Definition	Theorem
Level set of $f: \mathbb{R}^d \supset D \to \mathbb{R}$	Compact level sets and minima
Nonlinear optimisation	Nonlinear optimisation
Example	Theorem
Linear regression	For convex problems, local minima are global
NONLINEAR OPTIMISATION	Nonlinear optimisation
Theorem	Theorem
Strictly convex problems are uniquely solvable	Variational inequality for the directional derivative
NONLINEAR OPTIMISATION	Nonlinear optimisation
Theorem	Theorem
Second order necessary condition	Sufficient second order condition
Nonlinear optimisation	Nonlinear optimisation
Theorem	Theorem
Quadratic growth condition	Tangents of convex functions
Nonlinear optimisation	Nonlinear optimisation

Let  $f: \mathbb{R}^n \supset D \to \mathbb{R}$  by a continuous function and  $\Omega$  closed. If there exists a  $\omega \in \Omega$ , such that  $\mathcal{N}(f, f(\omega))$  is compact there exists a global minimum of f on  $\Omega$ .

Let  $a := \inf_{x \in \Omega} f(x) \leq f(\omega)$ . As  $\Omega$  is closed,  $N := \Omega \cap \mathcal{N}(f, f(\omega))$  is compact and we have  $a = \inf_{x \in N} f(x)$ . By the theorem of WEIERSTRASS there exists a  $\hat{x} \in \Omega$  with  $\inf_{x \in N} f(x) = f(\hat{x})$ .

Let f be convex,  $D \neq \emptyset$  open and  $\Omega \subset D$  convex. Any local minimum of f is global. The set of solutions is convex.

Let  $x \in \Omega$  be a local minima. Then  $\exists r > 0$  with  $f(x) \leq f(y)$ for all  $y \in \Omega \cap B(x, r)$ . Let  $y \in \Omega$  and t > 0 so small, that  $x_t := x + t(y - x) \in B(x, r)$ . Since  $\Omega$  is convex,  $x_t \in \Omega$  for all  $t \in [0, 1]$ . Since f is convex,  $f(x) \leq f(x_t) = f((1-t)x+ty) \leq$ (1-t)f(x) + tf(y), which yields  $f(x) \leq f(y)$ . If  $x, y \in \Omega$  are solutions, for all  $z \in \Omega$   $f((1-t)x + ty) \leq$  $(1-t)f(x) + tf(y) \leq (1-t)f(z) + tf(z) = f(z)$ , so (1-t)x + tyis a minimum, too.

If $x \in D$	is a local minimum $f$	and $f$ is dir	rectionally	differen-
tiable in	$f(x;h) \ge 0$ for all	$h \in \mathbb{R}^n$ .		

As *D* is open,  $\exists r > 0$  with  $f(y) \ge f(x)$  for all  $y \in B(x, r)$ . For  $h \in \mathbb{R}^n$  and small *t* we have  $x + th \in B(x, r)$  and thus  $f(x + th) - f(x) \ge 0$ , i.e.  $\frac{f(x+th) - f(x)}{t} \ge 0$ . We have  $f'(x; h) := \lim_{t \searrow 0} \frac{f(x+th) - f(x)}{t}$ . The absolute has a minimum in 0, but is not differentiable

there, but we have  $|\cdot|'(x;h) = |h| \ge 0$  for all h. If  $f \in C^1$  and x is a local min,  $f'(x;h) = \nabla f(x)^{\mathsf{T}}h \ge 0 \ \forall h \in \mathbb{R}^n$ (var. ineq.). Taking  $h = -\nabla f(x)^{\mathsf{T}}h$ , we get  $\nabla f(x) = 0$ .

Let f be  $C^2$  in a neighbourhood of  $x \in D$ ,  $\nabla f(x) = 0$  and f''(z) be **positive semidefinite for all**  $z \in B(x, \delta)$  with some  $\delta > 0$ . Then x is a local minimum of f. For  $y \in B(x, \delta)$  and  $\theta \in [0, 1]$ 

 $f(y) - f(x) = \underbrace{f'(x)}_{=0} (y - x) + \frac{1}{2} \underbrace{(y - x)}_{h}^{\mathsf{T}} f''(\underbrace{x + \theta(y - x)}_{=z \in B_{\mathfrak{T}}(x)})(y - x) \ge 0$ 

by TAYLOR's theorem.

x = 0 is a local minimum of  $f(x) \coloneqq x^{2p}$ , where  $p \in \mathbb{N}_{\geq 2}$ . We have f'(0) = f''(0) = 0, which is not positive definite.

Let f be differentiable on D. Then f is (strictly) convex on  $\Omega$  iff  $f(y) \stackrel{(>)}{\geq} f(x) + \nabla f(x)^{\mathsf{T}}(y-x)$  for all  $x, y \in \Omega$ .



For  $a \in \mathbb{R}$ ,  $\mathcal{N}(f, a) := \{x \in D : f(x) \leq a\}$  is the **level set** of f with respect to a.



$$f(a,b) := \sum_{k=1}^{m} (y_k - ax_k - b)^2 = \frac{1}{2} z^{\mathsf{T}} H z + b^{\mathsf{T}} z + c$$
 with

$$H := 2 \begin{pmatrix} \sum_{k=1}^{m} x_k^2 & \sum_{k=1}^{m} x_k \\ \sum_{k=1}^{m} x_k^2 & m \end{pmatrix}, b := -2 \begin{pmatrix} \sum_{k=1}^{m} x_k y_k \\ \sum_{k=1}^{m} y_k \end{pmatrix}$$

 $c := \sum_{k=1}^{m} y_k^2$ . If two  $x_k$  are different, H is positive definite, as all principal minors are positive (CS). Thus F is strictly convex and has a unique minimum.

Let f by strictly convex. If x is a minimum of f, it is unique and thus strict.

Let  $x \neq y \in \Omega$  be two (by the previous theorem, global) minima of f and  $a := \min_{x \in \Omega} f(x)$ . Then  $f\left(\frac{x+y}{2}\right) < \frac{f(x)+f(y)}{2} = a$ , which that only x and y are minima of f.

**Examples.** The exponential function is strictly convex (AM-GM), but has no minimum. If H is positive definite,  $\frac{1}{2}x^{\mathsf{T}}Hx + b^{\mathsf{T}}x$  is strictly convex.

Let f be  $C^2$  in a neighbourhood of  $x \in D$  and x a local minimum of f. Then we have  $\nabla f(x) = 0$  and that f''(x) is **positive semidefinite**.

For  $h \in \mathbb{R}^n$  let  $g(t) \coloneqq f(x + th)$ . Then  $g \in \mathcal{C}^2$  has a local minimum in t = 0. By TAYLOR  $\exists \theta \in [0, 1]$  with  $g(t) = g(0) + g'(0)t + \frac{t^2}{2}g''(\theta t)$ . As x is a local minimum of  $g, 0 \leq \frac{g(t) - g(0)}{t^2} = \frac{1}{2}g''(\theta t)$ . The continuity of g'' yields  $g''(0) = h^{\mathsf{T}}f''(x)h \ge 0$  for  $t \searrow 0$ .

 $f(x) = x^4$  has a global minimum in  $\tilde{x} = 0$ , but  $f''(\tilde{x}) = 0$ .

Let f be  $C^2$  in a neighbourhood of  $x \in D$ ,  $\nabla f(x) = 0$  and f''(x) positive definite. Then  $\exists r, a > 0$  such that  $f(y) \ge$   $f(x) + a \|y - x\|^2$  for all  $x \in B(x, r)$ , so x is a strict minimum. TAYLOR:  $f(y) = f(x) + \frac{1}{2}(y - x)f''(x + \theta(y - x))(y - x)$  and  $(y - x)f''(x + \theta(y - x))(y - x)$   $= \underbrace{(y - x)f''(x)(y - x)}_{\ge a \|y - x\|^2} + \underbrace{(y - x)\left[f''(x + \theta(y - x)) - f''(x)\right](y - x)}_{|\cdot| \le \frac{a}{2} \|y - x\|^2 \text{ for small } \|y - x\|^2, \text{ as } f \in C^2}$  $\ge \frac{a}{2} \|y - x\|^2.$ 

Theorem	Definition
Convex variational inequality	Descent direction
Nonlinear optimisation	Nonlinear optimisation
Examples + Lemma	Algorithm
Descent direction	General descent algorithm
NONLINEAR OPTIMISATION	Nonlinear optimisation
Assumptions	Definition
ALC and AFD	Efficient step size
Nonlinear optimisation	Nonlinear optimisation
Definition	Assumptions
(strictly) Gradient-related descent direction	(ALG) and (AHP)
Nonlinear optimisation	Nonlinear optimisation
Lemma	Тнеогем
Convergence results for general descent methods	Convergence of descent algorithms
Nonlinear optimisation	Nonlinear optimisation

$d \in \mathbb{R}^n$ with $\nabla f(x)^{T} d < 0$ is a descent direction of $f$ in $x$ . <b>For each of the second sec</b>	Let $f$ be differentiable in $D$ and convex in $\Omega \subset D$ . Then $x \in \Omega$ is a minimiser of $f$ if and only if $\nabla f(x)^{T}(y-x) \ge 0 \ \forall y \in \Omega$ . " $\implies$ ": If $x$ be a local solution, then $x + t(y - x) = (1 - t)x + ty \in \Omega$ for all $t \in [0, 1], y \in \Omega$ . For small $t > 0$ $\frac{f(x+t(y-x))-f(x)}{t} \ge 0$ . Take $t \searrow 0$ . (convexity of $f$ not needed) " $\Leftarrow$ ": As is $f$ convex and all tangents lie below the graph, we have $f(y) - f(x) \ge \nabla f(x)^{T}(y - x) \ge 0$ and by a previous theorem $x$ is a global minimum. For $x \in \operatorname{int}(\Omega)$ , we have $\nabla f(x)^{T} d \ge 0$ for all directions $d \in \mathbb{R}^n$ and thus $\nabla f(x) = 0$ .
<ol> <li>Choose x<sup>0</sup> ∈ ℝ<sup>n</sup> and set k := 0.</li> <li>If ∇f(x<sup>k</sup>) = 0 holds, stop.</li> <li>Compute a descent direction d<sup>k</sup> and a step size σ<sub>k</sub> such that f(x<sup>k</sup> + σ<sub>k</sub>d<sup>k</sup>) &lt; f(x<sup>k</sup>). Define x<sup>k+1</sup> = x<sup>k</sup> + σ<sub>k</sub>d<sup>k</sup>.</li> <li>Set k → k + 1 and return to step 2.</li> <li>Step 2 is only of academic nature, e.g. use  ∇f(x)  &lt; ε instead.</li> </ol>	For a descent direction $d \exists c > 0$ with $f(x + ad) < f(x)$ for all $a \in (0, c]$ : We have $\nabla f(x)^{T} d = \lim_{a \searrow 0} \frac{f(x+ad)-f(x)}{a} < 0$ and thus there exists a $c > 0$ such that $\frac{f(x+ad)-f(x)}{a} < 0$ for all $a \in (0, c]$ . The reverse direction of this lemma doesn't hold, take $x \mapsto -x^2$ , $\tilde{x} = 0$ , $d \coloneqq 1$ . The <b>antigradient / steepest descent</b> $d = -\nabla f(x) \neq 0$ and $-A^{-1}\nabla f(x)$ for positive definite A are descent directions.
Assume (ALG). A step size with $f(x^{(k)} + \sigma_k d^{(k)}) \leq f(x^{(k)}) - c \left(\frac{\nabla f(x^{(k)})^{T} d^{(k)}}{ d^{(k)} }\right)^2  (\text{ES})$ with a constant $c > 0$ independent of $k$ , is called <b>efficient</b> .	(ALC): for $x^{(0)} \in \mathbb{R}^d$ the level set $\mathcal{N}(f, f(x^{(0)}))$ is compact. (AFD): We have $f \in \mathcal{C}^1$ on an open, convex set $D_0 \supset \mathcal{N}(f, f(x^{(0)}))$ . In descent methods, $f(x^{k+1}) < f(x^{(k)})$ and thus $x^{(k)} \in \mathcal{N}(f, f(x^{(0)}))$ . If (ALC) holds, $(x^{(k)})_{k \in \mathbb{N}}$ and $(f(x^{(k)}))_{k \in \mathbb{N}}$ are bounded.
(AGL): $\nabla f$ is LIPSCHITZ continuous. (AHP): (uniformly positive definite) for $f \in C^2$ and $a > 0$ there holds that $h^{T} f''(x)h \ge a h ^2$ for all $h \in \mathbb{R}^n$ and for all $x \in D \subset \mathbb{R}^n$ (which is an open set). The function $x \mapsto e^x$ is not uniformly positive definite for $D = \mathbb{R}$ .	Let $x \in \mathcal{N}(f, f(x^{(0)}))$ . Then $d \in \mathbb{R}^n$ is (strictly) gradient- related if there exists a $c_3 > 0$ such that $-\nabla f(x)^{T} d \ge c_3  \nabla f(x)   d $ holds (and there exists a $c_4 > 0$ independent of $x$ and $d$ such that $c_4  \nabla f(x)  \ge  d  \ge \frac{1}{c_4}  \nabla f(x) $ ). The antigradient (and assuming (AHP), the NEWTON descent direction) is strictly gradient related ( $c_3 = c_4 = 1$ ).
Let $f: \mathbb{R}^n \supset D \to \mathbb{R}$ be a $\mathcal{C}^2$ function and $D$ be an <b>open con-</b> <b>vex subset</b> containing $N\left(f, f(x^{(0)})\right)$ and (AHP) be fulfilled. If $d^{(k)}$ is <b>gradient related</b> in $x^{(k)}$ and $(\sigma_k)_{k\in\mathbb{N}}$ are <b>efficient</b> , then $x^{(k)} \to \tilde{x}$ , which is the <b>unique</b> minimiser of $f$ . There exists a $q \in (0, 1)$ such that $f(x^{(k)}) - f(\tilde{x}) \leq q^k \left(f(x^{(0)}) - f(\tilde{x})\right)$ and $ x^{(k)} - \tilde{x} ^2 \leq \frac{2}{a}q^k \left(f(x^{(0)}) - f(\tilde{x})\right).$	<ul> <li>Let f: R<sup>n</sup> ⊃ D → R be a C<sup>2</sup> function and D be an open convex subset containing N (f, f(x<sup>(0)</sup>)) and (AHP) be fulfilled.</li> <li>① N (f, f(x<sup>(0)</sup>)) is convex and compact,</li> <li>② F has a unique minimiser x̃, which is the only stationary point of f,</li> <li>③ a/2  x - x̃ <sup>2</sup> ≤ f(x) - f(x̃) ≤ 1/2a  ∇f(x) <sup>2</sup> ∀x ∈ N(f, f(x<sup>(0)</sup>)).</li> </ul>

Definition	Definition
Exact step size	ARMIJO step size and algorithm
Nonlinear optimisation	Nonlinear optimisation
Definition	Assumptions
POWELL step size	R1 and R2: Sufficiently fast decay
Nonlinear optimisation	Nonlinear optimisation
Algorithm	Lemma
Powell Algorithm	Origin of the term steepest descent
NONLINEAR OPTIMISATION	Nonlinear optimisation
Algorithm + remark	Algorithm + interpretation
The gradient method (steepest descent)	Damped NEWTON method
Nonlinear optimisation	Nonlinear optimisation
Example	Algorithm
Steepest descent direction vs NEWTON direction for a quadratic function.	Variable metric method
NONLINEAR OPTIMISATION	Nonlinear optimisation

<ol> <li>Choose the flattening parameter δ ∈ (0, 1), efficiency parameter γ &gt; 0 and 0 &lt; β<sub>1</sub> ≤ β<sub>2</sub> &lt; 1.</li> <li>Initial step size. Take σ<sub>0</sub> ≥ −γ <sup>∇f(x)<sup>T</sup>d</sup>/<sub> d <sup>2</sup></sub>.</li> <li>If f(x + σ<sub>j</sub>d) ≤ f(x) + δσ<sub>j</sub>∇f(x)<sup>T</sup>d, then σ<sub>A</sub> = σ<sub>j</sub>.</li> <li>Else: reduce σ<sub>j</sub> such that σ̃<sub>j</sub> ∈ [β<sub>1</sub>σ<sub>j</sub>, β<sub>2</sub>σ<sub>j</sub>] and iterate j → j + 1 and return to step (3).</li> <li>Assuming (ALC), one can show that after finitely many steps, (R1) and (R2) are satisfied, so σ<sub>A</sub> is efficient</li> </ol>	We consider $\varphi(\sigma) \coloneqq f(x + \sigma d)$ . The exact step size $\sigma_E > 0$ is such that $\varphi'(\sigma_E) = 0$ and $\varphi'(s) < 0$ for $s \in [0, \sigma_E)$ . The exact step size is the "first" local minimum of $\varphi$ . If $\nabla f$ is L-cts (AGL), we have $\sigma_E \ge -\frac{\nabla f(x)^{T}d}{L d ^2}$ and $f(x + \sigma_E d) \le f(x) - \frac{1}{2L} \left(\frac{\nabla f(x)^{T}d}{ d }\right)^2$ , so $\sigma_E$ is <b>efficient</b> . If $f(x) = \frac{1}{2}x^{T}Hx + b^{T}x$ with positive definite $H$ , $\sigma_E = -\frac{\nabla f(x)^{T}d}{d^{T}Hd}$
(R1): There exists a constant $c_1 > 0$ independent of $k$ , such that $f(x^{(k)} + \sigma_k d^{(k)}) - f(x^{(k)}) \leq c_1 \sigma_k \nabla f(x^{(k)})^{T} d^{(k)} < 0$ . The sequence $(f(x^{(k)}))_{k \in \mathbb{N}}$ is bounded by (ALC) and monotone (by design of descent algorithm) and thus <b>convergent</b> . Then $\sigma_k \nabla f(x^{(k)}) d^{(k)} \to 0$ . (R2): There exists a constant $c_2 > 0$ independent of $k$ such that $\sigma_k \ge -c_2 \frac{\nabla f(x^{(k)})^{T} d^{(k)}}{ d^{(k)} ^2}$ . If (R1) and (R2) hold, then $\sigma_k$ satisfies the <b>sufficient decrease condition</b> : $f(x^k + \sigma_k d^k) \le f(x^k) - c_1 c_2 \frac{(\nabla f(x^k)^{T} d^k)^2}{ d^k ^2}$ . arg $\min_{d \in \mathbb{R}^n,  d =1} \nabla f(x)^{T} d = -\frac{\nabla f(x)}{ \nabla f(x) }$ . For $d \in \mathbb{R}^n$ with $ d  = 1$ , we have $\nabla f(x)^{T} d \stackrel{\mathrm{CS}}{\ge} - \nabla f(x)  d  = - \nabla f(x) $ . For $d = -\frac{\nabla f(x)}{ \nabla f(x) }$ we get $\nabla f(x)^{T} d = - \nabla f(x) ^2$ .	$\sigma_P \text{ should fulfil (R1) and } \nabla f(x + \sigma d)^{T} d \ge \beta \nabla f(x)^{T} d \text{ with } 0 < \delta < \beta < 1. \text{ The intersections } s_1 \text{ and } s_2 \text{ divide } [0, \infty) \text{ into } \text{three intervals } I_1 := [0, s_1), I_2 := [s_1, s_2] \text{ and } I_3 := (s_2, \infty).$ $G_1(\sigma) := \begin{cases} \frac{f(x + \sigma d) - f(x)}{\sigma \nabla f(x)^{T} d}, & \sigma > 0, \\ 1, & \sigma = 0, \end{cases} (\text{cts.}),  G_2(\sigma) := \\ \frac{\nabla f(x + \sigma d)^{T} d}{\nabla f(x)^{T} d}. \text{ From (R1) we get } G_1(\sigma) \ge \delta \text{ and from the second condition we get } G_2(\sigma) \le \beta. \text{ Moreover, } G_1(\sigma) \ge \delta \text{ and } G_2(\sigma) \le \beta \text{ holds only in } I_1 \text{ and } G_1(\sigma) \ge \delta \text{ and } G_2(\sigma) \le \beta \text{ holds only in } I_1 \text{ and } G_1(\sigma) \ge \delta \text{ only in } I_3. \\ (\text{ALC}), (\text{AGL}) \text{ imply that } \sigma_P \text{ is an efficient step size.} \end{cases}$ (1) Initialisation. Choose $\sigma_0 > 0$ and set $j := 0$ . (a) If $G_1(\sigma) \ge \delta$ and $G_2(\sigma) \le \beta$ , stop and let $\sigma_P := \sigma_0$ . (b) If $\sigma_0 \in I_1$ , define $a_0 := \sigma_0$ and $a_0 = 2^{-\ell} \sigma_0$ , where $\ell$ is chosen minimally, such that $G_1(b_0) < \delta$ . Go to step (2). (c) If $\sigma_0 \in I_3$ , define $b_0 := \sigma_0$ and $a_0 = 2^{-\ell} \sigma_0$ , where $\ell$ is chosen minimally, such that $G_2(a_0) > \beta$ . (2) Compute $\sigma_j := \frac{1}{2}(a_j + b_j)$ . (a) If $\sigma_j \in I_2$ , stop and set $\sigma_P := \sigma_j$ . (b) If $\sigma_j \in I_3$ , set $a_{j+1} := \sigma_j$ and $b_{j+1} := b_j$ . (c) If $\sigma_j \in I_3$ , set $a_{j+1} := a_j$ and $b_{j+1} := \sigma_j$ .
	(3) Set $j \to j + 1$ and go to step (2).
Like gradient method but instead $d^k = -f''(x^k)^{-1} \nabla f(x^k)$ . Let $A \coloneqq f''(x)$ be SPD and $\langle x, y \rangle_A \coloneqq x^{T} A y$ . We have $\tilde{d} \coloneqq -\frac{A^{-1} \nabla f(x)}{ A^{-1} \nabla f(x) _A} = \arg \min_{ d _A = 1} \nabla f(x)^{T} d$ . For $d \in \mathbb{R}^n$ with $ d _A = 1$ we have $\nabla f(x)^{T} d = \langle A^{-1} \nabla f(x), d \rangle_A \overset{\mathrm{CS}}{\geq} - A^{-1} \nabla f(x) _A  d _A$ $= - A^{-1} \nabla f(x) _A$ . We have $\nabla f(x)^{T} \tilde{d} = - A^{-1} \nabla f(x) _A$ .	<ol> <li>Initialise. Choose x<sup>0</sup> ∈ ℝ<sup>n</sup>, ε &gt; 0 and set k := 0.</li> <li>If  ∇f(x<sup>k</sup>)  &lt; ε, then stop.</li> <li>Compute d<sup>k</sup> := -∇f(x<sup>k</sup>) and choose an efficient step size σ<sub>k</sub>. Define x<sup>k+1</sup> = x<sup>k</sup> + σ<sub>k</sub>d<sup>k</sup>, k → k + 1, return to (2).</li> <li>After initially fast decrease, one observes slow convergence especially for functions with e.g. ellipse-shaped isolines. We e.g. have 0 = φ'(σ<sub>E</sub>) = ∇f(x<sup>k</sup> + σ<sub>E</sub>d<sup>k</sup>)<sup>T</sup>d<sup>k</sup> = d<sup>k+1</sup>d<sup>k</sup>, i.e. d<sup>k+1</sup> ⊥ d<sup>k</sup>, which leads to the slow convergence detailed above.</li> </ol>
We want to account for curvature information $(f'')$ with- out having to compute the second derivative.	f(x) = 1

- (1) Choose  $x^{(0)} \in \mathbb{R}^n$ ,  $\varepsilon > 0$  and set  $k \coloneqq 0$ .
- (2) If  $|\nabla f(x)| < \varepsilon$ , stop.
- (3) Compute the **positive definite matrix**  $A^{(k)}$  and the search direction  $d^{(k)} := -(A^{(k)})^{-1}\nabla f(x^{(k)})$  and an efficient step size  $\sigma_k$ . Set  $x^{(k+1)} = x^{(k)} + \sigma_k d^{(k)}$ , k = k + 1 and go to step (2).

The grey anti-gradient direction  $d_g = -\nabla f(x)$ , (orthogonal to isolines of f), is not optimal. The NEWTON direction is better:  $d_N = -f''(x)^{-1}\nabla f(x) = -H^{-1}Hx = -x$ .

vs.

Algorithm	Definition
BFGS-method	$H ext{-}Orthogonality$
Nonlinear optimisation	Nonlinear optimisation
Theorem	Theorem
BFGS method for quadratic problems	Properties of the CG method
NONLINEAR OPTIMISATION	Nonlinear optimisation
Algorithm	Algorithm
CG method	Trust region NEWTON method
Nonlinear optimisation	Nonlinear optimisation
Idea	Definition
Trust region method	Active / inactive inequality constraints and the active set
Nonlinear optimisation	Nonlinear optimisation
Definition	Definition
Tangent cone	Admissible approximation
Nonlinear optimisation	Nonlinear optimisation

Let $H \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix. Then directions $d^{(0)}, \ldots, d^{(k)}$ for $k < n$ are conjugate or $H$ -orthogonal if $d^{(i)} \neq 0$ and $(d^{(i)})^{T} H d^{(j)} = 0$ for all $0 \leq i < j \leq k$ .	<ol> <li>Choose x<sup>(0)</sup> ∈ ℝ<sup>n</sup>, A<sup>(0)</sup> ∈ ℝ<sup>n×n</sup> positive definite, ε &gt; 0 and set k := 0.</li> <li>If  ∇f(x)  &lt; ε, stop.</li> <li>Compute d<sup>(k)</sup> = -(A<sup>(k)</sup>)<sup>-1</sup>∇f(x<sup>(k)</sup>), an exact step size σ<sub>k</sub> and set x<sup>(k+1)</sup> = x<sup>(k)</sup> + σ<sub>k</sub>d<sup>k</sup>, s<sup>(k)</sup> = x<sup>(k+1)</sup> - x<sup>(k)</sup> and y<sup>(k+1)</sup> = ∇f(x<sup>(k+1)</sup>) - f(x<sup>(k)</sup>) and preform the rank-2- update A<sup>(k+1)</sup> = A<sup>(k)</sup> - A<sup>(k)</sup>s<sup>(k)</sup>(A<sup>(k)</sup>s<sup>(k)</sup>)<sup>T</sup> - y<sup>(k)</sup>(y<sup>(k)</sup>)<sup>T</sup> s<sup>(k)</sup>. Set k → k + 1 and go to step 2.</li> </ol>
As long as $\nabla f(x^{(k-1)}) \neq 0$ , we have (1) $d^{(k-1)} \neq 0$ and $d^{(0)}, \dots, d^{(k)}$ are <i>H</i> -orthogonal, (2) $V_k = \operatorname{span}(\nabla f(x^{(0)}), H \nabla f(x^{(0)}), \dots, H^{k-1} \nabla f(x^{(0)}))$ $= \operatorname{span}(\nabla f(x^{(0)}), \dots, \nabla f(x^{(k-1)})))$ $= \operatorname{span}(d^{(0)}, \dots, d^{(k-1)}),$ (3) $f(x^{(k)}) = \min_{z \in V_k} f(x^{(0)} + z).$	Let $H \in \mathbb{R}^{n \times n}$ be a symmetric positive definite matrix. Then, the BFGS-method generated $H$ -orthogonal search directions $d^{(k)}$ . The minimum is found in $m \leq n$ steps. If $m = n$ , then $A^{(n)} = H$ .
$ \begin{aligned} \text{Given: } 0 < \delta_1 < \delta_2 < 1,  \sigma_1 \in (0, 1),  \sigma_2 > 1,  \sigma_0 > 0,  x^{(0)} \in \mathbb{R}^n. \\ \hline (1)  d^{(k)} = \text{ solution of } \min_{ d  \leq \rho_k} f_k(d). \text{ If } f(x^{(k)}) = f_k(d^{(k)}), \text{ then stop.} \\ \hline (2)  r_k := \frac{f(x^{(k)}) - f(x^{(k)} + d^{(k)})}{f(x^{(k)}) - f_k(x^{(k)} + d^{(k)})}, \text{ If } r_k \geq \delta_1 \text{ (successful step), set } x^{(k+1)} = x^{(k)} + d^{(k)}, \text{ compute } \nabla f(x^{(k+1)}), f''(x^{(k+1)}) \text{ and update } \rho_k: \\ & \text{ if } r_k \begin{cases} \in [\delta_1, \delta_2), & \text{choose } \rho_{k+1} \in [\delta_1 \rho_k, \rho_k], \\ \geq \delta_2, & \text{choose } \rho_{k+1} \in [\rho_k, \delta_2 \rho_k], \end{cases} \\ & \text{ set } k \to k+1 \text{ and go to } (2). \end{aligned} $ $ \end{aligned} $	(1) choose $x^{(0)} \in \mathbb{R}^n, \varepsilon > 0$ and set $k := 0$ and $d^{(0)} = -H(x^{(0)} + b)$ . (2) If $ \nabla f(x^{(k)})  \leq \varepsilon$ , stop. (3) Compute $\sigma_k = \frac{ \nabla f(x^{(k)}) ^2}{ d^{(k)} ^2_H}$ and set $x^{(k+1)} = x^{(k)} + \sigma_k d^{(k)}$ . We have $\nabla f(x^{(k+1)}) = Hx^{(k+1)} + b = \nabla f(x^{(k)}) + \sigma_k Hd^k$ . Compute $\beta_k := \frac{ \nabla f(x^{(k+1)}) ^2}{ \nabla f(x^{(k)}) ^2}$ and set $d^{(k+1)} = -\nabla f(x^{(k+1)} + \beta_k d^{(k)}$ . Set $k \to k + 1$ and return to (2).
We consider $\min_{x \in \mathbb{R}^n} f(x)$ subject to $\begin{cases} c_i(x) = 0, & i \in E, \\ c_i(x) \ge 0, & i \in I \end{cases}$ where $I, E \subset \mathbb{N}$ are <b>disjoint</b> index sets. The constraints $c_i(x) \stackrel{(\geq)}{=} 0$ are called (in)equality constraints. The <b>admiss-able set</b> is $\Omega = \{x \in \mathbb{R}^n : c_i(x) = 0, i \in E, c_i(x) \ge 0, i \in I\},$ Let $x \in \Omega$ , then $c_i(x), i \in I$ is called <b>active</b> if $c_i(x) = 0$ and <b>inactive</b> if $c_i(x) > 0$ . The <b>active set</b> is $\mathcal{A}(x) := E \cup \{i \in I : c_i(x) = 0\}.$	<ul> <li>Up to now, we have computed a search direction d<sup>k</sup> and a step size σ<sub>k</sub> (line search) and we used the update x<sup>(k+1)</sup> = x<sup>(k)</sup> + σ<sub>k</sub>d<sup>(k)</sup>. The new idea is now to</li> <li>use a local model f<sub>k</sub> of f, e.g. f<sub>k</sub> = f(x<sup>(k)</sup>) + ∇f(x<sup>(k)</sup>)<sup>T</sup>d or f<sub>k</sub> = f(x<sup>(k)</sup>) + ∇f(x<sup>(k)</sup>)<sup>T</sup>d + ½d<sup>T</sup>f''(x<sup>(k)</sup>)d,</li> <li>take radius ρ<sub>k</sub> &gt; 0 and consider the trust region B<sub>ρ<sub>k</sub></sub>(x<sup>(k)</sup>),</li> <li>compute d<sup>(k)</sup> as a global solution to min<sub> d ≤ρ<sub>k</sub></sub> f<sub>k</sub>(d),</li> <li>update x<sup>(k+1)</sup> = x<sup>(k)</sup> + d<sup>(k)</sup>.</li> </ul>
Let $x \in \Omega$ . Then the sequence $(x^{(n)})_{n \in \mathbb{N}}$ is called <i>admissable</i> approximation of x if $x^{(n)} \to x$ and $x^{(n)} \in \Omega$ for almost all $n \in \mathbb{N}$ .	A direction $d \in \mathbb{R}^n$ is a <b>tangent to</b> $\Omega$ <b>in</b> $x \in \Omega$ if there exists an <b>admissable approximation</b> $(x^{(k)})_{k\in\mathbb{N}}$ of $x$ and a sequence $(t_k)_{k\in\mathbb{N}} \subset \mathbb{R}_+$ <b>converging to zero</b> such that $\lim_{k\to\infty} \frac{x^{(k)}-x}{t_k} = d$ . The <b>tangent cone</b> of $\Omega$ in $x$ is $T_{\Omega}(x) :=$ $\{d \in \mathbb{R}^n : d \text{ is tangent to } \Omega \text{ in } x\}.$ The tangent cone is a cone $(\tilde{t_k} := \frac{1}{a}t_k)$ . If $x \in \operatorname{int}(\Omega)$ , then $T_{\Omega}(x) = \mathbb{R}^n$ .

Theorem	Definition
Variational inequality - General case	Linearised cone
NONLINEAR OPTIMISATION	Nonlinear optimisation
Assumptions	Theorem
ACQ and LICQ	KKT conditions
NONLINEAR OPTIMISATION	Nonlinear optimisation
Theorem	Proof
$T_{\Omega}(\tilde{x}) \subset L_{\Omega}(\tilde{x})$	of the KKT theorem
Nonlinear optimisation	Nonlinear optimisation
Lemma	Definition
Farkas	Critical cone
Nonlinear optimisation	Nonlinear optimisation
Theorem	Theorem
Second order necessary condition for constraint problems	Second order sufficient optimality condition for constraint problems
Nonlinear optimisation	Nonlinear optimisation

For  $x \in \Omega$ , the **linearised cone of**  $\Omega$  **in**  $x \in \Omega$  is

$$L_{\Omega}(x) \coloneqq \left\{ d \in \mathbb{R}^{n} : \begin{array}{l} d^{\mathsf{T}} \nabla c_{i}(x) = 0 \ \forall i \in E, \\ d^{\mathsf{T}} \nabla c_{i}(x) \ge 0 \ \forall i \in I \cap \mathcal{A}(x) \end{array} \right\}$$

Thus  $L_{\Omega}(x)$  for  $\Omega := \{x \in \mathbb{R}^n : h(x) = 0, g(x) \leq 0\}$  depends on g and h, whereas  $N_{\Omega}(x)$  and  $T\Omega(x)$  don't. Let  $\hat{x} \in \Omega$  be a solution of the constrained problem and  $f \in \mathcal{C}^1$ . Then  $\nabla f(\hat{x})^{\mathsf{T}} d \ge 0$  holds for all  $d \in T_{\Omega}(\hat{x})$ .

For  $d \in T_{\Omega}(\hat{x})$  we have by TAYLORS theorem,

$$0 \leq \frac{f(x^{(k)}) - f(\hat{x})}{t_k} = \frac{1}{t_k} \left( f(\hat{x} + (x^{(k)} - \hat{x})) - f(\hat{x}) \right)$$
$$= \frac{1}{t_k} \left( f(\hat{x}) + \nabla f(\hat{x} + \xi(x^{(k)} - \hat{x}))^{\mathsf{T}}(x^{(k)} - \hat{x}) - f(\hat{x}) \right)$$
$$= \underbrace{\nabla f(\hat{x} + \xi(x^{(k)} - \hat{x}))}_{\to \nabla f(\hat{x})}^{\mathsf{T}} \underbrace{\frac{x^{(k)} - \hat{x}}{t_k}}_{\to d} \to \nabla f(\hat{x})^{\mathsf{T}} d.$$

Let  $\hat{x}$  be solution to the constrained problem, f and  $(c_i)_{i \in I \cup E}$ be  $\mathcal{C}^1$  functions such that (ACQ) is satisfied. Then there exists a vector  $(\tilde{\lambda}_i)_{i \in I \cup E}$  of LAGRANGE multipliers such that

(1) 
$$\nabla_x L(\tilde{x}, \tilde{\lambda}) = 0,$$

(2)  $c_i(\tilde{x}) = 0$  for all  $i \in E$ ,

(3)  $c_i(\tilde{x}) \ge 0$  for all  $i \in I$ ,

(4)  $\tilde{\lambda}_i \ge 0$  for all  $i \in I$ ,

(5)  $\tilde{\lambda}_i c_i(\hat{x}) = 0$  for all  $i \in E \cup I$  (complementarity).

Let  $N := \left\{ \sum_{i \in \mathcal{A}(\tilde{x})} \lambda_i \nabla c_i : \lambda \ge 0 \right\}$  and  $g := \nabla f(\tilde{x})$ . By FARKAS lemma either  $\nabla f(\tilde{x}) = \sum_{i \in \mathcal{A}(\tilde{x})} \lambda_i A^{\mathsf{T}}(\tilde{x}) \tilde{\lambda}$  with  $\tilde{\lambda}_i \ge 0$ for  $i \in \mathcal{A}(\tilde{x}) \cap I$  or there exists a  $d \in \mathbb{R}^n$  such that  $\nabla f(\tilde{x})^{\mathsf{T}} d < 0$ ,  $\nabla c_i^{\mathsf{T}} d = 0$  for  $i \in E$  and  $\nabla c_i^{\mathsf{T}} d \ge 0$  for  $i \in \mathcal{A}(\tilde{x}) \cap I$ . We can rewrite those three conditions as  $\nabla f(\tilde{x})^{\mathsf{T}} d < 0$  and  $d \in L_{\Omega}(\tilde{x})$ . By assumption  $\tilde{x} \in \Omega$  and (ACQ) hold. Thus we have  $\nabla f(\tilde{x})^{\mathsf{T}} d < 0$  for a  $d \in T_{\Omega}(\tilde{x})$ , which is a contradiction to the variational inequality, so the first option has to hold. Define  $\tilde{\lambda}_i = 0$  for  $i \notin \mathcal{A}(\tilde{x})$ , so the last condition (complementarity condition) holds.

If  $(\tilde{x}, \tilde{\lambda})$  satisfy the KKT conditions, the **critical cone** is

 $C(\tilde{x}, \tilde{\lambda}) = \{ w \in L_{\Omega}(\tilde{x}) : \nabla c_i(\tilde{x})^{\mathsf{T}} w = 0 \ \forall i \in \mathcal{A}(\tilde{x}) \cap I \text{ s.th. } \tilde{\lambda}_i > 0 \}$ 

We have  $w \in C(\tilde{x}, \tilde{\lambda})$  if and only if  $\nabla c_i(\tilde{x})^\mathsf{T} w = 0 \ \forall i \in E$ ,  $\forall i \in \mathcal{A}(\tilde{x}) \cap I \text{ s.th. } \tilde{\lambda}_i > 0 \text{ and } \nabla c_i(\tilde{x})^\mathsf{T} w = 0 \ \forall i \in \mathcal{A}(\tilde{x}) \cap I \text{ s.th. } \tilde{\lambda}_i = 0.$ For  $d \in C(\tilde{x}, \tilde{\lambda})$  we have  $\nabla f(\tilde{x})^\mathsf{T} d = \sum_{i \in \mathcal{A}(\tilde{x})} \tilde{\lambda}_i \nabla c_i(\tilde{x})^\mathsf{T} d = 0.$ Thus  $C(\tilde{x}, \tilde{\lambda})$  contains all directions where, based on first order information, we cannot decide if f decreases or increases.

Let  $\tilde{x} \in \Omega$  and  $\tilde{\lambda}$  such that  $(\tilde{x}, \tilde{\lambda})$  satisfies the KKT conditions. If there exists a  $\sigma > 0$  such that

$$w^{\mathsf{T}} \nabla_{xx}^2 L(\tilde{x}, \tilde{\lambda}) w \ge \sigma |w|^2$$

holds for all  $w \in C(\tilde{x}, \tilde{\lambda})$ , then  $\tilde{x}$  is a strict local solution to the constrained problem.

Let  $x \in \Omega$ .

ABADIE constraint qualification (ACQ):  $T_{\Omega}(x) = L_{\Omega}(x)$ . Linear independence constraint qualification (LICQ):  $\{\nabla c_i(x) : i \in \mathcal{A}(x)\}$  is **linearly independent**. LICQ  $\implies$  ACQ.

W.l.o.g. assume  $c_i(x), i \in \{1, ..., m\}$  be the **active** constraints in  $\tilde{x}$ . Let  $d \in T_{\Omega}(\tilde{x})$ . For k **sufficiently large** and  $c_i$  is an **equality constraint**, by TAYLOR expansion,  $\exists \alpha \in [0, 1]$ 

$$0 = \frac{1}{t_k} c_i(x^{(k)}) = \frac{1}{t_k} c_i(\tilde{x} + (x^{(k)} - \tilde{x}))$$
$$= \left[\underbrace{c_i(\tilde{x})}_{=0} + \underbrace{\nabla c_i(\tilde{x} + \alpha(x^{(k)} - \tilde{x}))^{\mathsf{T}}}_{\to \nabla c_i(\tilde{x})}\right] \underbrace{\frac{x^{(k)} - \tilde{x}}_{t_k}}_{\to d} \xrightarrow{k \to \infty} \nabla c_i(\tilde{x})^{\mathsf{T}} d.$$

Similarly we can show that  $\nabla c_i(\tilde{x})^{\mathsf{T}} d \ge 0$  for  $i \in I \cap \mathcal{A}(\tilde{x})$ . Thus  $d \in L_{\Omega}(\tilde{x})$ .

Let  $K := \{By + Cw : y \in \mathbb{R}^m, y \ge 0, w \in \mathbb{R}^p\}$  with  $B \in \mathbb{R}^{n \times m}$ and  $C \in \mathbb{R}^{n \times p}$ . For each  $g \in \mathbb{R}^n$  either  $g \in K$  or there exists  $d \in \mathbb{R}^n$  such that  $g^{\mathsf{T}}d < 0, B^{\mathsf{T}}d \ge 0$  and  $C^{\mathsf{T}}d = 0$ .

Let  $\tilde{x}$  be a local solution to the constrained problem, assume that (LICQ) holds and let  $\tilde{\lambda}$  be such that the KKT conditions are satisfied. Then

$$w^{\mathsf{T}} \nabla^2_{xx} L(\tilde{x}, \tilde{\lambda}) w \geqslant 0$$

holds for all  $w \in C(\tilde{x}, \tilde{\lambda})$ .

Lemma	Assumption
(ACQ) holds for affine linear constraints $c_i(x) = a_i^{T} x + b_i$ for $a_i \in \mathbb{R}^n$ and $b_i \in \mathbb{R}$ .	Mangasarian-Fromovitz
Nonlinear optimisation	Nonlinear optimisation
Lemma	WITHOUT PROOF
LICQ implies MFCQ	Implications between constraint qualifications
Nonlinear optimisation	Nonlinear optimisation
Definition	Theorem
Normal cone	Convergence of the NEWTON method
Nonlinear optimisation	Nonlinear optimisation
EXAMPLE	Theorem
Step size decreases to fast	TAYLOR in $\mathbb R$ with remainder
Nonlinear optimisation	Nonlinear optimisation
Definition	
Stationary points	Convergence analyisis steps for descent methods
Nonlinear optimisation	Nonlinear optimisation

(MFCQ) holds if there exists a $w \in \mathbb{R}^n$ such that $\nabla c_i(\tilde{x})^{T} w \begin{cases} > 0, & \forall i \in \mathcal{A}(\tilde{x}) \cap I \\ = 0, & \forall i \in E. \end{cases}$ and $\{\nabla c_i\}_{i \in E}$ is linearly independent.	We show $L_{\Omega}(\tilde{x}) \subset T_{\Omega}(\tilde{x})$ . Let $w \in L_{\Omega}(\tilde{x})$ . Then $a_i^{T}w = 0$ for $i \in E$ and $a_i^{T}w \ge 0$ for $i \in \mathcal{A}(\tilde{x}) \cap I$ , as $\nabla c_i = a_i$ . If $i \in I \setminus \mathcal{A}(\tilde{x})$ , then $c_i(\tilde{x}) > 0$ . Then $\exists t_0 > 0$ such that $c_i(\tilde{x} + tw) > 0$ $\forall t \in [0, t_0]$ , so $c_i$ "stays" inactive. Let $(x^{(k)} := \tilde{x} + \frac{t_0}{k}w)_{k \in \mathbb{N}}$ . For $i \in \mathcal{A}(\tilde{x}) \cap I$ we have $c_i(x^{(k)}) = c_i(x^{(k)}) - c_i(\tilde{x}) = a_i^T(x^{(k)} - \tilde{x}) = \frac{t_0}{k}a_i^Tw \ge 0$ since $c_i(\tilde{x}) = 0$ and $w \in L_{\Omega}(\tilde{x})$ , so $(x^{(k)})_{k \in \mathbb{N}}$ is an admissable approximation. For $i \in E$ we have $c_i(x^{(k)}) = c_i(x^{(k)}) - c_i(\tilde{x}) = \frac{t_0}{k}a_i^Tw \ge 0$ , by the same reasoning as above, so $(x^{(k)})_{k \in \mathbb{N}}$ is an admissable approximation. Moreover, $\lim_{k \to \infty} \frac{x^{(k)} - \tilde{x}}{\frac{t_0}{k}} = \lim_{k \to \infty} \frac{\frac{t_0}{k}w}{\frac{t_0}{k}} = w$ , so $w \in T_{\Omega}(\tilde{x})$ .
$(\text{LICQ}) \implies (\text{MFCQ}) \implies (\text{ACQ}).$ If (SQC) holds in $\tilde{x} \in \Omega$ , then (MFCQ) holds.	Let $G(\tilde{x}) := (\nabla c_i(\tilde{x})^{T})_{i \in \mathcal{A}(\tilde{x})}$ . By (LICQ) it has maximal rank. Then there exists a $w \in \mathbb{R}^n$ such that $\nabla c_i(\tilde{x})^{T}w =$ $\begin{cases} 1, & \forall i \in \mathcal{A}(\tilde{x}) \cap I, \\ 0, & \forall i \in E. \end{cases}$ This is because as $G(\tilde{x})$ has maximal $0, & \forall i \in E. \end{cases}$ rank, adding an additional column doesn't change the rank. A linear system $Ax = b$ is solvable if the rank of $A$ is equal to the rank of the extended matrix $A b$ . The system is solvable as A as maximal rank and thus we can append any $b$ , in particu- lar one with ones in first components for the active inequality constraints and zeros for all the equality constraints.
Let $f''$ be Lipschitz continuous in a neighbourhood of a local minimum $\hat{x}$ of $f$ and let $f''(\hat{x})$ be positive definite. Then the NEWTON method $x^{(k+1)} = x^{(k)} - f''(x^{(k)})^{-1} \nabla f(x^{(k)})$	For $x \in \Omega$ , $N_{\Omega}(x) := \{v \in \mathbb{R}^n : v^{T}w \leq 0 \ \forall w \in T_{\Omega}(x)\}$ is the <b>normal cone</b> to $T_{\Omega}(x)$ . The elements of $N_{\Omega}(x)$ are <b>normal vectors</b> .
converges locally quadratically to $\hat{x}$ . $d^{(k)} := f''(x^{(k)})^{-1} \cdot (-\nabla f(x^{(k)}))$ is the <b>Newton direction</b> . The <b>damped Newton method</b> is $x^{(k+1)} = x^{(k)} - \sigma_k f''(x^{(k)})^{-1} \nabla f(x^{(k)})$ with $\sigma_k < 1$ .	Let $\tilde{x}$ be a local solution to the constraint problem. Then $-\nabla f(\tilde{x}) \in N_{\Omega}(\tilde{x}).$ By the variational inequality $\nabla f(\tilde{x})d \ge 0$ , i.e. $-\nabla f(\tilde{x})d \le 0$ holds for all $d \in T_{\Omega}(\tilde{x}).$
Let $I \subset \mathbb{R}$ be an interval and $f: I \to \mathbb{R}$ in $\mathcal{C}^{n+1}(I)$ . Then there exits a $\theta \in [0,1]$ such that $f(x) = \sum_{k=0}^{n} \frac{f^{(k)}(a)}{k!} (x-a)^{k} + \frac{f^{(n+1)}(a+\theta(x-a))}{(n+1)!} (x-a)^{n+1}.$	Consider $f(x) := x^2$ , $d^{(k)} := -1$ and $\sigma_k := 2^{-k-2}$ for all $k \ge 0$ . The sequence $(x^{(k)})_{k\in\mathbb{N}}$ defined by $x^{(k+1)} = x^{(k)} + \sigma_k d^{(k)} = x^{(k)} - \frac{1}{2^{k+2}}$ and $x^{(0)} = 1$ converge to $\frac{1}{2}$ : $x^{(k+1)} = x^{(0)} - \sum_{j=0}^k \frac{1}{2^{k+2}} = 1 - \frac{1}{4} \frac{1 - \frac{1}{2^{k+1}}}{1 - \frac{1}{2}} = \frac{1}{2} + \frac{1}{2^{k+2}} \xrightarrow{k \to \infty} \frac{1}{2}$ .
We want to show $\nabla f(x^{(k)}) \to 0$ . We first show that $\frac{\nabla f(x^{(k)})^{T}d^{(k)}}{ d^{(k)} } \xrightarrow{k \to \infty} 0$ . If $\sigma_k$ is efficient, we have $0 \xleftarrow{k \to \infty} f(x^{(k+1)}) - f(x^{(k)}) \leqslant -c \left(\frac{\nabla f(x^{(k)})^{T}d^{(k)}}{ d^{(k)} }\right)^2 < 0$ . Then $\frac{\nabla f(x^{(k)})^{T}d^{(k)}}{ d^{(k)} } \xrightarrow{k \to \infty} 0$ , as we wanted. We have $\frac{\nabla f(x^{(k)})^{T}d^{(k)}}{ d^{(k)} } =  \nabla f(x^{(k)})  \cos\left(\triangleleft(\nabla f(x^{(k)}), d^{(k)})\right)$ , so to ensure that $\nabla f(x^{(k)}) \to 0$ we have to avoid $\nabla f(x^{(k)}) + d^{(k)}$	If $\nabla f(x) = 0$ holds, $x$ is a <b>stationary point</b> of $f$ . Stationary points need not be extrema, consider z.B. $f(x) := x^3$ and $x = 0$ .

for large k.

	Definition
Requirements for the search directions	Strict complementarity
Nonlinear optimisation	Nonlinear optimisation
	Definition
Problems with box constraints	Slater constraint qualification
NONLINEAR OPTIMISATION	NONLINEAR OPTIMISATION
Algorithm	
Solve $\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T Q x + q^T x$ subject to Ax = b	
Nonlinear optimisation	

A LAGRANGE multiplier  $\lambda$  satisfies strict complementarity Thus to ensure that  $\nabla f(x^{(k)}) \rightarrow 0$  we have to avoid if  $\lambda_i > 0$  for all  $i \in I \cap \mathcal{A}(\tilde{x})$ .  $\nabla f(x^{(k)}) \perp d^{(k)}$  for large k (this is slow convergence). We Then  $C(\tilde{x}, \tilde{\lambda}) = \{ d \in \mathbb{R}^n : \nabla c_i(\tilde{x})^\mathsf{T} d = 0 \ \forall i \in \mathcal{A}(\tilde{x}) \} =$ have  $\ker(G(\tilde{x}))$  for  $G(\tilde{x}) := (\nabla c_i(\tilde{x})^{\mathsf{T}})_{i \in \mathcal{A}(\tilde{x})}$ . Let  $(s_k)_{k=1}^{\ell}$  be a basis  $\cos(\sphericalangle(\nabla f(x^{(k)}), d^{(k)})) = \frac{\nabla f(x^{(k)})^{\mathsf{T}} d^{(k)}}{|d^{(k)}| |\nabla f(x^{(k)})|} =: \beta_k$ of ker $(G(\tilde{x}))$ . The second order optimality conditions reduce to  $Z^{\mathsf{T}} \nabla^2_{xx} L(\tilde{x}, \tilde{\lambda}) Z$  being positive definite on  $\mathbb{R}^{\ell}$ . Then  $\beta_k |\nabla f(x^{(k)})| = \frac{\nabla f(x^{(k)})^{\mathsf{T}} d^{(k)}}{|d^{(k)}|} \to 0$ . We can infer from this that  $\nabla f(x^{(k)}) \to 0$  if  $-\beta_k \ge c > 0$  is bounded away from zero for all  $k \in \mathbb{N}$ . Let  $D \subset \mathbb{R}^n$  be an open and *convex* subset such that  $-c_i$  is a Let  $\Omega := \{x \in \mathbb{R}^n : v_i \leq x_i \leq w_i \ \forall i \in \{1, \dots, n\}\}$  and for convex  $\mathcal{C}^1$  function on D for  $i \in I$  and  $c_i(x) \coloneqq a_i^\mathsf{T} x + b_i$  is an simplicity assume v < w componentwise. Then  $x \in \Omega$  can be affine linear function for  $i \in E$ . rewritten as  $Gx \ge r$ , where  $G = (I, -I)^{\mathsf{T}}$ , r = (v, -w). Then the global SLATER condition holds if the set  $(a_i)_{i \in E}$ At most one constraint can be active, so G(x):= is linearly independent and there exists a  $v \in \mathbb{R}^n$  such that  $(\nabla c_i(x))_{i \in \mathcal{A}(x)} = \operatorname{diag}((\pm 1)_{i=1}^n)$  and thus  $\{\nabla c_i : i \in \mathcal{A}(\tilde{x})\}$  $c_i(v) = 0$  for all  $i \in E$  and  $c_i(v) \ge 0$  for  $i \in I$ . is linearly independent and thus (ACQ) holds.  $L(x,\lambda) = f(x) - \sum_{j=1}^{n} \lambda_{j}^{(\ell)} (x_{j} - v_{j}) - \sum_{j=1}^{n} \lambda_{j}^{(u)} (-x_{j} + w_{j}),$ by KKT:  $\lambda_{i}^{(\ell)} = \left[\frac{\partial f(x)}{\partial x_{i}}\right]_{+}$  and  $\lambda_{i}^{(u)} = \left[\frac{\partial f(x)}{\partial x_{i}}\right]_{-}$  are unique. One can show that if (SQC) holds in  $\tilde{x} \in \Omega$ , then (MFCQ) holds.  $C(\tilde{x}, \tilde{y}) = \left\{ d \in L_{\Omega}(\tilde{x}) : d_i = 0 \text{ if } \frac{\partial f(\tilde{x})}{\partial x_i} \neq 0 \right\}.$  $Q \in \mathbb{R}^{n \times n}$  sym., PD on ker $(A), A \in \mathbb{R}^{m \times n}$ , rang $(A) = m \leq n$ . (1) Compute the QR-decomposition of  $A^{\mathsf{T}}$ : compute  $H \in$  $\mathbb{R}^{n \times n}$  and  $R \in \mathbb{R}^{m \times m}$ , such that  $HA^{\mathsf{T}} = \begin{pmatrix} R \\ 0 \end{pmatrix}$ . Define  $h \coloneqq -Hq = \begin{pmatrix} h_1 \\ h_2 \end{pmatrix}$  and  $B \coloneqq HQH^{\mathsf{T}} \eqqcolon \begin{pmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{pmatrix}$ , where  $h_1 \in \mathbb{R}^m$  and  $B_{11} \in \mathbb{R}^{m \times m}$ . (2) Solve  $R^{\mathsf{T}}\tilde{x}_y = b$  and  $B_{22}\tilde{x}_z = h_2 - B_{21}\tilde{x}_y$ .  $\tilde{x} \coloneqq H^{\mathsf{T}}\left(\frac{\tilde{x}_y}{\tilde{x}_z}\right)$ . (3) Solve  $R\lambda = B_{11}\tilde{x}_y + B_{12}\tilde{x}_z - h_1$  for  $\lambda$  via forward substitution.