f(\mathbf{x}_\lambda) \leqslant \lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{y})$. By Lemma 3, $\mathbf{x}_\lambda \in \operatorname{int} W$. By assumption, there exists $\boldsymbol{\xi} \in \mathbb{R}^n$ such that

$$f(\mathbf{x}) \geqslant f(\mathbf{x}_\lambda) + \boldsymbol{\xi}^\mathsf{T}(\mathbf{x} - \mathbf{x}_\lambda), \quad f(\mathbf{y}) \geqslant f(\mathbf{x}_\lambda) + \boldsymbol{\xi}^\mathsf{T}(\mathbf{y} - \mathbf{x}_\lambda).$$

Hence,

$$\lambda f(\mathbf{x}) + (1-\lambda)f(\mathbf{y}) \geqslant f(\mathbf{x}_\lambda) + \boldsymbol{\xi}^\mathsf{T}\left[\lambda(\mathbf{x} - \mathbf{x}_\lambda) + (1-\lambda)(\mathbf{y} - \mathbf{x}_\lambda)\right] = f(\mathbf{x}_\lambda),$$

as the terms in the brackets cancel each other out. The convexity of $f$ is proved. If the assumed inequalities are sharp, then also the inequalities in the calculation above are sharp and the function $f$ is strictly convex. $\square$

<u>Theorem 28</u> *Let $W \subset \mathbb{R}^n$ be nonempty, open and convex and let $f\colon W \to \mathbb{R}$ be twice differentiable. Then,*

*I) $f$ is convex if and only if the Hessian $\mathrm{D}^2 f(\mathbf{x})$ is nonnegative-definite for all $\mathbf{x} \in W$,*

*II) if the Hessian is positive-definite for all $\mathbf{x} \in W$, then $f$ is strictly convex (this is not a necessary condition).*

<u>Proof.</u> Suppose that the Hessian is nonnegative-definite for all $\mathbf{x} \in W$. Then, by Corollary 1, for all $\overline{\mathbf{x}}, \mathbf{x} \in W$ we have

$$f(\mathbf{x}) = f(\overline{\mathbf{x}}) + \mathrm{D}f(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) + \frac{1}{2}(\mathbf{x} - \overline{\mathbf{x}})^\mathsf{T}\mathrm{D}^2 f(\tilde{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}),$$

where $\tilde{\mathbf{x}}$ is a point of the line segment $\overline{\overline{\mathbf{x}}\mathbf{x}}$. As the Hessian is assumed to be nonnegative-definite, the last term above is nonnegative. Hence,

$$f(\mathbf{x}) \geqslant f(\overline{\mathbf{x}}) + \mathrm{D}f(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}).$$

This inequality holds for all $\overline{\mathbf{x}}, \mathbf{x} \in W$, the function $f$ is convex by Theorem 27.

If the Hessian is positive-definite in $W$, then for $\mathbf{x} \neq \overline{\mathbf{x}}$ the last inequality is sharp, and the function $f$ is strictly convex.

Now we prove that the convexity of $f$ implies that the Hessian is nonnegative-definite. Assume that $f$ is convex. Let $\overline{\mathbf{x}} \in W$ and $\mathbf{h} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ be fixed. As the set $W$ is open, there exists $\delta > 0$ such that $\overline{\mathbf{x}} + t\mathbf{h} \in W$ for all $t \in (-\delta, \delta)$. Let $g(t) \stackrel{\text{def}}{=} f(\overline{\mathbf{x}} + t\mathbf{h})$. It is a convex and twice differentiable function of one variable. By Theorem 26,

$$g(t) \geqslant g(0) + g'(0)t, \quad t \in (-\delta, \delta).$$

Using the Taylor's formula with the remainder in Peano form, we obtain

$$g(t) = g(0) + g'(0)t + \frac{1}{2}g''(0)t^2 + o(t^2), \quad t \in (-\delta, \delta).$$

The last inequality and the Taylor's formula give us the following estimation:

$$\frac{1}{2}g''(0)t^2 + o(t^2) \geqslant 0.$$

After dividing both sides by $t^2$ we get

$$\frac{1}{2}g''(0) + \frac{o(t^2)}{t^2} \geqslant 0.$$

With $t$ tending to $0$, the second term vanishes, which leaves us $g''(0) \geqslant 0$. With this we return to the function $f$:

$$g'(t) = \mathrm{D}f(\overline{\mathbf{x}} + t\mathbf{h})\mathbf{h}, \quad g''(t) = \mathbf{h}^\mathsf{T}\mathrm{D}^2 f(\overline{\mathbf{x}} + t\mathbf{h})\mathbf{h}.$$

Hence, $g''(0) = \mathbf{h}^\mathsf{T}\mathrm{D}^2 f(\overline{\mathbf{x}})\mathbf{h}$. As the vector $\mathbf{h}$ may be arbitrarily chosen, the Hessian at $\overline{\mathbf{x}}$ is nonnegative-definite. $\square$

## Subdifferential

We are going to generalise the notion of derivative to non-differentiable convex functions. Let $W \subset \mathbb{R}^n$ be a convex set and $f\colon W \to \mathbb{R}$ be a convex function.

<u>Definition 11</u>  *A vector $\boldsymbol{\xi} \in \mathbb{R}^n$ is called the <u>subgradient</u> of the function f at a point $\mathbf{x}_0 \in W$, if*

$$f(\boldsymbol{x}) \geqslant f(\boldsymbol{x}_0) + \boldsymbol{\xi}^\mathsf{T}(\boldsymbol{x} - \boldsymbol{x}_0) \quad \text{for all } \boldsymbol{x} \in W.$$

*The set of all subgradients of f at $\mathbf{x}_0$ is called the <u>subdifferential</u> and it is denoted by $\partial f(\mathbf{x}_0)$.*

<u>Corollary 4</u>  *If $W \subset \mathbb{R}^n$ is a convex set with the nonempty interior, then $f\colon W \to \mathbb{R}$ is a convex function if and only if at each point $\boldsymbol{x} \in \operatorname{int} W$ there exists a subgradient, i.e.,*

$$\partial f(\boldsymbol{x}) \neq \emptyset.$$

<u>Proof.</u>  By the supporting hyperplane theorem, the convexity implies the existence of the gradient at each point of $W$. By Theorem 27, it is a sufficient condition. $\square$

<u>Lemma 4</u>  *Let $W \subset \mathbb{R}^n$ be convex and let $f\colon W \to \mathbb{R}$ be a convex function. Then the subdifferential $\partial f(\boldsymbol{x})$ is a convex and closed set. If $\boldsymbol{x} \in \operatorname{int} W$, then the subdifferential is also bounded, and thus it is a compact set.*

<u>Proof.</u>  The proof of convexity and closedness is an exercise. Let $\overline{\boldsymbol{x}} \in \operatorname{int} W$ be fixed. There exists $\varepsilon > 0$ such that the ball $B(\overline{\boldsymbol{x}}, \varepsilon) \subset \operatorname{int} W$. For any $\boldsymbol{\xi} \in \partial f(\overline{\boldsymbol{x}})$ we have

$$f(\boldsymbol{x}) \geqslant f(\overline{\boldsymbol{x}}) + \boldsymbol{\xi}^\mathsf{T}(\boldsymbol{x} - \overline{\boldsymbol{x}}) \quad \text{for all } \boldsymbol{x} \in W.$$

Therefore,

$$\sup_{\boldsymbol{x} \in B(\overline{\boldsymbol{x}}, \varepsilon)} f(\boldsymbol{x}) \geqslant f(\overline{\boldsymbol{x}}) + \sup_{\boldsymbol{x} \in B(\overline{\boldsymbol{x}}, \varepsilon)} \boldsymbol{\xi}^\mathsf{T}(\boldsymbol{x} - \overline{\boldsymbol{x}}).$$

The left hand side does not depend on $\boldsymbol{\xi}$, and, by continuity of f in $\operatorname{int} W$, it is finite. The supremum on the right hand side is attended for $\boldsymbol{x} = \overline{\boldsymbol{x}} + \varepsilon \boldsymbol{\xi}/\|\boldsymbol{\xi}\|$, and it is equal to $\varepsilon\|\boldsymbol{\xi}\|$. Hence,

$$\varepsilon\|\boldsymbol{\xi}\| \leqslant \sup_{\boldsymbol{x} \in B(\overline{\boldsymbol{x}}, \varepsilon)} f(\boldsymbol{x}) - f(\overline{\boldsymbol{x}}),$$

which proves that the set $\partial f(\overline{\boldsymbol{x}})$ is bounded. $\square$

<u>Definition 12</u>  *The <u>directional derivative</u> of a function f at a point $\overline{\mathbf{x}}$ in the direction $\mathbf{d}$ is the limit*

$$f'(\overline{\mathbf{x}}; \mathbf{d}) = \lim_{\lambda \searrow 0} \frac{f(\overline{\mathbf{x}} + \lambda\mathbf{d}) - f(\overline{\mathbf{x}})}{\lambda}.$$

<u>Definition 13</u>  *A <u>divided difference of order 0</u> of a function $f\colon I \subset \mathbb{R} \to \mathbb{R}$ at a point $t_i$ is the number $f[t_i] = f(t_i)$. A <u>divided difference of order $k > 0$</u> at different points $t_i, \ldots, t_{i+k}$ is given by the recursive formula*

$$f[t_i, \ldots, t_{i+k}] = \frac{f[t_i, \ldots, t_{i+k-1}] - f[t_{i+1}, \ldots, t_{i+k}]}{t_i - t_{i+k}}.$$

<u>Lemma 5</u>  *I) Divided differences of any order are symmetric functions of the arguments $t_i, \ldots, t_{i+k}$, i.e., $f[t_i, \ldots, t_{i+k}] = f[t_{\sigma(i)}, \ldots, t_{\sigma(i+k)}]$ for any permutation $\sigma$ of the set $\{i, \ldots, i+k\}$.*

*II) If f is a convex function, then the first order divided difference $f[t_i, t_{i+1}]$ is a monotone (nondecreasing) function of the arguments $t_i, t_{i+1}$.*

<u>Proof.</u>  The proof of (I) is an exercise. To prove (II) we notice that

$$f[x_0, x_1, x_2] = c_0 f(x_0) + c_1 f(x_1) + c_2 f(x_2),$$
$$c_0 = \frac{1}{(x_1 - x_0)(x_2 - x_0)}, \quad c_1 = \frac{1}{(x_2 - x_1)(x_0 - x_1)}, \quad c_2 = \frac{1}{(x_0 - x_2)(x_1 - x_2)}.$$

Assume that $x_0 < x_1 < x_2$; then, $c_0, c_2 > 0$, $c_1 < 0$. Then,

$$\lambda = \frac{x_1 - x_0}{x_2 - x_0} \in (0, 1), \quad (1 - \lambda) = \frac{x_2 - x_1}{x_2 - x_0}.$$

We can check that if $f(x_1) = (1 - \lambda)f(x_0) + \lambda f(x_2)$, then $f[x_0, x_1, x_2] = 0$; as the value of a convex function at $x_1$ is less than or equal to this expression, due to $c_1 < 0$ there is $f[x_0, x_1, x_2] \geqslant 0$.

Now we use the symmetry of the divided differences; we choose the numbers $t_0, t_1, t_2$ such that $t_1 < t_2$. Then,

$$0 \leqslant f[t_2, t_0, t_1] = \frac{f[t_2, t_0] - f[t_0, t_1]}{t_2 - t_1} = \frac{f[t_2, t_0] - f[t_1, t_0]}{t_2 - t_1} = \frac{f[t_0, t_2] - f[t_0, t_1]}{t_2 - t_1}.$$

Hence, if f is convex and $t_2 > t_1$, then $f[t_2, t_0] \geqslant f[t_1, t_0]$ and $f[t_0, t_2] \geqslant f[t_0, t_1]$. $\square$

<u>Lemma 6</u> *Let $W \subset \mathbb{R}^n$ be a convex open set and let $f\colon W \to \mathbb{R}$ be a convex function. Then, for all $\mathbf{d} \in \mathbb{R}^n$ and $\mathbf{x} \in W$*

*I) there exists the directional derivative $f'(\overline{\mathbf{x}}; \mathbf{d})$,*

*II) $f'(\overline{\mathbf{x}}; \mathbf{d}) = \inf_{\lambda > 0} \frac{f(\overline{\mathbf{x}} + \lambda \mathbf{d}) - f(\overline{\mathbf{x}})}{\lambda}$,*

*III) $f'(\overline{\mathbf{x}}; \mathbf{d}) \geqslant -f'(\overline{\mathbf{x}}; -\mathbf{d})$.*

<u>Proof.</u> Let $g(t) \stackrel{\text{def}}{=} f(\mathbf{x} + t\mathbf{d})$ for $t$ such that $\mathbf{x} + t\mathbf{d} \in W$. As $W$ is open, the function $g$ is defined in an interval $(-\delta, \delta)$. This function is convex. By Lemma 5, its divided difference is monotone, i.e., for $t_1, t_2 \in (-\delta, \delta) \setminus \{0\}$, $t_1 < t_2$, we have

$$\frac{g(t_1) - g(0)}{t_1} \leqslant \frac{g(t_2) - g(0)}{t_2}. \tag{*}$$

By the monotonicity of the divided difference, there exists the left-side derivative, $g'(0^-)$, and the right-side derivative, $g'(0^+)$, such that $g'(0^-) \leqslant g'(0^+)$, and

$$g'(0^-) = \sup_{t<0} \frac{g(t) - g(0)}{t}, \quad g'(0^+) = \inf_{t>0} \frac{g(t) - g(0)}{t}.$$

It suffices to notice that $f'(\mathbf{x}; \mathbf{d}) = g'(0^+)$ and $f'(\mathbf{x}; -\mathbf{d}) = -g'(0^-)$. $\square$

<u>Lemma 7</u> *Let $W \subset \mathbb{R}^n$ be a convex open set and let $f\colon W \to \mathbb{R}$ be a convex function. Then a vector $\xi$ is a subgradient if and only if*

$$f'(\overline{\mathbf{x}}; \mathbf{d}) \geqslant \xi^\mathsf{T} \mathbf{d} \quad \text{for all } \mathbf{d} \in \mathbb{R}^n.$$

<u>Proof.</u> Let $\overline{\mathbf{x}} \in W$ and $\xi \in \partial f(\overline{\mathbf{x}})$. Then, for $\lambda > 0$ and $\mathbf{d} \in \mathbb{R}^n$ such that $\overline{\mathbf{x}} + \lambda \mathbf{d} \in W$, there is

$$f(\overline{\mathbf{x}} + \lambda \mathbf{d}) \geqslant f(\overline{\mathbf{x}}) + \lambda \xi^\mathsf{T} \mathbf{d}.$$

Hence,

$$\frac{f(\overline{\mathbf{x}} + \lambda \mathbf{d}) - f(\overline{\mathbf{x}})}{\lambda} \geqslant \xi^\mathsf{T} \mathbf{d},$$

i.e., $f'(\overline{\mathbf{x}}; \mathbf{d}) \geqslant \xi^\mathsf{T} \mathbf{d}$.

Now, let $\xi \in \mathbb{R}^n$ be a vector such that $f'(\overline{\mathbf{x}}; \mathbf{d}) \geqslant \xi^\mathsf{T} \mathbf{d}$ for all $\mathbf{d} \in \mathbb{R}^n$. By Lemma 6(II), for $\lambda > 0$ we obtain

$$f'(\overline{\mathbf{x}}; \mathbf{d}) \leqslant \frac{f(\overline{\mathbf{x}} + \lambda \mathbf{d}) - f(\overline{\mathbf{x}})}{\lambda}.$$

Hence,

$$f(\overline{\mathbf{x}} + \lambda \mathbf{d}) \geqslant f(\overline{\mathbf{x}}) + \lambda \xi^\mathsf{T} \mathbf{d}.$$

As $\lambda$ and $\mathbf{d}$ may be arbitrary (such that $\overline{\mathbf{x}} + \lambda \mathbf{d} \in W$), the vector $\xi$ is a subgradient. $\square$

<u>Theorem 29</u> *Let $f\colon W \to \mathbb{R}$ be a convex function in an open convex set $W \subset \mathbb{R}^n$. For each point $\overline{\mathbf{x}} \in W$ and vector $\mathbf{d} \in \mathbb{R}^n$ there is*

$$f'(\overline{\mathbf{x}}; \mathbf{d}) = \max_{\xi \in \partial f(\overline{\mathbf{x}})} \xi^\mathsf{T} \mathbf{d}.$$

*Moreover, the function $f$ is differentiable at $\overline{\mathbf{x}}$ if and only if the subdifferential $\partial f(\overline{\mathbf{x}})$ has only one element. This element is $Df(\overline{\mathbf{x}})^\mathsf{T}$.*

<u>Proof.</u> By Lemma 7, $f'(\overline{\mathbf{x}}; \mathbf{d}) \geqslant \xi^\mathsf{T} \mathbf{d}$ for all $\xi \in \partial f(\overline{\mathbf{x}})$. Hence,

$$f'(\overline{\mathbf{x}}; \mathbf{d}) \geqslant \max_{\xi \in \partial f(\overline{\mathbf{x}})} \xi^\mathsf{T} \mathbf{d}.$$

The opposite inequality may be proved using the weak separation theorem. Let

$$C_1 = \{ (\mathbf{x}, z) \in W \times \mathbb{R} \colon z > f(\mathbf{x}) \}$$
$$C_2 = \{ (\mathbf{x}, z) \in W \times \mathbb{R} \colon \mathbf{x} = \overline{\mathbf{x}} + \lambda \mathbf{d}, \, z = f(\overline{\mathbf{x}}) + \lambda f'(\overline{\mathbf{x}}; \mathbf{d}), \, \lambda \geqslant 0 \}.$$

Note that $C_1$ is the interior of the epigraph of $f$; hence, $C_1$ is a convex set. The set $C_2$ is a halfline with the origin at $(\overline{\mathbf{x}}, f(\overline{\mathbf{x}}))$ and the direction $(\mathbf{d}, f'(\overline{\mathbf{x}}; \mathbf{d}))$, which is also a convex set. This halfline is the graph of a linear approximation of $f$ along the line segment $\{ \overline{\mathbf{x}} + \lambda \mathbf{d} \colon \lambda \geqslant 0 \} \cap W$.

By Lemma 6, $f'(\mathbf{x}; \mathbf{d}) \leqslant \frac{f(\overline{\mathbf{x}} + \lambda \mathbf{d}) - f(\overline{\mathbf{x}})}{\lambda}$, i.e.,

$$f(\overline{\mathbf{x}} + \lambda \mathbf{d}) \geqslant f(\overline{\mathbf{x}}) + \lambda f'(\mathbf{x}; \mathbf{d}).$$

Hence, the sets $C_1$ and $C_2$ are disjoint. By the weak separation theorem, there exists a nonzero vector $(\mu, \gamma) \in \mathbb{R}^{n+1}$ such that

$$\mu^\mathsf{T} \mathbf{x} + \gamma z \geqslant \mu^\mathsf{T}(\overline{\mathbf{x}} + \lambda \mathbf{d}) + \gamma\big(f(\overline{\mathbf{x}}) + \lambda f'(\overline{\mathbf{x}}; \mathbf{d})\big), \quad \text{for all } (\mathbf{x}, z) \in C_1, \, \lambda \in [0, L),$$

where $L = \sup\{\lambda > 0 \colon \overline{x} + \lambda d \in W\}$. The number $\gamma$ cannot be negative, as the left hand side might be arbitrarily small (after choosing a large $z$). Also, $\gamma$ cannot be zero, as in that case the inequality $\mu^T(x - \overline{x}) \geqslant \lambda \mu^T d$ would have to hold for all $x \in W$, and this is possible only with $\mu = 0$. This inconsistency with $(\mu, \gamma) \neq 0$ proves that $\gamma > 0$.

By rescaling the vector $(\mu, \gamma)$, we can assume $\gamma = 1$. Then,

$$\mu^T x + z \geqslant \mu^T(\overline{x} + \lambda d) + \big(f(\overline{x}) + \lambda f'(\overline{x}; d)\big), \quad \text{for all } (x, z) \in C_1, \lambda \in [0, L],$$

With $z$ tending to $f(x)$ we obtain the following inequality, which holds for all $x \in W$ and $\lambda \in [0, L]$:

$$\mu^T x + f(x) \geqslant \mu^T(\overline{x} + \lambda d) + \big(f(\overline{x}) + \lambda f'(\overline{x}; d)\big) \tag{*}$$

With $\lambda = 0$, we obtain

$$\mu^T(x - \overline{x}) + f(x) \geqslant f(\overline{x}),$$

i.e.,

$$f(x) \geqslant f(\overline{x}) - \mu^T(x - \overline{x}),$$

Hence, $-\mu \in \partial f(\overline{x})$. Now, substituting $\lambda > 0$ and $x = \overline{x}$ in (*), we obtain

$$-\mu^T(\lambda d) \geqslant \lambda f'(\overline{x}; d),$$

i.e.,

$$\sup_{\xi \in \partial f(\overline{x})} \xi^T d \geqslant f'(\overline{x}; d).$$

The proof of the first claim is complete.

To prove the second claim, we notice that the function $f$ is differentiable at $\overline{x}$ if and only if there exists $\alpha \in \mathbb{R}^n$ such that $f'(\overline{x}; d) = \alpha^T d$ for all $d \in \mathbb{R}^n$ (then $\alpha = Df(\overline{x})^T$). Thus, if the set $\partial f(\overline{x})$ has only one element, then $f$ is differentiable at $\overline{x}$.

Suppose that $f$ is diferentiable at $\overline{x}$. Then, for sufficiently small $\lambda > 0$ and $d \in \mathbb{R}^n$ (without loss of generality we assume that $\|d\| = 1$), we have

$$f(\overline{x} + \lambda d) = f(\overline{x}) + \lambda Df(\overline{x})d + o(\lambda).$$

By definition of the subgradient, we have

$$f(\overline{x} + \lambda d) \geqslant f(\overline{x}) + \lambda \xi^T d,$$

where $\xi$ is a subgradient. By subtracting the sides of the above, we obtain

$$\lambda\big(\xi^T - Df(\overline{x})\big)d - o(\lambda) \leqslant 0.$$

After dividing both sides of this inequality by $\lambda$ and passing with $\lambda$ to $0$, we obtain

$$\big(\xi^T - Df(\overline{x})\big)d \leqslant 0.$$

The substitution $d = \pm\frac{\xi - Df(\overline{x})^T}{\|\xi - Df(\overline{x})^T\|}$ yields the equality

$$\xi^T = Df(\overline{x}),$$

which means that the subdifferential consists of one element. $\square$

<u>Theorem 30</u> *Let $W \subset \mathbb{R}^n$ be an open convex set and $f_1, f_2 \colon W \to \mathbb{R}$ be convex functions.*

*I) Let $f = f_1 + f_2$. Then, $\partial f(x) = \partial f_1(x) + \partial f_2(x)$, i.e.,*

$$\partial f(x) = \partial f_1(x) + \partial f_2(x) = \{\xi_1 + \xi_2 \colon \xi_1 \in \partial f_1(x), \xi_2 \in \partial f_2(x)\}.$$

*II) Let $f = \max(f_1, f_2)$. Then,*

$$\partial f(x) = \begin{cases} \partial f_1(x) & \text{if } f_1(x) > f_2(x), \\ \operatorname{conv}\big(\partial f_1(x) \cup \partial f_2(x)\big) & \text{if } f_1(x) = f_2(x), \\ \partial f_2(x) & \text{if } f_1(x) < f_2(x), \end{cases}$$

*where $\operatorname{conv}\big(\partial f_1(x) \cup \partial f_2(x)\big)$ is the convex hull of the union $\partial f_1(x) \cup \partial f_2(x)$, i.e., the set of all convex combinations of the subgradients in both subdifferentials.*

<u>Proof.</u> (I): Let $\overline{x} \in W$. Let $\xi_1 \in \partial f_1(\overline{x})$ and $\xi_2 \in \partial f_2(\overline{x})$. Then, for all $x \in W$ we have

$$f_1(x) \geqslant f_1(\overline{x}) + \xi_1^T(x - \overline{x}),$$
$$f_2(x) \geqslant f_2(\overline{x}) + \xi_2^T(x - \overline{x}).$$

By adding the above inequalities side by side we obtain

$$f(x) \geqslant f(\overline{x}) + (\xi_1 + \xi_2)^T(x - \overline{x}),$$

i.e., $\xi_1 + \xi_2 \in \partial f(\overline{x})$. Hence, $\partial f_1(\overline{x}) + \partial f_2(\overline{x}) \subset \partial f(\overline{x})$. Suppose that there exists $\xi \in \partial f(\overline{x})$ such that $\xi \notin \partial f_1(\overline{x}) + \partial f_2(\overline{x})$. By Lemma 4, the subdifferentials $\partial f_1(\overline{x})$

and $\partial f_2(\overline{x})$ are compact convex sets. Their algebraic sum is, therefore, also a compact convex set. By the strong separation theorem, applied to the sets $\{\xi\}$ and $\partial f_1(\overline{x}) + \partial f_2(\overline{x})$, there exists $\mu \in \mathbb{R}^n$, such that

$$\mu^\mathsf{T}\xi_1 + \mu^\mathsf{T}\xi_2 < \mu^\mathsf{T}\xi, \quad \text{for all } \xi_1 \in \partial f_1(\overline{x}) \text{ and } \xi_2 \in \partial f_2(\overline{x}).$$

We take $\xi_1, \xi_2$ to maximise the left hand side. By Theorem 29,

$$f_1'(\overline{x}; \mu) + f_2'(\overline{x}; \mu) < \xi^\mathsf{T}\mu \leqslant f'(\overline{x}; \mu).$$

On the other hand, by the properties of directional derivatives,

$$f_1'(\overline{x}; \mu) + f_2'(\overline{x}; \mu) = f'(\overline{x}; \mu).$$

This is an inconsistency; a vector $\xi$ with assumed properties cannot exist, which completes the proof of (I).

Now we prove (II). The form of the subdifferential $\partial f$ in the sets $W_1$ and $W_2$ defined as $W_i = \{x \in W: f_i(x) > f_{3-i}(x)\}$ is obvious, which leaves the set $W_0 = \{x \in W: f_1(x) = f_2(x)\}$ to investigate. Let $\overline{x} \in W$ and $f_1(\overline{x}) = f_2(\overline{x})$. Denote $A = \text{conv}\big(\partial f_1(\overline{x}) \cup \partial f_2(\overline{x})\big)$. For $i = 1, 2$ and $x \in W$ we have

$$f(x) - f(\overline{x}) \geqslant f_i(x) - f(\overline{x}) = f_i(x) - f_i(\overline{x}) \geqslant \xi_i^\mathsf{T}(x - \overline{x}), \quad \text{for all } \xi_i \in \partial f_i(\overline{x}).$$

From the above we obtain $\partial f_1(\overline{x}) \cup \partial f_2(\overline{x})$. By convexity of the subdifferential, $A \subset \partial f(\overline{x})$. Suppose that there exists $\xi \in \partial f(\overline{x}) \setminus A$. The set $A$ is convex and compact. By the strong separation theorem applied to the sets $\{\xi\}$ and $A$, there exists a vector $\mu \in \mathbb{R}^n$ and a constant $b$ such that

$$\mu^\mathsf{T}\tilde{\xi} < b < \mu^\mathsf{T}\xi \quad \text{for all } \tilde{\xi} \in A.$$

In particular, $\mu^\mathsf{T}\xi_i < b$ for $\xi_i \in \partial f_i(\overline{x})$, $i = 1, 2$. By Theorem 29,

$$\max\big\{f_1'(\overline{x}; \mu), f_2'(\overline{x}; \mu)\big\} \leqslant b.$$

Similarly, $b < \xi^\mathsf{T}\mu \leqslant f'(\overline{x}; \mu)$; hence,

$$\max\big\{f_1'(\overline{x}; \mu), f_2'(\overline{x}; \mu)\big\} < f'(\overline{x}; \mu). \tag{*}$$

On the other hand, by definition of the directional derivative, due to $f(\overline{x}) = f_1(\overline{x}) = f_2(\overline{x})$, we obtain the equality

$$\frac{f(\overline{x} + \lambda d) - f(\overline{x})}{\lambda} = \max\left\{\frac{f_1(\overline{x} + \lambda d) - f(\overline{x})}{\lambda}, \frac{f_2(\overline{x} + \lambda d) - f(\overline{x})}{\lambda}\right\}, \quad \lambda > 0.$$

Passing with $\lambda$ to $0$, we obtain

$$f'(\overline{x}; d) = \max\big\{f_1'(\overline{x}; d), f_2'(\overline{x}; d)\big\}.$$

With $d = \mu$ we obtain an inconsistency with (*). Hence, the set $\partial f(\overline{x}) \setminus A$ is empty. $\square$

<u>Theorem 31</u> *Let $W \subset \mathbb{R}^n$ be an open and convex set, let $f: W \to \mathbb{R}$ be a convex function and let $A$ be an $n \times m$ matrix. If $\tilde{W} = \{x \in \mathbb{R}^m: Ax \in W\}$, then $\tilde{W}$ is an open convex set and the function $F: \tilde{W} \to \mathbb{R}$ given by the formula $F(x) = f(Ax)$ at any point $x \in \tilde{W}$ has the subdifferential given by*

$$\partial F(x) = A^\mathsf{T}\partial f(Ax).$$

<u>Proof.</u> Let $\overline{x} \in \tilde{W}$ and let $\xi \in \partial f(A\overline{x})$. Then,

$$f(Ax) \geqslant f(A\overline{x}) + \xi^\mathsf{T}(Ax - A\overline{x}) = f(A\overline{x}) + (A^\mathsf{T}\xi)^\mathsf{T}(x - \overline{x}),$$

i.e., $A^\mathsf{T}\xi \in \partial F(\overline{x})$. Hence, $A^\mathsf{T}\partial f(A\overline{x}) \subset \partial F(A\overline{x})$. Suppose that there exists $\xi \in \partial F(A\overline{x}) \setminus A^\mathsf{T}\partial f(A\overline{x})$. The set $A^\mathsf{T}\partial f(A\overline{x})$ is convex and closed, as the image of a closed and convex set in a linear transformation. We apply the strong separation theorem to this set and $\{\xi\}$. There exists $\mu \in \mathbb{R}^m$ and $b \in \mathbb{R}$ such that

$$\mu^\mathsf{T}A^\mathsf{T}\tilde{\xi} < b < \mu^\mathsf{T}\xi \quad \text{for all } \tilde{\xi} \in \partial f(A\overline{x}).$$

By taking the supremum over $\tilde{\xi} \in \partial f(A\overline{x})$, and using Theorem 29, we obtain $f'(A\overline{x}; A\mu) \leqslant b$. The right hand side may be estimated by the directional derivative: $\mu^\mathsf{T}\xi \leqslant F'(\overline{x}; \mu)$. Hence,

$$f'(A\overline{x}; A\mu) \leqslant b < F'(\overline{x}; \mu).$$

But the directional derivatives satisfy the equality $F'(\overline{x}; d) = f'(A\overline{x}; Ad)$ for all $d \in \mathbb{R}^m$. Thus, we have an inconsistency, which proves that $\partial F(A\overline{x}) \setminus A^\mathsf{T}\partial f(A\overline{x})$ is the empty set. $\square$

# 4. Extrema of convex functions with constraints

We consider a convex function $f\colon W \to \mathbb{R}$ defined in a convex set $W \subset \mathbb{R}^n$ and the following problem:

$$\begin{cases} f(\boldsymbol{x}) \to \min, \\ \boldsymbol{x} \in W. \end{cases}$$

A <u>global solution</u> is a feasible point $\overline{\boldsymbol{x}}$ such that $f(\overline{\boldsymbol{x}}) \leqslant f(\boldsymbol{x})$ for all $\boldsymbol{x} \in W$.
A <u>local solution</u> is a point $\overline{\boldsymbol{x}} \in W$ such that there exists $\varepsilon > 0$ such that $f(\overline{\boldsymbol{x}}) \leqslant f(\boldsymbol{x})$ for all $\boldsymbol{x} \in W \cap B(\overline{\boldsymbol{x}}, \varepsilon)$, i.e., if the point $\overline{\boldsymbol{x}}$ is a minimum of $f$ in its neighbourhood. The local solution is <u>strict</u> if $f(\overline{\boldsymbol{x}}) < f(\boldsymbol{x})$ for $\boldsymbol{x} \neq \overline{\boldsymbol{x}}$.

<u>Theorem 32</u> *Let the set $W \subset \mathbb{R}^n$ be convex, and the function $f\colon W \to \mathbb{R}^n$ be convex. If $\overline{\boldsymbol{x}} \in W$ is a local solution of the problem above, then*

*I) $\overline{\boldsymbol{x}}$ is a global solution,*

*II) the set of global solutions is convex,*

*III) if $f$ is strictly convex, then $\overline{\boldsymbol{x}}$ is a strict local solution,*

*IV) if $\overline{\boldsymbol{x}}$ is a strict local solution, then it is the unique global solution.*

We do not assume the differentiability of $f$.

<u>Proof.</u> (I) is proved by contradiction. Suppose that there exists $\boldsymbol{x}^* \in W$ such that $f(\boldsymbol{x}^*) < f(\overline{\boldsymbol{x}})$. As $\overline{\boldsymbol{x}}$ is a local solution, $f(\overline{\boldsymbol{x}}) \leqslant f(\boldsymbol{x})$ for all $\boldsymbol{x} \in W \cap B(\overline{\boldsymbol{x}}, \varepsilon)$, for some $\varepsilon > 0$. By convexity of $W$, this set contains the line segment $\overline{\boldsymbol{x}}\boldsymbol{x}^*$. This line segment has a nonempty intersection with the ball $B(\overline{\boldsymbol{x}}, \varepsilon)$; there exists $\lambda \in (0, 1)$ such that $\lambda \overline{\boldsymbol{x}} + (1 - \lambda)\boldsymbol{x}^* \in B(\overline{\boldsymbol{x}}, \varepsilon)$. By convexity of $f$, we obtain

$$f\big(\lambda\overline{\boldsymbol{x}} + (1-\lambda)\boldsymbol{x}^*\big) \leqslant \lambda f(\overline{\boldsymbol{x}}) + (1-\lambda)f(\boldsymbol{x}^*) < f(\overline{\boldsymbol{x}}),$$

which contradicts $\overline{\boldsymbol{x}}$ being a local minimum.

The proofs of (II), (III) and (IV) are left as exercises. $\square$

So far, we have shown that a necessary and sufficient condition for a minimum of a differentiable convex function in an *open* set is the zero of the derivative or the gradient. This result may be generalised to arbitrary convex sets.

<u>Theorem 33</u> *Let the set $W \subset \mathbb{R}^n$ be convex and the function $f\colon W \to \mathbb{R}$ be convex. If $f$ is differentiable at $\overline{\boldsymbol{x}} \in W$, then there is the following equivalence: $\overline{\boldsymbol{x}}$ is a minimum if and only if $Df(\overline{\boldsymbol{x}})(\boldsymbol{x} - \overline{\boldsymbol{x}}) \geqslant 0$ for all $\boldsymbol{x} \in W$.*

<u>Remark.</u> To speak of differentiability of $f$ at a point $\overline{\boldsymbol{x}}$, this function must be defined in a neighbourhood of this point, i.e., in a ball $B(\overline{\boldsymbol{x}}, \varepsilon)$, $\varepsilon > 0$. If $\overline{\boldsymbol{x}}$ is at the boundary of $W$, then we assume that $f$ is defined in $W \cup B(\overline{\boldsymbol{x}}, \varepsilon)$, though we omit it in the theorem's assumptions.

<u>Remark.</u> If $\overline{\boldsymbol{x}} \in \operatorname{int} W$, then the condition in the theorem is equivalent to $Df(\overline{\boldsymbol{x}}) = \boldsymbol{0}^\mathsf{T}$.

<u>Proof.</u> Let $Df(\overline{\boldsymbol{x}})(\boldsymbol{x} - \overline{\boldsymbol{x}}) \geqslant 0$ for all $\boldsymbol{x} \in W$. Suppose that there is no minimum at $\overline{\boldsymbol{x}}$. Then, there exists a point $\boldsymbol{x}' \in W$ such that $f(\boldsymbol{x}') < f(\overline{\boldsymbol{x}})$. We construct a sequence $\boldsymbol{x}_k = (1 - \frac{1}{k})\overline{\boldsymbol{x}} + \frac{1}{k}\boldsymbol{x}'$. By convexity of $W$, this sequence is contained in $W$. We consider the directional derivative of $f$ in the direction of the vector $\boldsymbol{x}' - \overline{\boldsymbol{x}}$:

$$f'(\overline{\boldsymbol{x}}; \boldsymbol{x}' - \overline{\boldsymbol{x}}) = \lim_{k\to\infty} \frac{f\big(\overline{\boldsymbol{x}} + \frac{1}{k}(\boldsymbol{x}' - \overline{\boldsymbol{x}})\big) - f(\overline{\boldsymbol{x}})}{1/k} = \lim_{k\to\infty} \frac{f(\boldsymbol{x}_k) - f(\overline{\boldsymbol{x}})}{1/k}$$

$$\leqslant \lim_{k\to\infty} \frac{(1 - \frac{1}{k})f(\overline{\boldsymbol{x}}) + \frac{1}{k}f(\boldsymbol{x}') - f(\overline{\boldsymbol{x}})}{1/k} = f(\boldsymbol{x}') - f(\overline{\boldsymbol{x}}) < 0.$$

By the assumption, we have

$$f'(\overline{\boldsymbol{x}}; \boldsymbol{x}' - \overline{\boldsymbol{x}}) = Df(\overline{\boldsymbol{x}})(\boldsymbol{x}' - \overline{\boldsymbol{x}}) \geqslant 0$$

This inconsistency proves that $\overline{\boldsymbol{x}}$ is a minimum of $f$ in $W$.

Now suppose that $\overline{\boldsymbol{x}}$ is a solution. Let $\boldsymbol{x} \in W$. The convexity of $W$ implies that $\overline{\boldsymbol{x}} + \lambda(\boldsymbol{x} - \overline{\boldsymbol{x}}) = (1 - \lambda)\overline{\boldsymbol{x}} + \lambda\boldsymbol{x} \in W$ for all $\lambda \in [0, 1]$. By definition of the derivative,

$$Df(\overline{\boldsymbol{x}})(\boldsymbol{x} - \overline{\boldsymbol{x}}) = \lim_{\lambda\searrow 0, \lambda<1} \frac{f\big(\overline{\boldsymbol{x}} + \lambda(\boldsymbol{x} - \overline{\boldsymbol{x}})\big) - f(\overline{\boldsymbol{x}})}{\lambda}.$$

As $\overline{\boldsymbol{x}}$ is a minimum, $f\big(\overline{\boldsymbol{x}} + \lambda(\boldsymbol{x} - \overline{\boldsymbol{x}})\big) \geqslant f(\overline{\boldsymbol{x}})$. Hence, $Df(\overline{\boldsymbol{x}})(\boldsymbol{x} - \overline{\boldsymbol{x}}) \geqslant 0$. $\square$

<u>Corollary 5</u> *If $\overline{\boldsymbol{x}} \in W$, where $W \subset \mathbb{R}^n$ is convex, is a local minimum of a (not necessarily convex) function $f\colon W \to \mathbb{R}$, differentiable at $\overline{\boldsymbol{x}}$, then $Df(\overline{\boldsymbol{x}})(\boldsymbol{x} - \overline{\boldsymbol{x}}) \geqslant 0$ for all $\boldsymbol{x} \in W$.*

<u>Theorem 34</u> *Let $\mathbb{X} \subset \mathbb{R}^n$ be a convex and open set and $f\colon \mathbb{X} \to \mathbb{R}$ be a convex function. Suppose that the feasible set $W$ is a subset of $\mathbb{X}$. Then $\overline{x} \in W$ is a minimum if and only if there exists $\xi \in \partial f(\overline{x})$, such that $\xi^{\mathsf{T}}(x - \overline{x}) \geqslant 0$ for all $x \in W$.*

<u>Corollary 6</u> *If $\overline{x} \in \operatorname{int} W$, then $f$ has a global minimum at $\overline{x}$ if and only if $0 \in \partial f(\overline{x})$.*

<u>Proof.</u> Suppose that there exists $\xi \in \partial f(\overline{x})$ such that $\xi^{\mathsf{T}}(x - \overline{x}) \geqslant 0$ for all $x \in W$. As $\xi$ is a subgradient, it follows that

$$f(x) \geqslant f(\overline{x}) + \xi^{\mathsf{T}}(x - \overline{x}), \quad x \in W.$$

Now it suffices to use the assumption to notice that $f(x) \geqslant f(\overline{x})$, i.e., $\overline{x}$ is a minimum.

Now suppose that $\overline{x} \in W$ is a minimum. We define two sets:

$$C_1 = \{ (x, z) \in \mathbb{R}^{n+1} \colon x \in \mathbb{X},\ z > f(x) - f(\overline{x}) \},$$
$$C_2 = \{ (x, z) \in \mathbb{R}^{n+1} \colon x \in W,\ z \leqslant 0 \}.$$

Both sets are convex and the interior of $C_1$ is nonempty (the interior of $C_2$ may be empty if the interior of $W$ is empty). From $\overline{x}$ being a solution it follows that $C_1 \cap C_2 = \emptyset$. We use the weak separation theorem: there exists a nonzero vector $(\mu, \gamma) \in \mathbb{R}^{n+1}$ and a constant $b$ such that

$$\mu^{\mathsf{T}} x + \gamma z \leqslant b, \quad \text{for all } x \in \mathbb{X},\ z > f(x) - f(\overline{x}),$$
$$\mu^{\mathsf{T}} x + \gamma z \geqslant b, \quad \text{for all } x \in W,\ z \leqslant 0.$$

Before we proceed, let's take a look at Figure 3. As we can see, the two sets "touch" each other at the point $(\overline{x}, 0)$. The separating hyperplane must therefore contain this point. It is tangent to the graph of the function $x \to f(x) - f(\overline{x})$ and in fact it determines a subgradient of this function, which is also a subgradient of the function $f$ at $\overline{x}$.

Now we prove it analytically. We can write

$$\mu^{\mathsf{T}}(x - \overline{x}) + \gamma z \leqslant \tilde{b}, \quad \text{for all } x \in \mathbb{X},\ z > f(x) - f(\overline{x}), \tag{$*$}$$
$$\mu^{\mathsf{T}}(x - \overline{x}) + \gamma z \geqslant \tilde{b}, \quad \text{for all } x \in W,\ z \leqslant 0, \tag{$**$}$$



Figure 3: The sets $C_1$ and $C_2$

where $\tilde{b} = b - \mu^{\mathsf{T}}\overline{x}$. We notice that $\gamma$ cannot be positive, as we can take an arbitrarily small $z$ in $(**)$, which leads to inconsistency.

Taking $x = \overline{x}$ and $z = 0$ in $(**)$, we obtain $\tilde{b} \leqslant 0$. On the other hand, with $x = \overline{x}$, the inequality $(*)$ turns into $\gamma z \leqslant \tilde{b}$ for $z > 0$; hence, $\tilde{b} \geqslant 0$. Therefore, $\tilde{b} = 0$. Using this fact we show that $\gamma$ cannot be zero. From $(*)$, due to $\mathbb{X}$ being open, we would then obtain $\mu = 0$ which contradicts $(\mu, \gamma) \neq 0$. Thus we proved that $\gamma < 0$. Taking $z = f(x) - f(\overline{x})$ in $(*)$, we obtain

$$\mu^{\mathsf{T}}(x - \overline{x}) + \gamma\big(f(x) - f(\overline{x})\big) \leqslant 0.$$

After dividing the sides by $\gamma$, which is negative, we obtain

$$\frac{\mu^{\mathsf{T}}}{\gamma}(x - \overline{x}) + f(x) - f(\overline{x}) \geqslant 0,$$

which proves that $\xi = -\frac{\mu}{\gamma} \in \partial f(\overline{x})$. Taking $z = 0$ in $(**)$, we obtain $\mu^{\mathsf{T}}(x - \overline{x}) \geqslant \tilde{b} = 0$. After dividing the sides by $-\gamma$, we obtain $\xi^{\mathsf{T}}(x - \overline{x}) \geqslant 0$, and the proof is complete. $\square$

## Pseudoconvex functions

We introduce a family of functions such that

$$\mathrm{Df}(\overline{\mathbf{x}}) = \mathbf{0}^\mathsf{T} \quad \Leftrightarrow \quad \overline{\mathbf{x}} \text{ is a global minimum of } \mathsf{f}.$$

<u>Definition 14</u> *Let $W \subset \mathbb{R}^n$ be convex, open and nonempty and let $\mathsf{f}\colon W \to \mathbb{R}$. The function $\mathsf{f}$ is <u>pseudoconvex in W</u> if it is differentiable in $W$ and*

$$\mathrm{Df}(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geqslant 0 \quad \Rightarrow \quad \mathsf{f}(\mathbf{y}) \geqslant \mathsf{f}(\mathbf{x}) \qquad \text{for all } \mathbf{x}, \mathbf{y} \in W.$$

*A function $\mathsf{f}$ is <u>strictly pseudoconvex in W</u> if*

$$\mathrm{Df}(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geqslant 0 \quad \Rightarrow \quad \mathsf{f}(\mathbf{y}) > \mathsf{f}(\mathbf{x}) \qquad \text{for all } \mathbf{x}, \mathbf{y} \in W, \ \mathbf{x} \neq \mathbf{y}.$$

*A function $\mathsf{f}$ is <u>pseudoconvex at a point $\overline{\mathbf{x}} \in W$</u> if it is differentiable at $\overline{\mathbf{x}}$ and*

$$\mathrm{Df}(\overline{\mathbf{x}})(\mathbf{y} - \overline{\mathbf{x}}) \geqslant 0 \quad \Rightarrow \quad \mathsf{f}(\mathbf{y}) \geqslant \mathsf{f}(\overline{\mathbf{x}}) \qquad \text{for all } \mathbf{y} \in W.$$

*Similarly is defined a function <u>strictly pseudoconvex at a point $\overline{\mathbf{x}} \in W$</u>.*

*A function $\mathsf{f}$ is <u>(strictly) pseudoconcave</u> if $-\mathsf{f}$ is (strictly) pseudoconvex.*



Figure 4: Examples explaining the notion of pseudoconvexity

<u>Remark.</u> Pseudoconvexity at a point is a property of a function related with the entire set $W$, even if the differentiability is needed at that point only.

<u>Remark.</u> A condition equivalent to that in the definition is the following:

$$\mathsf{f}(\mathbf{y}) < \mathsf{f}(\mathbf{x}) \quad \Rightarrow \quad \mathrm{Df}(\mathbf{x})(\mathbf{y} - \mathbf{x}) < 0.$$

<u>Lemma 8</u> *Let $\mathsf{f}\colon W \to \mathbb{R}$, where $W \subset \mathbb{R}^n$ is nonempty, open and convex. If $\mathsf{f}$ is (strictly) convex and differentiable in $W$, then $\mathsf{f}$ is (strictly) pseudoconvex.*

<u>Proof.</u> Suppose that $\mathsf{f}$ is convex. By the supporting hyperplane theorem, for any $\overline{\mathbf{x}}, \mathbf{x} \in W$ we have

$$\mathsf{f}(\mathbf{x}) \geqslant \mathsf{f}(\overline{\mathbf{x}}) + \mathrm{Df}(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}).$$

Thus, if $\mathrm{Df}(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) \geqslant 0$, then $\mathsf{f}(\mathbf{x}) \geqslant \mathsf{f}(\overline{\mathbf{x}})$ and $\mathsf{f}$ is indeed pseudoconvex. The proof of strict pseudoconvexity of a strictly convex function is similar. $\square$

<u>Lemma 9</u> *Let $\mathsf{f}\colon W \to \mathbb{R}$, where $W \subset \mathbb{R}^n$ is nonempty, open and convex. If a function $\mathsf{f}$ is pseudoconvex at $\overline{\mathbf{x}} \in W$, then $\overline{\mathbf{x}}$ is a global minimum of $\mathsf{f}$ if and only if $\mathrm{Df}(\overline{\mathbf{x}}) = \mathbf{0}^\mathsf{T}$.*

<u>Proof.</u> Identical as the proof of Corollary 3. $\square$

<u>Lemma 10</u> *Let $W \subset \mathbb{R}^n$ be convex and let $\mathsf{f}\colon W \to \mathbb{R}$ be pseudoconvex. Then, $\overline{\mathbf{x}}$ is a solution of the minimization problem if and only if $\mathrm{Df}(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) \geqslant 0$ for all $\mathbf{x} \in W$.*

<u>Proof.</u> Identical as the proof of Theorem 33. $\square$

## Finding maxima of convex functions

<u>Definition 15</u> *An <u>extremal point</u> of a convex set $W \subset \mathbb{R}^n$ is such a point $\overline{\mathbf{x}} \in W$, which is not an internal point of any line segment contained in $W$, i.e., if $\overline{\mathbf{x}} = \lambda\mathbf{x}_1 + (1 - \lambda)\mathbf{x}_2$, where $\lambda \in (0, 1)$ and $\mathbf{x}_1, \mathbf{x}_2 \in W$, then $\mathbf{x}_1 = \mathbf{x}_1 = \overline{\mathbf{x}}$.*

<u>Definition 16</u> *The <u>convex hull</u> of points $A = \{\mathbf{x}_i \colon i \in I\}$ is the set of points being convex combinations of all finite subsets of the set $A$.*

Equivalently, the convex hull may be defined as the smallest convex set containing the set $A$.

Definition 17 *The <u>dimension</u> of a convex set $U \in \mathbb{R}^n$ is the dimension of the smallest affine subspace $A \in \mathbb{R}^n$ containing $U$ (<u>affine hull</u> of $U$), i.e., the set*

$$\text{aff } U = \Big\{ \sum_{i=1}^{k} \lambda_i x_i : x_1, \dots, x_k \in U \Big\}.$$

<u>Remark.</u> Any affine subspace of $\mathbb{R}^n$, subject to a translation, may become a linear subspace.

<u>Remark.</u> A set $U$ whose dimension is $m$ may be seen as a subset of $\mathbb{R}^m$.

<u>Remark.</u> A convex set $U \subset \mathbb{R}^n$ has a nonempty interior if and only if its dimension is $n$.

<u>Theorem 35</u> *(<u>Krein–Milman, finite-dimensional case</u>) Let $U \subset \mathbb{R}^n$ be convex and compact. It is then the convex hull of the set of its extremal points.*

<u>Lemma 11</u> *Let $U \subset \mathbb{R}^n$ be a convex set with a nonempty interior and let $\overline{x} \in \partial U$. The point $\overline{x}$ is an element of a hyperplane such that $U$ is contained in one of two halfspaces separated by this hyperplane, which we call a <u>supporting hyperplane</u>.*

<u>Proof.</u> By the weak separation theorem applied to int $U$ and $V = \{\overline{x}\}$ there exists $a \in \mathbb{R}^n \setminus \{0\}$ such that $a^\top x \leqslant a^\top \overline{x}$ for all $x \in U$. The hyperplane sought is

$$H = \{x \in \mathbb{R}^n : a^\top x = a^\top \overline{x}\}.$$

The set $U$ is contained in the halfspace $\{x \in \mathbb{R}^n : a^\top x \leqslant a^\top \overline{x}\}$. □

<u>First proof of the Krein–Milman theorem.</u> We use induction with respect to the dimension $m$ of the compact and convex set $U$. The cases $m = 0$ ($U$ consists of a single point) and $m = 1$ ($U$ is a line segment) are obvious. The induction step is the following: assume that each convex and compact set of dimension not greater than $m$ is the convex hull of the set of its extremal points. Let $U$ be a convex and

compact set of dimension $m + 1$. We look at $U$ as a subset of $\mathbb{R}^{m+1}$. Its interior is then nonempty. Let $\overline{x} \in U$.

First, let $\overline{x} \in \partial U$. By Lemma 11, there exists a supporting hyperplane $H$ of the set $U$. The set $U_{\overline{x}} = U \cap H$ is convex and compact and its dimension is at most $m$. By the inductive assumption, $\overline{x}$ is a convex combination of the extremal points of $U_{\overline{x}}$. It has to be shown that these points are also extremal points of $U$. But this is a consequence of the fact that no extremal point of $U_{\overline{x}}$ is in the interior of a line segment whose both end points are in $U$.

Let $\overline{x} \in \text{int } U$. We can take an arbitrary line passing through $\overline{x}$; its intersection with $U$ is a line segment whose end points $x_1$, $x_2$ are located on the boundary of $U$. Both these points are convex combinations of extremal points of $U$, and so is the point $\overline{x}$. □

<u>Second proof of the Krein–Milman theorem.</u> If $U$ is a subset of $\mathbb{R}^1$, then the claim is trivial. Assume that any convex and compact subset of $\mathbb{R}^m$ is the convex hull of the set of its extremal points. We shall prove the theorem for the subsets of $\mathbb{R}^{m+1}$. Let $W$ be the convex hull of the set of extremal points of $U$. Obviously, $W \subset U$. Suppose that there exists $\overline{x} \in U \setminus W$. Then, there exists a ball $B(\overline{x}, \varepsilon)$ disjoint with $W$. By the strong separation theorem, there exists a nonzero vector $a \in \mathbb{R}^{m+1}$ and a number $\alpha$ such that $a^\top x \leqslant \alpha$ for $x \in W$ and $a^\top \overline{x} > \alpha$. Let $\beta = \sup_{x \in U} a^\top x$. As $U$ is compact, $\beta$ is finite. The hyperplane $P = \{x \in \mathbb{R}^{m+1} : a^\top x = \beta\}$ does not intersect $W$, but it has a common point with $U$; indeed, $P_U \overset{\text{def}}{=} P \cap U$ is nonempty, as $U$ is compact; hence, the supremum $\beta$ of $a^\top x$ is taken at some point $x \in U$. We are going to show that in the set $P_U$ there is an extremal point of $U$, which is inconsistent with the definition of $W$. The set $P_U$ is a nonempty, compact and convex set of dimension $m$. It may be seen as a subset of $\mathbb{R}^m$; by the inductive assumption it is the convex hull of the set of its extremal points. Let $\overline{y}$ be one of extremal points of $P_U$; suppose that $\overline{y}$ is a convex combination of some points of $U$: $\overline{y} = \lambda y_1 + (1 - \lambda) y_2$, $y_1, y_2 \in U$, $\lambda \in (0, 1)$. Then, $\beta = a^\top \overline{y} = \lambda a^\top y_1 + (1 - \lambda) a^\top y_2$. By definition of $\beta$, both terms, $a^\top y_1$ and $a^\top y_2$ must be equal to $\beta$; hence, $y_1, y_2 \in P_U$. But $\overline{y}$ is an extremal point of $P_U$, therefore, $\overline{y} = y_1 = y_2$ and thus $\overline{y}$ is an extremal point of $U$. □

<u>Theorem 36</u> *Let $f : W \to \mathbb{R}$ be a convex and continuous function defined in a convex and compact set $W \subset \mathbb{R}^n$. Then at least one of global solutions of*

*the problem*

$$\begin{cases} f(\mathbf{x}) \to \max, \\ \mathbf{x} \in W. \end{cases}$$

*is an extremal point of the set $W$.*

Proof. A continuous function in a compact set achievess its extrema. Therefore, the problem formulated above has a solution $\overline{\mathbf{x}} \in W$. By Theorem 35, the point $\overline{\mathbf{x}}$ is a convex combination of a finite number of extremal points of the set $W$, i.e.,

$$\overline{\mathbf{x}} = a_1 \mathbf{x}_1 + \cdots + a_m \mathbf{x}_m,$$

where $\mathbf{x}_1, \ldots, \mathbf{x}_m$ are extremal points, $a_1, \ldots, a_m > 0$ and $a_1 + \cdots + a_m = 1$.
By convexity of $f$ we obtain

$$f(\overline{\mathbf{x}}) \leqslant a_1 f(\mathbf{x}_1) + \cdots + a_m f(\mathbf{x}_m).$$

As $\overline{\mathbf{x}}$ is a maximum of $f$ in $W$, there must be $f(\mathbf{x}_1) = \cdots = f(\mathbf{x}_m)$. $\square$

Definition 18  *A set $W \subset \mathbb{R}^n$ is called a underline{polyhedral set} if it is the intersection of a finite number of halfspaces, i.e.,*

$$W = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{p}_i^\mathsf{T} \mathbf{x} \leqslant \alpha_i, \, i = 1, \ldots, m\},$$

*where $\mathbf{p}_i \in \mathbb{R}^n \setminus \{\mathbf{0}\}$, $\alpha_i \in \mathbb{R}$.*

Lemma 12  *A polyhedral set is convex and closed.*

Proof. Obvious. $\square$

# 5. Necessary condition of the first order

We consider the optimization problem

$$\begin{cases} f(\mathbf{x}) \to \min, \\ \mathbf{x} \in W, \end{cases} \tag{*}$$

where $W \subset \mathbb{R}^n$ and $f\colon W \to \mathbb{R}$. Let $\bar{\mathbf{x}}$ be a local solution. We are going to connect the local geometry of the set $W$ at $\bar{\mathbf{x}}$ with the behaviour of the function $f$, i.e., with directions of descent of its values. By the local geometry we understand the set of directions at which we can move without leaving the set $W$.

<u>Definition 19</u> *The <u>cone of tangents</u> $T(\bar{\mathbf{x}})$ to the set $W$ at $\bar{\mathbf{x}} \in \operatorname{cl} W$ is the set of vectors $\mathbf{d} \in \mathbb{R}^n$ such that*

$$\mathbf{d} = \lim_{k \to \infty} \lambda_k(\mathbf{x}_k - \bar{\mathbf{x}})$$

*for some numbers $\lambda_k > 0$ and points $\mathbf{x}_k \in W$ such that $\mathbf{x}_k \to \bar{\mathbf{x}}$.*

According to the definition the vector $\mathbf{d}$ is an element of the cone of tangents $T(\bar{\mathbf{x}})$ if it is the limit of a sequence of vectors determined by a sequence of feasible points $(\mathbf{x}_k)_k$ tending to $\bar{\mathbf{x}}$. It may formally be described as follows

$$T(\bar{\mathbf{x}}) = \Big\{ \mathbf{d} \in \mathbb{R}^n \colon \mathbf{d} = \lambda \lim_{k \to \infty} \frac{\mathbf{x}_k - \bar{\mathbf{x}}}{\|\mathbf{x}_k - \bar{\mathbf{x}}\|} \\ \text{for some } (\mathbf{x}_k)_k \subset W,\, \mathbf{x}_k \to \bar{\mathbf{x}},\, \mathbf{x}_k \neq \bar{\mathbf{x}},\, \lambda \geqslant 0 \Big\},$$

which is to be proved as an exercise, as well as the lemma below.

<u>Lemma 13</u> *I) The set $T(\bar{\mathbf{x}})$ is a cone, i.e., $\lambda \mathbf{d} \in T(\bar{\mathbf{x}})$ for all $\mathbf{d} \in T(\bar{\mathbf{x}})$ and $\lambda \geqslant 0$. In particular, $\mathbf{0} \in T(\bar{\mathbf{x}})$.*

*II) If $\bar{\mathbf{x}} \in \operatorname{int} W$, then $T(\bar{\mathbf{x}}) = \mathbb{R}^n$.*

*III) The cone $T(\bar{\mathbf{x}})$ is closed.*

<u>Definition 20</u> *Let $f\colon \mathbb{X} \to \mathbb{R}$ be differentiable at $\bar{\mathbf{x}} \in \mathbb{X}$. The set of <u>descent directions</u>(or <u>improving directions</u>) of $f$ at $\bar{\mathbf{x}}$ is the set*

$$D(\bar{\mathbf{x}}) = \{ \mathbf{d} \in \mathbb{R}^n \colon Df(\bar{\mathbf{x}})\mathbf{d} < 0 \}.$$

<u>Theorem 37</u> *Let $\bar{\mathbf{x}}$ be a local solution of the problem (\*). If $f$ is differentiable at $\bar{\mathbf{x}}$, then*

$$D(\bar{\mathbf{x}}) \cap T(\bar{\mathbf{x}}) = \emptyset.$$

<u>Proof.</u> Let $\mathbf{d} \in T(\bar{\mathbf{x}})$. Then, $\mathbf{d} = \lim_{k \to \infty} \lambda_k(\mathbf{x}_k - \bar{\mathbf{x}})$ for some sequence of points $(\mathbf{x}_k)_k$ converging to $\bar{\mathbf{x}}$ and a sequence $(\lambda_k)_k \subset (0, \infty)$. As the function $f$ is differentiable at $\bar{\mathbf{x}}$, there is

$$f(\mathbf{x}_k) = f(\bar{\mathbf{x}}) + Df(\bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}}) + o(\|\mathbf{x}_k - \bar{\mathbf{x}}\|).$$

As $\bar{\mathbf{x}}$ is a local solution, $f(\mathbf{x}_k) \geqslant f(\bar{\mathbf{x}})$ for $k$ sufficiently large. Together with the formula above, we obtain the following estimation:

$$0 \leqslant f(\mathbf{x}_k) - f(\bar{\mathbf{x}}) = Df(\bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}}) + o(\|\mathbf{x}_k - \bar{\mathbf{x}}\|)$$

By multiplying the sides by $\lambda_k$, we obtain

$$0 \leqslant Df(\bar{\mathbf{x}})\big(\lambda_k(\mathbf{x}_k - \bar{\mathbf{x}})\big) + \lambda_k o(\|\mathbf{x}_k - \bar{\mathbf{x}}\|).$$

A little trick allows us to deal with the remainder; here we go with $k$ to infinity:

$$0 \leqslant Df(\bar{\mathbf{x}}) \underbrace{\big(\lambda_k(\mathbf{x}_k - \bar{\mathbf{x}})\big)}_{\to \mathbf{d}} + \underbrace{\lambda_k \|\mathbf{x}_k - \bar{\mathbf{x}}\|}_{\to \|\mathbf{d}\|} \underbrace{\frac{o(\|\mathbf{x}_k - \bar{\mathbf{x}}\|)}{\|\mathbf{x}_k - \bar{\mathbf{x}}\|}}_{\to 0}.$$

We proved that $Df(\bar{\mathbf{x}})\mathbf{d} \geqslant 0$, i.e., $\mathbf{d} \notin D(\bar{\mathbf{x}})$. $\square$

<u>Example.</u> Consider the following problem:

$$\begin{cases} x_1^2 + x_2^2 \to \min, \\ x_1 + x_2 \geqslant 1. \end{cases}$$

We denote $f(x_1, x_2) = x_1^2 + x_2^2$, $W = \{ \mathbf{x} \in \mathbb{R}^2 \colon x_1 + x_2 \geqslant 1 \}$.

We investigate the sets $T(\bar{\mathbf{x}})$ and $D(\bar{\mathbf{x}})$ at three points: $(1, 1)$, $(1, 0)$, $(\frac{1}{2}, \frac{1}{2})$.

The point $(1, 1)$ is located in the interior of $W$, i.e., $T(\bar{\mathbf{x}}) = \mathbb{R}^2$. The set of descent directions is

$$D(\bar{\mathbf{x}}) = \{ \mathbf{d} \in \mathbb{R}^2 \colon Df(\bar{\mathbf{x}})\mathbf{d} < 0 \} = \{ \mathbf{d} \in \mathbb{R}^2 \colon [2, 2]\mathbf{d} < 0 \}$$
$$= \{ (d_1, d_2) \in \mathbb{R}^2 \colon d_1 + d_2 < 0 \}.$$

Obviously, the two sets have a nonempty intersection; hence, there is no minimum at $(1, 1)$.

The point $(1,0)$ is located at the boundary of $W$, and we have

$$T(\overline{\mathbf{x}}) = \{\, \mathbf{d} \in \mathbb{R}^2 \colon d_1 + d_2 \geqslant 0 \,\}, \quad D(\overline{\mathbf{x}}) = \{\, \mathbf{d} \in \mathbb{R}^2 \colon d_1 < 0 \,\}.$$

The intersection of the two sets is nonempty and thus there is no minimum at $(1,0)$.

For the point $(\tfrac{1}{2}, \tfrac{1}{2})$ we find the sets

$$T(\overline{\mathbf{x}}) = \{\, \mathbf{d} \in \mathbb{R}^2 \colon d_1 + d_2 \geqslant 0 \,\},$$
$$D(\overline{\mathbf{x}}) = \{\, \mathbf{d} \in \mathbb{R}^2 \colon [1,1]\mathbf{d} < 0 \,\} = \{\, (d_1, d_2) \in \mathbb{R}^2 \colon d_1 + d_2 < 0 \,\}.$$

Their intersection is empty and thus *it is possible* (it is still to be verified) that there is a minimum at $(\tfrac{1}{2}, \tfrac{1}{2})$.

## Inequality constraints

Now we consider the problem given in the following form:

$$\begin{cases} f(\mathbf{x}) \to \min, \\ g_i(\mathbf{x}) \leqslant 0, \quad i = 1, \ldots, m, \\ \mathbf{x} \in \mathbb{X}, \end{cases} \tag{**}$$

where $\mathbb{X} \subset \mathbb{R}^n$ is an open set and $f, g_1, \ldots, g_m \colon \mathbb{X} \to \mathbb{R}$. Here the feasible set is

$$W = \{\, \mathbf{x} \in \mathbb{X} \colon g_1(\mathbf{x}) \leqslant 0, \ldots, g_m(\mathbf{x}) \leqslant 0 \,\}.$$

The functions $g_1, \ldots, g_m$ are called <u>inequality constraints</u> and the problem (**) is called the <u>optimization problem with inequality constraints</u>.

We assume that the functions $g_i$ are continuous; then the motion around a point $\overline{\mathbf{x}}$ is restricted by only those functions equal to $0$ at $\overline{\mathbf{x}}$. The others, due to their continuity, are less than zero in a neighbourhood of $\overline{\mathbf{x}}$.

<u>Definition 21</u> *The set of <u>active (or binding or tight) constraints at a point $\overline{\mathbf{x}}$</u> is the set*

$$I(\overline{\mathbf{x}}) = \{\, i \in \{1, \ldots, m\} \colon g_i(\overline{\mathbf{x}}) = 0 \,\}.$$

We are going to connect the properties of active constraints at a point $\overline{\mathbf{x}} \in W$ with the local geometry of the set $W$ around $\overline{\mathbf{x}}$. To do this we introduce the following definition:

<u>Definition 22</u> *Let $\overline{\mathbf{x}} \in W$ and let the functions $g_i$ which describe the constraints active at $\overline{\mathbf{x}}$ be differentiable at $\overline{\mathbf{x}}$. The <u>cone of tangents for active (binding) constraints</u> is the set*

$$T_{\text{lin}}(\overline{\mathbf{x}}) = \{\, \mathbf{d} \in \mathbb{R}^n \colon Dg_i(\overline{\mathbf{x}})\mathbf{d} \leqslant 0 \text{ for all } i \in I(\overline{\mathbf{x}}) \,\}.$$

The cone of tangents for active constraints is a polyhedral set, i.e., it is convex and closed.

<u>Lemma 14</u> *If $\overline{\mathbf{x}} \in W$, then $T(\overline{\mathbf{x}}) \subset T_{\text{lin}}(\overline{\mathbf{x}})$.*

<u>Proof.</u> The proof is similar to that of Theorem 37. Let $\mathbf{d} \in T(\overline{\mathbf{x}})$. Then, $\mathbf{d} = \lim_{k \to \infty} \lambda_k(\mathbf{x}_k - \overline{\mathbf{x}})$ for some sequence of points $(\mathbf{x}_k)_k \subset W$ converging to $\overline{\mathbf{x}}$ and a sequence of positive numbers $\lambda_k$. Let $i \in I(\overline{\mathbf{x}})$. As $g_i$ is differentiable at $\overline{\mathbf{x}}$, we have

$$g_i(\mathbf{x}_k) = g_i(\overline{\mathbf{x}}) + Dg_i(\overline{\mathbf{x}})(\mathbf{x}_k - \overline{\mathbf{x}}) + o(\|\mathbf{x}_k - \overline{\mathbf{x}}\|).$$

The $i$-th constraint is active at $\overline{\mathbf{x}}$; hence, $g_i(\overline{\mathbf{x}}) = 0$. Obviously, $g_i(\mathbf{x}_k) \leqslant 0$ due to $\mathbf{x}_k \in W$. From the above we have the estimation

$$0 \geqslant g_i(\mathbf{x}_k) - g_i(\overline{\mathbf{x}}) = Dg_i(\overline{\mathbf{x}})(\mathbf{x}_k - \overline{\mathbf{x}}) + o(\|\mathbf{x}_k - \overline{\mathbf{x}}\|).$$

We multiply the sides by $\lambda_k$ to obtain

$$0 \geqslant Dg_i(\overline{\mathbf{x}})\big(\lambda_k(\mathbf{x}_k - \overline{\mathbf{x}})\big) + \lambda_k o(\|\mathbf{x}_k - \overline{\mathbf{x}}\|).$$

The same trick as before gives us

$$0 \geqslant Dg_i(\overline{\mathbf{x}}) \underbrace{\big(\lambda_k(\mathbf{x}_k - \overline{\mathbf{x}})\big)}_{\to \mathbf{d}} + \underbrace{\lambda_k \|\mathbf{x}_k - \overline{\mathbf{x}}\|}_{\to \|\mathbf{d}\|} \underbrace{\frac{o(\|\mathbf{x}_k - \overline{\mathbf{x}}\|)}{\|\mathbf{x}_k - \overline{\mathbf{x}}\|}}_{\to 0}.$$

Here we proved that $Dg_i(\overline{\mathbf{x}})\mathbf{d} \leqslant 0$. This holds for all $i \in I(\overline{\mathbf{x}})$; hence, $\mathbf{d} \in T_{\text{lin}}(\overline{\mathbf{x}})$. $\square$

<u>Example.</u> Let $W = \{\, (x_1, x_2) \in \mathbb{R}^2 \colon x_1^2 + x_2^2 \leqslant 1, \, x_2 \geqslant 0 \,\}$. In the canonical form we have

$$W = \{\, (x_1, x_2) \in \mathbb{X} = \mathbb{R}^2 \colon x_1^2 + x_2^2 - 1 \leqslant 0, -x_2 \leqslant 0 \,\},$$

and $g_1(x_1, x_2) = x_1^2 + x_2^2 - 1$, $g_2 = -x_2$. We look at three points of the set $W$: $(\tfrac{1}{2}, \tfrac{1}{2})$, $(0,1)$, $(1,0)$.

No constraint is active at the point $(\frac{1}{2}, \frac{1}{2})$; there is $I\big((\frac{1}{2}, \frac{1}{2})\big) = \emptyset$ and $T\big((\frac{1}{2}, \frac{1}{2})\big) = T_{\text{lin}}\big((\frac{1}{2}, \frac{1}{2})\big) = \mathbb{R}^2$.

At the point $(0,1)$ we have $I\big((0,1)\big) = \{1\}$, $T\big((0,1)\big) = \{(d_1, d_2) \in \mathbb{R}^2 \colon d_2 \leqslant 0\}$ and

$$T_{\text{lin}}\big((0,1)\big) = \{\, \mathbf{d} \in \mathbb{R}^2 \colon Dg_1(0,1)\mathbf{d} \leqslant 0 \,\} = \{\, \mathbf{d} \in \mathbb{R}^2 \colon [0,2]\mathbf{d} \leqslant 0 \,\} = T\big((0,1)\big).$$

At the point $(1,0)$ there is $I\big((1,0)\big) = \{1,2\}$, $T\big((1,0)\big) = \{\, (d_1, d_2) \in \mathbb{R}^2 \colon d_1 \leqslant 0,\ d_2 \geqslant 0 \,\}$ and

$$\begin{aligned}
T_{\text{lin}}\big((1,0)\big) &= \{\, \mathbf{d} \in \mathbb{R}^2 \colon Dg_1(1,0)\mathbf{d} \leqslant 0,\ Dg_2(1,0)\mathbf{d} \leqslant 0 \,\} \\
&= \{\, \mathbf{d} \in \mathbb{R}^2 \colon [2,0]\mathbf{d} \leqslant 0,\ [0,-1]\mathbf{d} \leqslant 0 \,\} = T\big((1,0)\big).
\end{aligned}$$

Example. The same set $W$ may have another description: $W = \{\, (x_1, x_2) \in \mathbb{R}^2 \colon x_1^2 + x_2^2 - 1 \leqslant 0,\ -x_2^3 \leqslant 0 \,\}$. The second constraint is now described by the function $g_2(x_1, x_2) = -x_2^3$. The cone of tangent directions at $(1,0)$ is unchanged, but

$$\begin{aligned}
T_{\text{lin}}\big((1,0)\big) &= \{\, \mathbf{d} \in \mathbb{R}^2 \colon Dg_1(1,0)\mathbf{d} \leqslant 0,\ Dg_2(1,0)\mathbf{d} \leqslant 0 \,\} \\
&= \{\, \mathbf{d} \in \mathbb{R}^2 \colon [2,0]\mathbf{d} \leqslant 0,\ [0,0](1,0)\mathbf{d} \leqslant 0 \,\} \\
&= \{\, (d_1, d_2) \in \mathbb{R}^2 \colon d_1 \leqslant 0 \,\}
\end{aligned}$$

and in this case $T(\overline{\mathbf{x}}) \neq T_{\text{lin}}(\overline{\mathbf{x}})$.

## Necessary Kuhn–Tucker conditions

By Lemma 14, $T(\overline{\mathbf{x}}) \subset T_{\text{lin}}(\overline{\mathbf{x}})$. Often, but not always, the two sets are equal. It is an important property, being the starting point to the entire theory of nonlinear optimization by Kuhn and Tucker. We begin studying it with the lemma:

Lemma 15 (*G. Farkas, 1901*) *Let* $A$ *be an* $m \times n$ *real matrix and let* $\mathbf{d} \in \mathbb{R}^n$. *Then exactly one of the two following systems of equations and inequalities has a solution:*

$$(1)\ \begin{cases} A\mathbf{x} \leqslant 0, \\ \mathbf{d}^\mathsf{T}\mathbf{x} > 0, \\ \mathbf{x} \in \mathbb{R}^n \end{cases} \qquad (2)\ \begin{cases} A^\mathsf{T}\mathbf{y} = \mathbf{d}, \\ \mathbf{y} \geqslant 0, \\ \mathbf{y} \in \mathbb{R}^m. \end{cases}$$

Proof. First we show that if (2) has a solution, then (1) does not. Let $\mathbf{y}$ satisfy (2). Then, $\mathbf{d} = A^\mathsf{T}\mathbf{y}$. After substituting this to (1) we obtain

$$\begin{cases} A\mathbf{x} \leqslant 0, \\ \mathbf{y}^\mathsf{T}A\mathbf{x} > 0. \end{cases}$$

The first inequality means that each coordinate of the vector $A\mathbf{x}$ is nonpositive. As the coordinates of $\mathbf{y}$ are nonnegative, the scalar product $\mathbf{y}^\mathsf{T}A\mathbf{x}$ is not positive. This is inconsistent with the second inequality above; hence, (1) does not have a solution.

Now suppose that (2) does not have a solution. Let

$$V \overset{\text{def}}{=} \{\, \mathbf{x} \in \mathbb{R}^n \colon \mathbf{x} = A^\mathsf{T}\mathbf{y},\ \mathbf{y} \in \mathbb{R}^m,\ \mathbf{y} \geqslant 0 \,\}.$$

The set $V$ is a polyhedral set; hence, it is convex and closed. As (2) does not have a solution, $\mathbf{d} \notin V$. By the strong separation theorem applied to the sets $V$ and $U = \{\mathbf{d}\}$, there exists a vector $\mathbf{a} \in \mathbb{R}^n$ such that

$$\mathbf{a}^\mathsf{T}\mathbf{d} > \sup_{\mathbf{x} \in V} \mathbf{a}^\mathsf{T}\mathbf{x}.$$

Below we show that $\overline{\mathbf{x}} = \mathbf{a}$ is a solution of (1). Let $\alpha = \sup_{\mathbf{x} \in V} \overline{\mathbf{x}}^\mathsf{T}\mathbf{x}$. From $0 \in V$ it follows that $\alpha \geqslant 0$; hence, $\mathbf{a}^\mathsf{T}\mathbf{d} = \mathbf{d}^\mathsf{T}\overline{\mathbf{x}} > 0$. It remains to be proved that $A\overline{\mathbf{x}} \leqslant 0$. Suppose that the $i$-th coordinate of $A\overline{\mathbf{x}}$ is positive. By definition of $V$, for any $\mathbf{y} \geqslant 0$ there is $\alpha \geqslant \overline{\mathbf{x}}^\mathsf{T}A^\mathsf{T}\mathbf{y} = \mathbf{y}^\mathsf{T}A\overline{\mathbf{x}}$. Let $\mathbf{y}_k = k\mathbf{e}_i$ (i.e., the $i$-th coordinate of the vector $\mathbf{y}_k$ is $k$ and all other coordinates are zero). If the $i$-th coordinate $(A\overline{\mathbf{x}})_i$ of $A\overline{\mathbf{x}}$ is positive, then,

$$\lim_{k \to \infty} \mathbf{y}_k^\mathsf{T}A\overline{\mathbf{x}} = \lim_{k \to \infty} k(A\overline{\mathbf{x}})_i = \infty,$$

which is inconsistent with $\mathbf{y}^\mathsf{T}A\overline{\mathbf{x}} \leqslant \alpha$ for all $\mathbf{y} \geqslant 0$. Hence, $A\overline{\mathbf{x}} \leqslant 0$. $\square$

Remark. If a vector $\mathbf{x}$ is a solution of System (1), then for all $a > 0$ the vector $a\mathbf{x}$ is also a solution; hence, the set of solutions of (1) is either empty or infinite. On the other hand, if the rows of the matrix $A$ are linearly independent and System (2) has a solution $\mathbf{y}$, then this solution is unique.

An example is shown in Figure 5. The matrix $A$ is $2 \times 2$, its rows are the vectors $\mathbf{a}_1$ and $\mathbf{a}_2$ (they are row matrices $1 \times 2$). For System (1) from the lemma, the condition $A\mathbf{x} \leqslant 0$ is equivalent to two inequalities, $\mathbf{a}_1\mathbf{x} \leqslant 0$ and $\mathbf{a}_2\mathbf{x} \leqslant 0$. The sets of vectors $\mathbf{x}$ satisfying each of the two inequalities are hatched with dashed lines; their intersection is the set of vectors $\mathbf{x}$ such that $A\mathbf{x} \leqslant 0$. The set of vectors $\mathbf{x}$
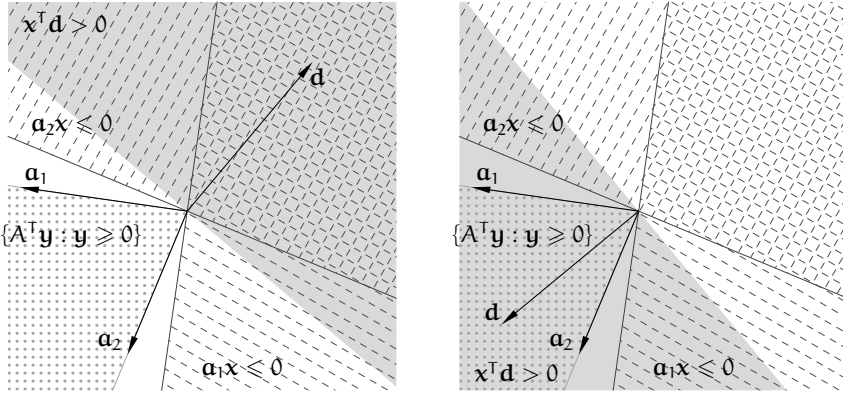
Figure 5: Farkas' lemma: on the left side (1) has a solution and on the right side (2) has a solution

such that $\mathbf{d}^\mathsf{T}\mathbf{x} > 0$, where $\mathbf{d}$ is given, is grey. Clearly, (1) has solutions if and only if the grey area and the doubly hatched area intersect. This is the case in the picture on the left side. On the right side, System (2) has no solution. To find a geometric interpretation of (2), we need to see how the set $V = \{A^\mathsf{T}\mathbf{y} : \mathbf{y} \geqslant \mathbf{0}\}$ looks like. Any vector $A^\mathsf{T}\mathbf{y}$ may be represented as a linear combination $y_1 \mathbf{a}_1^\mathsf{T} + y_2 \mathbf{a}_2^\mathsf{T}$, where $\mathbf{y} = (y_1, y_1)$. One can easily notice that the set $V$ is a cone spanned by the vectors $\mathbf{a}_1^\mathsf{T}$ and $\mathbf{a}_2^\mathsf{T}$; it is dotted on the picture. Equation $A^\mathsf{T}\mathbf{y} = \mathbf{d}$ has a solution when the vector $\mathbf{d}$ is inside the dotted part of the cone. This is the case in the picture on the right side. On the left side the vector $\mathbf{d}$ is outside of the dotted area, and System (2) does not have a solution.

Theorem 38 *(Kuhn–Tucker theorem) Let $\overline{x}$ be a local solution of the problem (**). If the functions $f$ and $g_i$, where $i \in I(\overline{x})$, are differentiable at $\overline{x}$ and $T(\overline{x}) = T_{\lin}(\overline{x})$, then there exists $\mu = (\mu_1, \dots, \mu_m) \in [0, \infty)^m$ such that*

$$\begin{cases} Df(\overline{x}) + \sum_{i \in I(\overline{x})} \mu_i Dg_i(\overline{x}) = \mathbf{0}^\mathsf{T}, \\ \mu_i g_i(\overline{x}) = 0, \quad i = 1, 2, \dots, m. \end{cases}$$

Remark. The second condition is called the complementary slackness condition. *Whenever the constraint $g_i \leqslant 0$ is slack (meaning $g_i < 0$), the constraint $\mu_i \geqslant 0$ must not be slack (meaning $\mu_i = 0$) and reverse.*

Remark. Often the system of Kuhn–Tucker equations is written with the sum over all indices $i = 1, \dots, m$:

$$\begin{cases} Df(\overline{x}) + \sum_{i=1}^{m} \mu_i Dg_i(\overline{x}) = \mathbf{0}^\mathsf{T}, \\ \mu_i g_i(\overline{x}) = 0, \quad i = 1, 2, \dots, m. \end{cases}$$

This notation is an abuse, as the functions $g_i$ which describe inactive constraints need not be differentiable at $\overline{x}$. On the other hand, the derivatives are multiplied by $\mu_i$, equal to $0$ for inactive constraints, which is a sort of justification for this notation.

Remark. The assumptions of the Kuhn–Tucker theorem are obviously satisfied if $\overline{x} \in \text{int } W$, i.e., $I(\overline{x})$ is the empty set. Then we have $Df(\overline{x}) = \mathbf{0}^\mathsf{T}$, $\mu = \mathbf{0}$.

Proof of the Kuhn–Tucker theorem. By Theorem 37 we have $D(\overline{x}) \cap T(\overline{x}) = \emptyset$. By assumption, $D(\overline{x}) \cap T_{\lin}(\overline{x}) = \emptyset$, which means that the system

$$\begin{cases} Df(\overline{x})z < 0, \\ Dg_i(\overline{x})z \leqslant 0, \quad i \in I(\overline{x}), \end{cases}$$

has no solution $z \in \mathbb{R}^n$. Let $\mathbf{d} = -\big(Df(\overline{x})\big)^\mathsf{T}$, let $k = |I(\overline{x})|$ and let $A$ be the $k \times n$ matrix whose rows are gradients of the active constraints $Dg_i(\overline{x})$. By the Farkas' lemma, there exists a vector $\mathbf{y} \in [0, \infty)^k$ such that $\mathbf{y}^\mathsf{T}A = \mathbf{d}$, i.e.,

$$Df(\overline{x}) + \mathbf{y}^\mathsf{T}A = \mathbf{0}^\mathsf{T}.$$

Let $\mu \in [0, \infty)^m$ be defined as follows: $(\mu_i)_{i \in I(\overline{x})} = \mathbf{y}$ and $(\mu_i)_{i \notin I(\overline{x})} = \mathbf{0}$. Then, the equality above is equivalent to

$$Df(\overline{x}) + \sum_{i \in I(\overline{x})} \mu_i Dg_i(\overline{x}) = \mathbf{0}^\mathsf{T}.$$

By definition of $\mu$ it is obvious that $\mu_i g_i(\overline{x}) = 0$ for all $i = 1, \dots, m$. $\square$

The assumptions of the Kuhn–Tucker theorem are called the necessary conditions of the first order. Due to the importance of the vector $\mu$ in what follows, we give it a name:

Definition 23 *The vector $\mu$ which appears in the necessary first-order condition is called the vector of Lagrange multipliers.*

# 6. Regularity conditions and examples

## Constraints qualifications

Three conditions defined below are sufficient for the equality $T(\overline{x}) = T_{lin}(\overline{x})$; we shall prove it.

Definition 24 *At a point $\overline{x} \in W \subset \mathbb{X}$*

- *the <u>linear independence condition</u> is satisfied if the functions $g_i$ are continuous at $\overline{x}$ for $i \notin I(\overline{x})$ and the vectors $Dg_i(\overline{x})$ are linearly independent for $i \in I(\overline{x})$,*

- *the <u>affine function condition</u> is satisfied if the functions $g_i$ are continuous at $\overline{x}$ for $i \notin I(\overline{x})$ and the functions $g_i$ are affine for $i \in I(\overline{x})$,*

- *the <u>Slater condition</u> is satisfied if the functions $g_i$ are continuous at $\overline{x}$ for $i \notin I(\overline{x})$, the functions $g_i$ are pseudoconvex at $\overline{x}$ for $i \in I(\overline{x})$ (i.e., $Df(\overline{x})(y - \overline{x}) \geqslant 0 \Rightarrow f(y) \geqslant f(\overline{x})$) and there exists a point $x^* \in \mathbb{X}$ such that $g_i(x^*) < 0$ for $i \in I(\overline{x})$.*

Note that the point $x$ considered in the Slater condition needs not satisfy the inactive constraints, i.e., it is not required that $x \in W$.

Theorem 39 *If the affine function condition is satisfied at $\overline{x} \in W$, then $T(\overline{x}) = T_{lin}(\overline{x})$.*

Proof. The set inclusion $T(\overline{x}) \subset T_{lin}(\overline{x})$ follows from Lemma 14. We need to prove the opposite inclusion. Let $d \in T_{lin}(\overline{x})$. We are going to prove the existence of $\lambda^* > 0$ such that the entire line segment $\overline{x} + \lambda d$, $\lambda \in [0, \lambda^*]$, is a subset of $W$.

We can notice that if $i \notin I(\overline{x})$, then $g_i(\overline{x}) < 0$. As these functions $g_i$ are continuous, there exists $\lambda^* > 0$ such that $g_i(\overline{x} + \lambda d) \leqslant 0$ for all $\lambda \in [0, \lambda^*]$. It remains to be proved that this inequality is satisfied also for the active constraints. Let $i \in I(\overline{x})$. By definition of $d$, $Dg_i(x)d \leqslant 0$. The active constraint $g_i$ is an affine function equal to $0$ at $\overline{x}$, i.e., it has the form $g_i(x) = a_i^T(x - \overline{x})$ for some nonzero vector $a_i \in \mathbb{R}^n$. There is $Dg_i(\overline{x})d = a_i^T d$. Therefore, for any $\lambda \geqslant 0$ we have

$$0 \geqslant \lambda a_i^T d = g_i(\overline{x} + \lambda d);$$

hence,

$$\{\overline{x} + \lambda d : \lambda \in [0, \lambda^*]\} \subset W.$$

It remains to construct appropriate sequences $(x_k)_k \subset W$, and $(\lambda_k)_k \subset (0, \infty)$. Let

$$x_k = \overline{x} + \frac{\lambda^*}{k}d, \quad \lambda_k = \frac{k}{\lambda^*}.$$

Then, $x_k \in W$, $x_k \to \overline{x}$ and $\lambda_k(x_k - \overline{x}) = d$ for all k, and thus $\lim_{k \to \infty} \lambda_k(x_k - \overline{x}) = d$. □

The proofs of the other two regularity conditions refer to the set

$$T_{int}(\overline{x}) \stackrel{\text{def}}{=} \{ d \in \mathbb{R}^n : Dg_i(\overline{x})d < 0 \text{ for all } i \in I(\overline{x}) \}.$$

Note that $T_{int}(\overline{x}) = \emptyset$ if $g_i(\overline{x}) = 0$ and $Dg_i(\overline{x}) = 0^T$ for some $i \in \{1, \ldots, m\}$.

Lemma 16 *If the functions $g_i$, $i \in I(\overline{x})$ are differentiable at $\overline{x} \in W$ and the other constraints $g_i$ are continuous, then from $d \in T_{int}(\overline{x})$ it follows that $\overline{x} + \lambda d \in \text{int } W$ for sufficiently small $\lambda > 0$.*

Proof. If $i \notin I(\overline{x})$, then $g_i(\overline{x}) < 0$. By the continuity of $g_i$, there is $g_i(\overline{x} + \lambda d) < 0$ for all $\lambda$ sufficiently small, say, $\lambda \in (0, \lambda^*]$. If $i \in I(\overline{x})$, then, by differentiability of $g_i$ at $\overline{x}$, we have

$$\lim_{\lambda \searrow 0} \frac{g_i(\overline{x} + \lambda d) - g_i(\overline{x})}{\lambda} = Dg_i(\overline{x})d < 0,$$

because $d \in T_{int}(\overline{x})$. This inequality holds also if $\lambda > 0$ is arbirarily small. Therefore,

$$\frac{g_i(\overline{x} + \lambda d) - g_i(\overline{x})}{\lambda} < 0,$$

and $g_i(\overline{x} + \lambda d) - g_i(\overline{x}) = g_i(\overline{x} + \lambda d) < 0$. □

Lemma 17 *Let $\overline{x} \in W$, let the functions $g_i$, where $i \in I(\overline{x})$, be differentiable at $\overline{x}$ and let the functions $g_i$, where $i \notin I(\overline{x})$ be continuous at $\overline{x}$. Then,*

*I) $T_{int} \subset T(\overline{x})$,*

*II) If $T_{int}(\overline{x}) \neq \emptyset$, then $\text{cl}\left(T_{int}(\overline{x})\right) = T_{lin}(\overline{x})$.*

Proof. (I) follows directly from Lemma 16. To prove (II) we notice that $T_{int}(\overline{x})$ is the interior of the set $T_{lin}(\overline{x})$. As $T_{lin}(\overline{x})$ is a polyhedral set, it is convex and closed. The claim follows from Lemma 3. $\square$

Theorem 40 *If the Slater condition is satisfied at $\overline{x} \in W$, then $T(\overline{x}) = T_{lin}(\overline{x})$.*

Proof. First we show that $T_{int}(\overline{x}) \neq \emptyset$. Let $x^* \in \mathbb{X}$ such that $g_i(x^*) < 0$ for $i \in I(\overline{x})$. By pseudoconvexity of $g_i$ at $\overline{x}$ we have $Dg_i(\overline{x})(x^* - \overline{x}) < 0$ for $i \in I(\overline{x})$, i.e., $x^* - \overline{x} \in T_{int}(\overline{x})$.

By Lemma 17, $\mathrm{cl}\left(T_{int}(\overline{x})\right) = T_{lin}(\overline{x})$. We have also proved that $T_{int}(\overline{x}) \subset T(\overline{x}) \subset T_{lin}(\overline{x})$ and that the set $T(\overline{x})$ is closed. Hence, $T(\overline{x}) = T_{lin}(\overline{x})$. $\square$

Lemma 18 *(P. Gordan, 1873) Let $A$ be an $m \times n$ matrix. Then exactly one of the following two systems,*

$$(1) \begin{cases} Ax < 0, \\ x \in \mathbb{R}^n, \end{cases} \qquad (2) \begin{cases} A^{\mathsf{T}}y = 0, \\ y \geqslant 0, \ y \neq 0, \\ y \in \mathbb{R}^m, \end{cases}$$

*has a solution.*

Proof. First we prove that the two systems cannot have solutions at the same time. Suppose that they have; let $x$ and $y$ satisfy (1) and (2), respectively. As $y$ satisfies (2), we have $y^{\mathsf{T}}Ax = 0$. On the other hand, $x$ satisfies (1), i.e., $(Ax)_i < 0$ for all $i = 1, \ldots, m$. Knowing that $y_i \geqslant 0$ for all $i$ and $y \neq 0$, we have $y^{\mathsf{T}}Ax < 0$. This inconsistency proves that at least one of the systems does not have a solution.

Now we prove that at least one of the systems has a solution, by proving that if (1) has no solution, then (2) has at least one. We define the convex sets

$$U \stackrel{\text{def}}{=} (-\infty, 0)^m, \quad V \stackrel{\text{def}}{=} \{z \in \mathbb{R}^m \colon z = Ax \text{ for some } x \in \mathbb{R}^n\}.$$

The set $U$ is the interior of a cone and $V$ is a linear subspace. If (1) has no solution, then the two sets are disjoint. By the weak separation theorem, there exists a nonzero vector $a \in \mathbb{R}^m$ such that

$$\sup_{z \in U} a^{\mathsf{T}} z \leqslant \inf_{z \in V} a^{\mathsf{T}} z.$$

We show that $a \geqslant 0$. On the contrary, suppose that one of the coordinates of $a$, say, $a_i < 0$. Consider the sequence of vectors $z_k = -\left(ke_i + \frac{1}{k}\sum_{j \neq i} e_j\right)$. Then, $z_k \in U$ and $\lim_{k \to \infty} a^{\mathsf{T}} z_k = \infty$, which is an inconsistency, as $\sup_{z \in U} a^{\mathsf{T}} z$ is finite (we have $0 \in V$; hence, $\inf_{z \in V} a^{\mathsf{T}} z \leqslant 0$).

We proved that all coordinates of $a$ are nonnegative. Hence, $\sup_{z \in U} a^{\mathsf{T}} z = 0$. Let $z = A(-A^{\mathsf{T}}a)$. Then, $z \in V$, i.e., $a^{\mathsf{T}} z \geqslant 0$. Therefore,

$$0 \leqslant a^{\mathsf{T}} z = -a^{\mathsf{T}} A A^{\mathsf{T}} a = -\|A^{\mathsf{T}} a\|^2,$$

which implies that $\|A^{\mathsf{T}} a\| = 0$, i.e., $A^{\mathsf{T}} a = 0$. The vector $y = a$ is, therefore, a solution of (2). $\square$
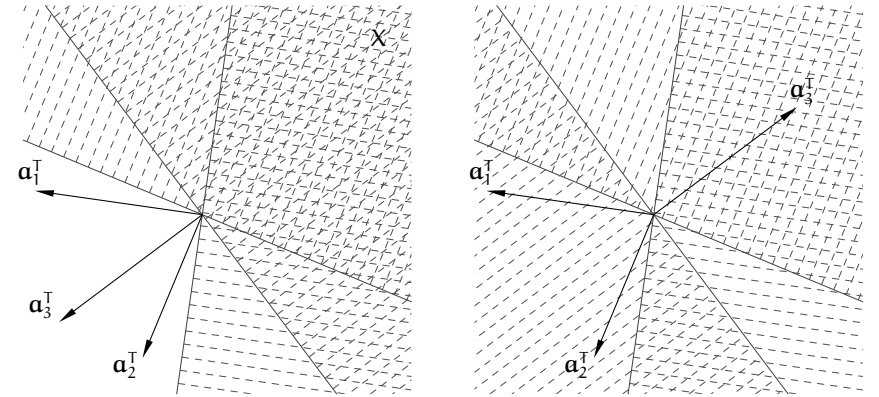


Figure 6: Gordan's lemma

Figure 6 illustrates the Gordan's lemma. The matrix $A$ is $3 \times 2$, its transposed rows are the vectors in $\mathbb{R}^2$. Obviously, in this example they are linearly dependent. The sets of solutions of the inequalities $a_i x < 0$ are the hatched halfplanes (without the boundaries).

On the left side there exist vectors $v \in \mathbb{R}^2$ (e.g., $a_3^{\mathsf{T}}$) such that for all $i$ the number $a_i v$ is positive; hence, the coefficients of the linear combination $a_1^{\mathsf{T}} y_1 + a_2^{\mathsf{T}} y_2 + a_3^{\mathsf{T}} y_3 = 0$, if not all equal to 0, must have both signs, i.e., at least one of them is positive and at least one is negative. Thus, the system (2) has no solution. The cone $X$ being the intersection of all hatched halfplanes is nonempty; the vectors $x$ satisfying (1) are elements of this cone.

On the right side we have vectors $a_1, a_2, a_3$, such that there exists a linear combination with positive coefficients, which is the zero vector. Therefore, there

exists a vector $\mathbf{y} \in \mathbb{R}^3$ satisfying (2). The set X—the intersection of the appropriate halfplanes—is empty.

<u>Theorem 41</u> *If the linear independence condition is satisfied at a point $\overline{\mathbf{x}} \in W$, then $T(\overline{\mathbf{x}}) = T_{\mathrm{lin}}(\overline{\mathbf{x}})$.*

<u>Proof.</u> Just as in the case of the Slater condition, it suffices to prove that $T_{\mathrm{int}}(\overline{\mathbf{x}}) \neq \emptyset$. Let A be the matrix whose rows are gradients of active constraints. By their linear independence, a nonzero vector $\boldsymbol{\mu} \in \mathbb{R}^{|I(\overline{\mathbf{x}})|}$ such that $A^{\mathsf{T}}\boldsymbol{\mu} = \mathbf{0}$ does not exist. In other words, System (2) of Lemma 18 has no solution. Hence, System (1) has a solution, i.e., there exists $\mathbf{d} \in \mathbb{R}^n$ such that $A\mathbf{d} < \mathbf{0}$, which means that

$$Dg_i(\overline{\mathbf{x}})\mathbf{d} < 0 \quad \text{for all } i \in I(\overline{\mathbf{x}}).$$

Hence, $\mathbf{d} \in T_{\mathrm{int}}(\overline{\mathbf{x}})$. □

## Examples

<u>Example.</u> Consider an optimization problem in the set

$$W = \{\mathbf{x} \in \mathbb{R}^2 \colon x_1^2 + x_2^2 \leqslant 1,\ x_1 + 2x_2 \leqslant 1,\ x_1 - 3x_2 \leqslant 1\}.$$

There is $\mathbb{X} = \mathbb{R}^2$. The constraints are described by the functions

$$g_1(x_1, x_2) = x_1^2 + x_2^2 - 1, \quad g_2(x_1, x_2) = x_1 + 2x_2 - 1, \quad g_3(x_1, x_2) = x_1 - 3x_2 - 1.$$

At the point $\overline{\mathbf{x}} = (1, 0)$ all three constraints are active and

$$Dg_1(\overline{\mathbf{x}}) = [2, 0], \quad Dg_2(\overline{\mathbf{x}}) = [1, 2], \quad Dg_3(\overline{\mathbf{x}}) = [1, -3].$$

The linear independence condition is not satisfied at $\overline{\mathbf{x}}$. Also, not all active constraints are described by affine functions. But all these functions are convex and differentiable, i.e., they are pseudoconvex. At $\mathbf{x} = (0, 0)$ we have $g_i(\mathbf{x}) = -1 < 0$ for $i = 1, 2, 3$. Thus, the Slater condition is satisfied.

<u>Example (Kuhn, Tucker, 1951).</u> Consider the optimization problem

$$\begin{cases} x_1 \to \min, \\ x_2 \leqslant x_1^3, \\ x_2 \geqslant 0. \end{cases}$$

The constraints are described by the functions

$$g_1(x_1, x_2) = -x_1^3 + x_2, \quad g_2(x_1, x_2) = -x_2.$$

At each feasible point except $(0, 0)$ the linear independence condition is satisfied. However, the solution of this problem is $\overline{\mathbf{x}} = (0, 0)$, but

$$T(\overline{\mathbf{x}}) = \{(d_1, d_2) \in \mathbb{R}^2 \colon d_1 \geqslant 0,\ d_2 = 0\} \neq T_{\mathrm{lin}}(\overline{\mathbf{x}}) = \{(d_1, d_2) \in \mathbb{R}^2 \colon d_2 = 0\},$$

as $Dg_1(\overline{\mathbf{x}}) = [0, 1]$, $Dg_2(\overline{\mathbf{x}}) = [0, -1]$, which implies $Dg_1(\overline{\mathbf{x}})\mathbf{d} \leqslant 0 \Rightarrow d_2 \leqslant 0$ and $Dg_2(\overline{\mathbf{x}})\mathbf{d} \leqslant 0 \Rightarrow d_2 \geqslant 0$. On the other hand,

$$D(\overline{\mathbf{x}}) = \{(d_1, d_2) \in \mathbb{R}^2 \colon d_1 < 0\},$$

because $Df(\overline{\mathbf{x}}) = [1, 0]$; hence, $Df(\overline{\mathbf{x}})\mathbf{d} < 0 \Rightarrow d_1 < 0$.

For the optimization problem

$$\begin{cases} x_2 \to \min, \\ x_2 \leqslant x_1^3, \\ x_2 \geqslant 0, \end{cases}$$

we still have $T(\overline{\mathbf{x}}) \neq T_{\mathrm{lin}}(\overline{\mathbf{x}})$. But for this new problem $D(\overline{\mathbf{x}}) = \{(d_1, d_2) \in \mathbb{R}^2 \colon d_2 < 0\}$. It follows that at $\overline{\mathbf{x}} = (0, 0)$

$$T(\overline{\mathbf{x}}) \cap D(\overline{\mathbf{x}}) = T_{\mathrm{lin}}(\overline{\mathbf{x}}) \cap D(\overline{\mathbf{x}}) = \emptyset.$$

The above condition is sufficient for the claim of the Kuhn–Tucker theorem to hold.

<u>Example.</u> Let A be a symmetric $n \times n$ matrix. Consider the optimization problem

$$\begin{cases} \mathbf{x}^{\mathsf{T}} A \mathbf{x} \to \max, \\ \|\mathbf{x}\| \leqslant 1, \\ \mathbf{x} \in \mathbb{R}^n. \end{cases}$$

First, we rewrite the problem in the canonical form; note that the constraint is equivalent to $\|\mathbf{x}\|^2 = \mathbf{x}^{\mathsf{T}}\mathbf{x} \leqslant 1$:

$$\begin{cases} -\mathbf{x}^{\mathsf{T}} A \mathbf{x} \to \max, \\ \mathbf{x}^{\mathsf{T}}\mathbf{x} - 1 \leqslant 0, \\ \mathbf{x} \in \mathbb{R}^n. \end{cases}$$

By $W$ we denote the feasible set. At each point the linear independence condition and the Slater condition are satisfied. Therefore, the solution of the problem

satisfies the necessary first order conditions (i.e., the Kuhn–Tucker conditions).
We are looking for all points "under suspicion".

The Kuhn–Tucker conditions have the form

$$\begin{cases} -2x^{\mathsf{T}}A + 2\mu x^{\mathsf{T}} = 0^{\mathsf{T}}, \\ \mu(x^{\mathsf{T}}x - 1) = 0, \\ \mu \geqslant 0, \; x \in \mathbb{R}^n. \end{cases}$$

First we check the case $x^{\mathsf{T}}x - 1 < 0$. Then (due to the second equation), $\mu = 0$ and the first equation is equivalent to $x^{\mathsf{T}}A = 0^{\mathsf{T}}$ or $Ax = 0$. This equation is satisfied by all vectors $x \in \ker A$ such that $\|x\| < 1$. In particular, it has at least one solution, $x = 0$.

Then we check $x^{\mathsf{T}}x - 1 = 0$; from the second equation nothing may be said about $\mu$. However, the first equation takes the form $\mu x^{\mathsf{T}} = x^{\mathsf{T}}A$, or $Ax = \mu x$. The solutions are the unit eigenvectors corresponding to nonnegative eigenvalues $\mu$ of the matrix $A$. The set of solutions may be empty if all the eigenvalues are negative. To conclude, the set of points "under suspicion", i.e., satisfying the necessary first order condition, is

$$\{\, x \in \mathbb{R}^n \colon (x \in \ker A \text{ and } \|x\| < 1) \text{ or}$$
$$(\|x\| = 1 \text{ and } x \text{ is an eigenvector of } A$$
$$\text{associated with a nonnegative eigenvalue}) \,\}.$$

At this point we do not have any technique of finding solutions, except for using the common sense. If $x \in \ker A$, then the target function $x^{\mathsf{T}}Ax$ is equal to $0$. For any eigenvector $x$ the value of the target function is the associated eigenvalue. If the greatest eigenvalue is positive, then each of its unit eigenvectors is a global solution. If the matrix $A$ has no positive eigenvalue, then any vector from $\ker A$ whose norm is not greater than $1$ is a solution.

Example. Consider the problem

$$\begin{cases} x_1 + x_2 \to \min, \\ x_2 \geqslant x_1^2, \\ x_2 \leqslant 0. \end{cases}$$

It is easy to notice that the solution is $\overline{x} = (0, 0)$, which is the only feasible point. On the other hand, $Df = [1, 1]$, $Dg_1 = [2x_1, -1]$ and $Dg_2 = [0, 1]$. Therefore, for any numbers $\mu_1$ and $\mu_2$, there is

$$Df(\overline{x}) + \mu_1 Dg_1(\overline{x}) + \mu_2 Dg_2(\overline{x}) \neq 0,$$

i.e., the Kuhn–Tucker conditions are not satisfied. But
$T(\overline{x}) = \{0\} \neq T_{\mathrm{lin}}(\overline{x}) = \{\, (d_1, d_2) \in \mathbb{R}^2 \colon d_2 = 0 \,\}$. Also,
$\emptyset = T(\overline{x}) \cap D(\overline{x}) \neq T_{\mathrm{lin}}(\overline{x}) \cap D(\overline{x})$, as $D(\overline{x}) = \{\, (d_1, d_2) \in \mathbb{R}^2 \colon d_1 + d_2 < 0 \,\}$. As we can see, the assumptions of the Kuhn–Tucker theorem are not satisfied at $\overline{x}$.

# 7. Quasi-convex functions and sufficient conditions

## Quasi-convexity

Below we extend the class of functions whose maxima are located at the extremal points of their domains.

<u>Definition 25</u> *Let $W \subset \mathbb{R}^n$ be a convex set and let $f \colon W \to \mathbb{R}$. The function $f$ is <u>quasi-convex</u> if for all $\mathbf{x}, \mathbf{y} \in W$ and $\lambda \in [0,1]$ there is*

$$f\bigl(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}\bigr) \leqslant \max\{f(\mathbf{x}), f(\mathbf{y})\}.$$

*The function $f$ is <u>quasi-concave</u> if $-f$ is quasi-convex, i.e., if*

$$f\bigl(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}\bigr) \geqslant \min\{f(\mathbf{x}), f(\mathbf{y})\} \quad \text{for all } \mathbf{x}, \mathbf{y} \in W, \ \lambda \in [0,1].$$

*The function $f$ is quasi-linear if it is both quasi-convex and quasi-concave.*

<u>Theorem 42</u> *Let $W \subset \mathbb{R}^n$ be a convex set and let $f \colon W \to \mathbb{R}$. The function $f$ is quasi-linear if and only if each sublevel set of this function is convex.*

<u>Proof.</u> Suppose that the function $f$ is quasi-convex and let $\alpha \in \mathbb{R}$ be fixed. Let $\mathbf{x}, \mathbf{y} \in W_\alpha(f)$. Then, $f(\mathbf{x}) \leqslant \alpha$ and $f(\mathbf{y}) \leqslant \alpha$. For any $\lambda \in (0,1)$ we obtain

$$f\bigl(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}\bigr) \leqslant \max\{f(\mathbf{x}), f(\mathbf{y})\} \leqslant \alpha.$$

Hence, $\lambda \mathbf{x} + (1-\lambda)\mathbf{y} \in W_\alpha(f)$, i.e., $W_\alpha(f)$ is a convex set.

Now suppose that $W_\alpha(f)$ is convex for all $\alpha \in \mathbb{R}$. We fix $\mathbf{x}, \mathbf{y} \in W$ and $\lambda \in (0,1)$. By assumption, the set $W_\alpha(f)$, where $\alpha = \max\{f(\mathbf{x}), f(\mathbf{y})\}$, is convex. It follows that $\lambda \mathbf{x} + (1-\lambda)\mathbf{y} \in W_\alpha(f)$; hence,

$$f\bigl(\lambda \mathbf{x} + (1-\lambda)\mathbf{y}\bigr) \leqslant \alpha = \max\{f(\mathbf{x}), f(\mathbf{y})\}$$

and the proof is complete. $\square$

<u>Corollary 7</u> *Any convex function is a quasi-convex function.*

<u>Example.</u> The function $f(x) = -e^x$ is quasi-convex even if it is strictly concave. For $\alpha \geqslant 0$ the set $W_\alpha(f) = \mathbb{R}$ and if $\alpha < 0$, then $W_\alpha(f) = [\ln(-\alpha), \infty)$. All these sets are convex; hence, by Theorem 42, the function $f$ is quasi-convex. It is also quasi-concave, i.e., it is quasi-linear.

<u>Example.</u> The function $f(x) = x^2$ is quasi-convex (as it is convex), but it is not quasi-concave. This may be checked by looking at the sublevel sets of $-f$: for $\alpha < 0$ we have $W_\alpha(f) = (-\infty, -\sqrt{-\alpha}] \cup [\sqrt{-\alpha}, \infty)$, which is not a convex set.

<u>Lemma 19</u> *If a set $W \subset \mathbb{R}^n$ is convex, then a function $f \colon W \to \mathbb{R}$ is quasi-linear if and only if its restriction to any interval is a monotone function.*

<u>Proof</u> is left as an exercise.

<u>Example.</u> The function $f(x) = -e^{-x^2}$ has the following sublevel sets:

$$W_\alpha(f) = \begin{cases} \emptyset & \text{if } \alpha \leqslant -1, \\ \left[-\sqrt{-\ln(-\alpha)}, \sqrt{-\ln(-\alpha)}\right] & \text{if } \alpha \in (-1, 0], \\ \mathbb{R} & \text{if } \alpha > 0. \end{cases}$$

All these sets are convex; hence, the function $f$ is quasi-convex.

<u>Example.</u> The function $f \colon [0, \infty)^2 \to \mathbb{R}$ given by the formula $f(x_1, x_2) = -x_1 x_2$ is quasi-convex. Its sublevel sets for $\alpha \geqslant 0$ are trivial, $W_\alpha(f) = [0, \infty)^2$, and if $\alpha < 0$, then the sublevel sets have a boundary being one branch of a hyperbola. This function is neither convex nor concave, as its Hessian has the eigenvalues $-1$ and $1$.

<u>Example.</u> Let $\mathbf{a}, \mathbf{c} \in \mathbb{R}^n$ and $b, d \in \mathbb{R}$. Let $D = \{\mathbf{x} \in \mathbb{R}^n \colon \mathbf{c}^\mathsf{T}\mathbf{x} + d > 0\}$. Then, the rational function $f \colon D \to \mathbb{R}$ given by

$$f(\mathbf{x}) = \frac{\mathbf{a}^\mathsf{T}\mathbf{x} + b}{\mathbf{c}^\mathsf{T}\mathbf{x} + d}$$

is quasi-linear. The proof is left as an exercise.

<u>Theorem 43</u> *Let $W \subset \mathbb{R}^n$ be a convex set and let $f \colon W \to \mathbb{R}$.*

*I) If the function $f$ is quasi-convex and differentiable at a point $\mathbf{y} \in W$, then*

$$f(\mathbf{x}) \leqslant f(\mathbf{y}) \Rightarrow Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) \leqslant 0 \text{ for all } \mathbf{x} \in W.$$

*II) Suppose that the function* $f$ *is differentiable in the entire set* $W$. *Then* $f$ *is quasi-convex if and only if*

$$f(\mathbf{x}) \leqslant f(\mathbf{y}) \Rightarrow Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) \leqslant 0 \text{ for all } \mathbf{x}, \mathbf{y} \in W.$$

<u>Remark.</u> The implication $f(\mathbf{x}) \leqslant f(\mathbf{y}) \Rightarrow Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) \leqslant 0$ is equivalent to

$$Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) > 0 \Rightarrow f(\mathbf{x}) > f(\mathbf{y}).$$

If the function $f$ is quasi-linear and $f(\mathbf{x}) = f(\mathbf{y})$, then $Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) = 0$.

<u>Proof.</u> (I) We fix $\mathbf{x}, \mathbf{y} \in W$ such that $f(\mathbf{x}) \leqslant f(\mathbf{y})$. For any $\lambda \in (0, 1)$ we have

$$f\big(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})\big) = f\big(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}\big) \leqslant \max\{f(\mathbf{x}), f(\mathbf{y})\} = f(\mathbf{y}).$$

Therefore,

$$\frac{f\big(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})\big) - f(\mathbf{y})}{\lambda} \leqslant 0.$$

By definition of the directional derivative, we obtain

$$Df(\mathbf{y})(\mathbf{x} - \mathbf{y}) = \lim_{\lambda \searrow 0} \frac{f\big(\mathbf{y} + \lambda(\mathbf{x} - \mathbf{y})\big) - f(\mathbf{y})}{\lambda} \leqslant 0$$

The proof of (II) is left as a (not very easy) exercise. □

<u>Theorem 44</u> *If a function* $f \colon W \to \mathbb{R}$ *defined in a convex set* $W \subset \mathbb{R}^n$ *is pseudoconvex, then it is quasi-convex.*

<u>Proof.</u> By assuming that the function $f$ is not quasi-convex we shall get an inconsistency with its pseudoconvexity. Let $\mathbf{x}, \mathbf{y} \in W$ and $\lambda \in (0, 1)$ such that

$$f\big(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}\big) > \max\{f(\mathbf{x}), f(\mathbf{y})\}.$$

Let $\mathbf{z} = \lambda \mathbf{x} + (1 - \lambda)\mathbf{y}$. By pseudoconvexity, we have

$$f(\mathbf{x}) < f(\mathbf{z}) \Rightarrow Df(\mathbf{z})(\mathbf{x} - \mathbf{z}) < 0,$$
$$f(\mathbf{y}) < f(\mathbf{z}) \Rightarrow Df(\mathbf{z})(\mathbf{y} - \mathbf{z}) < 0.$$

The vectors $\mathbf{x} - \mathbf{z}$ and $\mathbf{y} - \mathbf{z}$ have the same direction, but opposite orientations. The directional derivatives of $f$ at $\mathbf{z}$ cannot be negative in both directions. □

<u>Theorem 45</u> *Let a function* $f \colon W \to \mathbb{R}$ *defined in a convex open set* $W \subset \mathbb{R}^n$ *be quasi-convex and continuous. If the function* $f$ *is differentiable at a point* $\overline{\mathbf{x}} \in W$ *and* $Df(\overline{\mathbf{x}}) \neq \mathbf{0}^\mathsf{T}$, *then* $f$ *is pseudoconvex at* $\overline{\mathbf{x}}$.

<u>Proof.</u> We need to prove that for all $\mathbf{y} \in W$ the condition $Df(\overline{\mathbf{x}})(\mathbf{y} - \overline{\mathbf{x}}) \geqslant 0$ implies that $f(\mathbf{y}) \geqslant f(\overline{\mathbf{x}})$. Let $A$ denote the affine space perpendicular to $Df(\overline{\mathbf{x}})$ passing through $\overline{\mathbf{x}}$. As $Df(\overline{\mathbf{x}}) \neq \mathbf{0}^\mathsf{T}$, the dimension of $A$ is $n - 1$, i.e., $A$ is a hyperplane in $\mathbb{R}^n$.

We notice that if $\mathbf{y} \in W \setminus A$ and $Df(\overline{\mathbf{x}})(\mathbf{y} - \overline{\mathbf{x}}) \geqslant 0$, then the directional derivative is (strictly) positive: $Df(\overline{\mathbf{x}})(\mathbf{y} - \mathbf{x}) > 0$. It follows (see the remark above) that $f(\mathbf{y}) > f(\overline{\mathbf{x}})$, which is what was to be proved. Now we fix $\mathbf{y} \in W \cap A$. As $W$ is open and $A$ is a hyperplane, there exists a sequence of points $(\mathbf{y}_k)_k \subset W \setminus A$ converging to $\mathbf{y}$ and such that $Df(\overline{\mathbf{x}})(\mathbf{y}_k - \overline{\mathbf{x}}) > 0$, i.e., $f(\mathbf{y}_k) > f(\overline{\mathbf{x}})$. By continuity of $f$ we obtain $f(\mathbf{y}) \geqslant f(\overline{\mathbf{x}})$. □

## Finding maxima of quasi-convex functions

<u>Theorem 46</u> *Let* $f \colon W \to \mathbb{R}$ *be a quasi-convex and continuous function defined in a convex and compact set* $W \subset \mathbb{R}^n$. *Then at least one of the global solutions of the problem*

$$\begin{cases} f(\mathbf{x}) \to \max, \\ \mathbf{x} \in W \end{cases}$$

*is an extremal point of the set* $W$.

<u>Proof.</u> A convex function reaches its extremal values in a compact set. Therefore, a solution $\overline{\mathbf{x}}$ exists. By Theorem 35, the point $\overline{\mathbf{x}}$ is either an extremal point or a convex combination of a finite set of extremal points of $W$, i.e.,

$$\overline{\mathbf{x}} = a_1 \mathbf{x}_1 + \cdots + a_m \mathbf{x}_m,$$

where $a_1, \ldots, a_m > 0$, $a_1 + \cdots + a_m = 1$. By quasi-convexity of $f$,

$$f(\overline{\mathbf{x}}) \leqslant \max\{f(\mathbf{x}_1), \ldots, f(\mathbf{x}_m)\}.$$

As the function $f$ at $\overline{\mathbf{x}}$ takes its maximal value in $W$, the equality $f(\overline{\mathbf{x}}) = f(\mathbf{x}_i)$ has to hold for some $i \in \{1, \ldots, m\}$. □

## Sufficient conditions

Below we consider the optimization problem with both non-equality and equality constraints:

$$\begin{cases} f(\mathbf{x}) \to \min, \\ g_i(\mathbf{x}) \leqslant 0, \ i = 1, \ldots, m, \\ h_j(\mathbf{x}) = 0, \ j = 1, \ldots, l, \\ \mathbf{x} \in \mathbb{X}, \end{cases} \qquad (*)$$

where $\mathbb{X} \subset \mathbb{R}^n$ is an open set and $f, g_1, \ldots, g_m, h_1, \ldots, h_l \colon \mathbb{X} \to \mathbb{R}$. Thus, the feasible set is

$$W = \{\, \mathbf{x} \in \mathbb{X} \colon g_1(\mathbf{x}) \leqslant 0, \ \ldots, \ g_m(\mathbf{x}) \leqslant 0, \ h_1(\mathbf{x}) = \cdots = h_l(\mathbf{x}) = 0 \,\}.$$

The functions $g_i$ are called <u>non-equality constraints</u>, the functions $h_j$ are called <u>equality constraints</u> and the problem above is called an <u>optimization problem with mixed constraints</u>.

The theorem below describes sufficient conditions for a point satisfying the first order necessary conditions to be a global solution of $(*)$.

<u>Theorem 47</u> *Let $\overline{\mathbf{x}} \in W$. Assume that*

- *the functions $g_i$, where $i \notin I(\overline{\mathbf{x}})$, are continuous at $\overline{\mathbf{x}}$, the functions $g_i$, where $i \in I(\overline{\mathbf{x}})$, are differentiable at $\overline{\mathbf{x}}$ and quasi-convex,*

- *the functions $h_j$, where $j \in \{1, \ldots, l\}$, are quasi-linear and differentiable at $\overline{\mathbf{x}}$,*

- *the function $f$ is pseudoconvex at $\overline{\mathbf{x}}$.*

*If there exists $\boldsymbol{\mu} \in [0, \infty)^m$ and $\boldsymbol{\lambda} \in \mathbb{R}^l$ which satisfy the following first order condition:*

$$\begin{cases} Df(\overline{\mathbf{x}}) + \displaystyle\sum_{i \in I(\overline{\mathbf{x}})} \mu_i Dg_i(\overline{\mathbf{x}}) + \sum_{j=1}^{l} \lambda_j Dh_j(\overline{\mathbf{x}}) = \mathbf{0}^\mathsf{T}, \\ \mu_i g_i(\overline{\mathbf{x}}) = 0, \ i = 1, \ldots, m, \end{cases}$$

*then the point $\overline{\mathbf{x}}$ is a global solution.*

<u>Proof.</u> Let $\mathbf{x} \in W$. We multiply the sides of the first equation by $\mathbf{x} - \overline{\mathbf{x}}$:

$$Df(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) + \sum_{i \in I(\overline{\mathbf{x}})} \mu_i Dg_i(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) + \sum_{j=1}^{l} \lambda_j Dh_j(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) = 0$$

By Theorem 43, $Dh_j(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) = 0$ for all $j$, as $h_j(\mathbf{x}) = h_j(\overline{\mathbf{x}}) = 0$. By the same theorem, also $Dg_i(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) \leqslant 0$ for $i \in I(\overline{\mathbf{x}})$, as $0 = g_i(\overline{\mathbf{x}}) \geqslant g_i(\mathbf{x})$. From the above we conclude that

$$Df(\overline{\mathbf{x}})(\mathbf{x} - \overline{\mathbf{x}}) \geqslant 0.$$

By the definition of a pseudoconvex function, $f(\mathbf{x}) \geqslant f(\overline{\mathbf{x}})$. As the choice of $\mathbf{x} \in W$ was arbitrary, $\overline{\mathbf{x}}$ is a global solution. $\square$

<u>Remark.</u> If the assumptions of Theorem 47 are satisfied locally, in a neighbourhood of $\overline{\mathbf{x}}$, then $\overline{\mathbf{x}}$ is a local solution.

<u>Remark.</u> By Theorem 45, instead of the pseudoconvexity of $f$ at $\overline{\mathbf{x}}$, we can assume its continuity in $\mathbb{X}$, quasi-convexity and the condition $Df(\overline{\mathbf{x}}) \neq \mathbf{0}^\mathsf{T}$. It is one of necessary conditions, given by Arrow and Enthoven in 1961.

# 8. A necessary condition for mixed constraints

Below we are going to derive a necessary first-order condition for the optimization problem given as follows:

$$\begin{cases} f(\boldsymbol{x}) \to \min, \\ g_i(\boldsymbol{x}) \leqslant 0, \ i = 1, \ldots, m \\ h_j(\boldsymbol{x}) = 0, \ j = 1, \ldots, l, \\ \boldsymbol{x} \in \mathbb{X}, \end{cases} \qquad (*)$$

where $\mathbb{X} \subset \mathbb{R}^n$ is an open set and $f, g_1, \ldots, g_m, h_1, \ldots, h_l \colon \mathbb{X} \to \mathbb{R}$. The set of feasible points is thus

$$W = \{\, \boldsymbol{x} \in \mathbb{X} \colon g_1(\boldsymbol{x}) \leqslant 0, \, \ldots, \, g_m(\boldsymbol{x}) \leqslant 0, \, h_1(\boldsymbol{x}) = \cdots = h_l(\boldsymbol{x}) = 0 \,\}.$$

<u>Example.</u> Consider the following problem:

$$\begin{cases} f(\boldsymbol{x}) \to \min \\ \boldsymbol{a}^\mathsf{T} \boldsymbol{x} + b = 0, \\ \boldsymbol{x} \in \mathbb{R}^n, \end{cases}$$

for some $\boldsymbol{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$. The equality constraint may be replaced by two inequality constraints:

$$\begin{cases} f(\boldsymbol{x}) \to \min \\ \boldsymbol{a}^\mathsf{T} \boldsymbol{x} + b \leqslant 0, \\ -\boldsymbol{a}^\mathsf{T} \boldsymbol{x} - b \leqslant 0, \\ \boldsymbol{x} \in \mathbb{R}^n, \end{cases}$$

The constraints are affine functions; hence, the affine function condition is satisfied at each point. If $\overline{\boldsymbol{x}}$ is a local solution, then there exists a vector of Lagrange multipliers $\boldsymbol{\mu} = [\mu_1, \mu_2]^\mathsf{T}$ such that the Kuhn–Tucker conditions are satisfied:

$$\begin{cases} Df(\overline{\boldsymbol{x}}) + \mu_1 \boldsymbol{a}^\mathsf{T} + \mu_2(-\boldsymbol{a}^\mathsf{T}) = \boldsymbol{0}^\mathsf{T}, \\ \mu_1(\boldsymbol{a}^\mathsf{T} \boldsymbol{x} + b) = 0, \\ \mu_2(-\boldsymbol{a}^\mathsf{T} \boldsymbol{x} - b) = 0, \\ \mu_1, \mu_2 \geqslant 0. \end{cases}$$

The point $\overline{\boldsymbol{x}}$, being a solution, is feasible, i.e., it satisfies the constraints: $\boldsymbol{a}^\mathsf{T} \overline{\boldsymbol{x}} + b = 0$. Therefore, the second and third equalities are trivially satisfied. The conditions above are thus equivalent to

$$\begin{cases} Df(\overline{\boldsymbol{x}}) + (\mu_1 - \mu_2) \boldsymbol{a}^\mathsf{T} = \boldsymbol{0}^\mathsf{T}, \\ \mu_1, \mu_2 \geqslant 0. \end{cases}$$

Denote $\lambda = \mu_1 - \mu_2$. As these two numbers must only be nonnegative, $\lambda \in \mathbb{R}$. As a result we obtain

$$Df(\overline{\boldsymbol{x}}) + \lambda \boldsymbol{a}^\mathsf{T} = \boldsymbol{0}^\mathsf{T}, \quad \lambda \in \mathbb{R},$$

which is the Kuhn–Tucker condition for equality constraints.

The example above suggests that the theory for inequality constraints, developed earlier, is also suitable for dealing with equality constraints. Unfortunately, it is not the case. Affine constraints are very special. If any constraint is not affine and we split it to a pair of inequality constraints, as above, then neither the linear independence condition nor the Slater condition is satisfied at any point of the set $W$.

## Necessary first-order condition

The theory described below is a direct extension of what we have done for the optimization problems with inequality contraints. We begin with extending the definition of the set $T_{\text{lin}}(\overline{\boldsymbol{x}})$:

<u>Definition 26</u> *Let $\overline{\boldsymbol{x}} \in W$, let $g_i$, where $i \in I(\overline{\boldsymbol{x}})$ be functions differentiable at $\overline{\boldsymbol{x}}$ and let the functions $h_1, \ldots, h_l$ be differentiable at $\overline{\boldsymbol{x}}$. The set*

$$T_{\text{lin}}(\overline{\boldsymbol{x}}) = \{\, \boldsymbol{d} \in \mathbb{R}^n \colon Dg_i(\overline{\boldsymbol{x}}) \boldsymbol{d} \leqslant 0 \ \text{for } i \in I(\overline{\boldsymbol{x}}), Dh_j(\overline{\boldsymbol{x}}) \boldsymbol{d} = 0 \ \text{for } i = 1, \ldots, l \,\}$$

*is called the <u>cone of tangents for active (binding) constraints</u>.*

Just as in the case of inequality constraints, the cone of tangents for binding constraints is a polyhedral set, i.e., it is a convex and closed set. However, with at least one equality constraint the interior of this set is empty.

The necessary condition for the existence of a local solution of the optimization problem with mixed constraints given below assumes the equality of the cone of tangents to the set $W$ and the cone of tangents for binding constraints. Later we show generalised regularity conditions which imply this equality.

Theorem 48 *(Kuhn–Tucker theorem) Let $\overline{x}$ be a local solution of the problem (\*). If the functions $f, g_i, i \in I(\overline{x}), h_j, j = 1, \dots, l$ are differentiable at $\overline{x}$ and $T(\overline{x}) = T_{\text{lin}}(\overline{x})$, then there exist vectors $\mu \in [0, \infty)^m$ and $\lambda \in \mathbb{R}^l$ such that*

$$\begin{cases} Df(\overline{x}) + \sum_{i \in I(\overline{x})} \mu_i Dg_i(\overline{x}) + \sum_{j=1}^{l} \lambda_j h_j(\overline{x}) = 0^\mathsf{T}, \\ \mu_i g_i(\overline{x}) = 0, \ i = 1, \dots, m. \end{cases}$$

Proof. By Theorem 37, $D(\overline{x}) \cap T(\overline{x}) = \emptyset$. Then, by assumption, $D(\overline{x}) \cap T_{\text{lin}}(\overline{x}) = \emptyset$, which means that there is no solution $z$ of the system

$$\begin{cases} Df(\overline{x})z < 0, \\ Dg_i(\overline{x})z \leqslant 0, \ i \in I(\overline{x}), \\ Dh_j(\overline{x})z \leqslant 0, \ j = 1, \dots, l, \\ -Dh_j(\overline{x})z \leqslant 0, \ j = 1, \dots, l. \end{cases}$$

We use the Farkas' lemma with $\mathbf{d} = -Df(\overline{x})$ and with the following matrix $A$:

$$A = \begin{bmatrix} Dh_j(\overline{x}), & j = 1, \dots, l, \\ -Dh_j(\overline{x}), & j = 1, \dots, l, \\ Dg_i(\overline{x}), & i \in I(\overline{x}) \end{bmatrix}.$$

Hence, there exists $\mathbf{y} \in [0, \infty)^{|I(\overline{x})|+2l}$ such that $\mathbf{y}^\mathsf{T} A = -Df(\overline{x})$, or

$$Df(\overline{x}) + \mathbf{y}^\mathsf{T} A = 0^\mathsf{T}. \tag{$\diamond$}$$

Let $\lambda_j = y_j - y_{l+j}$ for $j = 1, \dots, l$. Let the coordinates of the vector $\mu \in \mathbb{R}^m$ corresponding to the active constraints be equal to the last $|I(\overline{x})|$ coordinates of $\mathbf{y}$ and let the other coordinates of $\mu$ be equal to $0$. Then, $(\diamond)$ is equivalent to

$$Df(\overline{x}) + \sum_{i \in I(\overline{x})} \mu_i Dg_i(\overline{x}) + \sum_{j=1}^{l} \lambda_j h_j(\overline{x}) = 0^\mathsf{T}.$$

By definition of $\mu$ it is obvious that $\mu_i g_i(\overline{x}) = 0$ for all $i$. $\square$

## Regularity constraints

Below three sufficient conditions for the equality $T(\overline{x}) = T_{\text{lin}}(\overline{x})$, called the regularity conditions, are defined.

Definition 27 *At a point $\overline{x} \in W$*

- *the linear independence condition is satisfied if the functions $g_i$, where $i \notin I(\overline{x})$, are continuous at $\overline{x}$ and all the other inequality and equality constraints are continuously differentiable in a neighbourhood of $\overline{x}$ and the vectors $Dg_i(\overline{x})$ for $i \in I(\overline{x})$ and $Dh_j$ for $j = 1, \dots, l$ are linearly independent,*

- *the affine function condition is satisfied if the functions $g_i(\overline{x})$ for $i \in I(\overline{x})$ and $h_j$ for $j = 1, \dots, l$ are affine, while $g_i$, for $i \notin I(\overline{x})$, are continuous at $\overline{x}$,*

- *the Slater condition is satisfied if*
  *—the functions $g_i, i \in I(\overline{x})$ are pseudoconvex at $\overline{x}$,*
  *—the functions $g_i, i \notin I(\overline{x})$ are continuous at $\overline{x}$,*
  *—the functions $h_j, j = 1, \dots, l$, are affine,*
  *—there exists $x^* \in \mathbb{X}$ such that $g_i(x^*) < 0$ for $i \in I(\overline{x})$ and $h_j(x^*) = 0$ for $l = 1, \dots, l$.*

Theorem 49 *If the affine function condition is satisfied at a point $\overline{x} \in W$, then $T(\overline{x}) = T_{\text{lin}}(\overline{x})$.*

Proof. Just as in the last example, we change the affine equality constraints into pairs of affine inequality constraints. The claim follows from Theorem 39. $\square$

Theorem 50 *If the Slater condition is satisfied at a point $\overline{x} \in W$, then $T(\overline{x}) = T_{\text{lin}}(\overline{x})$.*

Proof. The functions $h_j$, for $j = 1, \dots, l$, have the form

$$h_j(\mathbf{y}) = \mathbf{a}_j^\mathsf{T} \mathbf{y} + b_j, \quad \mathbf{a}_j \in \mathbb{R}^n, b_j \in \mathbb{R}.$$

We introduce the generalisation of set $T_{int}(\overline{x})$ for mixed constraints:

$$T_{int}(\overline{x}) = \{\, d \in \mathbb{R}^n \colon Dg_i(\overline{x})d < 0 \text{ for } i \in I(\overline{x}), Dh_j(\overline{x})d = 0 \text{ for } j = 1, \ldots, l \}.$$

(1) We prove that $T_{int}(\overline{x}) \neq \emptyset$. We take $x$ from the Slater condition. By pseudoconvexity of $g_i$, we have

$$Dg_i(\overline{x})(x - \overline{x}) < 0, \quad i \in I(\overline{x}).$$

For all $j$ we also have

$$a_j^T(x - \overline{x}) = a_j^T x + b_j - a_j^T \overline{x} - b_j = h_j(x) - h_j(\overline{x}) = 0.$$

Hence, $(x - \overline{x}) \in T_{int}(\overline{x})$.

(2) We prove that $T_{int}(\overline{x}) \subset T(\overline{x})$. We choose any $d \in T_{int}(\overline{x})$. It suffices to show that there exists a line segment contained in $W$ whose one end is $\overline{x}$ and whose direction is $d$. Let $y(\lambda) = \overline{x} + \lambda d$. As the inactive constraints are continuous functions, there exists $\varepsilon > 0$ such that $g_i(y(\lambda)) \leqslant 0$ for $\lambda \in [0, \varepsilon]$ and $i \notin I(\overline{x})$. From $d \in T_{int}(\overline{x})$ it follows also that $h_j(y(\lambda)) = 0$ for $j = 1, \ldots, l$ (because $h_j$ are affine functions). And as the functions $g_i$, where $i \in I(\overline{x})$, are differentiable at $\overline{x}$ and $d \in T_{int}(\overline{x})$, it follows that

$$\lim_{\lambda \searrow 0} \frac{g_i(y(\lambda)) - g_i(\overline{x})}{\lambda} = Dg_i(\overline{x})d < 0.$$

Hence, $g_i(y(\lambda)) - g_i(\overline{x}) < 0$ for $\lambda$ small enough.

(3) We prove that $\operatorname{cl} T_{int}(\overline{x}) = T_{lin}(\overline{x})$. The sets $T_{int}(\overline{x})$ and $T_{lin}(\overline{x})$ are contained in the affine subspace $H$ determined by the affine equality constraints. Therefore, there exists an affine mapping $P$, such that the image of the subspace $H$ is $\mathbb{R}^{n'}$, where $n' = \dim H$ (if the vectors $a_j$ are linearly independent, then $n' = n - l$). The mapping $P$ restricted to $H$ is a one-to-one function. The topologies of $H$ and $\mathbb{R}^{n'}$ are, therefore, identical, which means that it suffices to prove this claim for the images $T'_{int}(\overline{x})$ and $T'_{lin}(\overline{x})$ of the sets $T_{int}(\overline{x})$ and $T_{lin}(\overline{x})$ under the mapping $P$. We notice that $T'_{int}(\overline{x})$ is open in $\mathbb{R}^{n'}$ and nonempty. It is also the interior of $T'_{lin}(\overline{x})$. By Lemma 17 we have $T'_{int}(\overline{x}) = T'_{lin}(\overline{x})$.

(4) The proof that $T(\overline{x}) \subset T_{lin}(\overline{x})$ is identical to the proof of Lemma 14.

(5) It suffices to recall that $T(\overline{x})$ is a closed set. Therefore,

$$\operatorname{cl} T_{int}(\overline{x}) \subset T(\overline{x}) \subset T_{lin}(\overline{x}) = \operatorname{cl} T_{int}(\overline{x}). \quad \square$$

Before tackling the third regularity condition, we recall the implicit function theorem, which we need to describe the cone of tangents to a surface determined by equality constraints.

<u>Theorem 51</u> *(implicit function theorem) Let* $f \colon \mathbb{X} \to \mathbb{R}^n$, *where* $\mathbb{X} \subset \mathbb{R}^{n+m}$ *is an open set, be a function of class* $C^k$. *Assume that* $f(\overline{x}, \overline{y}) = 0$, *where* $\overline{x} \in \mathbb{R}^n$, $\overline{y} \in \mathbb{R}^m$, $(\overline{x}, \overline{y}) \in \mathbb{X}$. *By* $A_x$ *we denote the matrix of partial derivatives of* $f$ *with respect to the first* $n$ *variables at* $(\overline{x}, \overline{y})$: $A_x \in \mathbb{R}^{n \times n}$, $(A_x)_{ij} = \frac{\partial f_i}{\partial u_j}(\overline{x}, \overline{y})$.

*If the matrix* $A_x$ *is nonsingular, then there exists an open set* $W \subset \mathbb{R}^m$, *whose element is* $\overline{y}$, *and a function* $g \colon W \to \mathbb{R}^n$ *of class* $C^k$, *such that* $g(\overline{y}) = \overline{x}$ *and* $(g(y), y) \in \mathbb{X}$ *and* $f(g(y), y) = 0$ *for all* $y \in W$. *Moreover,* $Dg(\overline{y}) = -(A_x)^{-1} A_y$, *where the matrix* $A_y \in \mathbb{R}^{n \times m}$ *consists of the derivatives of* $f$ *at* $(\overline{x}, \overline{y})$ *with respect to the last* $m$ *variables:* $(A_y)_{ij} = \frac{\partial f_i}{\partial u_{n+j}}(\overline{x}, \overline{y})$.

We consider the surface $S$ described by a system of $m^*$ equations:

$$S = \{\, x \in \mathbb{X} \colon c_i(x) = 0, \, i = 1, \ldots, m^* \},$$

where $\mathbb{X}$ is an open set. By $T^S(\overline{x})$ we denote the cone of tangents to $S$ at a point $\overline{x} \in S$.

<u>Theorem 52</u> *Let* $k \geqslant 1$. *Assume that the functions* $c_1, \ldots, c_{m^*}$, *are of class* $C^k$ *in a neighbourhood of* $\overline{x}$ *and assume that their gradients* $Dc_i(\overline{x})$ *for* $i = 1, \ldots, m^*$ *are linearly independent. Then,*

$$T^S(\overline{x}) = T_{lin}^S(\overline{x}) \stackrel{\text{def}}{=} \{\, d \in \mathbb{R}^n \colon Dc_i(\overline{x})d = 0, \, i = 1, \ldots, m^* \}.$$

*Moreover, for any* $d \in T^S(\overline{x})$ *there exists an* $\varepsilon > 0$ *and a parametric curve* $y \colon (-\varepsilon, \varepsilon) \to S$ *of class* $C^k$ *such that* $y(0) = \overline{x}$ *and* $y'(0) = d$.

<u>Proof.</u> First we show that $T^S(\overline{x}) \subset T_{lin}^S(\overline{x})$. Let $d \in T^S(\overline{x})$. Then, $d = \lim_{k \to \infty} \lambda_k(x_k - \overline{x})$ for some sequences $(x_k)_k \subset S$, $x_k \neq \overline{x}$, and $(\lambda_k)_k \subset (0, \infty)$. By definition of the directional derivative, for each $i = 1, \ldots, m^*$ we obtain

$$\underbrace{c_i(x_k)}_{=0} = \underbrace{c_i(\overline{x})}_{=0} + Dc_i(\overline{x}) \underbrace{\lambda_k(x_k - \overline{x})}_{\to d} + \underbrace{\lambda_k \|x_k - \overline{x}\|}_{\to \|d\|} \underbrace{\frac{o(\|x_k - \overline{x}\|)}{\|x_k - \overline{x}\|}}_{\to 0},$$

i.e., $Dc_i(\overline{x})d = 0$. Hence, $d \in T_{lin}^S(\overline{x})$.

It remains to be proved that $T^S_{lin}(\overline{x}) \subset T^S(\overline{x})$, which is more difficult. We fix $d \in T^S_{lin}(\overline{x})$. We are going to construct a curve located on S and passing through $\overline{x}$, whose derivative at $\overline{x}$ is $d$. We denote $c(x) = (c_1(x), \ldots, c_{m^*}(x))$ and we define the function $\Phi \colon \mathbb{R}^{m^*} \times \mathbb{R} \to \mathbb{R}^{m^*}$ by the formula

$$\Phi(u, t) = c\big(\overline{x} + td + (Dc(\overline{x}))^\mathsf{T} u\big).$$

We can see that $\Phi(0, 0) = 0$. By $D_u \Phi$ we denote the matrix of partial derivatives with respect to the coordinates of the vector $u$: $D_u \Phi = \big(\frac{\partial \Phi_i}{\partial u_j}\big)_{i,j=1,\ldots,m^*}$. At the point $(0, 0)$ we have $D_u \Phi(0, 0) = Dc(\overline{x})(Dc(\overline{x}))^\mathsf{T}$. Recall that due to the assumption, the matrix $Dc(\overline{x})$ is of full rank (equal to $m^*$), i.e., $D_u \Phi(0, 0)$ is nonsingular. By the implicit function theorem, there exists $\varepsilon > 0$ and a function $u \colon (-\varepsilon, \varepsilon) \to \mathbb{R}^{m^*}$ of class $C^k$ such that $\Phi\big(u(t), t\big) = 0$ and $u(0) = 0$. Let

$$y(t) = \overline{x} + td + \big(Dc(\overline{x})\big)^\mathsf{T} u(t).$$

This curve, according to the construction, is located on the surface S, as $c\big(y(t)\big) = \Phi\big(u(t), t\big) = 0$ for $t \in (-\varepsilon, \varepsilon)$ and $y(0) = \overline{x}$. The derivative of the composition of functions $c \circ y$ is

$$\frac{d}{dt} c\big(y(t)\big) = Dc\big(y(t)\big)\big(d + (Dc(\overline{x}))^\mathsf{T} u'(t)\big),$$

and at $t = 0$ we have

$$\frac{d}{dt} c\big(y(t)\big)\Big|_{t=0} = Dc\big(y(0)\big)\big(d + (Dc(\overline{x}))^\mathsf{T} u'(0)\big).$$

On the other hand, we know that $c\big(y(t)\big) = 0$, i.e., the derivative above is equal to 0: $Dc\big(y(0)\big)\big(d + (Dc(\overline{x}))^\mathsf{T} u'(0)\big) = 0$. But we have chosen $d \in T^S_{lin}(\overline{x})$, which in our notation means that $Dc(\overline{x})d = 0$. Hence, $Dc(\overline{x})\big(Dc(\overline{x})\big)^\mathsf{T} u'(0) = 0$. As the rank of the matrix $Dc(\overline{x})$ is equal to $m^*$, we obtain $u'(0) = 0$. To complete the proof we compute the derivative of $y$:

$$y'(t) = d + \big(Dx(\overline{x})\big)^\mathsf{T} u'(t),$$

and the above, for $t = 0$, gives us $y'(0) = d$. This equality allows us to conclude that $d \in T^S(\overline{x})$. □

Remark. The theorem above establishes a well known fact about spaces tangent to manifolds. From its assumptions it follows that S locally (in a neighbourhood of $\overline{x}$) is a differential manifold of class $C^k$. The tangent space at $\overline{x}$ is defined as the set of vectors being derivatives (at $\overline{x}$) of curves contained in this manifold and passing through $\overline{x}$ (which is equivalent to the definition of $T(\overline{x})$). The equality

$T_{lin}(\overline{x}) = T(\overline{x})$ means that the tangent space is the kernel of the linear mapping $Dc(\overline{x})$.

Theorem 52 will often be used in what follows; it is the main tool used in the proof of the sufficient second-order condition. Using this theorem, we can easily prove the equality $T(\overline{x}) = T_{lin}(\overline{x})$ assuming that the linear independence condition is satisfied:

<u>Theorem 53</u> *If the linear independence condition is satisfied at a point $\overline{x} \in W$, then $T(\overline{x}) = T_{lin}(\overline{x})$.*

<u>Proof.</u> Let $d \in T_{lin}(\overline{x})$. Let $\hat{I}(\overline{x}) = \{i \in I(\overline{x}) \colon Dg_i(\overline{x})d = 0\}$. We define the surface

$$S = \{x \in \mathbb{X} \colon c_k(x) = 0, \text{ where } c_k = g_i \text{ if } i \in I(\overline{x}) \text{ or } c_k = h_j, j = 1, \ldots, l\}.$$

Then,

$$T^S(\overline{x}) = \{d \in \mathbb{R}^n \colon Dg_i(\overline{x})d = 0, i \in \hat{I}(\overline{x}), Dh_j(\overline{x})d = 0, j = 1, \ldots, l\}.$$

By Theorem 52, there exists a curve $y \colon (-\varepsilon, \varepsilon) \to \mathbb{R}^n$ such that $y(0) = \overline{x}$, $y'(0) = d$ and $g_i\big(y(t)\big) = 0$, $i \in \hat{I}(\overline{x})$, $h_j\big(y(t)\big) = 0$, $j = 1, \ldots, l$. For $i \in I(\overline{x}) \setminus \hat{I}(\overline{x})$ let $\hat{g}_i(t) = g_i\big(y(t)\big)$, $t \in (-\varepsilon, \varepsilon)$. Then, $\hat{g}_i'(0) = Dg_i(\overline{x})d < 0$, i.e., there exists $\varepsilon_i > 0$ such that $\hat{g}_i(t) < 0$ for $t \in [0, \varepsilon_i)$. For $i \notin I(\overline{x})$, by continuity of $g_i$, there is $g_i\big(y(t)\big) < 0$ in a neighbourhood of 0. Therefore, there exists $\overline{\varepsilon} > 0$ such that $y(t) \in W$ for $t \in [0, \overline{\varepsilon})$. Hence, trivially, $d \in T(\overline{x})$.

The proof of the set inclusion $T(\overline{x}) \subset T_{lin}(\overline{x})$ is identical to the proof of Lemma 14. □

# 9. Second-order conditions

We consider the optimization problem with mixed constraints. We assume that at a point $\overline{x} \in W$ there is $T(\overline{x}) = T_{lin}(\overline{x})$ and the first-order necessary condition is satisfied, i.e., there exist vectors $\mu \in [0, \infty)^m$ and $\lambda \in \mathbb{R}^l$ such that

$$
\begin{cases}
Df(\overline{x}) + \sum_{i \in I(\overline{x})} \mu_i Dg_i(\overline{x}) + \sum_{j=1}^{l} \lambda_j Dh_j(\overline{x}) = 0^T, \\
\mu_i g_i(\overline{x}) = 0, \quad i = 1, \ldots, m.
\end{cases}
$$

<u>Definition 28</u> *The <u>Lagrange function</u> is the function given by the formula*

$$
L(x, \mu, \lambda) = f(x) + \sum_{i \in I(\overline{x})} \mu_i g_i(x) + \sum_{j=1}^{l} \lambda_j h_j(x)
$$

The first-order condition may be written in a shorter form

$$
\begin{cases}
D_x L(x, \mu, \lambda) = 0^T, \\
\mu_i g_i(\overline{x}) = 0, \quad i = 1, \ldots, m,
\end{cases}
$$

where $D_x$ denotes the derivative with respect to $x$.

<u>Definition 29</u> *The set*

$$
I^*(\overline{x}) = \{\, i \in I(\overline{x}) : \mu_i > 0 \,\}
$$

*is called the set of <u>strongly binding inequality constraints</u>. The set*

$$
I^0(\overline{x}) = I(\overline{x}) \setminus I^*(\overline{x})
$$

*is called the set of <u>weakly binding inequality constraints</u>.*

<u>Theorem 54</u> *(<u>necessary second-order condition</u>) Suppose that $\overline{x}$ is a local solution of the problem with mixed constraints and the linear independence condition is satisfied at $\overline{x}$. Let $\mu$ and $\lambda$ be vectors of Lagrange multipliers of the first-order condition. If the functions $f$, $g_i$ for $i \in I(\overline{x})$ and $h_1, \ldots, h_l$ are twice differentiable in a neighbourhood of $\overline{x}$, then*

$$
d^T D_x^2 L(x, \mu, \lambda) d \geqslant 0
$$

*for all $d \in \mathbb{R}^n$ such that*

$$
\begin{aligned}
Dg_i(\overline{x}) d &= 0, \quad i \in I(\overline{x}), \\
Dh_j(\overline{x}) d &= 0, \quad j = 1, \ldots, l.
\end{aligned}
$$

<u>Proof.</u> Let $d \in \mathbb{R}^n$ be a vector satisfying the conditions given above. By Theorem 52 there exists $\varepsilon > 0$ and a parametric curve $y : (-\varepsilon, \varepsilon) \to \mathbb{R}^n$ of class $C^2$ having the following properties: $y(0) = \overline{x}$, $y'(0) = d$ and for $t \in (-\varepsilon, \varepsilon)$ there is $h_j(y(t)) = 0$, $j = 1, \ldots, l$ and $g_i(y(t)) = 0$, $i \in I(\overline{x})$. Hence, the function $F(t) \stackrel{\text{def}}{=} L(y(t), \mu, \lambda)$ is equal to $f(y(t))$ for $t \in (-\varepsilon, \varepsilon)$. The continuity of the inactive constraints allows us to conclude that $y(t) \in W$ for $t$ in a neighbourhood of $0$. Hence, $F$ has a local minimum at $0$, as $y(0)$ is a local minimum of $f$ in $W$. From the assumptions it follows that $F$ is of class $C^2$. The existence of its minimum at $0$ implies that $F''(0) \geqslant 0$, i.e.,

$$
0 \leqslant d^T D_x^2 L(\overline{x}, \mu, \lambda) d + D_x L(\overline{x}, \mu, \lambda) y''(0).
$$

The above completes the proof, as at the point $\overline{x}$ the necessary conditions are satisfied, in particular

$$
D_x L(\overline{x}, \mu, \lambda) = 0^T. \quad \square
$$

Theorem 54 may be generalised as follows:

<u>Theorem 55</u> *If the assumptions of Theorem 54 are satisfied, then the inequality*

$$
d^T D_x^2 L(x, \mu, \lambda) d \geqslant 0
$$

*is satisfied for all $d \in \mathbb{R}^n$ such that*

$$
\begin{aligned}
Dg_i(\overline{x}) d &= 0, \quad i \in I^*(\overline{x}), \\
Dg_i(\overline{x}) d &\leqslant 0, \quad i \in I^0(\overline{x}), \\
Dh_j(\overline{x}) d &= 0, \quad j = 1, \ldots, l.
\end{aligned}
$$

The next theorem describes a sufficient condition for a local solution. Note that if this condition is satisfied then the solution is strict. This leaves us a "grey zone", where the necessary condition is satisfied, but the sufficient condition is not, just like in the case of unconstrained optimization.

<u>Theorem 56</u> *(<u>sufficient second-order condition</u>) Suppose that the first-order necessary condition is satisfied at a point $\overline{x} \in W$ and the functions $g_i$ for $i \in I^*(\overline{x})$ and $h_1 \ldots, h_l$ are twice differentiable at $\overline{x}$. If*

$$
d^T D_x^2 L(x, \mu, \lambda) d > 0
$$

*for all* $\mathbf{d} \in \mathbb{R}^n \setminus \{0\}$ *such that*

$$Dg_i(\overline{\mathbf{x}})\mathbf{d} = 0, \quad i \in I^*(\overline{\mathbf{x}}),$$
$$Dg_i(\overline{\mathbf{x}})\mathbf{d} \leqslant 0, \quad i \in I^0(\overline{\mathbf{x}}), \qquad\qquad (\otimes)$$
$$Dh_j(\overline{\mathbf{x}})\mathbf{d} = 0, \quad j = 1, \dots, l,$$

*then the point* $\overline{\mathbf{x}}$ *is a strict local solution.*

Note that no regularity condition is assumed in this theorem.

<u>Proof.</u> The proof is done by contradiction. Suppose that $\overline{\mathbf{x}}$ is not a strict local minimum. Then, there exists a sequence of feasible points, $(\mathbf{x}_k)_k \in W$, convergent to $\overline{\mathbf{x}}$ and such that $\mathbf{x}_k \neq \overline{\mathbf{x}}$ and $f(\mathbf{x}_k) \leqslant f(\overline{\mathbf{x}})$ for all $k$. Let

$$\mathbf{d}_k = \frac{\mathbf{x}_k - \overline{\mathbf{x}}}{\|\mathbf{x}_k - \overline{\mathbf{x}}\|} \quad \text{and} \quad s_k = \|\mathbf{x}_k - \overline{\mathbf{x}}\|.$$

Then, $\mathbf{x}_k = \overline{\mathbf{x}} + s_k \mathbf{d}_k$. As $\|\mathbf{d}_k\| = 1$ for all $k$, there exists a subsequence $(\mathbf{d}_{k_p})_p$ convergent to a unit vector $\mathbf{d}$. To simplify the notation, assume that the original sequence $(\mathbf{d}_k)_k$ converges to $\mathbf{d}$. It is obvious that $\lim_{k \to \infty} s_k = 0$.

In what follows we are going to prove two properties of the vector $\mathbf{d}$:
(a) $\mathbf{d}^\mathsf{T} D_x^2 L(\overline{\mathbf{x}}, \mu, \lambda)\mathbf{d} \leqslant 0$ and (b), the vector $\mathbf{d}$ satisfies the conditions $(\otimes)$. Which contradicts the assumptions.

(a) By definition of the second order derivative,

$$L(\mathbf{x}, \mu, \lambda) = L(\overline{\mathbf{x}}, \mu, \lambda) + D_x L(\overline{\mathbf{x}}, \mu, \lambda)(\mathbf{x} - \overline{\mathbf{x}}) +$$
$$\frac{1}{2}(\mathbf{x} - \overline{\mathbf{x}})^\mathsf{T} D_x^2 L(\overline{\mathbf{x}}, \mu, \lambda)(\mathbf{x} - \overline{\mathbf{x}}) + o(\|\mathbf{x} - \overline{\mathbf{x}}\|)$$

From $f(\mathbf{x}_k) \leqslant f(\overline{\mathbf{x}})$, $g_i(\mathbf{x}_k) \leqslant g_i(\overline{\mathbf{x}})$ for $i \in I(\overline{\mathbf{x}})$ and $h(\mathbf{x}_k) = h(\overline{\mathbf{x}}) = 0$ it follows that $L(\mathbf{x}_k, \mu, \lambda) \leqslant L(\overline{\mathbf{x}}, \mu, \lambda)$. From the first-order condition it follows that $D_x L(\overline{\mathbf{x}}, \mu, \lambda) = \mathbf{0}^\mathsf{T}$. Hence,

$$(\mathbf{x} - \overline{\mathbf{x}})^\mathsf{T} D_x^2 L(\overline{\mathbf{x}}, \mu, \lambda)(\mathbf{x} - \overline{\mathbf{x}}) + o(\|\mathbf{x} - \overline{\mathbf{x}}\|) \leqslant 0.$$

To the above we substitute $\mathbf{x}_k = \overline{\mathbf{x}} + s_k \mathbf{d}_k$ and we obtain

$$s_k^2 \mathbf{d}_k^\mathsf{T} D_x^2 L(\overline{\mathbf{x}}, \mu, \lambda)\mathbf{d}_k + o(s_k^2 \|\mathbf{d}_k\|^2) \leqslant 0.$$

After dividing both sides of the above by $s_k^2$ and recalling that $\|\mathbf{d}_k\| = 1$ we obtain

$$\mathbf{d}_k^\mathsf{T} D_x^2 L(\overline{\mathbf{x}}, \mu, \lambda)\mathbf{d}_k + \frac{o(s_k^2)}{s_k^2} \leqslant 0.$$

With $k$ tending to $\infty$ the second term above tends to $0$ and $\mathbf{d}_k$ tends to $\mathbf{d}$; hence,

$$\mathbf{d}^\mathsf{T} D_x^2 L(\overline{\mathbf{x}}, \mu, \lambda)\mathbf{d} \leqslant 0.$$

(b) is proved in a similar way. For a function $f$ differentiable at $\overline{\mathbf{x}}$ we have

$$f(\mathbf{x}_k) = f(\overline{\mathbf{x}}) + Df(\overline{\mathbf{x}})(\mathbf{x}_k - \overline{\mathbf{x}}) + o(\|\mathbf{x}_k - \overline{\mathbf{x}}\|).$$

With $f(\mathbf{x}_k) \leqslant f(\overline{\mathbf{x}})$ it follows that

$$Df(\overline{\mathbf{x}})(\mathbf{x}_k - \overline{\mathbf{x}}) + o(\|\mathbf{x}_k - \overline{\mathbf{x}}\|) \leqslant 0.$$

Using again $\mathbf{x}_k = \overline{\mathbf{x}} + s_k \mathbf{d}_k$ we obtain

$$Df(\overline{\mathbf{x}})\mathbf{d}_k + \frac{o(s_k)}{s_k} \leqslant 0.$$

At the limit for $k \to \infty$ we obtain $Df(\overline{\mathbf{x}})\mathbf{d} \leqslant 0$. Taking into account that $g_i(\mathbf{x}_k) \leqslant g_i(\overline{\mathbf{x}}) = 0$ for $i \in I(\overline{\mathbf{x}})$ and proceeding as previously, we obtain $Dg_i(\overline{\mathbf{x}})\mathbf{d} \leqslant 0$, $i \in I(\overline{\mathbf{x}})$. Similarly we prove that $Dh_j(\overline{\mathbf{x}})\mathbf{d} = 0$ for $j = 1, \dots, l$.

We multiply the first equation of the necessary condition by $\mathbf{d}$:

$$Df(\overline{\mathbf{x}})\mathbf{d} + \sum_{i \in I(\overline{\mathbf{x}})} \mu_i Dg_i(\overline{\mathbf{x}})\mathbf{d} + \sum_{j=1}^l \lambda_j Dh_j(\overline{\mathbf{x}})\mathbf{d} = 0.$$

All terms of the sum above are less than or equal to $0$. As their sun is $0$, all of them must be zero. In particular,

$$Df(\overline{\mathbf{x}})\mathbf{d} = 0, \quad \text{and} \quad Dg_i(\overline{\mathbf{x}})\mathbf{d} = 0, \; i \in I^*(\overline{\mathbf{x}}).$$

Thus we proved that the vector $\mathbf{d}$ satisfies the conditions $(\otimes)$. □

A <u>general procedure</u> for optimization problems with mixed constraints is the following:

<u>Step 1.</u> We are looking for candidates for solutions, which form two sets:

$$A_1 = \{\mathbf{x} \in \mathbb{X} : \textit{no regularity condition} \text{ is satisfied at } \mathbf{x}\},$$
$$A_2 = \{\mathbf{x} \in \mathbb{X} : \text{a regularity condition}$$
$$\text{and the necessary first-order condition are satisfied at } \mathbf{x}\}.$$

<u>Step 2.</u> We check, whether the assumptions of Theorem 47, i.e., sufficient first-order conditions, are satisfied at each point of the sets $A_1$ and $A_2$. If they are, then we obtain global solutions.

The steps that follow should be done if we have not found any global solution <u>or</u> we need to find *all* global solutions <u>or</u> we need to find *all* local solutions. We remove all global solutions from the sets $A_1$ and $A_2$; the sets of remaining points we denote by $A_1'$ and $A_2'$.

<u>Step 3.</u> From $A_2'$ we exclude all points not satisfying the necessary condition of the second order, thus obtaining the set $A_2''$.

<u>Step 4.</u> At each point of $A_1' \cup A_2''$ we check the sufficient condition of the second order. The points satisfying it are local solutions.

<u>Step 5.</u> The analysis of the other points of $A_1' \cup A_2''$ must be done using other methods.

<u>Example.</u> Below is a detailed description of the analysis, step by step, of the following problem:

$$\begin{cases} (x_1 - 1)^2 + x_2^2 \to \min, \\ 2kx_1 - x_2^2 \leqslant 0, \\ \mathbf{x} = (x_1, x_2) \in \mathbb{R}^2, \end{cases}$$

where $k > 0$ is a parameter.

We notice that at each point, where the constraint is active, the linear independence condition is satisfied: the gradient of the constraint function is $[2k, 2x_2] \neq \mathbf{0}^{\mathsf{T}}$. Hence, $A_1 = \emptyset$.

The Lagrange function for this problem is

$$L(x_1, x_2; \mu) = (x_1 - 1)^2 + x_2^2 + \mu(2kx_1 - x_2^2).$$

The first-order necessary condition is thus

$$\begin{cases} [2(x_1 - 1), 2x_2] + \mu[2k, -2x_2] = [0, 0], \\ \mu(2kx_1 - x_2^2) = 0, \\ \mu \geqslant 0. \end{cases}$$

Is it possible that $\mu = 0$? By the first equation we obtain $2(x_1 - 1) = 0$, $2x_2 = 0$, i.e., $x_1 = 1$, $x_2 = 0$. This is not a feasible point for any parameter $k > 0$. Therefore $\mu$ must be positive and the necessary first-order condition may be satisfied only at the boundary of the feasible set, where $2kx_1 - x_2^2 = 0$.

Again, by the first equation,

$$\begin{cases} x_1 - 1 + \mu k = 0, \\ x_2 - \mu x_2 = 0. \end{cases}$$

There is either $\mu_2 = 1$ or $x_2 = 0$. If $x_2 = 0$, then due to the constraint we have $x_1 = 0$. The point $(0, 0)$ together with the Lagrange multiplier $\mu = 1/k$ satisfies the first-order necessary constraint.

Now we consider $\mu = 1$. Then, from $x_1 - 1 - \mu k = 0$ we obtain $x_1 = 1 - k$. If $k > 1$, then $x_1 < 0$ and no point with such a first coordinate is feasible. If $k = 1$, then we get the point $(0, 0)$ and for $k \in (0, 1)$ we have two points:

$$x_1 = 1 - k, \quad x_2 = \pm\sqrt{2k(1 - k)}.$$

We conclude that $A_1 = \emptyset$ and

$$A_2 = \{(0, 0)\} \quad \text{if } k \geqslant 1,$$
$$A_2 = \left\{ (0, 0), \left(1 - k, \sqrt{2k(1 - k)}\right), \left(1 - k, -\sqrt{2k(k - 1)}\right) \right\} \quad \text{if } 0 < k < 1.$$

The function $g_1(x_1, x_2) = 2kx_1 - x_2^2$ is not quasi-convex, so we cannot use Theorem 47 describing a sufficient first-order condition. Thus we have $A_2' = A_2$.

Now we pass to Step 3. We check the sufficient second-order condition for points of $A_2'$. We take a look at the point $(0, 0)$; with this point we have the Lagrange multiplier $\mu = 1/k$. The gradient and the Hessian of the Lagrange function are

$$D_x L(x_1, x_2; 1/k) = [2(x_1 - 1) + 2, 2x_2(1 - 1/k)]$$
$$D_x^2 L(x_1, x_2; 1/k) = \begin{bmatrix} 2 & 0 \\ 0 & 2(1 - 1/k) \end{bmatrix}$$

As the theorem states, it suffices to check if

$$\mathbf{d}^{\mathsf{T}} \begin{bmatrix} 2 & 0 \\ 0 & 2(1 - 1/k) \end{bmatrix} \mathbf{d} \geqslant 0$$

for all vectors $\mathbf{d} \in \mathbb{R}^2 \setminus \{\mathbf{0}\}$ such that $Dg_1(0, 0)\mathbf{d} = [2k, 0]\mathbf{d} = 0$, i.e., the vectors with the coordinate $d_1 = 0$. The last inequality is $(1 - 1/k)d_2^2 \geqslant 0$; it is satisfied if $k \geqslant 1$. Thus, if $k \geqslant 1$ it *may be* a local solution. If $k < 1$, then the inequality is not satisfied and there is no solution at $(0, 0)$.

Now we assume that $0 < k < 1$ and we consider the other two points. In both cases $\mu = 1$. Then we have

$$D_x L(x_1, x_2; 1/k) = [2(x_1 - 1) + 2k, 0]$$

$$D_x^2 L(x_1, x_2; 1/k) = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$

The Hessian is nonnegative-definite; hence, the necessary second-order condition is satisfied at both points. Thus, we have

$$A_2'' = \{(0,0)\} \quad \text{if } k \geqslant 1,$$
$$A_2'' = \left\{ \left(1 - k, \sqrt{2k(1-k)}\right), \left(1 - k, -\sqrt{2k(k-1)}\right) \right\} \quad \text{if } 0 < k < 1.$$

Now we check the sufficient second-order condition. If $k > 1$ the Hessian is positive-definite and by Theorem 56 there is a strict local minimum at $(0,0)$. We cannot make this conclusion if $k = 1$.

We assume $0 < k < 1$ and we analyse the point $\overline{x} = (1 - k, \sqrt{2k(k-1)})$. We need to check the assumptions of Theorem 56:

$$\mathbf{d}^{\mathsf{T}} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \mathbf{d} = 2d_1^2 > 0$$

for $\mathbf{d} \in \mathbb{R}^2 \setminus \{\mathbf{0}\}$ such that $[2k, -2\sqrt{2k(1-k)}]\mathbf{d} = 0$, i.e., $d_1 = \frac{1}{k}\sqrt{2k(1-k)}d_2$. If $\mathbf{d} \neq \mathbf{0}$, then $d_1 \neq 0$ and at the point $\overline{x}$ the sufficient condition of the second order is satisfied: this point is a strict local solution. Similarly we prove that the point $(1 - k, -\sqrt{2k(k-1)})$ is also a strict local solution.

It remains to analyse the case $k = 1$. Though the point $(0,0)$ is a global minimum, we have not managed to prove it using the Kuhn–Tucker theory. On the other hand, it is easy to do it in an elementary way.

# 10. Dual problems

Described below are so-called dual problems, i.e., another view at optimization problems with inequality constraints. In contrast to the Kuhn–Tucker approach, the dual setting does not assume differentiability of the target function $f$ nor of the constraint functions $g_i$. Moreover, with suitable assumptions, the solution of the original problem may easily be obtained from the solution of the dual problem. Solving the dual problem is finding the maximum of a concave target function in the nonnegative octant. As we shall see, the concavity may lead to a good convergence of numerical optimization algorithms. The simplicity of the feasible set is another feature contributing to the simplicity of implementation and to speed of numerical algorithms. Sometimes the dual problem is easier to solve analytically; examples will be given.

## Sufficient condition

<u>Definition 30</u> *Let* $A$, $B$ *be arbitrary sets and let* $h\colon A \times B \to \mathbb{R}$ *be a function. The point* $(\overline{x}, \overline{\mu}) \in A \times B$ *is called a* <u>saddle point</u> *of the function* $h$ *if*

$$h(\overline{x}, \mu) \leqslant h(\overline{x}, \overline{\mu}) \leqslant h(x, \overline{\mu}) \quad \text{for all } x \in A, \ \mu \in B.$$

<u>Example.</u> The simplest example is the "centre of the saddle", $A = B = \mathbb{R}$, $h(x, \mu) = x^2 - \mu^2$. The function $h$ has the saddle point at $(0, 0)$, which is the minimum with respect to $x$ and the maximum with respect to $\mu$.

It turns out that the global solution of the following problem:

$$\begin{cases} f(x) \to \min, \\ g_i(x) \leqslant 0, \ i = 1, \ldots, m, \\ x \in \mathbb{X}, \end{cases} \qquad (*)$$

is related to the saddle point of the Lagrange function defined for this problem. Recall that $W$ denotes the feasible set, i.e.,

$$W = \{x \in \mathbb{X}\colon g_i(x) \leqslant 0, i = 1, \ldots, m\}.$$

<u>Theorem 57</u> *If* $(\overline{x}, \overline{\mu}) \in W \times [0, \infty)^m$ *is a saddle point of the Lagrange function*

$$L(x, \mu) = f(x) + \sum_{i=0}^{m} \mu_i g_i(x)$$

*in the set* $W \times [0, \infty)^m$, *i.e.,*

$$L(\overline{x}, \mu) \leqslant L(\overline{x}, \overline{\mu}) \leqslant L(x, \overline{\mu}) \quad \text{for all } x \in W, \ \mu \in [0, \infty)^m,$$

*then* $\overline{x}$ *is a global solution of the problem (\*) and* $\overline{\mu}_i g_i(\overline{x}) = 0$ *for* $i = 1, \ldots, m$.

<u>Proof.</u> First we prove that $\overline{\mu}_i g_i(\overline{x}) = 0$ for $i = 1, \ldots, m$. The inequality $L(\overline{x}, \mu) \leqslant L(\overline{x}, \overline{\mu})$ may be expanded as follows:

$$f(\overline{x}) + \sum_{i=1}^{m} \mu_i g_i(\overline{x}) \leqslant f(\overline{x}) + \sum_{i=1}^{m} \overline{\mu}_i g_i(\overline{x}).$$

Hence, for all $\mu \in [0, \infty)^m$ we have

$$\sum_{i=1}^{m} \mu_i g_i(\overline{x}) \leqslant \sum_{i=1}^{m} \overline{\mu}_i g_i(\overline{x}).$$

By substituting $\mu = \frac{1}{2}\overline{\mu}$, we obtain

$$\sum_{i=1}^{m} \overline{\mu}_i g_i(\overline{x}) \geqslant 0.$$

The point $\overline{x}$ is feasible, i.e., $g_i(\overline{x}) \leqslant 0$ for all $i$. Taking into acount that all coordinates of $\overline{\mu}$ are nonnegative, we conclude that $\sum_{i=1}^{m} \overline{\mu}_i g_i(\overline{x}) = 0$ and each term is nonnegative. From the above it follows directly that $\overline{\mu}_i g_i(\overline{x}) = 0$ for all $i$.

Now we use the second inequality, $L(\overline{x}, \overline{\mu}) \leqslant L(x, \overline{\mu})$ for all $x \in W$, to prove that $\overline{x}$ is a global solution. We expand the inequality:

$$f(\overline{x}) + \sum_{i=1}^{m} \overline{\mu}_i g_i(\overline{x}) \leqslant f(x) + \sum_{i=1}^{m} \overline{\mu}_i g_i(x).$$

Earlier we have proved that $\overline{\mu}_i g_i(\overline{x}) = 0$ for all $i$. Due to $x \in W$, there is $\overline{\mu}_i g_i(x) \leqslant 0$. Hence,

$$f(\overline{x}) \leqslant f(x) \quad \text{for all } x \in W. \ \square$$

<u>Remarks.</u>

- We do not assume that $\mathbb{X}$ is open, also the functions $f$ and $g_i$ do not have to be continuous.

- The feasible set does not have to be convex.

- No regularity conditions are necessary.

- The theorem does not offer any methods of searching the saddle point. It might be found using the necessary first-order conditions, and then Theorem 57 may be used as the sufficient condition.

- Theorem 57 is the base for the dual approach and despite the previous remark it is useful in the development of numerical optimization algorithms.

## Necessary condition for convex programming

Now we assume that the set $\mathbb{X} \subset \mathbb{R}^n$ is convex and the functions $f$ and $g_i$, $i = 1, \ldots, m$, are convex. For such an optimization problem the saddle point of the Lagrange function is a necessary condition for the global solution. We begin the analysis with the simpler case, where all functions are differentiable, and later we prove the theorem which does not assume differentiability. As it was mentioned before, the lack of the requirement of differentiability distinguishes the saddle point method from the Kuhn–Tucker method.

<u>Lemma 20</u> *Suppose that the set $\mathbb{X}$ in the convex programming problem is open and the functions $f, g_1, \ldots, g_m$ are differentiable at a point $\overline{x}$. If $\overline{x}$ is a local solution of the problem (\*) and one of the regularity conditions: linear independence, affine function or Slater condition, is satisfied, then there exists $\overline{\mu} \in [0, \infty)^m$, such that $(\overline{x}, \overline{\mu})$ is a saddle point of the Lagrange function in the set $\mathbb{X} \times [0, \infty)^m$.*

<u>Proof.</u> By the Kuhn–Tucker theorem (Theorem 38), there exists a vector of Lagrange multipliers $\overline{\mu} \in [0, \infty)^m$ such that the first-order condition is satisfied (the assumptions of this theorem are satisfied due to regularity of $\overline{x}$ and Theorems 39–41). The function

$$L(x, \mu) = f(x) + \sum_{i=1}^{m} \mu_i g_i(x),$$

being the linear combination of convex functions with nonnegative coefficients, is convex. Therefore,

$$L(x, \overline{\mu}) \geqslant L(\overline{x}, \overline{\mu}) + D_x L(\overline{x}, \overline{\mu})(x - \overline{x}).$$

By the Kuhn–Tucker theorem,

$$D_x L(\overline{x}, \overline{\mu}) = Df(\overline{x}) + \sum_{i=1}^{m} \overline{\mu}_i Dg_i(\overline{x}) = 0^\mathsf{T},$$

i.e., $D_x L(\overline{x}, \overline{\mu})(x - \overline{x}) = 0$. Hence, $L(x, \overline{\mu}) \geqslant L(\overline{x}, \overline{\mu})$.

To prove the inequality $L(\overline{x}, \overline{\mu}) \geqslant L(\overline{x}, \mu)$, we notice that

$$\sum_{i=1}^{m} \mu_i g_i(\overline{x}) \leqslant 0 = \sum_{i=1}^{m} \overline{\mu}_i g_i(\overline{x}),$$

because $\mu_i \geqslant 0$ and $g_i(\overline{x}) \leqslant 0$. The last equality is the claim of Theorem 38. □

<u>Remark.</u> By Theorem 47, each point satisfying the first-order condition is a global solution of the convex programming problem. Therefore, we do not need to distinguish global and local solutions.

<u>Theorem 58</u> *Let $\overline{x} \in \mathbb{X}$ be a global solution of the convex programming problem (\*) and let there exist a point $x^* \in \mathbb{X}$ such that $g_i(x^*) < 0$ for all $i = 1, \ldots, m$. Then, there exists $\overline{\mu} \in [0, \infty)^m$ such that $(\overline{x}, \overline{\mu})$ is a saddle point of the Lagrange function in the space $\mathbb{X} \times [0, \infty)^m$, i.e.,*

$$L(\overline{x}, \mu) \leqslant L(\overline{x}, \overline{\mu}) \leqslant L(x, \overline{\mu}) \quad \text{for all } x \in \mathbb{X} \text{ and } \mu \in [0, \infty)^m.$$

*Moreover, $\overline{\mu}_i g_i(\overline{x}) = 0$ for $i = 1, \ldots, m$.*

<u>Remark.</u> In Theorems 57 and 58 the saddle point of the Lagrange function is considered in different spaces. In the second theorem the space is wider, as the first variable goes through the entire set $\mathbb{X}$, not just the feasible set $W$. Now we obtain the equivalence of the existence of a solution at the saddle point of the Lagrange function and the existence of the global solution of the convex programming problem.

<u>Proof.</u> Like in the proof of the necessary first-order condition (Theorem 38), the crucial role is played by the separation theorem. It will indicate the Lagrange multiplier vector $\overline{\mu}$.

Denote $g(x) = \big(g_1(x), \ldots, g_m(x)\big)$. We define the following subsets of $\mathbb{R}^{m+1}$:

$$A = \{\overline{y} = (y_0, y) \in \mathbb{R} \times \mathbb{R}^n \colon y_0 \geqslant f(x), \ y \geqslant g(x), \ x \in \mathbb{X}\},$$
$$B = \{\overline{y} = (y_0, y) \in \mathbb{R} \times \mathbb{R}^n \colon y_0 = f(x), \ y = g(x), \ x \in \mathbb{X}\},$$
$$C = \{\overline{y} = (y_0, y) \in \mathbb{R} \times \mathbb{R}^n \colon y_0 < f(\overline{x}), \ y < 0\}.$$

The "inequality between vectors" notation is to be understood as the inequality between their corresponding coordinates.

The set C is the Cartesian product of the interval $(-\infty, f(\overline{x})]$ and the cone $\{y < 0\}$. Obviously, this set is convex. As $\overline{x}$ is the minimum of f, there is $B \cap C = \emptyset$. The set B is not convex; therefore we cannot use the separation theorem for the sets B and C. A remedy is to take the set A, whose subset is B. Suppose that there exists $\overline{y} = (y_0, y) \in A \cap C$. It follows that there exists $x' \in \mathbb{X}$ such that

$$y_0 \geqslant f(x'), \ y \geqslant g(x'), \ y_0 < f(\overline{x}), \ y < 0.$$

From the inequalities above we conclude that $f(x') < f(\overline{x})$ and $g(x') < 0$. Thus, $x'$ is a feasible point at which the value of f is less than $f(\overline{x})$. It contradicts $\overline{x}$ being a solution.

<u>Example.</u> Before proceeding with the proof, we can take a look at the sets A, B, C for the following problem:

$$\begin{cases} -x \to \min, \\ x^2 - 1 \leqslant 0, \\ x \in \mathbb{X} = \mathbb{R}. \end{cases}$$

Its solution is $\overline{x} = 1$. There is only one constraint, and thus the sets are subsets of $\mathbb{R}^2$. As we can see in Figure 7, the part of the set B for $y_0 \leqslant 0$ is a part of the boundary of A, and the rest of it (for $y_0 > 0$) is located in the interior of A.



Figure 7: The sets A, B, C in the example

Back to the proof. The convexity of C is already established. The convexity of A may be proved directly. Let $\overline{y}', \overline{y}'' \in A$ and $\lambda \in (0, 1)$. Then, there exist points $x', x'' \in \mathbb{X}$ such that

$$y_0' \geqslant f(x'), \ y' \geqslant g(x'), \qquad y_0'' \geqslant f(x''), \ y'' \geqslant g(x'').$$

Let $x = \lambda x' + (1 - \lambda)x''$. By convexity of $\mathbb{X}$, $x \in \mathbb{X}$. We also have

$$\lambda y_0' + (1 - \lambda)y_0'' \geqslant \lambda f(x') + (1 - \lambda)f(x'') \geqslant f(\lambda x' + (1 - \lambda)x'') = f(x).$$

The first inequality above results from the assumed properties of $x'$, $x''$, and the second inequality is a consequence of the convexity of f. In a similar way, using the convexity of the components of g, we show that

$$\lambda y' + (1 - \lambda)y'' \geqslant g(x).$$

Hence, $\overline{y} = \lambda \overline{y}' + (1 - \lambda)\overline{y}'' \in A$, as with the point $x$ defined above,

$$y_0 \geqslant f(x), \ y \geqslant g(x).$$

By the weak separation theorem, there exists a nonzero vector $\tilde{\mu} \in \mathbb{R}^{m+1}$ such that

$$\tilde{\mu}^\top \overline{y} \geqslant \tilde{\mu}^\top \overline{z}, \quad \text{for all } \overline{y} \in A, \ \overline{z} \in C.$$

From $\sup_{\overline{z} \in C} \tilde{\mu}^\top \overline{z} < \infty$ it follows that $\tilde{\mu} \geqslant 0$. Due to the continuity of linear functions, we can take $\overline{z}$ from the closure of C:

$$\tilde{\mu}^\top \overline{y} \geqslant \tilde{\mu}^\top \overline{z}, \quad \text{for all } \overline{y} \in A, \ \overline{z} \in \text{cl } C.$$

Hence, for $\overline{z} = (f(\overline{x}), 0)$ we have

$$\tilde{\mu}_0 y_0 + \sum_{i=1}^m \tilde{\mu}_i y_i \geqslant \tilde{\mu}_0 f(\overline{x}) \quad \text{for all } (y_0, y) \in A. \tag{$\otimes$}$$

In particular, this inequality holds for $y_0 = f(x)$ and $y = g(x)$, where $x \in \mathbb{X}$:

$$\tilde{\mu}_0 f(x) + \sum_{i=1}^m \tilde{\mu}_i g_i(x) \geqslant \tilde{\mu}_0 f(\overline{x}).$$

Now we prove that $\tilde{\mu}_0 \neq 0$, which together with the observation $\tilde{\mu} \geqslant 0$ implies $\tilde{\mu}_0 > 0$. The proof is done by contradiction. Suppose that $\tilde{\mu}_0 = 0$. Then, by the last inequality, we have

$$\sum_{i=1}^m \tilde{\mu}_i g_i(x) \geqslant 0 \quad \text{for all } (y_0, y) \in A.$$

In particular, the above holds for the point $x^*$ from the assumptions of the theorem. However, at this point we have $g_i(x^*) < 0$ for $i = 1, \dots, m$. Together with the fact that $\tilde{\mu} \geqslant 0$ this implies $\tilde{\mu}_1 = \cdots = \tilde{\mu}_m = 0$. It follows that $\tilde{\mu} = 0$, and this is inconsistent with the choice of $\tilde{\mu}$ based on the separation theorem.

As we now know, $\tilde{\mu}_0 > 0$. We define

$$\overline{\mu} = \left(\frac{\tilde{\mu}_1}{\tilde{\mu}_0}, \ldots, \frac{\tilde{\mu}_m}{\tilde{\mu}_0}\right).$$

Obviously, $\overline{\mu} \in [0, \infty)^m$. As $\overline{x}$, being a solution, is a feasible point, there is $g_i(\overline{x}) \leqslant 0$ for $i = 1, \ldots, m$ and $\sum_{i=1}^m \overline{\mu}_i g_i(\overline{x}) \leqslant 0$. We add this sum to the right-hand side of the inequality $(\otimes)$ divided by $\tilde{\mu}_0$:

$$f(x) + \sum_{i=1}^m \overline{\mu}_i g_i(x) \geqslant f(\overline{x}) + \sum_{i=1}^m \overline{\mu}_i g_i(\overline{x}) \quad \text{for all } x \in \mathbb{X}.$$

In other words,

$$L(x, \overline{\mu}) \geqslant L(\overline{x}, \overline{\mu}) \quad \text{for all } x \in \mathbb{X}.$$

It remains to prove the other inequality of the saddle point. By taking $x = \overline{x}$ and dividing both sides of $(\otimes)$ by $\tilde{\mu}_0$, we obtain $\sum_{i=1}^m \overline{\mu}_i g_i(\overline{x}) \geqslant 0$. On the other hand, the point $\overline{x}$ is feasible, i.e., $g_i(\overline{x}) \leqslant 0$. Recalling that $\overline{\mu} \geqslant 0$, we conclude that each term of this sum is nonpositive. Hence, we obtain

$$\overline{\mu}_i g_i(\overline{x}) = 0, \quad i = 1, \ldots, m.$$

For any other $\mu$ we have $\sum_{i=1}^m \mu_i g_i(\overline{x}) = 0$, i.e.,

$$\sum_{i=1}^m \mu_i g_i(\overline{x}) \leqslant \sum_{i=1}^m \overline{\mu}_i g_i(\overline{x}) \quad \text{for all } \mu \in [0, \infty)^m.$$

This is equivalent to

$$L(\overline{x}, \mu) \leqslant L(\overline{x}, \overline{\mu}) \quad \text{for all } \mu \in [0, \infty)^m. \quad \square$$

## Primal and dual problems

The theory of saddle points leads to formulating primal and dual problems. Consider the optimization problem $(*)$ and its Lagrange function $L(x, \mu)$. We define the function $L_P \colon W \to (-\infty, \infty]$:

$$L_P(x) = \sup_{\mu \in [0, \infty)^m} L(x, \mu).$$

As we can notice,

$$L_P(x) = \begin{cases} f(x) & \text{if } g(x) \leqslant 0, \\ \infty & \text{else.} \end{cases}$$

Therefore the problem $(*)$ may be rewritten in what seems to be a simpler form

$$L_P(x) \to \min, \quad x \in \mathbb{X}.$$

Alas, the problem above reduces to the original problem, so it does not give any "added value", but soon it will. Before we reveal this value, we define another function, $L_D(\mu) \colon [0, \infty)^m \to [-\infty, \infty)$:

$$L_D(\mu) = \inf_{x \in \mathbb{X}} L(x, \mu).$$

Remarks. (1) for any $x \in \mathbb{X}$ and $\mu \in [0, \infty)^m$ there is $L_P(x) \geqslant L(x, \mu) \geqslant L_D(\mu)$.
(2) If $(\overline{x}, \overline{\mu})$ is a saddle point of the Lagrange function in $\mathbb{X} \times [0, \infty)^m$, then $L_P(\overline{x}) = L_D(\overline{\mu})$.

If $(\overline{x}, \overline{\mu})$ is a saddle point, then $L(\overline{x}, \overline{\mu}) \leqslant L(x, \overline{\mu})$. In view of the first remark above it gives us $L(\overline{x}, \overline{\mu}) = L_D(\overline{\mu})$. Similarly we prove the equality $L(\overline{x}, \overline{\mu}) = L_P(\overline{x})$. These observations lead us in the right direction. We are going to use the functions $L_P$ and $L_D$ in the search of saddle points.

Definition 31 *The underline{primal problem} is the optimization problem*

$$L_P(x) \to \min, \quad x \in \mathbb{X}.$$

*Its underline{dual problem} is the optimization problem*

$$L_D(\mu) \to \max, \quad \mu \in [0, \infty)^m.$$

By the properties mentioned in the last remark, the value of the target function of the primal problem at the solution is greater than or equal to the value of the target function of the dual problem at the solution:

$$\inf_{x \in \mathbb{X}} L_P(x) \geqslant \sup_{\mu \in [0, \infty)^m} L_D(\mu).$$

Moreover, the solution of the dual problem gives us an estimation from below of the function $f$.

Lemma 21 *(underline{weak duality theorem}) For any feasible point $x \in W$ and any vector $\mu \in [0, \infty)^m$ there is*

$$f(x) \geqslant L_D(\mu).$$

*Hence,*

$$f(x) \geqslant \sup_{\mu \in [0, \infty)^m} L_D(\mu).$$

Proof. We have

$$f(\boldsymbol{x}) \geqslant L(\boldsymbol{x}, \boldsymbol{\mu}) \geqslant L_D(\boldsymbol{\mu}).$$

The first inequality holds for $\boldsymbol{x} \in W$ because $g_i(\boldsymbol{x}) \leqslant 0$ for all $i$. The second inequality is a consequence of Remark (1) above. $\square$

Definition 32 *The underline{duality gap} is the difference between the values of target functions at the solutions of the primal and dual problem:*

$$\inf_{\boldsymbol{x} \in \mathbb{X}} L_P(\boldsymbol{x}) - \sup_{\boldsymbol{\mu} \in [0, \infty)^m} L_D(\boldsymbol{\mu}).$$

The saddle point condition written in terms of the primal and dual functions is the following: $(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\mu}})$ is a saddle point if

$$L_P(\overline{\boldsymbol{x}}) = L(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\mu}}) = L_D(\overline{\boldsymbol{\mu}}).$$

In other words, if the Lagrange function has a saddle point, then the duality gap is zero. This is the case if, for example, the assumptions of Theorem 58 are satisfied.



Figure 8: A scheme of solutions of a primal and dual problem

Figure 8 shows the solution of a primal and dual problem of an optimization problem with only one inequality constraint; $m = 1$. The set $G$ is the set of pairs

of values $\big(f(\boldsymbol{x}), g(\boldsymbol{x})\big)$ for $\boldsymbol{x} \in \mathbb{X}$. The lines $z + \mu y = \alpha$ show the values of the Lagrange function $L(\boldsymbol{x}, \mu) = f(\boldsymbol{x}) + \mu g(\boldsymbol{x})$. From the definition of the primal function $L_P$ it follows that $L_P(\boldsymbol{x}) = f(\boldsymbol{x})$ if $g(\boldsymbol{x}) \leqslant 0$, as $\sup_{\mu \geqslant 0}\{\, \alpha \colon \alpha = z + \mu y \,\}$ is taken for $\mu = 0$. This is the case of the point $\boldsymbol{x}_2$. If $g(\boldsymbol{x}) > 0$, then $L_P(\boldsymbol{x}) = +\infty$. This is the case of the point $\boldsymbol{x}_1$. The dual function $L_D(\mu)$ may be found by considering the lines $z + \mu y = \alpha$, looking for $\inf_{\boldsymbol{x} \in \mathbb{X}} \alpha$ with a fixed $\mu$. It may be seen in the picture that the minimum is taken for the line tangent to the boundary of the area $G$ and the value of this function is the intersection point of the line with the line $y = 0$. One can also see that the line which yields the greatest value of $L_D$ is the line $z + \overline{\mu} y = \overline{\alpha}$ whose inclination (tangent of the angle between the $y$ axis and the line) is $-\overline{\mu}$, tangent to the boundary of $G$ at the point $(\overline{y}, \overline{z})$. This point coresponds to the solution of the primal problem, because $\overline{z} = \inf_{\boldsymbol{x} \in \mathbb{X}, g(\boldsymbol{x}) \leqslant 0} f(\boldsymbol{x})$.

Theorem 59 *(underline{strong duality theorem}) Let $\mathbb{X}$ be a nonempty convex subset of $\mathbb{R}^n$ and let the functions $f$ and $g_i$, $i = 1, \ldots, m$, be convex in $\mathbb{X}$. Assume in addition that there exists a point $\boldsymbol{x}^* \in \mathbb{X}$ such that $g_i(\boldsymbol{x}^*) < 0$ for all $i$. Then,*

$$\inf_{\boldsymbol{x} \in \mathbb{X}} L_P(\boldsymbol{x}) = \sup_{\boldsymbol{\mu} \in [0, \infty)^m} L_D(\boldsymbol{\mu}).$$

*If $\inf_{\boldsymbol{x} \in \mathbb{X}} L_P(\boldsymbol{x})$ is finite, then $\sup_{\boldsymbol{\mu} \in [0, \infty)^m} L_D(\boldsymbol{\mu})$ is taken at a point $\overline{\boldsymbol{\mu}}$ such that $\overline{\boldsymbol{\mu}} \geqslant \boldsymbol{0}$. If $\inf_{\boldsymbol{x} \in \mathbb{X}} L_P(\boldsymbol{x})$ is taken at a point $\overline{\boldsymbol{x}}$, then $\overline{\mu}_i g_i(\overline{\boldsymbol{x}}) = 0$ for $i = 1, \ldots, m$.*

Proof. Let $\gamma = \inf_{\boldsymbol{x} \in \mathbb{X}} L_P(\boldsymbol{x})$. If $\gamma = -\infty$, then (by Lemma 21) $\sup_{\boldsymbol{\mu} \in [0, \infty)^m} L_D(\boldsymbol{\mu}) = -\infty$ and the claim is true. Now assume that $\gamma > -\infty$.

From the proof of Theorem 58 applied to the sets

$$A = \{\, \overline{\boldsymbol{y}} = (y_0, \boldsymbol{y}) \in \mathbb{R} \times \mathbb{R}^n \colon y_0 \geqslant f(\boldsymbol{x}),\ \boldsymbol{y} \geqslant g(\boldsymbol{x}),\ \boldsymbol{x} \in \mathbb{X} \,\},$$
$$C = \{\, \overline{\boldsymbol{y}} = (y_0, \boldsymbol{y}) \in \mathbb{R} \times \mathbb{R}^n \colon y_0 < f(\overline{\boldsymbol{x}}),\ \boldsymbol{y} < \boldsymbol{0} \,\},$$

we conclude that

$$f(\boldsymbol{x}) + \sum_{i=1}^m \overline{\mu}_i g_i(\boldsymbol{x}) \geqslant \gamma \quad \text{for all } \boldsymbol{x} \in \mathbb{X}. \tag{$\oplus$}$$

As a consequence, we obtain

$$L_D(\overline{\boldsymbol{\mu}}) = \inf_{\boldsymbol{x} \in \mathbb{X}} \left( f(\boldsymbol{x}) + \sum_{i=1}^m \overline{\mu}_i g_i(\boldsymbol{x}) \right) \geqslant \gamma.$$

From the remark made before Definition 31 it follows that

$$\gamma = \inf_{x \in \mathbb{X}} L_P(x) \geqslant \sup_{\mu \in [0,\infty)^m} L_D(\mu) = L_D(\overline{\mu}) \geqslant \gamma.$$

Hence, $L_D(\overline{\mu}) = \gamma$ and $\overline{\mu}$ is a solution of the dual problem.

If $\inf_{x \in \mathbb{X}} L_P(x)$ is taken at a point $\overline{x}$, then due to the definition of $L_P$ we have $L_P(\overline{x}) = f(\overline{x})$. The point $\overline{x}$ is a solution of the primal problem; there is $g_i(\overline{x}) \leqslant 0$ for $i = 1, \ldots, m$ and $f(\overline{x}) = \gamma$. By substituting $x = \overline{x}$ in the inequality ($\oplus$), we get $\sum_i \overline{\mu}_i g_i(\overline{x}) \geqslant 0$. Because $\mu_i \geqslant 0$ and $g_i(\overline{x}) \leqslant 0$, the equality $\overline{\mu}_i g_i(\overline{x}) = 0$ for $i = 1, \ldots, m$ follows. $\square$

The following algorithm of solving the optimization problem (*) using dual methods may be proposed:

1. Solve the dual problem. The solution gives a lower bound for the solution of the primal problem, due to Lemma 21.

2. Suppose that there exists a finite solution $\overline{\mu} \in [0, \infty)^m$ of the dual problem and a point $\overline{x} \in \mathbb{X}$ such that $L_D(\overline{\mu}) = L_P(\overline{x})$. If $\overline{x}$ is a feasible point and $f(\overline{x}) = L_D(\overline{\mu})$, then $(\overline{x}, \overline{\mu})$ is a saddle point of the Lagrange function and by Theorem 57 the point $\overline{x}$ is a solution of the problem (*).

Let's explain the conditions in the second step. From $L_D(\overline{\mu}) = L(\overline{x}, \overline{\mu})$ it follows that $L(\overline{x}, \overline{\mu}) \leqslant L(x, \overline{\mu})$ for all $x \in \mathbb{X}$. This is the second inequality of the saddle point. To verify the first inequality, recall that $L_P(x) = f(x)$ for any feasible point $x$ and $\inf_{x \in \mathbb{X}} L_P(x) \geqslant L_D(\mu)$ for all $\mu \in [0, \infty)^m$. In the second step we assume that $f(\overline{x}) = L_D(\overline{\mu})$, which implies that

$$L_P(\overline{x}) = f(\overline{x}) = L_D(\overline{\mu})$$

and, thus, $(\overline{x}, \overline{\mu})$ is indeed a saddle point.

# 11. Sensitivity theory

So far, the Lagrange multipliers seemed to be just a technical trick useful in finding solutions of optimization problems with constraints. Below we show that they represent costs of changing the constraints. The equality and inequality constraints are dealt with separately.

## Equality constraints

Consider a problem with equality constraints:

$$\begin{cases} f(\boldsymbol{x}) \to \min, \\ h_j(\boldsymbol{x}) = 0, \ j = 1, \ldots, l, \\ \boldsymbol{x} \in \mathbb{X}, \end{cases} \tag{*}$$

where $\mathbb{X} \subset \mathbb{R}^n$ is an open set and $f, h_1, \ldots, h_l \colon \mathbb{X} \to \mathbb{R}$. To simplify the notation, we denote $\boldsymbol{h}(\boldsymbol{x}) = \big(h_1(\boldsymbol{x}), \ldots, h_l(\boldsymbol{x})\big)$. Now we introduce a perturbation:

$$\begin{cases} f(\boldsymbol{x}) \to \min, \\ \boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{z}, \\ \boldsymbol{x} \in \mathbb{X}, \end{cases} \tag{**}$$

where $\boldsymbol{z} \in \mathbb{R}^l$.

__Theorem 60__ *Let $\overline{\boldsymbol{x}}$ be a solution of the problem (\*) and let $\overline{\boldsymbol{\lambda}}$ be the corresponding vector of Lagrange multipliers. Assume that the functions $f, h_1, \ldots, h_l$ are of class $C^2$ in a neighbourhood of $\overline{\boldsymbol{x}}$, the gradients of the constraint functions are linearly independent and*

$$\boldsymbol{d}^\mathsf{T} D_x^2 L(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) \boldsymbol{d} > 0 \tag{$\odot$}$$

*for all nonzero vectors $\boldsymbol{d} \in \mathbb{R}^n$ such that $Dh_j(\overline{\boldsymbol{x}})\boldsymbol{d} = 0$, $j = 1, \ldots, l$. Then, there exists a neighbourhood $\tilde{O}$ of the point $\boldsymbol{0} \in \mathbb{R}^l$ and a function $\boldsymbol{x} \colon \tilde{O} \to \mathbb{X}$ of class $C^1$ such that $\boldsymbol{x}(\boldsymbol{0}) = \overline{\boldsymbol{x}}$ and $\boldsymbol{x}(\boldsymbol{z})$ is a strict local solution of the modified problem (\*\*). Moreover,*

$$D(f \circ \boldsymbol{x})(\boldsymbol{0}) = -\overline{\boldsymbol{\lambda}}^\mathsf{T}.$$

__Proof.__ By Theorem 48, the point $\overline{\boldsymbol{x}}$ is a solution of the system of equations

$$\begin{cases} D_x L(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) = \boldsymbol{0}^\mathsf{T}, \\ \boldsymbol{h}(\overline{\boldsymbol{x}}) = \boldsymbol{0}, \end{cases}$$

where

$$D_x L(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) = Df(\overline{\boldsymbol{x}}) + \overline{\boldsymbol{\lambda}}^\mathsf{T} D\boldsymbol{h}(\overline{\boldsymbol{x}}) = Df(\overline{\boldsymbol{x}}) + \sum_{j=1}^{l} \overline{\lambda}_j Dh_j(\overline{\boldsymbol{x}}).$$

After adding the perturbation $\boldsymbol{z}$ to the right-hand side of the second equality we are going to show that there exists a solution being a function of class $C^1$ of the perturbation. Consider the system

$$\begin{cases} D_x L(\boldsymbol{x}, \boldsymbol{\lambda}) = \boldsymbol{0}^\mathsf{T}, \\ \boldsymbol{h}(\boldsymbol{x}) = \boldsymbol{z}, \end{cases}$$

with unknown $\boldsymbol{\lambda}$ and $\boldsymbol{x}$. We define the function $\mathbf{G} \colon \mathbb{R}^n \times \mathbb{R}^l \times \mathbb{R}^l \to \mathbb{R}^n \times \mathbb{R}^l$ by the formula

$$\mathbf{G}(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{z}) = \left[ \begin{array}{c} \big(D_x L(\boldsymbol{x}, \boldsymbol{\lambda})\big)^\mathsf{T} \\ \boldsymbol{h}(\boldsymbol{x}) - \boldsymbol{z} \end{array} \right].$$

The modified system may be rewritten as $\mathbf{G}(\boldsymbol{x}, \boldsymbol{\lambda}, \boldsymbol{z}) = \boldsymbol{0}$.

We know that $\mathbf{G}(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}, \boldsymbol{0}) = \boldsymbol{0}$. Using the implicit function theorem, we weave the first two variables as functions of the third one. To do that, we consider the matrix of derivatives of $\mathbf{G}$ (the blocks $\boldsymbol{0}$ and $-I$ are $l \times l$):

$$D\mathbf{G}(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}, \boldsymbol{0}) = \left[ \begin{array}{ccc} D^2 L(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) & \big(D\boldsymbol{h}(\overline{\boldsymbol{x}})\big)^\mathsf{T} & 0 \\ D\boldsymbol{h}(\overline{\boldsymbol{x}}) & 0 & -I \end{array} \right].$$

The linear independence of the gradients of the constraints implies that the submatrix

$$D\mathbf{G}(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}, \boldsymbol{0}) = \left[ \begin{array}{cc} D^2 L(\overline{\boldsymbol{x}}, \overline{\boldsymbol{\lambda}}) & \big(D\boldsymbol{h}(\overline{\boldsymbol{x}})\big)^\mathsf{T} \\ D\boldsymbol{h}(\overline{\boldsymbol{x}}) & 0 \end{array} \right].$$

is nonsingular (proof is an exercise). The assumptions of the implicit function theorem (Theorem 51) are, therefore, satisfied and there exists a neighbourhood $O$ of the point $\boldsymbol{0} \in \mathbb{R}^l$ and functions $\boldsymbol{x} \colon O \to \mathbb{X}$ and $\boldsymbol{\lambda} \colon O \to \mathbb{R}^l$ of class $C^1$ such that for all $\boldsymbol{z} \in O$ there is $\mathbf{G}\big(\boldsymbol{x}(\boldsymbol{z}), \boldsymbol{\lambda}(\boldsymbol{z})\big) = \boldsymbol{0}$, i.e.,

$$D_x L\big(\boldsymbol{x}(\boldsymbol{z}), \boldsymbol{\lambda}(\boldsymbol{z})\big) = \boldsymbol{0}^\mathsf{T}, \quad \boldsymbol{h}\big(\boldsymbol{x}(\boldsymbol{z})\big) = \boldsymbol{z}.$$

Using the fact that the functions $D_x^2 L$, $D\boldsymbol{h}$, $\boldsymbol{x}$, $\boldsymbol{\lambda}$ are continuous and the inequality $(\odot)$ is satisfied for the original problem, we conclude that there exists a (possibly smaller) neighbourhood $\tilde{O}$ of $\boldsymbol{0} \in \mathbb{R}^l$ such that

$$\boldsymbol{d}^\mathsf{T} D_x^2 L\big(\boldsymbol{x}(\boldsymbol{z}), \boldsymbol{\lambda}(\boldsymbol{z})\big) \boldsymbol{d} > 0$$

for all $z \in \tilde{O}$ and nonzero vectors $\mathbf{d} \in \mathbb{R}^n$ such that $Dh_j(\mathbf{x}(z))\mathbf{d} = 0$, $j = 1, \ldots, l$. The key to this result is that $(\odot)$ is a sharp inequality. By Theorem 56 the point $\mathbf{x}(z)$ is therefore a strict solution of the modified problem (**). Recall that $\mathbf{x}$ is a function of class $C^1$. Hence, we can define the derivative of the composition

$$D(f \circ \mathbf{x})(\mathbf{0}) = Df(\overline{\mathbf{x}})D\mathbf{x}(\mathbf{0}).$$

To complete the proof we need two observations. First, after multiplying the sides of the necessary condition for the original problem,

$$Df(\overline{\mathbf{x}}) + \overline{\boldsymbol{\lambda}}^{\mathsf{T}}D\mathbf{h}(\overline{\mathbf{x}}) = \mathbf{0}^{\mathsf{T}}$$

by $D\mathbf{x}(\mathbf{0})$ we obtain

$$Df(\overline{\mathbf{x}})D\mathbf{x}(\mathbf{0}) + \overline{\boldsymbol{\lambda}}^{\mathsf{T}}D\mathbf{h}(\overline{\mathbf{x}})D\mathbf{x}(\mathbf{0}) = \mathbf{0}^{\mathsf{T}}.$$

Second, by differentiating $\mathbf{h}(\mathbf{x}(z))$ with respect to $z$, at $z = \mathbf{0}$ we obtain the following derivative: $D(\mathbf{h} \circ z)(\mathbf{0}) = D\mathbf{h}(\overline{\mathbf{x}})D\mathbf{x}(\mathbf{0}) = I$. The equation above may, therefore, be simplified to the following:

$$Df(\overline{\mathbf{x}})D\mathbf{x}(\mathbf{0}) + \overline{\boldsymbol{\lambda}}^{\mathsf{T}} = \mathbf{0}^{\mathsf{T}}.$$

The claim follows immediately. $\square$

Theorem 60 may be understood as follows: a small change of the j-th constraint from $0$ to $\varepsilon$ causes the change of the local minimum of f by $-\overline{\lambda}_j\varepsilon + o(\varepsilon) \approx -\overline{\lambda}_j\varepsilon$.

## Inequality constraints

We use a different approach for inequality constraints. We focus on a convex optimization problem:

$$\begin{cases} f(\mathbf{x}) \to \min, \\ g_i(\mathbf{x}) \leqslant 0, \quad i = 1, \ldots, m, \\ \mathbf{x} \in \mathbb{X}, \end{cases} \qquad \overset{(**)}{\phantom{.}}$$

where $\mathbb{X} \subset \mathbb{R}^n$ is a convex set and the functions $f, g_1, \ldots, g_n \colon \mathbb{X} \to \mathbb{R}$ are convex. To simplify the notation we denote $\mathbf{g}(\mathbf{x}) = (g_1(\mathbf{x}), \ldots, g_m(\mathbf{x}))$. The problem may be rewritten as

$$\begin{cases} f(\mathbf{x}) \to \min, \\ \mathbf{g}(\mathbf{x}) \leqslant \mathbf{0}, \\ \mathbf{x} \in \mathbb{X}, \end{cases}$$

Consider the modified problem

$$\begin{cases} f(\mathbf{x}) \to \min, \\ \mathbf{g}(\mathbf{x}) \leqslant \mathbf{z}, \\ \mathbf{x} \in \mathbb{X}, \end{cases}$$

<u>Definition 33</u> *Let $D_M$ denote the set of vectors $\mathbf{z} \in \mathbb{R}^m$ such that the feasible set for the modified problem, $W_z = \{\mathbf{x} \in \mathbb{X} \colon \mathbf{g}(z) \leqslant \mathbf{z}\}$, is nonempty. The function*

$$M(\mathbf{z}) = \inf_{\mathbf{x} \in \mathbb{X}, \, \mathbf{g}(\mathbf{x}) \leqslant \mathbf{z}} f(\mathbf{x}),$$

*defined for all $\mathbf{z} \in D_M$, is called the <u>perturbation function</u> and the set $D_M$ is called the <u>perturbation function domain</u>.*

Note that $M(\mathbf{z}) < \infty$ for all $\mathbf{z} \in D_M$, but it is possible that $M(\mathbf{z}) = -\infty$.

The graph of the perturbation function in Figure 8 is the curve $M(y)$; as we can see it is a convex function. Note that this function is well defined between the points $A = (y_A, z_A)$ and $B = (y_B, z_B)$, because there exists $\mathbf{x} \in \mathbb{X}$ such that $(g(\mathbf{x}), f(\mathbf{x})) \in G$, and then $y = g(\mathbf{x})$ is an element of the domain of the function $M(y)$. If $y < y_A$, then the feasible set is empty and such points $y$ are outside the domain of the perturbation function. For $y > y_B$ the perturbation function is a constant.

<u>Theorem 61</u>     *1. The set $D_M$ is convex.*

   *2. The function $M \colon D_M \to \mathbb{R} \cup \{-\infty\}$ is convex.*

   *3. If there exists a point $\mathbf{x}^* \in \mathbb{X}$ such that $\mathbf{g}(\mathbf{x}^*) < \mathbf{0}$, then $\operatorname{int} D_M \neq \emptyset$ and $\mathbf{0} \in \operatorname{int} D_M$.*

<u>Proof.</u> From the convexity of each component of $\mathbf{g}$ it follows that

$$\begin{aligned} \mathbf{g}(\mathbf{x}_1) \leqslant \mathbf{z}_1, \ \mathbf{g}(\mathbf{x}_2) \leqslant \mathbf{z}_2 \ &\Rightarrow \\ \mathbf{g}(\lambda \mathbf{x}_1 + (1-\lambda)\mathbf{x}_2) &\leqslant \lambda \mathbf{z}_1 + (1-\lambda)\mathbf{z}_2 \quad \text{for all } \lambda \in [0, 1] \end{aligned} \qquad (\otimes)$$

(the inequalities are to be understood componentwise). This observation will be used a number of times in the proof below.

(1) Let $z_1, z_2 \in D_M$ and let $\lambda \in (0, 1)$. There exist $x_1, x_2 \in \mathbb{X}$ such that $g(x_1) \leqslant z_1$ and $g(x_2) \leqslant z_2$. From $(\otimes)$ we obtain $g(\lambda x_1 + (1 - \lambda) x_2) \leqslant \lambda z_1 + (1 - \lambda) z_2$; hence, $\lambda z_1 + (1 - \lambda) z_2 \in D_M$.

(2) Let $z_1, z_2 \in D_M$ and let $\lambda \in (0, 1)$. Then,

$$
\begin{aligned}
\lambda M(z_1) + (1 - \lambda) M(z_2) &= \inf_{x_1 \in \mathbb{X}, g(x_1) \leqslant z_1} \big(\lambda f(x_1)\big) + \inf_{x_2 \in \mathbb{X}, g(x_2) \leqslant z_2} \big((1 - \lambda) f(x_1)\big) \\
&= \inf_{\substack{x_1 \in \mathbb{X}, g(x_1) \leqslant z_1 \\ x_2 \in \mathbb{X}, g(x_2) \leqslant z_2}} \big(\lambda f(x_1) + (1 - \lambda) f(x_2)\big) \\
&\geqslant \inf_{\substack{x_1 \in \mathbb{X}, g(x_1) \leqslant z_1 \\ x_2 \in \mathbb{X}, g(x_2) \leqslant z_2}} f\big(\lambda x_1 + (1 - \lambda) x_2\big) \\
&\geqslant \inf_{x \in \mathbb{X}, g(x) \leqslant \lambda z_1 + (1 - \lambda) z_2} f(x).
\end{aligned}
$$

The first inequality above is a consequence of the convexity of $f$ and the second is obtained from $(\otimes)$, which implies the following:

$$
\begin{aligned}
\{\lambda x_1 + (1 - \lambda) x_2 \colon x_1, x_2 \in \mathbb{X}, \ &g(x_1) \leqslant z_1, \ g(x_2) \leqslant z_2\} \subset \\
&\{x \in \mathbb{X} \colon g(x) \leqslant \lambda z_1 + (1 - \lambda) z_2\}.
\end{aligned}
$$

(3) We need to prove that the feasible set $W_z$ is nonempty for any $z$ from some neighbourhood of $0 \in \mathbb{R}^m$. By assumption, there exists $x^* \in \mathbb{X}$ such that $g_i(x^*) < 0$ for $i = 1, \ldots, m$. Let $a = \min\{-g_1(x^*), \ldots, -g_m(x^*)\}$. Then, for all $z \in [-a, a]^m$ we have $x^* \in W_z$; hence, $[-a, a]^m \subset D_M$. $\square$

Remark. (1) If $M_{\overline{z}} = -\infty$ for some $\overline{z} \in D_M$, then by convexity of $M$ for any $z \in D_M$ and $\lambda \in (0, 1)$ we have $M(\lambda \overline{z} + (1 - \lambda) z) = -\infty$.
(2) As a consequence, if $M(\overline{z}) = -\infty$ for some $\overline{z} \in D_M$, then $M(z) = -\infty$ for all $z \in \operatorname{int} D_M$.
(3) Moreover, if there exists $\overline{z} \in \operatorname{int} D_M$ such that $M(\overline{z}) > -\infty$, then $M(z) > -\infty$ for all $z \in D_M$.

Theorem 62 *If in a problem of convex optimization there exists a point $x^* \in \mathbb{X}$ such that $g_i(x^*) < 0$ for $i = 1, \ldots, m$ and $M(0) > -\infty$, then $M(z) > -\infty$ for all $z \in D_M$ and there exists a vector $\mu \in [0, \infty)^m$ which determines a supporting hyperplane of $M$:*

$$
M(z) \geqslant M(0) - \mu^\mathsf{T} z, \quad z \in D_M.
$$

Proof. By Theorem 61, $M$ is a convex function and $0 \in \operatorname{int} D_M$. Hence, by Remark (3) above, we have $M(z) > -\infty$ for all $z \in D_M$. The existence of the supporting plane is a consequence of the supporting plane theorem (Theorem 26):

$$
M(z) \geqslant M(0) - \mu^\mathsf{T} z, \quad z \in D_M,
$$

for some $\mu \in \mathbb{R}^m$. We have to prove that all coordinates of $\mu$ are nonnegative. On the contrary, suppose that $\mu_i < 0$ for some $i \in \{1, \ldots, m\}$. As $0 \in \operatorname{int} D_M$, if a number $a$ is sufficiently small, the point $\overline{z} = a e_i$ is an element of $D_M$. Due to $\mu_i$ being negative, we obtain

$$
M(\overline{z}) \geqslant M(0) - \mu_i a > M(0).
$$

On the other hand, $W_0 \subset W_{\overline{z}}$ (due to $\overline{z} \geqslant 0$, i.e., $M(\overline{z}) \leqslant M(0)$. This inconsistency proves that $\mu \in [0, \infty)^m$. $\square$

The vector $\mu$ is called the sensitivity vector for the problem $\binom{**}{*}$. By Theorem 26, if the function $M$ is differentiable at $0$, then $\mu = -\big(DM(0)\big)^\mathsf{T}$. Therefore, $\mu$ denotes the speed and the direction of changes of the minimal value of $f$ caused by a perturbation of constraints, just like in the case of equality constraints discussed earlier.

Now we take a look at the relation of the sensitivity vector with the saddle point and the first-order condition. Note that the connection of the saddle point and the sensitivity vector does not assume the convexity of the optimization problem.

Theorem 63    *1. If $(\overline{x}, \overline{\mu})$ is a saddle point of the Lagrange function in the set $\mathbb{X} \times [0, \infty)^m$, then $\overline{\mu}$ is a sensitivity vector (i.e., it determines a supporting plane). This claim does not assume convexity of the problem.*

     *2. Suppose that the functions $f, g_1, \ldots, g_m$ are differentiable at $\overline{x}$ and convex. If the first-order condition is satisfied at $\overline{x}$ with a Lagrange multipliers vector $\overline{\mu} \in [0, \infty)^m$, then $\overline{\mu}$ is the sensitivity vector.*

Proof. (1) Let $L_z(x, \mu)$ denote the Lagrange function for a modified problem. Then,

$$
L_z(x, \mu) = f(x) + \sum_{i=1}^{m} \mu_i \big(g_i(x) - z_i\big) = L(x, \mu) - \mu^\mathsf{T} z.
$$

The point $(\overline{\mathbf{x}}, \overline{\boldsymbol{\mu}})$ is a saddle point; hence,

$$M(\mathbf{0}) = L(\overline{\mathbf{x}}, \overline{\boldsymbol{\mu}}) = \inf_{\mathbf{x} \in \mathbb{X}} L(\mathbf{x}, \overline{\boldsymbol{\mu}}).$$

Therefore,

$$M(\mathbf{0}) = \inf_{\mathbf{x} \in \mathbb{X}} L(\mathbf{x}, \overline{\boldsymbol{\mu}}) = \inf_{\mathbf{x} \in \mathbb{X}} \left( L_z(\mathbf{x}, \overline{\boldsymbol{\mu}}) + \overline{\boldsymbol{\mu}}^{\mathsf{T}} \mathbf{z} \right) = \inf_{\mathbf{x} \in \mathbb{X}} L_z(\mathbf{x}, \overline{\boldsymbol{\mu}}) + \overline{\boldsymbol{\mu}}^{\mathsf{T}} \mathbf{z}. \qquad (\otimes)$$

Note that for any $\mathbf{x} \in W_z$ and $\boldsymbol{\mu} \in [0, \infty)^m$ we have $f(\mathbf{x}) \geqslant L_z(\mathbf{x}, \boldsymbol{\mu})$; in particular,

$$M(\mathbf{z}) = \inf_{\mathbf{x} \in W_z} f(\mathbf{x}) \geqslant \inf_{\mathbf{x} \in W_z} L_z(\mathbf{x}, \overline{\boldsymbol{\mu}}) \geqslant \inf_{\mathbf{x} \in \mathbb{X}} L_z(\mathbf{x}, \overline{\boldsymbol{\mu}}).$$

By substituting the above to $(\otimes)$, we obtain

$$M(\mathbf{0}) \leqslant M(\mathbf{z}) + \overline{\boldsymbol{\mu}}^{\mathsf{T}} \mathbf{z}.$$


(2) By Lemma 20, the point $(\overline{\mathbf{x}}, \overline{\boldsymbol{\mu}})$ is a saddle point of the Lagrange function. The claim may thus be obtained from the first claim. $\square$

# 12. Introduction to numerical methods of optimization

So far, we studied the theory of optimization problems; more precisely, we proved theorems which *enable* finding candidates for solutions and recognising solutions among the candidates. Unfortunately, in most cases we need to solve systems of nonlinear equations. In practice, these systems may not have solutions possible to describe by algebraic formulae (and also the systems may be too large to be dealt with analytically). What may be done in such cases, is using numerical methods, which can find approximations of solutions with given accuracy. Some of those methods are described below. Many algorithms do not look for points satisfying the first-order necessary conditions; instead, they construct, step by step, sequences of points convergent to solutions. It does not mean that Lagrange multipliers and dual methods are useless in numerical computations. On the contrary, the classical approach with equality constraints uses extensively Lagrange multipliers, even if below we survey other methods.

Definition 34 *An <u>iterative process</u> is a four-tuple* $(Q, I, \Omega, h)$, *where $Q$ is a set,* $I \subset Q$, $\Omega \subset Q$ *and* $h \colon Q \to Q$ *is a mapping in the set $Q$, which is an identity mapping of the subset $\Omega$. This four-tuple represents a computation process; $I$ is the set of initial data, $\Omega$ is the set of solutions and the function $h$ describes the computation. Given an initial point $x \in I$, the process generates the sequence*

$$x_0 = x, \qquad x_{k+1} = h(x_k), \quad k = 0, 1, \dots$$

*The iterative process terminates after $n$ steps if $x_n \in \Omega$ (in accordance to the definition of $h$, as in that case $x_{n+1} = x_n \in \Omega$). An <u>algorithm</u> is an iterative process which terminates after a finite number of steps.*

In optimization we are interested in using algorithms to solve problems which have solutions. The algorithms ought to have the following properties:

1. <u>Correctness</u> of the algorithm, i.e., the property that for any admissible initial point $x \in I$ we get the correct result. For our purposes we assume that this property ensures the convergence of the sequence to a solution.

2. <u>Stop property</u>, which involves a condition for the iterative process to terminate. This condition is satisfied when an element of the sequence of points is in the set $\Omega$, which means that it is an accurate enough approximation of the solution.

3. <u>Effectiveness</u>, related with the rate of convergence of the sequence to the solution.

## General properties of optimization algorithms

Before discussing the algorithms we formulate the problem:

$$\begin{cases} f(x) \to \min, \\ x \in \mathbb{X} \subset \mathbb{R}^n. \end{cases}$$

<u>Remark.</u> For the algorithm to work it is necessary to provide procedures computing at given points $x_k$ the function values $f(x_k)$ and sometimes its derivatives, $Df(x_k)$, $D^2 f(x_k)$ etc. The evaluation of the function $f$ and its derivatives is not a part of the algorithm; it is assumed that those procedures compute their results with some required accuracy. Sometimes they are seen as an "oracle" yielding the function values.

Definition 35 *Let $x^*$ be a solution of an optimization problem and let $f^* = f(x^*)$. The <u>stop criteria</u> using an <u>absolute tolerance at a level</u> $\varepsilon > 0$ may be the following:*

1. $|f(x_k) - f^*| \leqslant \varepsilon$,

2. $\|x_k - x^*\| \leqslant \varepsilon$,

3. $\|Df(x_k)\| \leqslant \varepsilon$,

4. $|f(x_{k+1}) - f(x_k)| \leqslant \varepsilon$,

5. $\|x_{k+1} - x_k\| \leqslant \varepsilon$.

*We can also use the stop criteria with a <u>relative tolerance at a level</u> $\varepsilon > 0$:*

1. $|f(x_k) - f^*|/|f^*| \leqslant \varepsilon$,

2. $\|x_k - x^*\|/\|x^*\| \leqslant \varepsilon$,

3. $|f(x_{k+1}) - f(x_k)|/|f(x_k)| \leqslant \varepsilon$,

4. $\|x_{k+1} - x_k\|/\|x_k\| \leqslant \varepsilon$.

In both cases the first two criteria are the most natural, but they are the least practical, as we do not know the point $x^*$ nor the minimal function value $f^*$. In practice we use the other criteria, even if they do not guarantee the termination of the process close to the actual solution of the problem.

<u>Definition 36</u> *The greatest number* $p$ *such that the inequality*

$$\|x_{k+1} - x^*\| \leqslant c\|x_k - x^*\|^p$$

*is satisfied for all* $k > K$ *(with some* $K \in \mathbb{N}$*), where* $c > 0$ *is a constant, is called the* <u>*exponent of convergence*</u> *of the algorithm.*

<u>Remark.</u> This definition is not easy to use in practice, because we do not know $x^*$ (but it is useful in theoretical analysis of algorithms). The exponent of convergence may be measured *a posteriori* using the approximate criterion; $p$ is the greatest number such that

$$\limsup_{k\to\infty} \frac{\|x_{k+1} - x_k\|}{\|x_k - x_{k-1}\|^p} < \infty.$$

Note that a numerical computation does not produce the infinite sequence $(x_k)_{k\in\mathbb{N}}$, but only its finite initial subsequence.

## Optimization of strictly quasi-convex functions

<u>Definition 37</u> *Let* $W \subset \mathbb{R}^n$ *be a convex set. A function* $f\colon W \to \mathbb{R}$ *is* <u>*strictly quasi-convex*</u> *if for any* $x, y \in W$, $x \neq y$, *and* $\lambda \in (0,1)$ *there is*

$$f(\lambda x + (1-\lambda)y) < \max\{f(x), f(y)\}.$$

<u>Lemma 22</u>    *1. A strictly quasi-convex function has at most one minimum (local and global).*

  *2. A strictly convex function is strictly quasi-convex.*

  *3. A function defined on a line or on an* open *interval is strictly quasi-convex if it is strictly increasing, strictly decreasing, or there exists a point* $\overline{x}$ *in its domain such that this function is strictly decreasing for* $x < \overline{x}$ *and strictly increasing for* $x > \overline{x}$.

<u>Proof</u> is left as an exercise.

The strict quasi-convexity is a property making it possible to find a minimum in a closed interval without using derivatives. The main observation is made in the following lemma:

<u>Lemma 23</u> *Let* $f\colon [a,b] \to \mathbb{R}$ *be a strictly quasi-convex function and let* $a \leqslant x \leqslant y \leqslant b$.

  *1. If* $f(x) \geqslant f(y)$, *then* $f(z) > f(y)$ *for all* $z \in [a, x)$.

  *2. If* $f(x) \leqslant f(y)$, *then* $f(z) > f(x)$ *for all* $z \in (y, b]$.

<u>Proof.</u> We prove the claim (1) by contradiction. Suppose that there exists $z \in [a, x)$ such that $f(z) \leqslant f(y)$. Then, by quasi-convexity of $f$, it follows that $f(x) < \max\{f(z), f(y)\} = f(y)$, which is inconsistent with the assumption that $f(x) \geqslant f(y)$. The claim (2) is proved in a similar way. $\square$

Based on Lemma 23, we can construct many algorithms finding minima of strictly quasi-convex functions of one variable. First we take a look at the <u>dichotomic subdivision algorithm</u>. Its idea is quite simple: to find the minimum of a strictly quasi-convex function $f\colon [a,b] \to \mathbb{R}$, we choose two points, $\lambda < \mu$, in the interior of the interval $[a,b]$. By Lemma 23, we notice that if $f(\lambda) < f(\mu)$, then the interval with the minimum of $f$ may be restricted to $[a, \mu]$ and if $f(\lambda) > f(\mu)$, then the interval may be restricted fo $[\lambda, b]$. The best strategy of choosing the two points is to obtain the next interval as short as possible. We do not know *a priori*, at which of the two points the function $f$ takes a greater value. Therefore, to obtain the fastest convergence we should take into account the worst case and minimise the greater of the two numbers $\mu - a$, $b - \lambda$. This minimum is obtained with $\mu = \lambda = (a+b)/2$. As this solution does not yield two *different* points, we choose a small $\varepsilon_k > 0$ and we take $\lambda = (a+b)/2 - \varepsilon_k$ and $\mu = (a+b)/2 + \varepsilon_k$. The algorithm is the following:

<u>Preparation:</u> choose $\varepsilon_1 \in \big(0, (a+b)/2\big)$. Take $x_1 = a$, $y_1 = b$.

<u>k-th step, repeated in a loop:</u>

  1. Compute $\lambda_k = (x_k + y_k)/2 - \varepsilon_k$ and $\mu_k = (x_k + y_k)/2 + \varepsilon_k$.

2. If $f(\lambda_k) < f(\mu_k)$, then take $x_{k+1} = x_k$ and $y_{k+1} = \mu_k$.

3. If $f(\lambda) \geqslant f(\mu_k)$, then take $x_{k+1} = \lambda_k$ and $y_{k+1} = y_k$.

4. Take $\varepsilon_{k+1} = \varepsilon_k/2$.

Stop condition: $y_{k+1} - x_{k+1} \leqslant 2\varepsilon$, where $\varepsilon > 0$ is the required accuracy (we can take the solution $\tilde{x} = (x_{k+1} + y_{k+1})/2)$.

If a minimum $\bar{x}$ of $f$ exists in $[a, b]$ (which *may not* be the case if $f$ is not continuous), then, by Lemma 23, this minimum is located in the interval $[x_k, y_k]$ for all $k$. With $k \to \infty$ the lengths of these intervals tend to 0; note that $0 < \varepsilon_k < (y_{k-1} - x_{k-1})/2$ for all $k$. Hence, both sequences, $(x_k)_k$ and $(y_k)_k$, converge to the solution $\bar{x}$. This proves that the algorithm is correct.

If the iterations are terminated after the k-th step, then the approximation error of the exact solution by the midpoint $\tilde{x}$ of the interval $[x_{k+1}, y_{k+1}]$ is bounded by the half of length of this interval, which justifies the choice of the stop criterion.

We can prove that the lengths of the intervals $[x_k, y_k]$ are bounded by elements of a geometric sequence; there is

$$y_{k+1} - x_{k+1} \leqslant (b - a)c^k,$$

where $c = \frac{1}{2} + \frac{\varepsilon_1}{b-a}$. The exponent of convergence is $p = 1$; the convergence with such an exponent is called underline{linear}.

The dichotomic subdivision algorithm has to compute two function values in each step. The golden ratio algorithm computes just one function value in each step. The choice of points $\lambda_k, \mu_k$ in the interval $[x_k, y_k]$ is done as follows:

$$y_k - \lambda_k = \mu_k - x_k; \quad \text{hence,} \quad \lambda_k = \tau x_k + (1-\tau)y_k, \ \mu_k = (1-\tau)x_k + \tau y_k,$$

where $\tau \in (0, 1)$ is such that

$$\lambda_k = (1-\tau)x_k + \tau\mu_k \quad \text{and} \quad \mu_k = \tau\lambda_k + (1-\tau)y_k.$$

The above determines the number $\tau$; there is

$$\lambda_k = \tau x_k + (1-\tau)y_k = (1-\tau)x_k + \tau\mu_k = (1-\tau)x_k + \tau\big((1-\tau)x_k + \tau y_k\big)$$
$$= \tau x_k + y_k - \tau y_k = x_k - \tau^2 x_k + \tau^2 y_k.$$

The last equality must hold for any $x_k, y_k$; hence,

$$\tau x_k = x_k - \tau^2 x_k, \quad y_k - \tau y_k = \tau^2 y_k,$$
$$\text{i.e., } (\tau^2 + \tau - 1)x_k = (\tau^2 + \tau - 1)y_k = 0.$$

The positive zero of the polynomial $p(\tau) = \tau^2 + \tau - 1$ is $\tau = (\sqrt{5} - 1)/2 \approx 0.618$; this number describes the golden ratio proportion. The golden ratio algorithm is the following:

Preparation: Take $x_1 = a$, $y_1 = b$, $\lambda_1 = \tau x_1 + (1 - \tau)y_1$, $\mu_1 = (1 - \tau)x_1 + \tau y_1$ and compute $f(\lambda_1)$ and $f(\mu_1)$.

k-th step, repeated in a loop:

1. If $f(\lambda_k) < f(\mu_k)$, then take $x_{k+1} = x_k$, $y_{k+1} = \mu_k$, $\lambda_{k+1} = \tau x_k + (1 - \tau)y_k$, $\mu_{k+1} = \lambda_k$ and compute $f(\lambda_{k+1})$.

2. If $f(\lambda_k) \geqslant f(\mu_k)$, then take $y_{k+1} = y_k$, $x_{k+1} = \lambda_k$, $\mu_{k+1} = (1 - \tau)x_k + \tau y_k$, $\lambda_{k+1} = \mu_k$ and compute $f(\mu_{k+1})$.

Stop condition: $y_{k+1} - x_{k+1} \leqslant 2\varepsilon$, where $\varepsilon > 0$ is the required accuracy.

The convergence of the sequences $(x_k)_k$, $(y_k)_k$ to the minimum $\bar{x}$ is slower than that of the sequences from the dichotomic subdivision algorithm; there is

$$y_{k+1} - x_{k+1} = (b - a)\tau^k.$$

However, due to the twice lower computational cost of one iteration, the golden ratio algorithm achieves the same accuracy of the result in a shorter time.

The slightly more complicated underline{Fibonacci algorithm} guarantees finding the minimum with the assumed accuracy after evaluating the function at the minimal number of points sufficient to minimise *any* strictly quasi-convex function.

Remark. (1) One can prove that the golden ratio algorithm may need to evaluate the function at most at one point more than the Fibonacci algorithm.
(2) There exist algorithms working faster for functions $f$ having some properties in addition to the strict quasi-convexity, like differentiablility. For example, having a procedure of computing the derivative, we can use the secant method to find the zero of $f'$; if $f''$ satisfies the Lipschitz condition, then the exponent of convergence is $p = 1 + \tau \approx 1.618$.

The Fibonacci sequence is defined as follows:

$$F_0 = F_1 = 1,$$
$$F_{k+1} = F_k + F_{k-1}, \quad k = 1, 2, 3, \ldots$$

To use this algorithm we fix *a priori* the number of iterations, which is related with the desired accuracy in the way explained later. The numbers $\lambda_k$ and $\mu_k$ are obtained from the formulae

$$\begin{cases} \lambda_k = x_k + \dfrac{F_{n-k-1}}{F_{n-k+1}}(y_k - x_k), \\[2mm] \mu_k = x_k + \dfrac{F_{n-k}}{F_{n-k+1}}(y_k - x_k), \end{cases} \qquad k = 1, \ldots, n-1.$$

If $f(\lambda_k) \geqslant f(\mu_k)$, then we take $x_{k+1} = \lambda_k$ and $y_{k+1} = y_k$. The rules written above produce

$$\begin{aligned}
\lambda_{k+1} &= x_{k+1} + \frac{F_{n-k-2}}{F_{n-k}}(y_{k+1} - x_{k+1}) \\
&= x_k + \frac{F_{n-k-1}}{F_{n-k+1}}(y_k - x_k) + \frac{F_{n-k-2}}{F_{n-k}}\left(y_k - x_k - \frac{F_{n-k-1}}{F_{n-k+1}}(y_k - x_k)\right) \\
&= x_k + \left(\frac{F_{n-k-1}}{F_{n-k+1}} + \frac{F_{n-k-2}}{F_{n-k}} - \frac{F_{n-k-2}}{F_{n-k}}\frac{F_{n-k-1}}{F_{n-k+1}}\right)(y_k - x_k) \\
&= x_k + \frac{F_{n-k-1}F_{n-k} + F_{n-k-2}F_{n-k+1} - F_{n-k-2}F_{n-k-1}}{F_{n-k+1}F_{n-k}}(y_k - x_k) \\
&= x_k + \frac{F_{n-k-1}F_{n-k} + F_{n-k-2}F_{n-k}}{F_{n-k+1}F_{n-k}}(y_k - x_k) \\
&= x_k + \frac{F_{n-k}}{F_{n-k+1}}(y_k - x_k) = \mu_k.
\end{aligned}$$

If $f(\lambda_k) < f(\mu_k)$, then we take $x_{k+1} = x_k$ and $y_{k+1} = \mu_k$. A similar calculation proves that in this case

$$\mu_{k+1} = \lambda_k,$$

which means that in both cases we choose only one point in the interval $[x_k, y_k]$ at which the function value has not yet been computed.

Suppose that we need to find the minimum in the interval $[a, b]$ with the accuracy not worse than some $\varepsilon > 0$. The number of steps is $n - 1$, where $n$ is the smallest number such that $F_n \geqslant (b - a)/\varepsilon$. Note that $\lambda_{n-1} = \mu_{n-1} = \frac{1}{2}(x_{n-1} + y_{n-1})$; the point obtained in last step is the midpoint of the interval $[x_{n-1}, y_{n-1}]$, whose length is less than or equal to $F_2\varepsilon$.

The Fibonacci algorithm is the following:

Preparation: Find the smallest $n$ such that $F_n \geqslant (b - a)/\varepsilon$;
take $x_1 = a$, $y_1 = b$, $\lambda_1 = \big(F_{n-1}x_1 + (F_n - F_{n-1})y_1\big)/F_n$,
$\mu_1 = \big((F_n - F_{n-1})x_1 + F_{n-1}y_1\big)/F_n$ and compute $f(\lambda_1)$ and $f(\mu_1)$.

k-th step, repeated in a loop for $k = 1, \ldots, n-2$:

1. If $f(\lambda_k) < f(\mu_k)$, then take $x_{k+1} = x_k$, $y_{k+1} = \mu_k$, $\mu_{k+1} = \lambda_k$,
$\lambda_{k+1} = \big(F_{n-k}x_k + (F_{n-k+1} - F_{n-k})y_k\big)/F_{n-k+1}$ and compute $f(\lambda_{k+1})$.

2. If $f(\lambda_k) \geqslant f(\mu_k)$, then take $y_{k+1} = y_k$, $x_{k+1} = \lambda_k$, $\lambda_{k+1} = \mu_k$,
$\mu_{k+1} = \big((F_{n-k+1} - F_{n-k})x_k + F_{n-k}y_k\big)/F_{n-k+1}$ and compute $f(\mu_{k+1})$.

The last step: If $f(\lambda_{n-1}) < f(\mu_{n-1})$, then take $\tilde{x} = \lambda_{n-1}$, else take $\tilde{x} = \mu_{n-1}$. The point $\tilde{x}$ is an approximation of the minimum with the required accuracy.
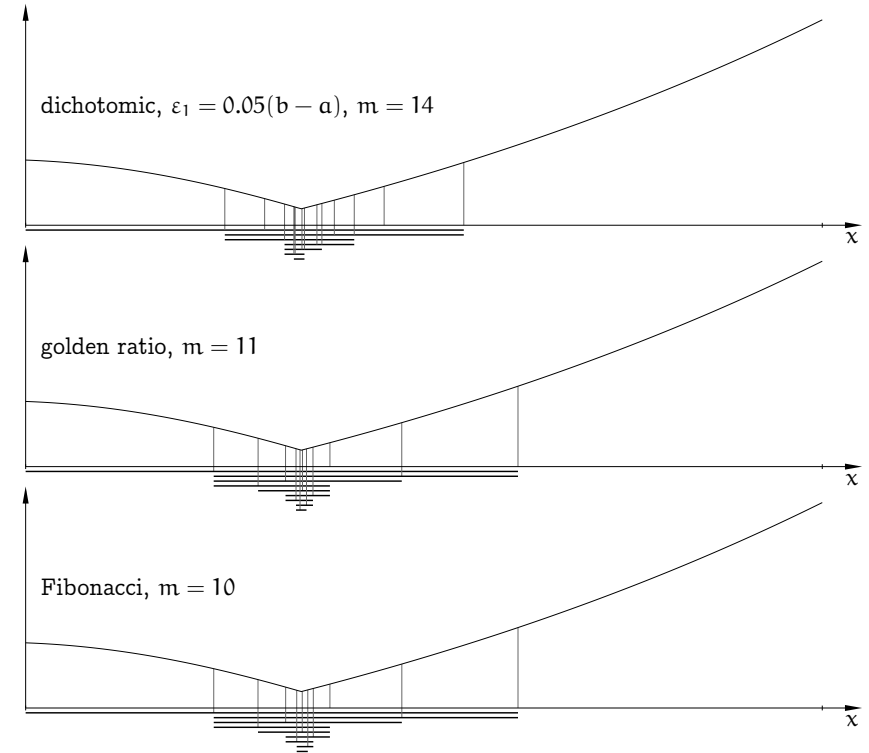


Figure 9: Three algorithms in action

An example is shown in Figure 9; the interval $[a, b] = [0, 1]$ and the assumed accuracy is 0.01. Below the $x$ axis of each graph the subsequent intervals $[x_k, y_k]$ are drawn; $m$ is the number of points at which the function value was computed.

# 13. Algorithms for unconstrained optimization

Now we take a look at multidimensional optimization methods for problems without constraints. We assume that the target function $f\colon \mathbb{R}^n \to \mathbb{R}$ is of class $C^1$ and its minimum is to be found. All methods considered here are based on the following scheme: beginning at a point $x_1$, which we believe to be close enough to the solution, we construct a sequence of points $x_1, x_2, \ldots$ so as to obtain $f(x_{k+1}) < f(x_k)$ for all $k$. We expect to find a minimum of $f$ in this way. However, it may turn out that the concentration points of the sequence $(x_k)_k$ are not solutions. We survey various methods of constructing sequences $(x_k)_k$, paying attention to their convergence. All results described below apply also to functions defined in open subsets of $\mathbb{R}^n$. Descent methods will play first fiddle in optimization problems with constraints, considered later.

We consider the problem

$$\begin{cases} f(x) \to \min, \\ x \in \mathbb{R}^n. \end{cases}$$

Descent methods are algorithms constructing consecutive points according to the formula

$$x_{k+1} = x_k + \alpha_k d_k,$$

where $\alpha_k > 0$ and the vector $d_k$ has a descent direction, i.e.,

$$\begin{aligned} Df(x_k)d_k &< 0 \quad \text{if } Df(x_k) \neq 0^\mathsf{T}, \\ d_k &= 0 \quad \text{if } Df(x_k) = 0^\mathsf{T}. \end{aligned}$$

With $d_k \neq 0$ and $\alpha_k$ sufficiently small we have $f(x_{k+1}) < f(x_k)$. We might expect the sequence $(x_k)_k$ to converge to a minimum. But without additional information about the problem we have no guarantee of finding a global minimum. We are going to look for methods of making the sequence converge to a local minimum.

## Steepest descent methods

The steepest descent methods take the vectors $d_k$ parallel to $\big(Df(x_k)\big)^\mathsf{T}$. The idea is simple:

Preparation: Choose an initial point $x_1$.

k-th step:

1. Choose $\alpha_k$,
2. Take $x_{k+1} = x_k - \alpha_k\big(Df(x_k)\big)^\mathsf{T}$.

Stop condition: $\|Df(x_{k+1})\| \leqslant \varepsilon$.

It remains to select rules of choosing $\alpha_k$. The rules below will also be used with any vectors $d_k$ having descent directions, not necessarily $-\big(Df(x_k)\big)^\mathsf{T}$.

- Exact minimization rule: choose $\alpha_k$ such that

$$f(x_k + \alpha_k d_k) = \min_{\alpha \geqslant 0} f(x_k + \alpha d_k).$$

- Limited minimization rule: with a fixed $A > 0$ choose $\alpha_k$ such that

$$f(x_k + \alpha_k d_k) = \min_{\alpha \in [0,A]} f(x_k + \alpha d_k).$$

- Armijo rule: with fixed $s > 0$, $\beta, \sigma \in (0,1)$ we take $\alpha_k = \beta^{m_k} s$, where $m_k$ is the smallest integer $m$ such that

$$f(x_k) - f(x_k + \beta^m s d_k) \geqslant -\sigma\beta^m s Df(x_k)d_k.$$

According to this rule, the following inequality holds:

$$f(x_k) - f(x_k + \beta^{m-1} s d_k) < -\sigma\beta^{m-1} s Df(x_k)d_k.$$

The constant $s$ is called a stride, $\beta$ controls the rate of decreasing or increasing the stride (the smaller it is, the faster the stride changes) and $\sigma$ influences the choice of $\alpha_k$ as follows: the smaller it is, the smaller $m_k$ satisfies the conditions above, resulting in the greater value of $\alpha_k$.

In the steepest descent methods, where $d_k = -\big(Df(x_k)\big)^\mathsf{T}$, the choice of $m_k$ is slightly simpler due to $Df(x_k)d_k = -\|Df(x_k)\|^2$: $m_k$ is the smallest integer $m$ such that

$$f(x_k) - f\big(x_k - \beta^m s\big(Df(x_k)\big)^\mathsf{T}\big) \geqslant \sigma\beta^m s\|Df(x_k)\|^2. \tag{*}$$

Figure 10 shows the idea of the Armijo rule, which is easy to implement. We assume that $Df(x_k) \neq 0^\mathsf{T}$. Starting from the point $y = x_k - s\big(Df(x_k)\big)^\mathsf{T}$ we check the condition (*). If it does not hold, then we check the points $x_k - s\beta\big(Df(x_k)\big)^\mathsf{T}$,
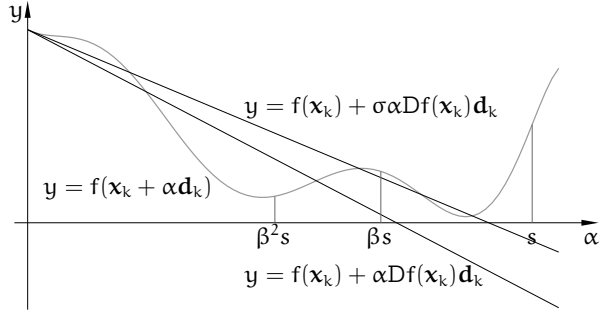
Figure 10: The Armijo rule in action; $\sigma = 0.8$, $\beta = 0.7$, $\mathbf{x}_{k+1} = \mathbf{x}_k + \beta^2 s \mathbf{d}_k$

$\mathbf{x}_k - s\beta^2 \big(Df(\mathbf{x}_k)\big)^\top$ etc. If (*) is satisfied with $m = 0$, corresponding to the point $\mathbf{y}$, we try the points $\mathbf{x}_k - s\beta^{-1}\big(Df(\mathbf{x}_k)\big)^\top$, $\mathbf{x}_k - s\beta^{-2}\big(Df(\mathbf{x}_k)\big)^\top$ etc. After a finite number of steps we find the number $m$ satisfying (*) such that $m-1$ does not satisfy this condition.

The exact minimization rule is well defined if the minimum on the right-hand side exists. The interval in which the minimum is to be found is unbounded; hence, the exact minimization rule does not have to be well defined. Later we show assumptions that guarantee its correctness. Restricting the interval to $[0, A]$ has two advantages: first, the problem always has a solution, as we minimise a continuous function in a compact set. Second, we can use faster methods of searching minima.

Both rules, exact and limited, find steps such that $f(\mathbf{x}_{k+1}) \leqslant f(\mathbf{x}_k)$; the proof of this fact is an exercise. By Inequality (*), also the Armijo rule produces such steps. There are, however, two questions. Is the stop criterion correct? Does the sequence $(\mathbf{x}_k)_k$ converge to a minimum?

Below we prove that any concentration point of the sequence $(\mathbf{x}_k)$ obtained by the steepest descent method is a critical point, i.e., it is a zero of the gradient of $f$. Such a point does not have to be a local minimum. However, if the function $f$ is pseudoconvex, then we can be sure that it is a minimum, moreover, a global minimum.

Theorem 64 *Let $(\mathbf{x}_k)_k$ be a sequence of points obtained using the steepest descent method, with any of the three rules described above. Then, any concentration point of this sequence is a critical point.*

Proof. Let $\overline{\mathbf{x}}$ be a concentration point and let $(\mathbf{x}_{k_n})_n$ be a subsequence convergent to $\overline{\mathbf{x}}$. The proof is done by contradiction; we suppose that $Df(\overline{\mathbf{x}}) \neq \mathbf{0}^\top$.

The main idea for the first two rules (exact and limited) is to show the existence of a constant $\gamma > 0$ such that with $n$ big enough, i.e., with $\mathbf{x}_{k_n}$ close enough to $\overline{\mathbf{x}}$, there is

$$f(\mathbf{x}_{k_n+1}) \leqslant f(\mathbf{x}_{k_n}) - \gamma,$$

which means that it is possible to decrease the value of $f$ in the step $k_n$ at least by $\gamma$. Bearing in mind that the sequence $\big(f(\mathbf{x}_k)\big)_k$ is decreasing, we have

$$f(\mathbf{x}_{k_{n+1}}) \leqslant f(\mathbf{x}_{k_n+1}) \leqslant f(\mathbf{x}_{k_n}) - \gamma.$$

By passing with $n$ to $\infty$ and using the fact that $\overline{\mathbf{x}}$ is a concentration point and $f$ is continuous, we obtain the inconsistency with the inequality $\gamma > 0$.

Therefore we need to prove the existence of $\gamma > 0$. The derivative of $f$ is continuous by assumption; hence, there exists a neighbourhood $V$ of $\overline{\mathbf{x}}$ such that

$$\frac{\|Df(\mathbf{x}) - Df(\mathbf{y})\|}{\|Df(\mathbf{x})\|} \leqslant \frac{1}{2} \quad \text{for all } \mathbf{x}, \mathbf{y} \in V. \tag{**}$$

The main observation to justify the inequality above is that in a neighbourhood of $\overline{\mathbf{x}}$ the norm of gradient of $f$ is strictly separated from $0$.

Let $\delta > 0$ be such that $B(\overline{\mathbf{x}}, 2\delta) \subset V$ and $\delta \leqslant A \inf_{\mathbf{x} \in V} \|Df(\mathbf{x})\|$, where $A$ is the constant of the definition of the limited minimization rule (from (**), we conclude that $\inf_{\mathbf{x} \in V} \|Df(\mathbf{x})\| > 0$). For a point $\mathbf{x} \in B(\overline{\mathbf{x}}, \delta)$, the point $\mathbf{x} - \frac{\delta}{\|Df(\mathbf{x})\|}\big(Df(\mathbf{x})\big)^\top$ is an element of $B(\overline{\mathbf{x}}, 2\delta)$ and, therefore, also of $V$. Moreover, with the limited minimization rule we have

$$\frac{\delta}{\|Df(\mathbf{x})\|} \leqslant \frac{A \inf_{\mathbf{x} \in V} \|Df(\mathbf{x})\|}{\|Df(\mathbf{x})\|} \leqslant A,$$

i.e., $\frac{\delta}{\|Df(\mathbf{x})\|} \in [0, A]$. By the mean value theorem,

$$f(\mathbf{x}) - f\Big(\mathbf{x} - \frac{\delta}{\|Df(\mathbf{x})\|}\big(Df(\mathbf{x})\big)^\top\Big) = Df(\theta)\Big(\frac{\delta}{\|Df(\mathbf{x})\|}\big(Df(\mathbf{x})\big)^\top\Big)$$
$$= \frac{\delta}{\|Df(\mathbf{x})\|} Df(\theta)\big(Df(\mathbf{x})\big)^\top,$$

where $\theta$ is an intermediate point; hence, it is an element of $V$. Now we focus on the product of the gradients:

$$Df(\theta)\big(Df(\mathbf{x})\big)^\top = \big(Df(\mathbf{x}) + Df(\theta) - Df(\mathbf{x})\big)\big(Df(\mathbf{x})\big)^\top$$
$$= \|Df(\mathbf{x})\|^2 + \big(Df(\theta) - Df(\mathbf{x})\big)\big(Df(\mathbf{x})\big)^\top$$
$$\geqslant \|Df(\mathbf{x})\|^2 - \|Df(\theta) - Df(\mathbf{x})\|\|Df(\mathbf{x})\|.$$

We substitute this estimation to the previous formula:

$$f(\mathbf{x}) - f\Big(\mathbf{x} - \frac{\delta}{\|Df(\mathbf{x})\|}\big(Df(\mathbf{x})\big)^{\mathsf{T}}\Big) \geqslant \delta\|Df(\mathbf{x})\| - \delta\|Df(\theta) - Df(\mathbf{x})\|$$

$$= \delta\|Df(\mathbf{x})\|\Big(1 - \frac{\|Df(\theta) - Df(\mathbf{x})\|}{\|Df(\mathbf{x})\|}\Big) \geqslant \frac{\delta}{2}\|Df(\mathbf{x})\|;$$

the last inequality follows from (**). The estimation above is true for all $\mathbf{x} \in B(\overline{\mathbf{x}}, \delta)$; in particular for all $\mathbf{x}_{k_n}$ with $n$ big enough we have

$$f(\mathbf{x}_{k_n}) - f(\mathbf{x}_{k_n+1}) \geqslant f(\mathbf{x}_{k_n}) - f\Big(\mathbf{x}_{k_n} - \frac{\delta}{\|Df(\mathbf{x}_{k_n})\|}\big(Df(\mathbf{x}_{k_n})\big)^{\mathsf{T}}\Big) \geqslant \frac{\delta}{2}\|Df(\mathbf{x}_{k_n})\|.$$

Therefore, we can take

$$\gamma = \frac{\delta}{2}\inf_{\mathbf{x} \in B(\overline{\mathbf{x}}, \delta)}\|Df(\mathbf{x})\|.$$

For the steepest descent method with the Armijo rule we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \geqslant \sigma\alpha_k\|Df(\mathbf{x}_k)\|^2,$$

which implies that the sequence $\big(f(\mathbf{x}_k)\big)_k$ is monotonically decreasing, i.e., it is either convergent, or it diverges to $-\infty$. Due to $f(\mathbf{x}_{k_n}) \to f(\overline{\mathbf{x}})$ there is $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \to 0$. It follows that also $\alpha_k\|Df(\mathbf{x}_k)\|^2 \to 0$. Due to $Df(\mathbf{x}_{k_n}) \to Df(\overline{\mathbf{x}}) \neq \mathbf{0}^{\mathsf{T}}$, there must be $\alpha_k \to 0$.

On the other hand, the numbers $\alpha_k$ in the Armijo rule are chosen in the optimal way, i.e.,

$$f(\mathbf{x}_k) - f\big(\mathbf{x}_k - \alpha_k/\beta\big(Df(\mathbf{x}_k)\big)^{\mathsf{T}}\big) < \sigma\alpha_k/\beta\|Df(\mathbf{x}_k)\|^2.$$

By applying the intermediate value theorem to the left-hand side of this inequality, we get

$$f(\mathbf{x}_k) - f\big(\mathbf{x}_k - \alpha_k/\beta\big(Df(\mathbf{x}_k)\big)^{\mathsf{T}}\big) = Df\big(\mathbf{x}_k - \tilde{\alpha}_k/\beta\big(Df(\mathbf{x}_{k_n})\big)^{\mathsf{T}}\big)\alpha_k/\beta\big(Df(\mathbf{x}_{k_n})\big)^{\mathsf{T}}.$$

Therefore, the previous inequality takes the form

$$Df\big(\mathbf{x}_k - \tilde{\alpha}_k/\beta\big(Df(\mathbf{x}_k)\big)^{\mathsf{T}}\big)\big(Df(\mathbf{x}_k)\big)^{\mathsf{T}} < \sigma\|Df(\mathbf{x}_k)\|^2$$

Taking the subsequence $(\mathbf{x}_{k_n})_n$, we obtain

$$Df\big(\mathbf{x}_{k_n} - \tilde{\alpha}_{k_n}/\beta\big(Df(\mathbf{x}_{k_n})\big)^{\mathsf{T}}\big)\big(Df(\mathbf{x}_{k_n})\big)^{\mathsf{T}} < \sigma\|Df(\mathbf{x}_{k_n})\|^2$$

If we pass with this inequality to the limit, then we see that $\tilde{\alpha}_{k_n} \in [0, \alpha_{k_n}]$; $\tilde{\alpha}_{k_n} \to 0$ because $\alpha_{k_n} \to 0$. The last inequality at the limit gives us $\|Df(\overline{\mathbf{x}})\|^2 \leqslant \sigma\|Df(\overline{\mathbf{x}})\|^2$, i.e., $(1 - \sigma)\|Df(\overline{\mathbf{x}})\|^2 \leqslant 0$. But $1 - \sigma \in (0, 1)$, which gives us an inconsistency with the supposition that $Df(\overline{\mathbf{x}}) \neq \mathbf{0}^{\mathsf{T}}$. $\square$

It turns out that slightly stronger assumptions ensure that $Df(\mathbf{x}_k) \to \mathbf{0}^{\mathsf{T}}$ for the entire sequence obtained by the steepest descent method, not just a subsequence convergent to the concentration point.

<u>Lemma 24</u> *Let $f$ be a function of class $C^1$ bounded from below. Let $(\mathbf{x}_k)_k$ be a sequence of points obtained by a steepert descent method. If there exists a constant $c > 0$ (independent of $k$) such that*

$$f(\mathbf{x}_k + \alpha_k\mathbf{d}_k) < f(\mathbf{x}_k) - c\|Df(\mathbf{x}_k)\|^2, \quad k = 1, 2, \ldots, \tag{$\overset{**}{*}$}$$

*then either there exists a number $K$ such that $Df(\mathbf{x}_K) = \mathbf{0}^{\mathsf{T}}$, or the sequence $\big(Df(\mathbf{x}_k)\big)_k$ converges to $\mathbf{0}^{\mathsf{T}}$.*

<u>Proof.</u> Suppose that there exists an infinite sequence $(\mathbf{x}_k)_k$ obtained by a steepest descent method (i.e., a number $K$ as described, does not exist). By $(\overset{**}{*})$, the sequence $\big(f(\mathbf{x}_k)\big)_k$ is monotonically decreasing. Being bounded from below, this sequence is convergent; hence, $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \to 0$. Due to $(\overset{**}{*})$, we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) > c\|Df(\mathbf{x}_k)\|^2,$$

therefore, $\|Df(\mathbf{x}_k)\| \to 0$. $\square$

<u>Theorem 65</u> *Let $f$ be a function bounded from below and let its gradient satisfy the Lipschitz condition with a constant $L$ in a sublevel set $S = W_{f(\mathbf{x}_1)}(f)$ for some $\mathbf{x}_1$. Let $(\mathbf{x}_k)_k$ be a sequence obtained using the steepest descent method with the exact rule (which is correct if the set $S$ is compact), limited rule with $A > \frac{1}{2L}$ or the Armijo rule. Then, either there exists a number $K$ such that $Df(\mathbf{x}_K) = \mathbf{0}^{\mathsf{T}}$, or the sequence $\big(Df(\mathbf{x}_k)\big)_k$ converges to $\mathbf{0}^{\mathsf{T}}$.*

<u>Proof.</u> The proof is done by showing that the assumptions of the theorem imply the assumptions of Lemma 24.

At first we consider the Armijo rule. Let $\alpha_k = s\beta^{m_k}$ be the step chosen by this rule in the $k$-th step. Then,

$$f(\mathbf{x}_k + \alpha_k\mathbf{d}_k) \leqslant f(\mathbf{x}_k) + \sigma\alpha_k Df(\mathbf{x}_k)\mathbf{d}_k. \tag{$\oplus$}$$

The number $\alpha_k$ has been chosen in the optimal way, i.e., with the lowest power of $\beta$ to satisfy the last inequality. Therefore, for an even lower power of $\beta$ we have

$$f(\mathbf{x}_k + \beta^{-1}\alpha_k\mathbf{d}_k) > f(\mathbf{x}_k) + \sigma\alpha_k Df(\mathbf{x}_k)\mathbf{d}_k. \qquad (\otimes)$$

Due to $Df(\mathbf{x}_k)\mathbf{d}_k < 0$ and $\mathbf{x}_k \in S$, using $(\oplus)$ we obtain $\mathbf{x}_{k+1} \in S$. By the intermediate value theorem, we obtain

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) = f(\mathbf{x}_k + \alpha_k\mathbf{d}_k) - f(\mathbf{x}_k) = \alpha_k Df(\tilde{\mathbf{x}}_k)\mathbf{d}_k,$$

where $\tilde{\mathbf{x}}_k$ is a point of the line segment $\overline{\mathbf{x}_k\mathbf{x}_{k+1}}$. This line segment is a subset of $S$; in particular, $\tilde{\mathbf{x}}_k \in S$. Taking into account that $\mathbf{d}_k = -\big(Df(\mathbf{x}_k)\big)^\top$, we have

$$\begin{aligned}
f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) &= \alpha Df(\tilde{\mathbf{x}}_k)\mathbf{d}_k = -\alpha_k\big(Df(\mathbf{x}_k) - Df(\mathbf{x}_k) + Df(\tilde{\mathbf{x}}_k)\big)\big(Df(\tilde{\mathbf{x}}_k)\big)^\top \\
&= -\alpha_k\|Df(\mathbf{x}_k)\|^2 + \alpha_k\big(Df(\mathbf{x}_k) - Df(\tilde{\mathbf{x}}_k)\big)\big(Df(\mathbf{x}_k)\big)^\top \\
&\leqslant -\alpha_k\|Df(\mathbf{x}_k)\|^2 + \alpha_k\|Df(\mathbf{x}_k) - Df(\tilde{\mathbf{x}}_k)\|\|Df(\mathbf{x}_k)\|.
\end{aligned}$$

There is $\mathbf{x}_k \in S$ and also $\tilde{\mathbf{x}}_k \in S$ for all $k$. From the Lipschitz condition satisfied by the gradient of $f$ it follows that

$$\|Df(\mathbf{x}_k) - Df(\tilde{\mathbf{x}}_k)\| \leqslant L\|\mathbf{x}_k - \tilde{\mathbf{x}}_k\| \leqslant L\|\mathbf{x}_k - \mathbf{x}_{k+1}\| = L\alpha_k\|Df(\mathbf{x}_k)\|.$$

After substituting this estimation to the previous inequality, we obtain

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leqslant -\alpha_k\|Df(\mathbf{x}_k)\|^2(1 - \alpha_k L).$$

It follows from the above that $(\oplus)$ is satisfied if $1 - \alpha_k L \geqslant \sigma$. Let $m_\sigma$ be such a number that $1 - s\beta^{m_\sigma}L \geqslant \sigma > 1 - s\beta^{m_\sigma-1}$. It follows that

$$\sigma s\beta^{m_\sigma} > \frac{\sigma(1-\sigma)\beta}{L}.$$

On the other hand, $m_\sigma$ does not have to be optimal, satisfying $(\oplus)$ and $(\otimes)$. However, the Armijo rule guarantees that $m_k \leqslant m_\sigma$, because $m_k$ is the smallest number such that $(\oplus)$ is satisfied. Hence,

$$\sigma s\beta^{m_k} \geqslant \sigma s\beta^{m_\sigma} > \frac{\sigma(1-\sigma)\beta}{L}.$$

After substituting the above to $(\oplus)$ and using the fact that $\mathbf{d}_k$ has the steepest descent direction, we obtain the following inequality:

$$\begin{aligned}
f(\mathbf{x}_k - \alpha_k\mathbf{d}_k) - f(\mathbf{x}_k) &\leqslant -\sigma s\beta^{m_k}\|Df(\mathbf{x}_k)\|^2 \leqslant -\sigma s\beta^{m_\sigma}\|Df(\mathbf{x}_k)\|^2 \\
&< -\frac{\sigma(1-\sigma)\beta}{L}\|Df(\mathbf{x}_k)\|^2,
\end{aligned}$$

i.e., the inequality $(\overset{*}{\ast})$ of Lemma 24.

With the exact minimization rule, we assume that it is well-defined in each step, i.e., there exists a finite number $\alpha_k \geqslant 0$ to minimise $f(\mathbf{x}_k + \alpha\mathbf{d}_k)$ for $\alpha \geqslant 0$. A calculation similar to that in the proof for the Armijo rule yields the estimation

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leqslant -\alpha_k\|Df(\mathbf{x}_k)\|^2(1 - \alpha_k L). \qquad (\overset{**}{\ast\ast})$$

The right-hand side of this inequality is minimal for $\alpha_k = \frac{1}{2L}$. From that we obtain

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leqslant -\frac{1}{4L}\|Df(\mathbf{x}_k)\|^2,$$

i.e., the inequality $(\overset{*}{\ast})$ from Lemma 24.

The proof for the restricted minimization rule is similar to the above. We need to assume that the interval $[0, A]$ in which we search for $\alpha_k$ contains the point $\frac{1}{2L}$ to minimise $(\overset{**}{\ast\ast})$. $\square$

Theorem 65 justifies correctness of the steepest descent methods, but it does not refer to the rate of their convergence. After further strengthtening its assumptions, we can prove that the convergence is at least linear and the stop criterion is correct.

We define the set $S = \{\mathbf{x} \in \mathbb{R}^n \colon f(\mathbf{x}) \leqslant f(\mathbf{x}_1)\}$, where $\mathbf{x}_1$ is the initial point. By $m(\mathbf{x})$ and $M(\mathbf{x})$ we denote respectively the smallest and the greatest eigenvalue of the Hessian matrix of $f$ at $\mathbf{x}$.

<u>Lemma 25</u> *Suppose that the set $S$ is convex and compact and the function $f$ is twice continuously differentiable in $S$. Let $m = \inf_{\mathbf{x}\in S} m(\mathbf{x}) > 0$. Then the point $\overline{\mathbf{x}}$ being the limit of the sequence $(\mathbf{x}_k)_k$ obtained by a steepest descent method is in $S$. It is the minimum of $f$ in $S$ and for all $\mathbf{x} \in S$ there is*

$$\|\mathbf{x} - \overline{\mathbf{x}}\| \leqslant \frac{1}{m}\|Df(\mathbf{x})\|, \quad f(\mathbf{x}) - f(\overline{\mathbf{x}}) \leqslant \frac{1}{m}\|Df(\mathbf{x})\|^2.$$

<u>Proof.</u> Due to $m = \inf_{\mathbf{x}\in S} m(\mathbf{x}) > 0$, the Hessian $D^2f(\mathbf{x})$ is positive-definite at all points of $S$ and it follows that $f$ is a convex function in $S$.

The set $S$ is compact; the function $f$ is bounded from below and its gradient satisfies the Lipschitz condition. By Theorem 65, $Df(\mathbf{x}_k) \to 0$, i.e., $\overline{\mathbf{x}} \in S$. The critical point of a convex function is its minimum.

The Taylor formula for $\boldsymbol{x} \in S$ gives us

$$f(\boldsymbol{x}) = f(\overline{\boldsymbol{x}}) + \frac{1}{2}(\boldsymbol{x} - \overline{\boldsymbol{x}})^{\mathsf{T}} D^2(\tilde{\boldsymbol{x}})(\boldsymbol{x} - \overline{\boldsymbol{x}}),$$

$$f(\overline{\boldsymbol{x}}) = f(\boldsymbol{x}) + Df(\boldsymbol{x})(\overline{\boldsymbol{x}} - \boldsymbol{x}) + \frac{1}{2}(\overline{\boldsymbol{x}} - \boldsymbol{x})^{\mathsf{T}} D^2(\hat{\boldsymbol{x}})(\overline{\boldsymbol{x}} - \boldsymbol{x}).$$

After substituting the first equality to the second, we obtain

$$Df(\boldsymbol{x})(\overline{\boldsymbol{x}} - \boldsymbol{x}) + (\boldsymbol{x} - \overline{\boldsymbol{x}})^{\mathsf{T}} \frac{D^2 f(\hat{\boldsymbol{x}}) - D^2 f(\tilde{\boldsymbol{x}})}{2}(\boldsymbol{x} - \overline{\boldsymbol{x}}) = 0.$$

We can, therefore, estimate

$$Df(\boldsymbol{x})(\boldsymbol{x} - \overline{\boldsymbol{x}}) = (\boldsymbol{x} - \overline{\boldsymbol{x}})^{\mathsf{T}} \frac{D^2 f(\hat{\boldsymbol{x}}) - D^2 f(\tilde{\boldsymbol{x}})}{2}(\boldsymbol{x} - \overline{\boldsymbol{x}}) \geqslant m\|\boldsymbol{x} - \overline{\boldsymbol{x}}\|^2,$$

and then

$$\|Df(\boldsymbol{x})\|\|\boldsymbol{x} - \overline{\boldsymbol{x}}\| \geqslant (\boldsymbol{x} - \overline{\boldsymbol{x}})^{\mathsf{T}} \frac{D^2 f(\hat{\boldsymbol{x}}) - D^2 f(\tilde{\boldsymbol{x}})}{2}(\boldsymbol{x} - \overline{\boldsymbol{x}}) \geqslant m\|\boldsymbol{x} - \overline{\boldsymbol{x}}\|^2.$$

After dividing the sides of the last inequality by $m\|\boldsymbol{x} - \overline{\boldsymbol{x}}\|$, we obtain

$$\|\boldsymbol{x} - \overline{\boldsymbol{x}}\| \leqslant \frac{1}{m}\|Df(\boldsymbol{x})\|.$$

Using the Taylor formula again, due to the convexity of $f$ in $S$, we obtain

$$f(\overline{\boldsymbol{x}}) = f(\boldsymbol{x}) + Df(\boldsymbol{x})(\overline{\boldsymbol{x}} - \boldsymbol{x}) + \frac{1}{2}(\overline{\boldsymbol{x}} - \boldsymbol{x})^{\mathsf{T}} D^2 f(\hat{\boldsymbol{x}})(\overline{\boldsymbol{x}} - \boldsymbol{x}) \geqslant f(\boldsymbol{x}) + Df(\boldsymbol{x})(\overline{\boldsymbol{x}} - \boldsymbol{x}).$$

Hence,

$$f(\boldsymbol{x}) - f(\overline{\boldsymbol{x}}) \leqslant Df(\boldsymbol{x})(\boldsymbol{x} - \overline{\boldsymbol{x}}),$$

i.e.,

$$f(\boldsymbol{x}) - f(\overline{\boldsymbol{x}}) = |f(\boldsymbol{x}) - f(\overline{\boldsymbol{x}})| \leqslant \|Df(\boldsymbol{x})\|\|\boldsymbol{x} - \overline{\boldsymbol{x}}\| \leqslant \frac{1}{m}\|Df(\boldsymbol{x})\|^2. \quad \square$$

**Lemma 26** *Let $S$ be convex and compact. Let $f$ be of class $C^2$ in $s$ and $m = \inf_{x \in S} m(\boldsymbol{x})$. Denote $M = \sup_{x \in S} M(\boldsymbol{x})$. Then, $M < +\infty$ and the sequence $(\boldsymbol{x}_k)_k$ obtained with the steepest descent method with the exact minimization rule there is*

$$f(\boldsymbol{x}_{k+1}) - f(\overline{\boldsymbol{x}}) \leqslant \left(1 - \frac{m}{2M}\right)\left(f(\boldsymbol{x}_k) - f(\overline{\boldsymbol{x}})\right).$$

*With the restricted minimization rule there is*

$$f(\boldsymbol{x}_{k+1}) - f(\overline{\boldsymbol{x}}) \leqslant \left(1 - m\gamma + \frac{mM\gamma^2}{2}\right)\left(f(\boldsymbol{x}_k) - f(\overline{\boldsymbol{x}})\right),$$

*where $\gamma = \min\{\frac{1}{M}, A\}$.*

**Proof.** Consider the exact rule. By the Taylor formula, for $\delta \geqslant 0$, we have

$$f\left(\boldsymbol{x}_k - \delta\left(Df(\boldsymbol{x}_k)\right)^{\mathsf{T}}\right) \leqslant f(\boldsymbol{x}_k) + Df(\boldsymbol{x}_k)\left(-\delta\left(Df(\boldsymbol{x}_k)\right)^{\mathsf{T}}\right) + \delta^2 \frac{M}{2}\|Df(\boldsymbol{x}_k)\|^2$$

$$= f(\boldsymbol{x}_k) - \delta\|Df(\boldsymbol{x}_k)\|^2 + \delta^2 \frac{M}{2}\|Df(\boldsymbol{x}_k)\|^2. \qquad (\odot)$$

The minimum of the right-hand side is taken at $\delta = \frac{1}{M}$. Recall that $\boldsymbol{x}_{k+1}$ is the minimum for $\alpha \geqslant 0$:

$$f(\boldsymbol{x}_{k+1}) = \inf_{\alpha \geqslant 0} f\left(\boldsymbol{x}_k - \alpha\left(Df(\boldsymbol{x}_k)\right)^{\mathsf{T}}\right).$$

Therefore,

$$f(\boldsymbol{x}_{k+1}) \leqslant f\left(\boldsymbol{x}_k - \frac{1}{M}\left(Df(\boldsymbol{x}_k)\right)^{\mathsf{T}}\right) \leqslant f(\boldsymbol{x}_k) - \frac{1}{2M}\|Df(\boldsymbol{x}_k)\|^2.$$

Now we subtract $f(\overline{\boldsymbol{x}})$ from both sides and we apply the inequality $\|Df(\boldsymbol{x}_k)\|^2 \geqslant m\left(f(\boldsymbol{x}_k) - f(\overline{\boldsymbol{x}})\right)$, being a consequence of Lemma 25:

$$f(\boldsymbol{x}_{k+1}) - f(\overline{\boldsymbol{x}}) \leqslant f(\boldsymbol{x}_k) - f(\overline{\boldsymbol{x}}) - \frac{m}{2M}\left(f(\boldsymbol{x}_k) - f(\overline{\boldsymbol{x}})\right).$$

A simple calculation yields the claim.

Now we consider the restricted rule. The point $\boldsymbol{x}_{k+1}$ is the minimum for $\alpha \in [0, A)$. The minimum on the right-hand side of $(\odot)$ is obtained with $\delta = \gamma$. The rest of the proof is identical to the proof for the exact rule. $\square$

The above considerations may be concluded as follows: by Lemma 25, the stop condition based on the norm of the gradient of $f$ is correct; it gives us estimations of the approximation error of $\overline{\boldsymbol{x}}$ by $\boldsymbol{x}_{k+1}$ as well as the approximation error of the minimal function value $f(\overline{\boldsymbol{x}})$ by $f(\boldsymbol{x}_{k+1})$. Note that "the more convex" is the function $f$ in a neighbourhood of $\overline{\boldsymbol{x}}$, i.e., the greater is $m$, the sharper is the dependency between the norm of the gradient and the distance between $\boldsymbol{x}_{k+1}$ and $\overline{\boldsymbol{x}}$. Lemma 26 suggests that the convergence is fastest if the function $f$ behaves similarly in all directions, i.e., the eigenvalues of the Hessian matrix are close to each other. Then, the quotient $\frac{m}{M}$ is great (close to 1), which decreases the contraction factor $1 - \frac{m}{2M}$. Thus, the steepest descent algorithm works best for functions such that their corresponding numbers $m$ and $M$ are of the same order of magnitude.

The correctness of the stop criterion for the Armijo rule still holds, as Lemma 25 does not depend on the step length. Lemma 26 for this case has to be modified; however, the (at least) linear rate of convergence remains.

What is the advantage of the Armijo rule? It is simple to implement, as it does not need methods of searching minima of functions of one variable. The actual rate of convergence depends on the parameters $s, \beta, \delta$. Alas, there is no general rule of choosing these parameters—user's experience and intuition must be employed.

## Newton's method

The descent methods described above choose the directions for the next point based on the Taylor expansion up to the first-order term:

$$f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + Df(\mathbf{x})\mathbf{d}.$$

The Newton method uses the expansion up to the second-order term of the function $f \colon \mathbb{R}^n \to \mathbb{R}$; assume that it is twice differentiable. Then, we use the approximation

$$f(\mathbf{x} + \mathbf{d}) \approx f(\mathbf{x}) + Df(\mathbf{x})\mathbf{d} + \frac{1}{2}\mathbf{d}^\mathsf{T} D^2 f(\mathbf{x})\mathbf{d}.$$

Instead of finding a minimum of $f$, we minimise the expression on the right-hand side above. For this to make sense we have to assume that the Hessian matrix $D^2(\mathbf{x})$ is positive-definite. With $f(\mathbf{x})$ fixed, the problem reduces to the following:

$$\begin{cases} h(\mathbf{d}) = \frac{1}{2}\mathbf{d}^\mathsf{T} D^2 f(\mathbf{x})\mathbf{d} + Df(\mathbf{x})\mathbf{d} \to \min, \\ \mathbf{d} \in \mathbb{R}^n. \end{cases}$$

As the Hessian is, by assumption, positive-definite, the problem above has the unique solution

$$\mathbf{d} = -\left(D^2 f(\mathbf{x})\right)^{-1}\left(Df(\mathbf{x})\right)^\mathsf{T}.$$

Note that if the gradient of $f$ at $\mathbf{x}$ is zero, then $\mathbf{d} = \mathbf{0}$ and we shall not leave the critical point. The Newton algorithm is

<u>Preparation:</u> Choose the initial point $\mathbf{x}_1$ and the parameter $\varepsilon > 0$.

<u>k-th step:</u> Take $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$, where $D^2 f(\mathbf{x}_k)\mathbf{d}_k = -\left(Df(\mathbf{x}_k)\right)^\mathsf{T}$.

<u>Stop condition:</u> $\|Df(\mathbf{x}_{k+1})\| \leqslant \varepsilon$.

There are a number of doubts about this algorithm. It is not clear whether the sequence $(\mathbf{x}_k)_k$ is convergent; moreover, it is easy to find a regular function, having a minimum, for which the method generates a divergent (unbounded) sequence. The theorem below describes a sufficient condition for the convergence.

<u>Theorem 66</u> *Let $f$ be a function of class $C^3$ in a neighbourhood of a local minimum $\overline{\mathbf{x}}$ and let the Hessian $D^2 f(\overline{\mathbf{x}})$ be positive-definite. Then, there exists $\delta > 0$ and $c > 0$ such that for any $\mathbf{x}_k \in B(\overline{\mathbf{x}}, \delta)$ there is*

$$\|\mathbf{x}_{k+1} - \overline{\mathbf{x}}\| \leqslant c\|\mathbf{x}_k - \overline{\mathbf{x}}\|^2.$$

<u>Proof.</u> The continuity of $D^2 f$ implies the existence of $\delta > 0$ such that for $\mathbf{x} \in B(\overline{\mathbf{x}}, \delta)$ the norms of $D^2 f(\mathbf{x})$ and $\left(D^2 f(\mathbf{x})\right)^{-1}$ are bounded and greater than some $r > 0$.

Expanding the gradient of $f$ at $\mathbf{x}_k$, we get

$$\left(Df(\mathbf{x}_k + \mathbf{h})\right)^\mathsf{T} = \left(Df(\mathbf{x}_k)\right)^\mathsf{T} + D^2 f(\mathbf{x}_k)\mathbf{h} + O(\|\mathbf{h}\|^2).$$

Substituting to the above $\mathbf{h} = -\mathbf{h}_k = -(\mathbf{x}_k - \overline{\mathbf{x}})$, we obtain

$$\left(Df(\mathbf{x}_k)\right)^\mathsf{T} - D^2 f(\mathbf{x}_k)\mathbf{h}_k + O(\|\mathbf{h}_k\|^2) = \left(Df(\overline{\mathbf{x}})\right)^\mathsf{T} = \mathbf{0}.$$

Let $\mathbf{x}_k \in B(\overline{\mathbf{x}}, \delta)$. Multiplying the sides of the above by $\left(D^2 f(\mathbf{x}_k)\right)^{-1}$, we obtain

$$\mathbf{0} = \left(D^2 f(\mathbf{x})\right)^{-1}\left(Df(\mathbf{x}_k)\right)^\mathsf{T} - \mathbf{h}_k + O(\|\mathbf{h}_k\|^2 = -\mathbf{d}_k - \mathbf{h}_k + O(\|\mathbf{h}_k\|^2) \qquad (*)$$

There is

$$-\mathbf{d}_k - \mathbf{h}_k = \mathbf{x}_k - \mathbf{x}_{k+1} - (\mathbf{x}_k - \overline{\mathbf{x}}) = -(\mathbf{x}_{k+1} - \overline{\mathbf{x}}) = -\mathbf{h}_{k+1}.$$

Therefore, by (*), there exists a constant $c > 0$ such that

$$\|\mathbf{h}_{k+1}\| \leqslant c\|\mathbf{h}_k\|^2;$$

hence, the rate of convergence is quadratic.

If $\mathbf{x}_k \in B(\overline{\mathbf{x}}, \frac{\alpha}{c})$, where $\alpha \in (0, 1)$ and $\frac{\alpha}{c} \leqslant \delta$, then the last inequality gives us also

$$\|\mathbf{h}_{k+1}\| \leqslant c\frac{\alpha}{c}\|\mathbf{h}_k\| = \alpha\|\mathbf{h}_k\|,$$

which means that the sequence of vectors $(\mathbf{h}_k)_k$, being differences of the points $\mathbf{x}_k$ and the solution $\overline{\mathbf{x}}$ tends to $\mathbf{0}$. $\square$

<u>Remark.</u> (1) If the point $\mathbf{x}_k$ is distant from the solution $\overline{\mathbf{x}}$, then the Hessian $D^2 f(\mathbf{x}_k)$ does not have to be positive-definite.
(2) If the Hessian is positive-definite at $\mathbf{x}_k$, then the vector $\mathbf{d}_k$ constructed in the k-th step has a descent direction for the function $f$:

$$Df(\mathbf{x}_k)\mathbf{d}_k = -Df(\mathbf{x}_k)\left(D^2 f(\mathbf{x}_k)^{-1}\right)^{-1}\left(Df(\mathbf{x}_k)\right)^\mathsf{T} < 0.$$

(3) Even if the Hessian is positive-definite at $\mathbf{x}_k$, there is no guarantee that $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$, because there is no minimization along the line having the descent direction. The length of the vector $\mathbf{d}_k$ may simply be too big.

(4) A drawback of the Newton method is that it may converge to a critical point being a local maximum or a saddle point (this may be the case if the Hessian is not positive-definite). The assumption of pseudoconvexity of $f$ guarantees that the critical point is a minimum.

A remedy to the above may be <u>minimization along the Levenberg–Marquardt trajectory</u>. The system of linear equations $D^2 f(\mathbf{x}_k)\mathbf{d} = -\left(Df(\mathbf{x}_k)\right)^\mathsf{T}$ may be replaced by the following:

$$\left(D^2 f(\mathbf{x}_k) + \nu I\right)\mathbf{d} = -\left(Df(\mathbf{x}_k)\right)^\mathsf{T},$$

where $\nu$ is a parameter and $I$ is the identity matrix. We can define a function of one variable,

$$g(\nu) \stackrel{\text{def}}{=} f\left(\mathbf{x}_k - \left(D^2 f(\mathbf{x}_k) + \nu I\right)^{-1}\left(Df(\mathbf{x}_k)\right)^\mathsf{T}\right).$$

Then, we can chose $\nu$ so as to minimise the function $g$; the point $\mathbf{x}_{k+1}$ is then the argument of $f$ corresponding to that $\nu$. Note that if $\nu = 0$, then the point $\mathbf{x}_{k+1}$ is the one obtained with the Newton method. The eigenvalues of the matrix $\left(D^2 f(\mathbf{x}_k) + \nu I\right)$ are the eigenvalues of $D^2 f(\mathbf{x}_k)$ increased by $\nu$; hence, there exists $\nu_0$ such that the matrix $\left(D^2 f(\mathbf{x}_k) + \nu I\right)$ is positive-definite for all $\nu > \nu_0$.

The Levenberg–Marquardt trajectory is the parametric curve made of the points $\mathbf{x}_k - \left(D^2 f(\mathbf{x}_k) + \nu I\right)^{-1}\left(Df(\mathbf{x}_k)\right)^\mathsf{T}$, where $\nu > \nu_0$. Note that by increasing $\nu$ we obtain the vectors $-\left(D^2 f(\mathbf{x}_k) + \nu I\right)^{-1}\left(Df(\mathbf{x}_k)\right)^\mathsf{T}$ whose directions tend to the direction of $-Df(\mathbf{x}_k)$, i.e., to the steepest descent direction, and their lengths tend to 0. Therefore, if the gradient of $f$ at $\mathbf{x}_k$ is nonzero, we have a guarantee of finding on the Levenberg–Marquardt trajectory a point $\mathbf{x}_{k+1}$ such that $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$. Minimization along Levenberg–Marquardt trajectories may also produce a divergent sequence of points or a sequence convergent to a saddle point.

Another drawback of the Newton method are the difficulty of finding a good initial point $\mathbf{x}_1$ and the parameter $\varepsilon$ for the stop condition. The Levenberg–Marquardt approach makes the choice of initial points considerably easier.

If the dimension $n$ of the problem is large, each iteration, involving the computation of $n^2$ coefficients of the Hessian and solving a system of $n$ equations, is costly.

## Conjugate directions and conjugate gradient methods

<u>Definition 38</u> *Let* $H$ *be a symmetric and positive-definite* $n \times n$ *matrix. Nonzero vectors* $\mathbf{d}_1, \ldots, \mathbf{d}_n$ *are called* <u>*conjugate with respect to the matrix*</u> $H$ *if*

$$\mathbf{d}_i^\mathsf{T} H \mathbf{d}_j = 0, \quad i, j \in \{1, \ldots, n\}, \quad i \neq j.$$

Note that if the vectors $\mathbf{d}_1, \ldots, \mathbf{d}_n$ are conjugate with respect to a matrix, then they are linearly independent.

Consider a quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\mathsf{T} H \mathbf{x} + \mathbf{b}^\mathsf{T}\mathbf{x} + c$ with a symmetric and positive-definite matrix $H$. The minimum of this function may be found using a <u>conjugate directions method</u>: subsequent points approximating the minimum are searched in the directions of vectors conjugate with respect to the matrix $H$. The algorithm is the following:

<u>Preparation:</u> Choose the initial point $\mathbf{x}_1$.

<u>k-th step:</u>

1. Choose the vector $\mathbf{d}_k$.
2. Take $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$, where $t_k$ is chosen by the exact minimization rule, $t_k = \arg\min\{f(\mathbf{x}_k + t\mathbf{d}_k): t \geqslant 0\}$

<u>Stop condition:</u> $\|Df(\mathbf{x}_{k+1})\| = 0$.

<u>Theorem 67</u> *A conjugate directions method with the exact minimization rule finds the minimum of a quadratic function* $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\mathsf{T} H \mathbf{x} + \mathbf{b}^\mathsf{T}\mathbf{x} + c$ *with a positive-definite matrix* $H$ *after at most* $n$ *steps.*

<u>Proof.</u> Denote $\mathbf{g}_k = Df(\mathbf{x}_k)$. It is easy to notice that

$$\mathbf{g}_{k+1}\mathbf{d}_k = 0, \quad k = 1, \ldots, n-1.$$

To this end, consider the function $h(t) = f(\mathbf{x}_k + t\mathbf{d}_k)$; it has the minimum at $t_k$, therefore $h'(t_k) = 0$. On the other hand, $h'(t) = D(\mathbf{x}_k + t\mathbf{d}_k)\mathbf{d}_k$ and $h'(t_k) = \mathbf{g}_{k+1}\mathbf{d}_k$.

For the quadratic function $f$ we obtain

$$g_{k+1} - g_k = Df(x_{k+1}) - Df(x_k) = (Hx_{k+1})^\top - (Hx_k)^\top = t_k d_k^\top H.$$

Using this equality we prove by induction that

$$g_{k+1} d_j = 0, \quad k = 1, \ldots, n-1, \ j = 1 \ldots, k. \tag{$\square$}$$

For $k = 1$ it is already proved. Suppose that $g_k d_j = 0$ for $j = 1, \ldots, k-1$. Then, for $j = 1, \ldots, k-1$

$$g_{k+1} d_j = (g_k + t_k d_k^\top H) d_j = g_k d_j + t_k d_k^\top H d_j = 0$$

because $g_k d_j = 0$ due to the inductive assumption and $d_k^\top H d_j = 0$ because the vectors $d_k$ and $d_j$ are conjugate with respect to $H$. Also, the equality $g_{k+1} d_k$ is already proved, which completes the proof of $(\square)$.

By $(\square)$, the derivatives of the function $f$ at $x_{k+1}$ in the directions of $d_1, \ldots, d_k$ are zero; hence,

$$Df(x_{k+1}) d = 0, \quad d \in \mathrm{lin}\{d_1, \ldots, d_k\}.$$

Let

$$K_k = x_{k+1} + \mathrm{lin}\{d_1, \ldots, d_k\} = x_1 + \mathrm{lin}\{d_1, \ldots, d_k\}.$$

Let $F = f|_{K_k}$. The function $F$ is convex, because $f$ is convex and $K_k$ is a convex set. All directional derivatives equal to $0$ is a sufficient condition for a minimum of a convex function. Therefore,

$$x_{k+1} = \arg\min\{ F(x) : x \in K_k \} = \arg\min\{ f(x) : x \in x_1 + \mathrm{lin}\{d_1, \ldots, d_k\} \}.$$

To complete the proof we notice that $K_n = \mathbb{R}^n$. $\square$

The method analysed in the last theorem is very effective, as we can find the minimum in at most $n$ steps. Its drawback is the necessity of finding the entire set of conjugate vectors of the matrix $H$ (these may be the eigenvectors of $H$). The underline{conujgate gradient method} described below finds conjugate vectors in consecutive steps of the algorithm. This method is also known as the Fletcher–Reeves algorithm.

Preparation: Choose the initial point $x_1$.

1-st step: Take $d_1 = -g_1^\top$ (the vector with the steepest descent direction) and

$$x_2 = x_1 + t_1 d_1, \quad t_1 = \arg\min\{f(x_1 + t d_1)\}.$$

k-th step, $k > 1$: (we already know the vectors $d_1, \ldots, d_{k-1}$)

1. Take

$$\beta_{k-1} = \frac{g_k g_k^\top}{g_{k-1} g_{k-1}^\top}, \quad d_k = -g_k^\top + \beta_k d_{k-1}.$$

2. Take $x_{k+1} = x_k + t_k d_k$, where $t_k = \arg\min\{f(x_k + t d_k)\}$.

Stop condition: $\|Df(x_{k+1})\| = 0$.

Theorem 68 *The Fletcher–Reeves algorithm with the exact minimization rule applied to a quadratic function $f \colon \mathbb{R}^n \to \mathbb{R}$ with a positive Hessian matrix $H$ constructs vectors $d_1, d_2, \ldots$ conjugate with respect to the matrix $H$. Moreover, for $i \leqslant m = \max\{i \colon g_i \neq 0^\top\}$ there is*

$$g_i g_j^\top = 0, \quad j = 1, \ldots, i-1, \tag{$\oplus$}$$

*and*

$$g_i d_i = -g_i g_i^\top. \tag{$\otimes$}$$

Proof. If $m = 0$, then the initial point is the solution and there is nothing to prove. Let $m \geqslant 1$. The proof is done by the induction with respect to $i$.

If $i = 1$, then we only need to prove that $g_i d_i = -g_i g_i^\top$, which is obvious due to $d_1 = -g_1^\top$.

Suppose that the vectors $d_1, \ldots, d_i$ are conjugate with respect to $H$ and the equalities $(\oplus)$ and $(\otimes)$ are satisfied for some $i < m$. From the proof of Theorem 67 we know that for a quadratic function $f$

$$g_{i+1} - g_i = t_i d_i^\top H. \tag{$\ominus$}$$

Using this and $g_{i+1} d_i = 0$, we obtain

$$0 = g_{i+1} d_i = g_i d_i + t_i d_i^\top H d_i.$$

From the above we obtain

$$t_i = -\frac{g_i d_i}{d_i^\mathsf{T} H d_i} = \frac{g_i g_i^\mathsf{T}}{d_i^\mathsf{T} H d_i}; \qquad (\odot)$$

the last equality is based on the inductive assumption.

For $j < i$, by the inductive assumption and Step 2 of the algorithm we have

$$g_{i+1} g_j^\mathsf{T} = g_i g_j^\mathsf{T} + t_i d_i^\mathsf{T} H g_j = g_i g_j^\mathsf{T} + t_i d_i^\mathsf{T} H(d_j - \beta_{j-1} d_{j-1})$$
$$= g_i g_j^\mathsf{T} - t_i d_i^\mathsf{T} H d_j + t_i \beta_{j-1} d_i^\mathsf{T} H d_{j-1} = 0,$$

because the first term of the last sum is $0$ by the inductive assumption (equality $(\oplus)$) and the other two terms are $0$ because the vectors $d_i$, $d_j$ and $d_{j-1}$ are conjugate with respect to $H$.

A similar calculation with $i = j$, using the inductive assumption and $(\odot)$ gives us

$$g_{i+1} g_i^\mathsf{T} = g_i g_i^\mathsf{T} - t_i d_i^\mathsf{T} H d_i + t_i \beta_{i-1} d_i^\mathsf{T} H d_{i-1} = g_i g_i^\mathsf{T} - \frac{g_i g_i^\mathsf{T}}{d_i^\mathsf{T} H d_i} d_i^\mathsf{T} H d_i = 0.$$

We proved that $g_{i+1} g_j^\mathsf{T} = 0$ for $j = 1, \dots, i$, which is the inductive step of the proof of $(\oplus)$.

It remains to be proved that the vectors $d_1, \dots, d_m$ are conjugate with respect to $H$. Using $(\ominus)$ and the formula for $d_{i+1}$ we obtain

$$d_{i+1}^\mathsf{T} H d_j = -g_{i+1} H d_j + \beta_i d_i^\mathsf{T} H d_j = -\frac{1}{t_j} g_{i+1}(g_j - g_{j+1}) + \beta_i d_i^\mathsf{T} H d_j$$

$$= -\frac{1}{t_j} g_{i+1} g_{j+1}^\mathsf{T} + \beta_i d_i^\mathsf{T} H d_j.$$

For $j < i$ we obtain at once $d_{i+1}^\mathsf{T} H d_j = 0$, because, by inductive assumption, $g_{i+1} g_{j+1}^\mathsf{T} = 0$ and $d_i^\mathsf{T} H d_j = 0$.

For $j = i$, using the formulae for $t_i$ and $\beta_i$, we obtain

$$d_{i+1}^\mathsf{T} H d_j = -\frac{1}{t_j} g_{i+1} g_{j+1}^\mathsf{T} + \frac{g_{i+1} g_{i+1}^\mathsf{T}}{g_i g_i^\mathsf{T}} d_i^\mathsf{T} H d_j$$

$$= -\frac{1}{t_j} g_{i+1} g_{j+1}^\mathsf{T} + \frac{g_{i+1} g_{i+1}^\mathsf{T}}{g_i g_i^\mathsf{T}} \frac{g_i g_i^\mathsf{T}}{t_i} = 0.$$

We proved that

$$d_{i+1}^\mathsf{T} H d_j = 0, \quad j = 1, \dots, i.$$

The last thing to be proved is the equality $(\otimes)$. Using the formula for $d_{i+1}$ and $(\square)$, we obtain

$$g_{i+1} d_{i+1} = g_{i+1}(-g_{i+1}^\mathsf{T} + \beta_i d_i) = -g_{i+1} g_{i+1}^\mathsf{T} + \beta_i g_{i+1} d_i = -g_{i+1} g_{i+1}^\mathsf{T}.$$

This completes the entire inductive step. □

The conjugate gradient method is a powerful method of solving systems of linear equations with a symmetric and positive-definite $n \times n$ matrix, where $n$ is large. It is most often used to solve systems with millions of unknown variables obtained by discretization of partial differential equations, which is beyond the scope of this lecture. The idea is to find the minimum of the quadratic function $\frac{1}{2} x^\mathsf{T} A x - b^\mathsf{T} x$, whose gradient is $x^\mathsf{T} A - b^\mathsf{T}$; clearly, at the minimum $\overline{x}$ the gradient is zero, i.e., $A \overline{x} = b$. It turns out that the points obtained in consecutive iterations initially approach the solution, but the rounding errors (always present if floating point arithmetic is used) destroy the convergence; from a certain step the distance between the consecutive points and the solution may (rapidly) increase. Therefore, for huge systems of equations the conjugate gradient method is used as an iterative method: the computations are broken after $m$ iterations, where $m$ is much smaller than $n$.

If $f$ is not a quadratic function, then the Fletcher–Reeves conjugate gradient method has to be modified. Even without rounding errors we cannot expect the method to find the minimal point of a function of $n$ variables after $n$ iterations. Therefore the stop condition has to be based on a test, e.g., comparing the length of the gradient with some tolerance threshold:

Preparation: Choose the initial point $x_1$ and the parameter $\varepsilon > 0$.

1-st step: Take $d_1 = -g_1^\mathsf{T}$ (the vector with the steepest descent direction) and

$$x_2 = x_1 + t_1 d_1, \quad t_1 = \arg\min\{f(x_1 + t d_1)\}.$$

k-th step, k > 1: (we already know the vectors $d_1, \dots, d_{k-1}$)

1. Take

$$\beta_{k-1} = \frac{g_k g_k^\mathsf{T}}{g_{k-1} g_{k-1}^\mathsf{T}}, \quad d_k = -g_k^\mathsf{T} + \beta_k d_{k-1}.$$

2. Take $x_{k+1} = x_k + t_k d_k$, where $t_k = \arg\min\{f(x_k + t d_k) : t \geqslant 0\}$.

Stop condition: $\|Df(x_{k+1})\| \leqslant \varepsilon$.

Note that making more than $n$ iterations makes no sense for a quadratic function $f$. Therefore many implementations of the method make a "reset" every $m$ iterations, where $m \leqslant n$; the vector $\mathbf{d}_k$ after the reset is $-\mathbf{g}_k^\mathsf{T}$, i.e., it has the steepest descent direction.

In general, the choice of $\beta_{k-1}$ as above does not guarantee the vector $\mathbf{d}_k$ to have a descent direction, which is why the formula for $\beta_{k-1}$ is subject to various modifications. A modification often giving better results than the original Fletcher–Reeves method is the following:

$$\beta_{k-1} = \frac{\mathbf{g}_k(\mathbf{g}_k - \mathbf{g}_{k-1})^\mathsf{T}}{\mathbf{g}_{k-1}\mathbf{g}_{k-1}^\mathsf{T}}.$$

If $f$ is a quadratic function, then $\mathbf{g}_k\mathbf{g}_{k-1}^\mathsf{T} = 0$, and so this formula is equivalent to the original one.

Example. Figure 11 shows level sets of the Rosenbrock "banana valley" function,

$$f(x, y) = (x - 1)^2 + 100(x^2 - y)^2,$$

whose unique minimum is $\overline{x} = (1, 1)$. This function is known as troublesome in numerical optimization, which is why it is one of popular, or even classical tests for various algorithms. The initial point $\mathbf{x}_1 = (-0.5, 0.5)$ was taken for three algorithms discussed above.

The steepest descent method with restricted minimization rule with $A = 0.5$ yields the sequence of points shown in figure (a); note the very slow convergence to the minimum; the first 150 points have been plotted. One can see that improving the accuracy of computing minima in the steepest descent directions does not accelerate the convergence.

Figure (b) shows the sequence constructed with the Newton method augmented with minimization along Levenberg–Marquardt trajectories; three such trajectories were necessary, and the method has dealt with the problem pretty well.

In Figure (c) we can see the sequence constructed using the modified Fletcher–Reeves method, with reset every 2 iterations. After the initial quick progress the algorithm got stuck at a considerable distance from the minimal point.
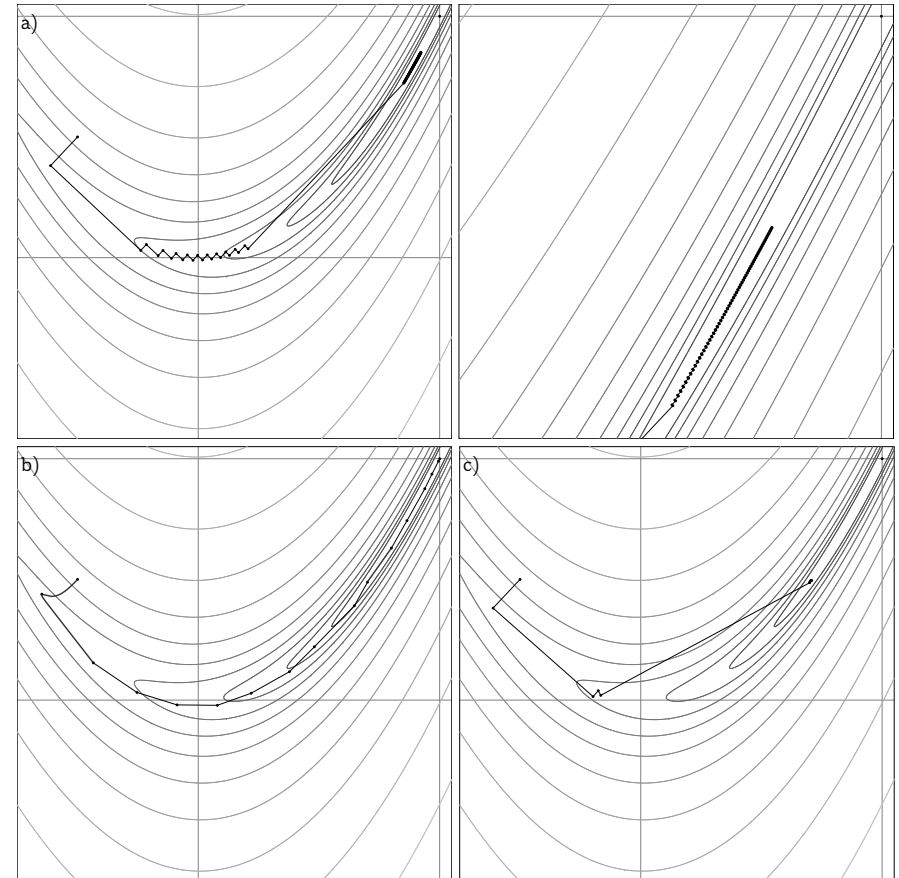


Figure 11: Searching the minimum of the Rosenbrock function

# 14. Algorithms for optimization with constraints

Now we focus on numerical methods for optimization problems with inequality constraints. We shal see problems arising when the steepest descent methods and their naive modifications are used and we shall see a more efective, though more complicated approach.

The problem considered here is

$$\begin{cases} f(\boldsymbol{x}) \to \min, \\ g_i(\boldsymbol{x}) \leqslant 0, \quad i = 1, \ldots, m, \\ \boldsymbol{x} \in \mathbb{R}^n, \end{cases}$$

where $f, g_1, \ldots, g_m \colon \mathbb{R}^n \to \mathbb{R}$. Thus, the feasible set is

$$W = \{\, \boldsymbol{x} \in \mathbb{R}^n \colon g_1(\boldsymbol{x}) \leqslant 0, \ldots, g_m(\boldsymbol{x}) \leqslant 0 \,\}.$$

In our considerations we use the notion of <u>feasible directions</u>, defined as follows:

$$F(\boldsymbol{x}) = \{\, \boldsymbol{d} \in \mathbb{R}^n \colon \boldsymbol{d} \neq \boldsymbol{0} \text{ and there exists } \lambda^* > 0$$
$$\text{such that } \boldsymbol{x} + \lambda \boldsymbol{d} \in W \text{ for all } \lambda \in [0, \lambda^*] \,\}.$$

## Zoutendijk algorithm for affine constraints

Consider a simple modification of the steepest descent method. If a point $\boldsymbol{x}$ is in the interior of $W$, then we can move in the steepest descent direction until we hit the boundary. If the point is already at the boundary, then it is natural to choose a direction of the possibly steep descent being feasible. Such a direction is called a <u>feasible descent direction</u> at the point $\boldsymbol{x}$; a vector $\boldsymbol{d}$ having such a direction satisfies the inequality $Df(\boldsymbol{x})\boldsymbol{d} < 0$.

It turns out that if the constraint functions are affine, then this idea works quite well; it is named the <u>Zoutendijk algorithm</u>. Recall that problems with affine constraints are simpler than the general case. This simplicity makes it possible to extend the analysis to problems with mixed affine constraints, i.e.,

$$\begin{cases} f(\boldsymbol{x}) \to \min, \\ A\boldsymbol{x} \leqslant \boldsymbol{b}, \\ Q\boldsymbol{x} = \boldsymbol{a}, \\ \boldsymbol{x} \in \mathbb{R}^n. \end{cases} \tag{*}$$

Here $A$ is an $m \times n$ matrix, $Q$ is $l \times n$, $\boldsymbol{b} \in \mathbb{R}^m$, $\boldsymbol{a} \in \mathbb{R}^l$. The following lemma characterises the set of feasible descent directions. Its proof is left as an exercise.

**Lemma 27** *Let $\boldsymbol{x}$ be a feasible point for the problem (\*). Assume that the matrix $A$ and the vector $\boldsymbol{b}$ may be divided into blocks $A_1$, $A_2$ and $\boldsymbol{b}_1$, $\boldsymbol{b}_2$ such that $A\boldsymbol{x} \leqslant \boldsymbol{b}$ is the conjunction of $A_1\boldsymbol{x} = \boldsymbol{b}_1$, $A_2\boldsymbol{x} < \boldsymbol{b}_2$ (depending on the set of active constraints at $\boldsymbol{x}$, this may require renumbering of the constraints). The vector $\boldsymbol{d} \in \mathbb{R}^n$ has a feasible direction at $\boldsymbol{x}$ if $A_1\boldsymbol{d} \leqslant \boldsymbol{0}$ and $Q\boldsymbol{d} = \boldsymbol{0}$. If in addition $Df(\boldsymbol{x})\boldsymbol{d} < 0$, then $\boldsymbol{d}$ has a feasible descent direction.*

How to choose the best descent direction at $\boldsymbol{x}$? It would be the simplest to solve the problem

$$Df(\boldsymbol{x})\boldsymbol{d} \to \min, \quad \boldsymbol{d} \in F(\boldsymbol{x}), \quad \|\boldsymbol{d}\| \leqslant 1. \tag{**}$$

The restriction for the norm of $\boldsymbol{d}$ is indispensable. Without it, for any vector $\boldsymbol{d}$ having a feasible descent direction there is $\lim_{\lambda \to \infty} Df(\boldsymbol{x})\lambda\boldsymbol{d} = -\infty$ and the problem above has no solution.

Using the block $A_1$ of the matrix $A$ as in Lemma 27, we can rewrite the problem (\*\*) in the form

$$\begin{cases} Df(\boldsymbol{x})\boldsymbol{d} \to \min, \\ A_1\boldsymbol{d} \leqslant \boldsymbol{0}, \\ Q\boldsymbol{d} = \boldsymbol{0}, \\ \boldsymbol{d}^\mathsf{T}\boldsymbol{d} \leqslant 1. \end{cases}$$

Note that the only nonlinear part above is the norm restriction. In practice, without loss of the algorithm quality it is replaced by linear restrictions, which make it possible to use fast methods of linear optimization (e.g. the simplex algorithm). The most popular replacements for the Euclidean norm $\|\boldsymbol{d}\|_2$ are

- the maximum norm, $\|\boldsymbol{d}\|_\infty = \max_j |d_j|$, which gives us

$$-1 \leqslant d_j \leqslant 1, \quad j = 1, \ldots, n,$$

- the first norm, $\|\boldsymbol{d}\|_1 = \sum_j |d_j|$; in this case we have the inequalities

$$\begin{cases} \sum_{j=1}^n \eta_j \leqslant 1, \\ -\eta_j \leqslant d_j \leqslant \eta_j, \quad j = 1, \ldots, n, \end{cases}$$

where $\eta_1, \ldots, \eta_n$ are new, auxiliary variables.

Below we consider the algorithm with the maximum norm restriction:

$$\begin{cases} Df(\mathbf{x})\mathbf{d} \to \min, \\ A_1\mathbf{d} \leqslant \mathbf{0}, \\ Q\mathbf{d} = \mathbf{0}, \\ -1 \leqslant d_j \leqslant 1, \quad j = 1, \dots, n. \end{cases}$$

The algorithm is the following:

Preparation: Choose the initial point $\mathbf{x}_1$.

k-th step:

1. Given a point $\mathbf{x}_k$, find the blocks $A_1$, $A_2$ of the matrix $A$ and the blocks $\mathbf{b}_1$, $\mathbf{b}_2$ of $\mathbf{b}$ so as to obtain $A_1\mathbf{x} = \mathbf{b}_1$ and $A_2\mathbf{x} < \mathbf{b}_2$ (like in Lemma 27).

2. Choose the vector $\mathbf{d}_k$ by solving the problem

$$\begin{cases} Df(\mathbf{x}_k)\mathbf{d} \to \min, \\ A_1\mathbf{d} \leqslant \mathbf{0}, \\ Q\mathbf{d} = \mathbf{0}, \\ -1 \leqslant d_j \leqslant 1, \quad j = 1, \dots, n. \end{cases} \qquad (\overset{**}{*})$$

3. If $Df(\mathbf{x}_k)\mathbf{d}_k = 0$, then stop, as the point $\mathbf{x}_k$ satisfies the necessary first-order condition. Else continue.

4. Take $\alpha_k = \arg\min_{\alpha \in [0, A_k]} f(\mathbf{x}_k + \alpha\mathbf{d}_k)$, where $A_k$ is the greatest number such that the line segment $\overline{\mathbf{x}_k, \mathbf{x}_k + A_k\mathbf{d}_k}$ is contained in the feasible set $W$.

5. Take $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k\mathbf{d}_k$.

Let's take a look at the choice of the vector $\mathbf{d}_k$. The vector $\mathbf{d}_k = \mathbf{0}$ satisfies all restrictions, therefore the minimal value of the function $Df(\mathbf{x}_k)\mathbf{d}_k$ is less than or equal to $0$.

Lemma 28 *The necessary first-order condition is satisfied at $\mathbf{x}_k$ if and only if the solution $\mathbf{d}_k$ of $(\overset{**}{*})$ satisfies the equality $Df(\mathbf{x}_k)\mathbf{d}_k = 0$*

Proof. Recall that the necessary first-order condition for the problem (*) is satisfied at $\mathbf{x}_k$ if and only if there exist vectors $\boldsymbol{\mu} \in [0, \infty)^{m_1}$ and $\boldsymbol{\lambda} \in \mathbb{R}^l$ such that

$$Df(\mathbf{x}_k) + \boldsymbol{\mu}^\mathsf{T}A_1 + \boldsymbol{\lambda}^\mathsf{T}Q = \mathbf{0}^\mathsf{T},$$

where the $m_1 \times n$ matrix $A_1$ corresponds to the constraints active at $\mathbf{x}_k$.

By the Farkas lemma (Lemma 15), if this system has a solution, then the system

$$\begin{cases} Df(\mathbf{x}_k)\mathbf{d} < 0, \\ A_1\mathbf{d} \leqslant \mathbf{0}, \\ Q\mathbf{d} = \mathbf{0} \end{cases} \qquad (\overset{**}{**})$$

is inconsistent. Due to the observation that $\mathbf{d} = \mathbf{0}$ is a solution of the system above with the inequality $Df(\mathbf{x}_k)\mathbf{d} < 0$ replaced by the equality $Df(\mathbf{x}_k)\mathbf{d} = 0$, the right implication is proved.

To prove the left implication, we notice that if $Df(\mathbf{x}_k)\mathbf{d} = 0$, then $\mathbf{d}_k$ does not satisfy $(\overset{**}{**})$. Using again the Farkas lemma we notice that the necessary first-order condition is then satisfied. $\square$

## Zoutendijk algorithm for nonlinear constraints

One can wonder if the Zoutendijk algorithm works just as well for nonlinear constraints:

Preparation: Choose the initial point $\mathbf{x}_1$.

k-th step:

1. Given a point $\mathbf{x}_k$, choose the vector $\mathbf{d}_k$ by solving the problem

$$Df(\mathbf{x}_k)\mathbf{d} \to \min, \quad \mathbf{d} \in F(\mathbf{x}_k), \quad \|\mathbf{d}\| \leqslant 1.$$

2. If $Df(\mathbf{x}_k)\mathbf{d}_k = 0$, then stop, as the point $\mathbf{x}_k$ satisfies the necessary first-order condition. Else continue.

3. Take $\alpha_k = \arg\min_{\alpha \in [0, A_k]} f(\mathbf{x}_k + \alpha\mathbf{d}_k)$, where $A_k$ is the greatest number such that the line segment $\overline{\mathbf{x}_k, \mathbf{x}_k + A_k\mathbf{d}_k}$ is contained in the feasible set $W$.

4. Take $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k\mathbf{d}_k$.

Figure 12 shows the steepest descent methods used to find the minimum of the function $f(x_1, x_2) = -2x_1 - x_2$ in two different sets; the minimal point is denoted by $\overline{\mathbf{x}}$. In the case (a) the minimum is found in the second step. If the set $W$ is not convex, as in the case (b), the algorithm may go to a dead end. This is, alas,
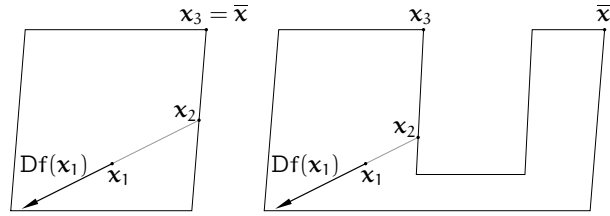
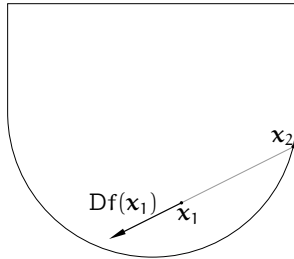Figure 12: Steepest descent method for two different feasible sets



Figure 13: Failure of the descent method in a convex feasible set

a property of all descent algorithms. Therefore, we have to require that the feasible set be convex. Does it suffice? No. In Figure 13 we can see that the algorithm may fail even with a convex feasible set. The steepest descent direction cannot be found, because the set of feasible descent directions $F(\boldsymbol{x}_2)$ is not closed. It is easy to find a way out: we need to choose such a vector $\boldsymbol{d}_k$ that the descent is steep and also a long part (line segment) of the halfline having the direction of $\boldsymbol{d}_k$ be contained in the set $W$. We still pursue the simplicity, i.e., linearity, of the optimization problem posed in order to choose $\boldsymbol{d}_k$. The solution is prompted by the following lemma:

<u>Lemma 29</u> *Let $\boldsymbol{x}$ be a feasible point. If the functions $f$ and $g_i$ for $i \in I(\boldsymbol{x})$ are differentiable at $\boldsymbol{x}$ and the functions $g_i$ for $i \notin I(\boldsymbol{x})$ are continuous, then any vector $\boldsymbol{d}$ such that $Df(\boldsymbol{x})\boldsymbol{d} < 0$ and $Dg_i(\boldsymbol{x})\boldsymbol{d} < 0$ for $i \in I(\boldsymbol{x})$ has a feasible descent direction.*

<u>Proof.</u> First we prove that $\boldsymbol{d}$ has a feasible direction, i.e., for a sufficiently small $\lambda > 0$ there is $\boldsymbol{x} + \lambda\boldsymbol{d} \in W$. For $i \notin I(\boldsymbol{x})$, due to the continuity of $g_i$ we have

$g_i(\boldsymbol{x} + \lambda\boldsymbol{d}) < 0$. For $i \in I(\boldsymbol{x})$ we have

$$g_i(\boldsymbol{x} + \lambda\boldsymbol{d}) = g_i(\boldsymbol{x}) + \lambda Dg_i(\boldsymbol{x})\boldsymbol{d} + o(\lambda\boldsymbol{d}).$$

There is $g_i(\boldsymbol{x}) = 0$, $Dg_i(\boldsymbol{x}) < 0$ and $o(\lambda\boldsymbol{d}) \to 0$ for $\lambda \to 0$; hence, $g_i(\boldsymbol{x} + \lambda\boldsymbol{d}) < 0$ for all sufficiently small positive $\lambda$.

The vector $\boldsymbol{d}$ has a feasible direction. By the asumptions of the lemma it has also a descent direction. There is

$$f(\boldsymbol{x} + \lambda\boldsymbol{d}) = f(\boldsymbol{x}) + \lambda Df(\boldsymbol{x})\boldsymbol{d} + o(\lambda\boldsymbol{d}).$$

Hence, for $\lambda$ sufficiently small there is $f(\boldsymbol{x} + \lambda\boldsymbol{d}) < f(\boldsymbol{x})$. □

Lemma 29 gives us only a *sufficient* condition. There are optimization problems with inequality constraints such that one of feasible descent directions does not satisfy the lemma's assumptions.

Choosing a vector $\boldsymbol{d} \in \mathbb{R}^n$ such that $Df(\boldsymbol{x})\boldsymbol{d} < 0$ and $Dg_i(\boldsymbol{x})\boldsymbol{d} < 0$ for $i \in I(\boldsymbol{x})$ may be done by solving the following problem:

$$\begin{cases} \max\{Df(\boldsymbol{x})\boldsymbol{d}, Dg_i(\boldsymbol{x})\boldsymbol{d}, i \in I(\boldsymbol{x})\} \longrightarrow \min, \\ -1 \leqslant d_j \leqslant 1, \quad j = 1, \ldots, n. \end{cases}$$

The target function above is tough to implement; we can reduce its minimization to the much simpler linear optimization problem

$$\begin{cases} \eta \to \min, \\ Df(\boldsymbol{x})\boldsymbol{d} \leqslant \eta, \\ Dg_i(\boldsymbol{x})\boldsymbol{d} \leqslant \eta, \quad i \in I(\boldsymbol{x}), \\ -1 \leqslant d_j \leqslant 1, \quad j = 1, \ldots, n. \end{cases} \qquad (\otimes)$$

Here the optimization is done with respect to *two* variables: $\boldsymbol{d} \in \mathbb{R}^n$ and $\eta \in \mathbb{R}$. Note that $\eta \leqslant 0$, because the pair $(\boldsymbol{d}, \eta) = (0, 0)$ satisfiest the constraints above. If the target function has a negative value, then by Lemma 29 the vector $\boldsymbol{d}$ has a feasible descent direction. If $\eta = 0$ is a solution and the linear independence condition is satisfied by the constraints, then the necessary first-order condition is satisfied at $\boldsymbol{x}$. The inverse implication is also true.

<u>Lemma 30</u> *If the linear independence condition is satisfied by the constraints at a feasible point $\boldsymbol{x}$ and $\eta = 0$ is the solution of the problem $(\otimes)$, then the necessary first-order condition is satisfied at $\boldsymbol{x}$. Also, if the first-order condition is satisfied at $\boldsymbol{x}$, then $\eta = 0$ is the solution of $(\otimes)$ (here the regularity of constraints at $\boldsymbol{x}$ needs not be assumed).*

Proof. If $\eta = 0$ is the solution, then the system $A\mathbf{d} < \mathbf{0}$, where $A$ is the matrix whose rows are gradients of $f$ and $g_i$, $i \in I(\mathbf{x})$, has no solution. By Gordan's lemma (Lemma 18), there exists $\mathbf{y} \geqslant \mathbf{0}$, $\mathbf{y} \neq \mathbf{0}$ such that $A^{\mathsf{T}}\mathbf{y} = \mathbf{0}$. Let $\mathbf{y} = \big(\hat{\mu}_0, \hat{\mu}_i, i \in I(\mathbf{x})\big)$ and let $\hat{\mu}_i = 0$ for $i \notin I(\mathbf{x})$. The equality $A^{\mathsf{T}}\mathbf{y} = \mathbf{0}$ may be rewritten as follows:

$$\hat{\mu}_0 Df(\mathbf{x}) + \sum_{i \in I(\mathbf{x})} \hat{\mu}_i Dg_i(\mathbf{x}) = \mathbf{0}^{\mathsf{T}}.$$

From the assumption of linear independence of gradients of the active constraints we conclude that $\hat{\mu}_0 \neq 0$. Taking $\mu_i = \hat{\mu}_i / \hat{\mu}_0$ for $i = 1, \ldots, m$, we obtain the Lagrange multipliers for the necessary first-order condition.

To prove the inverse implication we notice that if the necessary first-order condition is satisfied at $\mathbf{x}$, then the vector $\mathbf{y} = \big(1, \mu_i, i \in I(\mathbf{x})\big)$ satisfies the following: $\mathbf{y} \geqslant \mathbf{0}$, $\mathbf{y} \neq \mathbf{0}$ and $A^{\mathsf{T}}\mathbf{y} = \mathbf{0}$. By Lemma 18, there is no $\mathbf{d} \in \mathbb{R}^n$ such that $A\mathbf{d} < \mathbf{0}$. Then, $\eta = 0$ is the solution of $(\otimes)$. $\square$

The complete Zoutendijk algorithm for nonlinear problems with nonlinear constraints is the following:

Preparation: Choose the initial point $\mathbf{x}_1$.

k-th step:

1. Given a point $\mathbf{x}_k$, choose the vector $\mathbf{d}_k$ by solving the problem
$$\begin{cases} \eta \to \min, \\ Df(\mathbf{x}_k)\mathbf{d} \leqslant \eta, \\ Dg_i(\mathbf{x}_k)\mathbf{d} \leqslant \eta, \quad i \in I(\mathbf{x}_k), \\ -1 \leqslant d_j \leqslant 1, \quad j = 1, \ldots, n. \end{cases}$$

2. If $\eta = 0$, then stop, as the point $\mathbf{x}_k$ satisfies the necessary first-order condition. Else continue.

3. Take $\alpha_k = \arg\min_{\alpha \in [0, A_k]} f(\mathbf{x}_k + \alpha\mathbf{d}_k)$, where $A_k$ is the greatest number such that the line segment $\overline{\mathbf{x}_k, \mathbf{x}_k + A_k\mathbf{d}_k}$ is contained in the set $W$.

4. Take $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k\mathbf{d}_k$.

Example. Consider the problem
$$\begin{cases} 2x_1^2 + 2x_2^2 - 2x_1 x_2 - 4x_1 - 6x_2 \to \min, \\ x_1 + 5x_2 \leqslant 5, \\ 2x_1^2 - x_2 \leqslant 0, \\ x_1 \geqslant 0, \ x_2 \geqslant 0. \end{cases}$$

The Zoutendijk algorithm with the initial point $\mathbf{x}_1 = (0, 0.75)$ generates the following sequence of points:

$$\mathbf{x}_2 = (0.2803, 0.5477),$$
$$\mathbf{x}_3 = (0.5555, 0.8889),$$
$$\mathbf{x}_4 = (0.6479, 0.8397),$$
$$\mathbf{x}_5 = (0.6302, 0.8740).$$

As we can see, this sequence shows considerable oscillations in the feasible set, see Figure 14. This is a typical behaviour of descent methods for problems with constraints.
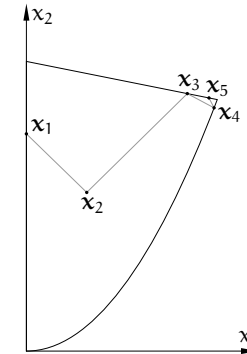


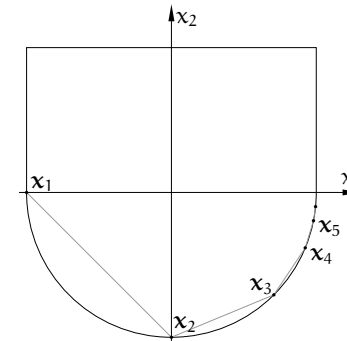Figure 14: Oscillations of the sequence generated by the Zoutendijk algorithm



Figure 15: Convergence to a point not being a solution

The Zoutendijk algorithm may fail even for rather simple problems. Consider

searching the minimum of the linear function $f(x_1, x_2) = -2x_1 - x_2$ in the set shown in Figure 15; the minimum is at $\overline{x} = (1, 1)$. Beginning at $x_1 = (-1, 0)$, we obtain a sequence of points converging to $(1, 0)$; thus, the algorithm will not approach the solution $\overline{x}$. Moreover, the value of $f$ at $(1, 0)$ is $-2$, while $f(\overline{x}) = -3$. However, a slight modification may improve the Zoutendijk algorithm.

## Topkis–Veinott modification

In 1967 Topkis and Veinott suggested a modification of the method of choosing $d_k$ in the Zoutendijk algorithm:

$$\begin{cases} \eta \to \min, \\ Df(x_k)d \leqslant \eta, \\ Dg_i(x_k)d \leqslant \eta - g_i(x_k), \quad i = 1, \ldots, m, \\ -1 \leqslant d_j \leqslant 1, \quad j = 1, \ldots, n. \end{cases} \qquad (\oplus)$$

The inequalities imposed for gradients of constraints include all constraints; for active constraints, $i \in I(x_k)$, we have $g_i(x_k) = 0$ and thus the conditions have not been changed. For inactive constraints, $g_i(x) < 0$ and the right-hand sides of the inequalities are greater than $\eta$. If the value of $g_i$ at $x_k$ is great, then the inequality is almost irrelevant. If the value of $g_i$ is close to $0$, i.e., the constraint is "almost active", then the corresponding inequality has a significant influence on the choice of $d_k$. Moreover, in the implementation this modification helps finding the active constraints, as due to the inexact representation of real numbers (the floating point representation), we usually cannot obtain $g_i(x) = 0$.

The theorem, which characterises the effectiveness of this modification is given without proof:

<u>Theorem 69</u> *Assume that $f, g_1, \ldots, g_m$ are of class $C^1$. If the sequence $(x_k)_k$ constructed by the Zoutendijk algorithm with the Topkis–Veinott modification has a concentration point $x$, at which the linear independence regularity condition is satisfied, then the necessary first-order condition is satisfied at $x$.*

## Conclusion

The numerical methods described above make it possible to approximate points satisfying the necessary first-order condition. Then, Theorem 47 may be used to guarantee the optimality of those points. In particular, if the constraints are linear, it suffices to assume pseudoconvexity of the function $f$. Note that similar

assumptions were needed for problems without constraints. The assumption about convexity is a natural and often necessary condition for the numerical methods to work well.