

# Nonlinear Optimization

## 0. Basics

Let  $f: \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ ,  $\overline{\mathbb{R}} := \mathbb{R} \cup \{\infty\}$

$$U \subset \mathbb{R}^n$$

look for  $x_* \in U: f(x_*) \leq f(x) \quad \forall x \in U$

## 0.1 Definition

(i)  $x_*$  is called minimizer if  $\forall x \in U: f(x_*) \leq f(x)$

(ii) strict minimizer " " <

(iii) the epigraph of  $f$  is

$$\text{epi } f := \{ (x, y) \in \mathbb{R}^n \times \mathbb{R} \mid y \geq f(x) \}$$



(iv) the level set  $L_\alpha$  for  $\alpha \in \mathbb{R}$  is

$$L_\alpha f := \{x \in U \mid f(x) \leq \alpha\}$$



(v) the domain of  $f$  is

$$\text{dom } f := \bigcup_{\alpha \in \mathbb{R}} L_\alpha f = \{x \in U \mid f(x) < \infty\}$$

(vi) the indicator function of  $U$  is

$$I_U := \begin{cases} 0, & x \in U \\ \infty, & \text{otherwise} \end{cases}$$



0.2 Lemma Let  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuous.

Then  $L_\alpha$  is closed for all  $\alpha \in \mathbb{R}$ .

Proof  $x_k \rightarrow \bar{x}$  with  $x_k \in L_\alpha f \Rightarrow f(x_k) \leq \alpha$  f.a.  $k \in \mathbb{N}$

$\Rightarrow f(\bar{x}) \leq \alpha \Rightarrow \bar{x} \in L_\alpha f$

continuity

□

### 0.3 Theorem (Existence of minimizers)

Let  $U \subset \mathbb{R}^n$  be closed, bounded, nonempty,

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  be continuous,

Then, a minimizer  $x_*$  exists.

Proof: Let  $m := \inf_{x \in U} f(x)$ . Then there is  $(x_k)_k: \lim_{k \rightarrow \infty} f(x_k) = m$ .

Due to  $m < \infty$ , the Heine-Borel  $\Rightarrow U$  compact.

$(x_k)_k$  contains a convergent subsequence. Wlog  $x_k \rightarrow \bar{x}$ .

Continuity of  $f$  implies  $f(\bar{x}) = f(\lim_{k \rightarrow \infty} x_k) = \lim_{k \rightarrow \infty} f(x_k) = m$ .

$\Rightarrow \bar{x}$  is a minimizer.  $\square$

### 0.4 Definition $f$ is lower semicontinuous,

if  $L_x f$  is closed for all  $x \in \mathbb{R}$ .

Generalization of 0.3: Assume  $f$  is only lower semicont.

A minimizer exists.

Proof: Assume  $f(\bar{x}) = m + \varepsilon$ ,  $\varepsilon > 0$ . There is a subsequence

$(x_{k_j})$  with  $f(x_{k_j}) \leq m + \frac{\varepsilon}{2} \Rightarrow x_{k_j} \in L_{m+\frac{\varepsilon}{2}} f$ .

Closedness of  $L_{m+\frac{\varepsilon}{2}} f$  implies  $\bar{x} \in L_{m+\frac{\varepsilon}{2}} f \Rightarrow f(\bar{x}) \leq m + \frac{\varepsilon}{2}$ .  $\square$

## 0.5 Definition (Convexity)

(i)  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex, if  $f(\lambda x + (1-\lambda)y)$

$$\leq \lambda f(x) + (1-\lambda)f(y)$$

for all  $x, y \in U$ ,  $\lambda \in [0, 1]$

strictly convex " < " for all  $x \neq y$ ,  $\lambda \in ]0, 1[$

(ii)  $U \subset \mathbb{R}^n$  is convex if  $x, y \in U \Rightarrow \lambda x + (1-\lambda)y \in U \quad \forall \lambda \in [0, 1]$

## 0.6 Theorem (uniqueness)

Let  $f: U \rightarrow \mathbb{R}$  be strictly convex and  $U \subset \mathbb{R}^n$  convex.

Then,  $f$  has at most one minimizers.

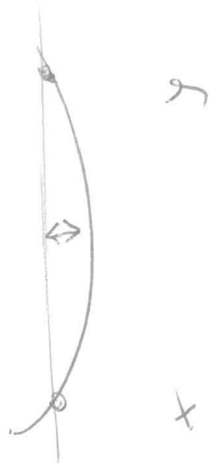
Proof: Let  $x, y$  be minimizers. Then  $\frac{x+y}{2} \in U$  and

$$f\left(\frac{x+y}{2}\right) < \frac{1}{2}(f(x) + f(y)) = f(x) \Rightarrow x = y \quad \square$$

## 0.7 Definition (local minimizer)

$x_0 \in U$  is a local minimizer if there is  $\varepsilon > 0$  such that  $x_0$  is a (global) minimizer of  $f$  in

$$U \cap B_\varepsilon(x_0).$$



## 0.8 Theorem (necessary conditions)

Let  $x_*$  be a local minimizer of  $f$  and  $x_* \in \text{int } U$ .

Then (i) if  $f \in C^1(U)$ :  $f'(x_*) = 0$

(ii) if  $f \in C^2(U)$ :  $f''(x_*)$  is positive semidefinite

If  $f$  and  $U$  are convex,  $f'(x_*) = 0 \Rightarrow x_*$  is a minimizer.

Proof (i)  $0 \leq f(x_* + h) - f(x_*) = \varepsilon f'(x_*)h + o(\varepsilon \|h\|)$

$$\Rightarrow f'(x_*)h \geq -\frac{1}{\varepsilon} o(\varepsilon \|h\|) \xrightarrow{\varepsilon \rightarrow 0} 0 \Rightarrow f'(x_*)h \geq 0 \quad \forall h \Rightarrow f'(x_*) = 0$$

$$(ii) \quad 0 \leq f(x_* + \varepsilon h) - f(x_*) = \varepsilon h^T f''(x_*) \varepsilon h + o(\varepsilon^2 \|h\|^2)$$

$$\Rightarrow h^T f''(x_*) h \geq 0.$$

$$f \text{ convex} \Rightarrow (1-\varepsilon)f(x_*) + \varepsilon f(x_* + \varepsilon h) \leq f(x_* + \varepsilon h) = f(x_*) + o(\varepsilon \|h\|)$$

$$\Rightarrow f(x_* + h) \geq \frac{1}{\varepsilon} ( \varepsilon f(x_*) + o(\varepsilon \|h\|) ) \xrightarrow{\varepsilon \rightarrow 0} f(x_*)$$

$$\Rightarrow f(x_* + h) \geq f(x_*) \Rightarrow x_* \text{ is minimizer.}$$

□

## 0.9 Theorem (sufficient conditions)

$$x_* \in \text{int } U, f \in C^2(U), f'(x_*) = 0, f''(x_*) \text{ spd } (v^T f''(x_*) v \geq \delta \|v\|^2)$$

$\Rightarrow x_*$  is local minimizer

$$\begin{aligned} \text{Proof} \quad f(x) - f(x_*) &= (x - x_*)^T f''(x_*) (x - x_*) + o(\|x - x_*\|^2) \\ &\geq \delta \|x - x_*\|^2 + o(\|x - x_*\|^2) \end{aligned}$$

$\rightarrow 0$  faster than  $\|x - x_*\|^2$

$\geq 0$  for  $\|x - x_*\|$  sufficiently small  $\square$

## I. Unconstrained Optimization

$U = \mathbb{R}^n$ ,  $f \in C^1(\mathbb{R}^n)$ ,  $f$  bounded  $\Rightarrow f$  is bounded from below

### I.1 Descent methods

#### (I.1.1) Algorithm

input:  $x_0 \in \mathbb{R}^n$   
while "not converged"

choose  $s_k$  with  $f'(x_k) s_k < 0$

choose  $\alpha_k > 0$  such that  $f(x_k + \alpha_k s_k) < f(x_k)$

$$x_{k+1} = x_k + \alpha_k s_k$$



#### (I.1.2) Lemma

- (i)  $f(x_k)$  is monotonically decreasing
- (ii)  $x_k$  has (at least) one accumulation point

Proof (i) by construction

- (ii)  $(x_k)_k \subset \mathbb{R}^n$  which is compact  
Bolzano-Weierstrass  $\Rightarrow$  convergent subsequence  $\square$

## I.2 Line search

Properties of "good" step sizes:

$$(i) \quad f(x + \alpha s) \leq f(x) + c_1 \alpha f'(x)s$$

$$0 < c_1 < 1$$

$$(ii) \quad |f'(x + \alpha s)s| \leq -c_2 f'(x)s$$

$$0 < c_2 < 1$$

Wolfe - conditions

(I.2.1) Lemma (existence of Wolfe step size)

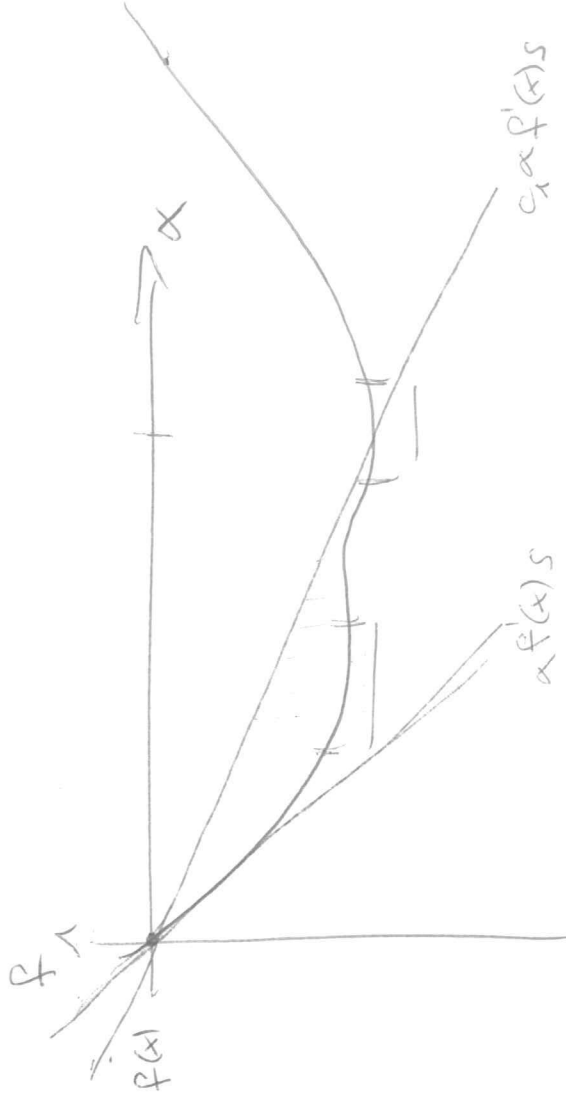
If  $0 < c_1 \leq c_2 < 1$ , then there is  $\alpha > 0$  satisfying conditions (i) and (ii), unless  $f'(x) \neq 0$ .

Proof.  $\varphi(\alpha) := f(x + \alpha s) \in C^1(\mathbb{R})$ ,  $\varphi'(0) < 0$

$L := L_{\varphi(0)} \quad \varphi = \{ \alpha \in \mathbb{R}_+ \mid \varphi(x + \alpha s) \leq f(x) \}$  is closed & bounded

$\Rightarrow \varphi$  assumes a minimum in  $L$ . Let  $\alpha_0$  be the minimum.

If  $\alpha \in \partial L \Rightarrow \varphi(\alpha) = \varphi(0)$ , but due to  $\varphi'(0) < 0$  there are  $\alpha > 0$  with  $\varphi(\alpha) < \varphi(0) \Rightarrow \alpha_0 \in \text{int } L \Rightarrow \varphi'(\alpha_0) = 0$ .





Let  $\underline{\alpha} := \min \{ \alpha \in \mathbb{R}_+ \mid |\varphi'(\alpha)| \leq -c_2 \varphi(0) \}$

$\underline{\alpha}$  satisfies (ii).

$$\begin{aligned} \text{Then } f(x+\alpha s) &= f(x) + \int_0^{\underline{\alpha}} f'(x+\alpha s) s \, d\alpha \\ &< f(x) + \int_0^{\underline{\alpha}} c_2 f'(x) s \, d\alpha \end{aligned}$$

$$\leq f(x) + \underline{\alpha} c_1 f'(x) s \Rightarrow \underline{\alpha} \text{ satisfies (i). } \square$$

### Practical line search

#### (I.2.2) Algorithm of Armijo linesearch

input  $\alpha_0 > 0$

while (i) violated

$$\alpha_{k+1} = \frac{\alpha_k}{2}$$

output  $\alpha_k$

#### (I.2.3) Wolfe line search

(a) find interval containing a "good" step size

$$\underline{\alpha} = 0, \bar{\alpha} > 0$$

while (i) satisfied

$$\bar{\alpha} \leftarrow 2\bar{\alpha}$$

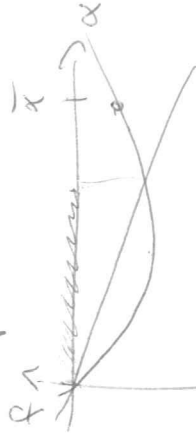
(b) find a Wolfe-admissible point

while true  
 $p = \frac{1}{2}(\underline{\alpha} + \bar{\alpha})$

if (i) violated:  $\bar{\alpha} \leftarrow p$

if  $\varphi'(p) < c_2 \varphi'(0)$ :  $\underline{\alpha} \leftarrow p$

if  $\varphi'(p) > -c_2 \varphi'(0)$ :  $\bar{\alpha} \leftarrow p$   
 else stop



### I.3. Steepest descent (Gradient method)

locally optimal choice of  $s_k$ :

$$s_k \in \arg \min \left\{ \frac{f'(x_k) s_k}{\|s_k\|} \right\} \Rightarrow s_k = \text{gradient}$$

for Euclidean norm:  $s_k = -\nabla f(x_k) = -f'(x)^T$

#### (I.3.1) Theorem

Let  $f'$  be Lipschitz continuous, i.e.  $\|f'(a) - f'(b)\| \leq L \cdot \|a - b\|$ .  
Then the gradient method (Alg. I.1.1. with linesearch I.2.1 and  $s_k = -\nabla f(x_k)$ ) creates a sequence  $(x_k)_k$ . All accumulation points  $\bar{x}$  are stationary:  $f'(\bar{x}) = 0$ .

$$\text{Proof} \quad (f'(x_{k+1}) - f'(x_k)) s_{k+1} \geq (c_2 - 1) f'(x_k) s_k \quad (\text{Wolfe (ii)})$$

$$\Rightarrow \alpha_k L \|s_k\|^2 \geq (c_2 - 1) f''(x_k) s_k \quad \text{Lipschitz}$$

$$f(x_{k+1}) = f(x_k) + c_1 x_k \underbrace{f'(x_k)}_{< 0} s_k \quad (\text{Wolfe (i)})$$

$$\leq f(x_k) + c_1 \frac{c_2 - 1}{L \|s_k\|^2} f'(x_k) s_k \quad f'(x_k) s_k$$

$$= f(x_k) - c_1 \frac{1 - c_2}{L} \| \nabla f(x_k) \|^2 \quad s_k = - \nabla f(x_k)$$

$$= f(x_0) - c_1 \frac{1 - c_2}{L} \sum_{i=0}^k \| \nabla f(x_i) \|^2$$

$$f \text{ bounded from below} \Rightarrow \sum_{i=0}^k \| \nabla f(x_i) \|^2 < \infty$$

$$\Rightarrow \lim_{k \rightarrow \infty} \| \nabla f(x_k) \| = 0$$

Let  $\bar{x}$  be an accumulation point of  $(x_k)_k$ . Then in any neighborhood of  $\bar{x}$  there are points with arbitrarily small gradient. Lipschitz continuity  $\Rightarrow f'(\bar{x}) = 0 \quad \square$

Example:  $f(x) = x^3$  with  $c_1 \geq \frac{1}{3}$ .

## Convergence speed

model problem  $f(x) = \frac{1}{2} x^T A x - b^T x$ ,  $A$  spd  $\Rightarrow$  unique minimize  $A^{-1}b$

gradient:  $\nabla f(x) = Ax - b$

exact linesearch:  $f(x_k - \alpha \nabla f(x_k)) = \min$

$$\frac{1}{2} (x - \alpha \nabla f)^T A (x - \alpha \nabla f) - b^T (x - \alpha \nabla f)$$

differentiate w.r.t  $\alpha$ :  $-\nabla f^T A (x - \alpha \nabla f) + b^T \nabla f = 0$

$$\Rightarrow \alpha = \frac{-b^T \nabla f + \nabla f^T A x}{\nabla f^T A \nabla f} = \frac{\nabla f^T \nabla f}{\nabla f^T A \nabla f}$$

(1.32) Theorem For the model problem, gradient method

with exact linesearch converges linearly

$$: \|x_{k+1} - x_*\|_A \leq \Theta \|x_k - x_*\|_A, \quad \Theta = 1 - \frac{1}{2\kappa}$$

$$\|x\|_A^2 := x^T A x, \quad \kappa = \|A\| \cdot \|A^{-1}\| = \frac{\lambda_{\max}}{\lambda_{\min}} \text{ condition of } A$$

$$\overline{\text{Proof}} \quad f(x_{k+1}) - f(x_k + \alpha_k s_k) = f(x_k) + f'(x_k) \alpha_k s_k + \frac{1}{2} f''(x_k) \alpha_k^2 s_k^2$$

$$= f(x_k) + \alpha_k \nabla f^T \nabla f + \frac{\alpha_k^2}{2} \nabla f^T A \nabla f$$

$$= f(x_k) - \frac{(\nabla f^T \nabla f)^2}{\nabla f^T A \nabla f} + \frac{1}{2} \frac{(\nabla f^T \nabla f)^2}{\nabla f^T A \nabla f}$$

$$= f(x_k) - \frac{1}{2} \frac{\|\nabla f\|_A^2}{\|\nabla f\|^2}$$

$$f(x) = f(x_*) + \underbrace{f'(x_*)}_{=0} (x - x_*) + \frac{1}{2} f''(x_*) (x - x_*)^2$$

$$= f(x_*) + \frac{1}{2} (x - x_*)^T A (x - x_*)$$

$$= f(x_*) + \frac{1}{2} \|x - x_*\|_A^2$$

$$\Rightarrow \|x_{k+1} - x_*\|_A^2 = 2 (f(x_{k+1}) - f(x_*))$$

$$= 2 (f(x_k) - \frac{1}{2} \frac{\|\nabla f\|_A^2}{\|\nabla f\|^2} - f(x_*))$$

$$= \|x_k - x_*\|_A^2 - \frac{\|\nabla f\|_A^2}{\|\nabla f\|^2}$$

Moreover  $\nabla f = Ax_k - b = Ax_k - Ax_* = A(x_k - x_*)$

$$\Rightarrow \|x_k - x_*\|_A^2 = (A^{-1} \nabla f)^T A (A^{-1} \nabla f) = \nabla f^T A^{-1} \nabla f$$

$$\Rightarrow \|x_{k+1} - x_*\|_A^2 = \|x_k - x_*\|_A^2 - \frac{\|\nabla f\|_A^2}{\|\nabla f\|_A^2} = \|x_k - x_*\|_A^2 \left( 1 - \underbrace{\frac{\|\nabla f\|_A^2}{\nabla f^T A^{-1} \nabla f}}_{\Theta_k^2} \right)$$

$$\Theta_k^2 \leq 1 - \frac{1}{d_{\max}^{-1}} = 1 - \frac{1}{\kappa}$$

$$\Rightarrow \Theta_k \leq 1 - \frac{1}{2\kappa} \quad \text{for } \kappa \gg 1$$

□

Remark: Improvement possible:  $\Theta = \frac{\kappa-1}{\kappa+1}$ .

## I.4 Termination criteria

$$f(x_k) - f(x_*) \leq \text{TOL}^2$$

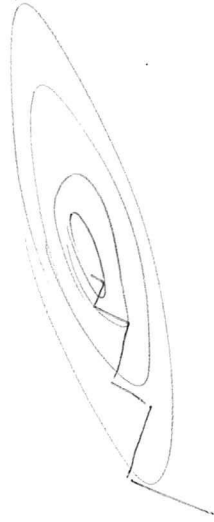
||

$$\frac{1}{2} \|x_k - x_*\|_A^2 \geq \frac{1}{2} \ominus^{-2m} \|x_{k+m} - x_*\|_A^2 = \ominus^{-2m} (f(x_{k+m}) - f(x_*))$$

$$\text{Thus } f(x_k) - f(x_*) = f(x_k) - f(x_{k+m}) + f(x_{k+m}) - f(x_*) \\ \leq f(x_k) - f(x_{k+m}) + \ominus^{2m} (f(x_k) - f(x_*))$$

$$f(x_k) - f(x_*) \leq \frac{f(x_k) - f(x_{k+m})}{1 - \ominus^{2m}} \quad ! \leq \text{TOL}^2$$

$\kappa(f''(x_*)) \gg 1 \Rightarrow$  steepest descent is slow



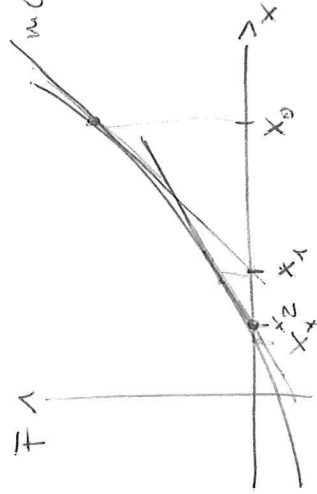
$$\nabla f(x_*) = 0$$

### I.5 Newton's method

ordinary Newton method for  $F(x) = 0$ ,  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$

$$w(x) = F(x_0) + F'(x_0)(x - x_0) = 0$$

$$w(x_1) = 0 \Rightarrow x_1 = x_0 - \underbrace{F'(x_0)^{-1} F(x_0)}_{\Delta x_0}$$



I.5.1 Theorem Let  $F \in C^2(\mathbb{R}^n)$ ,  $F(x_*) = 0$ ,  $F'(x)$  nonsingular

Then, Newton's method converges locally quadratically.



Proof

$$\begin{aligned}
 x_{k+1} - x_* &= x_k - \overbrace{F'(x_k)^{-1} F(x_k) - x_*}^{=0} + F'(x_k) F(x_*) \\
 &= x_k - x_* - F'(x_k)^{-1} (F(x_k) - F(x_*)) \\
 &= x_k - x_* - F'(x_k)^{-1} \int_{t=0}^1 F'(x_* + t(x_k - x_*)) (x_k - x_*) dt \\
 &= x_k - x_* - F'(x_k)^{-1} \int_0^1 [F'(x_k) + F'(x_* + t(x_k - x_*)) - F'(x_k)] (x_k - x_*) dt \\
 &= -F'(x_k)^{-1} \int_{t=0}^1 \underbrace{(F'(x_* + t(x_k - x_*)) - F'(x_k))}_{= O(\|x_k - x_*\|)} dt (x_k - x_*)
 \end{aligned}$$

$$\Rightarrow \|x_{k+1} - x_*\| = O(\|x_k - x_*\|^2) \quad \square$$

## I.6. Affine-conjugate Newton methods

$f(x) \Rightarrow \min$ ,  $F = \nabla f$ ,  $f \in C^2(\mathbb{R}^n)$  strictly convex

$$\Delta x = -F'(x)^{-1} F(x)$$

descent direction:  $-F(x)^T F(x) = F'(x) \Delta x < 0 \Leftrightarrow F'(x)^{-1} \text{ spd} \Leftrightarrow F'(x) \text{ spd}$

choice of stepsize  $\alpha$

$$\begin{aligned} f(x + \alpha \Delta x) &= f(x) + \int_{s=0}^1 f'(x + s\alpha \Delta x) \alpha \Delta x \, ds \\ &= f(x) + \alpha f'(x) \Delta x + \int_{s=0}^1 [f'(x + s\alpha \Delta x) - f'(x)] \alpha \Delta x \, ds \\ &= f(x) + \alpha f'(x) \Delta x + \int_{s=0}^1 \int_{t=0}^1 f''(x + t s \alpha \Delta x) s \alpha \Delta x \, dt \alpha \Delta x \, ds \\ &= f(x) + \alpha f'(x) \Delta x + \frac{\alpha^2}{2} \Delta x^T f''(x) \Delta x \\ &\quad + \int_{s=0}^1 \int_{t=0}^1 (f''(x + t s \alpha \Delta x) - f''(x)) s \alpha^2 \Delta x^2 \, dt \, ds \\ &= f(x) - \alpha \Delta x^T F'(x) \Delta x + \frac{\alpha^2}{2} \Delta x^T F''(x) \Delta x \\ &\quad + \int_{s=0}^1 \int_{t=0}^1 (F'(x + t s \alpha \Delta x) - F'(x)) s \alpha^2 \Delta x^2 \, ds \, dt \end{aligned}$$

## (I.6.1) affine conjugate Lipschitz condition

$$\Delta x^T (\bar{F}'(x + \Delta x) - \bar{F}'(x)) \Delta x \leq \omega (\Delta x^T \bar{F}'(x) \Delta x)^{\frac{3}{2}} \quad \text{f.o.a. } x, \Delta x \in \mathbb{R}^n$$

$$\Rightarrow f(x + \alpha \Delta x) \leq f(x) + (-\alpha + \frac{\alpha^2}{2} + \frac{\omega}{6} \alpha^3 \sqrt{\Delta x^T \bar{F}'(x) \Delta x}) \Delta x^T \bar{F}'(x) \Delta x$$

$$\text{energy norm } \|y\|_x^2 := y^T \bar{F}'(x) y$$

optimal step size

$$t(\alpha) = -\alpha + \frac{\alpha^2}{2} + \frac{\omega \alpha^3}{6} \|\Delta x\|_x$$

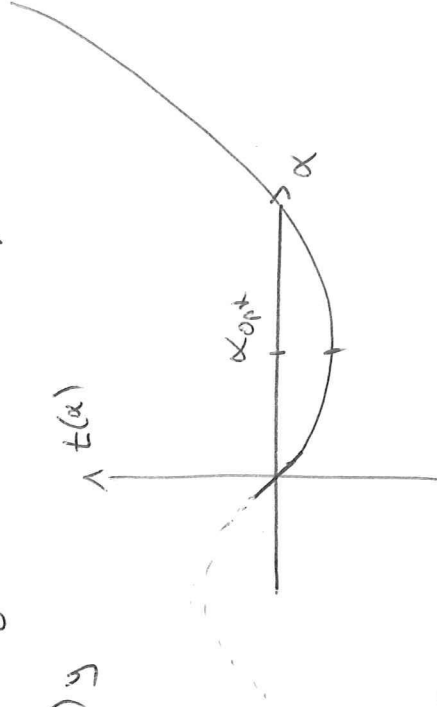
$$\alpha_{\text{opt}} = \arg \min_{\alpha > 0} t(\alpha)$$

$$t'(\alpha) = -1 + \alpha + \frac{\omega}{2} \alpha^2 \|\Delta x\|_x \stackrel{!}{=} 0$$

$$\Rightarrow \alpha_{\text{opt}} = \frac{-1 + \sqrt{1 + 2\omega \|\Delta x\|_x}}{\omega \|\Delta x\|_x}$$

$$= \frac{2}{1 + \sqrt{1 + 2\omega \|\Delta x\|_x}} \in ]0, 1]$$

$\rightarrow$  "damped" Newton method



I.6.4 Theorem  $\alpha_{\text{opt}} = \alpha$  satisfies the Armijo rule I.2.(i), i.e.

$$f(x + \alpha \Delta x) - f(x) \leq c_1 \alpha f'(x) \Delta x \quad \text{with } c_1 = \frac{1}{6}.$$

Proof: Let  $h = \omega \|\Delta x\|_x$ ,  $\alpha = \frac{2}{1 + \sqrt{1 + 2h}} \leq 1$ .

$$\text{Then } h \alpha^2 = \frac{4h}{1 + 2\sqrt{1 + 2h} + 1 + 2h} \leq \frac{4h}{2h} = 2.$$

$$-\frac{t(\alpha)}{\alpha} = 1 - \frac{\alpha}{2} - \frac{h}{6} \alpha^2 \geq 1 - \frac{1}{2} - \frac{2}{6} = \frac{1}{6}.$$

$$\text{Therefore } f(x + \alpha \Delta x) - f(x) \leq t(\alpha) (-f'(x))^T \Delta x = -\frac{t(\alpha)}{\alpha} \alpha f'(x) \Delta x \leq \frac{1}{6} \alpha f'(x) \Delta x \quad \square$$

I.6.5 Theorem If  $\omega \|\Delta x\|_x \leq \frac{2}{3}$ , then  $\alpha$  satisfies the curvature condition

$$\text{I.2.(ii), i.e. } |f'(x + \alpha \Delta x) \Delta x| \leq -c_2 f'(x) \Delta x$$

$$\text{with } c_2 < \frac{2}{3}.$$

$$\text{Proof: } |f'(x + \alpha \Delta x) \Delta x| = \left| f'(x)^T \Delta x + \int_{s=0}^1 f''(x + s\alpha \Delta x) \alpha \Delta x^2 ds \right|$$

$$= \left| f'(x)^T \Delta x + \alpha \Delta x^T f''(x) \Delta x + \int_{s=0}^1 (f''(x + s\alpha \Delta x) - f''(x)) \alpha \Delta x^2 ds \right|$$

$$\begin{aligned}
&\leq \left| -\Delta x^T \bar{F}'(x) \Delta x + \alpha \Delta x^T \bar{F}'(x) \Delta x \right| + \int_{s=0}^1 s \|\Delta x\|_x^3 ds \\
&= \left( \underbrace{-1 + \alpha}_{\leq 0} + \frac{\alpha}{2} \omega \|\Delta x\|_x \right) \|\Delta x\|_x^2 \\
&= \left| -1 + \alpha - \frac{\alpha}{2} \omega \|\Delta x\|_x \right| \|\Delta x\|_x^2 \\
&= \underbrace{\left| -1 + \alpha + \frac{\alpha}{2} \omega \|\Delta x\|_x - \alpha \omega \|\Delta x\|_x \right|}_{= F'(x) = 0} \|\Delta x\|_x^2 \\
&= \underbrace{\alpha \omega \|\Delta x\|_x}_{\leq \frac{3}{2}} (-f'(x) \Delta x)
\end{aligned}$$

For  $h = \omega \|\Delta x\|_x \leq \frac{3}{2}$  it is easily verified that  $\alpha h \leq \frac{2}{3}$  □

## I.7 Practical step size control

depends on (unknown)  $\omega$ !

$$\alpha_{\text{opt}} = \frac{2}{1 + \sqrt{1 + 2\omega \|\Delta x\|_x}}$$

From (I.6.2),  $f(x + \alpha \Delta x) \leq f(x) - (\alpha - \frac{\alpha^2}{2}) \|\Delta x\|_x^2 + \frac{\alpha^3}{6} \omega \|\Delta x\|_x^3$

we obtain  $\frac{\alpha^3}{6} \omega \|\Delta x\|_x^3 \geq \left| f(x + \alpha \Delta x) - f(x) + (\alpha - \frac{\alpha^2}{2}) \|\Delta x\|_x^2 \right| =: |\Xi|$

$$\Rightarrow \omega \geq \frac{6|\Xi|}{\alpha^3 \|\Delta x\|_x^3} =: [\omega]$$

$$\rightarrow \text{computable step size } [\alpha_{\text{opt}}] = \frac{2}{1 + \sqrt{1 + 2[\omega] \|\Delta x\|_x}} \geq \alpha_{\text{opt}}$$

### I.7.1 Lemma (bitcounting lemma)

Let  $\omega \leq 2[\omega]$ . Then

$$(I.7.2) \quad f(x + [\alpha_{\text{opt}}] \Delta x) \leq f(x) - \frac{1}{6} [\alpha_{\text{opt}}] ([\alpha_{\text{opt}}] + 2) \|\Delta x\|_x^2$$

Proof Let  $\alpha := [\alpha_{\text{opt}}]$ .

$$\begin{aligned} f(x + \alpha \Delta x) - f(x) &\leq -\alpha + \frac{\alpha^2}{2} + \frac{\alpha^3}{6} \omega \|\Delta x\|_x^2 \\ &\leq -\alpha + \frac{\alpha^2}{2} + \frac{\alpha^3}{6} 2[\omega] \|\Delta x\|_x^2 \end{aligned}$$

$\alpha$  is obtained by minimizing  $[t] = -\alpha + \frac{\alpha^2}{2} + \frac{\alpha^3}{6} [\omega] \|\Delta x\|_x^2$   
 $\Rightarrow 1 - \alpha - \frac{\alpha^2}{2} [\omega] \|\Delta x\|_x^2 = 0 \Rightarrow \frac{\alpha^2}{2} [\omega] \|\Delta x\|_x^2 = 1 - \alpha$

$$\begin{aligned} \Rightarrow f(x + \alpha \Delta x) - f(x) &\leq \left(-\alpha + \frac{\alpha^2}{2} + \frac{2}{3} \alpha (1-\alpha)\right) \|\Delta x\|_x^2 \\ &\leq \alpha \left(-\frac{1}{3} - \frac{\alpha}{6}\right) \|\Delta x\|_x^2 = -\frac{\alpha}{6} (1+\alpha) \|\Delta x\|_x^2 \end{aligned}$$

□

### I.7.3 Damped Newton iteration

choose  $[\omega]_{-1} > 0, x_0 \in \mathbb{R}^n$

for  $k = 0, \dots$

$$\bar{F}'(x_k) \Delta x_k = -\bar{F}(x_k)$$

$$[\omega]_k \leftarrow [\omega]_{k-1}$$

while (I.7.2) violated

$$[\omega]_k \leftarrow [\omega]_k (f(x_k + [\alpha_{\text{opt}}] \Delta x_k)) > [\omega]_k$$

$$x_{k+1} \leftarrow x_k + [\alpha_{\text{opt}}] \Delta x_k$$

$$[\omega]_k \leftarrow [\omega]_k (f(x_k + [\alpha_{\text{opt}}] \Delta x_k)) \quad // \text{ decrease in } [\omega]_k \text{ possible}$$

## Termination criterion

require (i) quadratic model is good approximation,  
ordinary Newton method converges to  
a unique minimizer

$$\Leftrightarrow [\omega] \|\Delta x\|_x \leq 1$$

$$(ii) \quad f(x) \leq f(x_*) + \text{TOL}^2$$

$$\begin{aligned} f(x_k) - f(x_*) &= \sum_{i=k}^{\infty} (f(x_i) - f(x_{i+1})) \\ &\leq \sum_{i=k}^{\infty} \left( \alpha_i - \frac{\alpha_i^2}{2} + \frac{\omega^3}{6} \alpha_i^3 \|\Delta x_i\|_{x_i} \right) \|\Delta x_i\|_{x_i}^2 \\ &= \sum_{i=k}^{\infty} \left( \frac{1}{2} + \frac{1}{6} \right) \|\Delta x_i\|_{x_i}^2 \\ &\leq \frac{5}{6} \sum_{i=k}^{\infty} \ominus^{2(i-k)} \|\Delta x_k\|_{x_k}^2 \\ &= \frac{5}{6} \frac{1}{1-\ominus^2} \|\Delta x_k\|_{x_k}^2 \end{aligned}$$

$$\ominus \approx \frac{\|\Delta x_k\|_{x_k}}{\|\Delta x_{k-1}\|_{x_{k-1}}}$$



## II.8 Nonconvex Problems

### II.8.1 Hessian-modifications

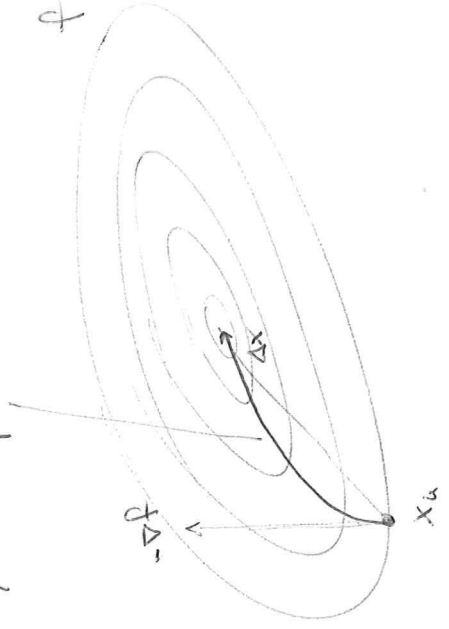
$Sx = -B^{-1}F(x)$  with  $B$  spd  $\Rightarrow Sx$  is descent direction  
simple choice  $B = F'(x) + \eta M$ ,  $M$  spd  
( $M = I$ : Levenberg-Marquardt method)

#### Choice of $\eta$

$\eta \rightarrow 0$ : ordinary Newton method

$\eta \rightarrow \infty$ : gradient method (preconditioned)  
with small step size

$\eta \mapsto Sx(\eta)$  for  $\eta \gg \eta_0$  defines a  
Levenberg-Marquardt path



II.8.2 Theorem Let  $\eta_k \geq 0$  be minimal such that

$B \geq \varepsilon I$  for  $\varepsilon > 0$ , and Armijo condition holds, then all accumulation points of  $(x_k)_k$  are stationary.

Proof as for Thm. I.3.1  $\square$

### II.8.3 Trust Region methods

quadratic model  $f(x) + f'(x) \delta x + \frac{1}{2} \delta x^T f''(x) \delta x$   
good for "small"  $\delta x$ , i.e.  $\|\delta x\|_H \leq \rho$

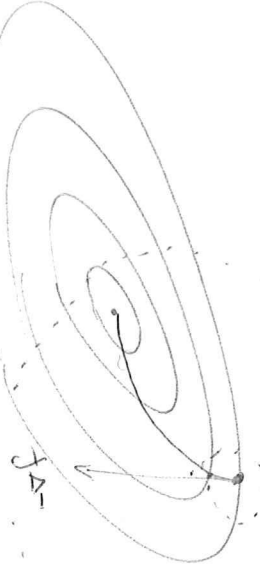
$$\delta x \in \arg \min_{\|\delta x\|_H \leq \rho} f(x) + f'(x) \delta x + \frac{1}{2} \delta x^T f''(x) \delta x$$

#### Choice of $\rho$

$\rho \rightarrow 0$ : preconditioned gradient method  
with very small step size

$\rho \rightarrow \infty$ : if  $f''(x)$  spd  $\rightarrow$  Newton

$\delta x(\rho)$  "defines" the trust-region path  $f$



II 8.4. Theorem Let  $\rho_k$  be chosen maximal such that  $Sx_k$  satisfy the Armijo condition. Then, all accumulation points  $\bar{z}$  of  $(x_k)_k$  are stationary and satisfy the necessary second order optimality conditions  $f''(\bar{z}) \succeq 0$ .

Proof

Stationarity of  $\bar{z}$  as in I.3.1.

(i) Due to  $f'' \in C(\mathbb{R}^n)$ ,  $f''$  is bounded in a neighborhood of  $\bar{z}$ :

$\|f''(x)\| \leq R$ . Then with  $\rho = \frac{\epsilon}{2R} (1 - c_1)$ ,  $s = -\nabla f(x_k)$

$$f(x_k + \rho s) \leq f(x_k) + \rho f'(x_k)s + \rho^2 \frac{R}{2} \|s\|^2$$

$$\leq f(x_k) + \rho(1 - \rho \frac{R}{2}) f'(x_k)s$$

$$= f(x_k) + c_1 \rho f'(x_k)s$$

$\Rightarrow x_{k+1}(\rho)$  satisfies Armijo  $\Rightarrow \rho_k \geq \rho$  bounded from below

(ii) wlog.  $x_k \rightarrow \bar{z}$ .  $f(\bar{z}) \leq f(x_k)$ .

Assume  $\exists p \in \mathbb{R}^n, \|p\|_1 = 1: p^T f''(\bar{z})p = -\epsilon < 0$ .

Then  $f(x_{k+1}) \leq f(x_k + \rho p) \leq \underbrace{f(x_k) + f'(x_k)\rho p}_{\rightarrow f(\bar{z})} \underbrace{\rightarrow 0}_{\rightarrow \frac{\epsilon}{2}\rho^2}$

$$\rightarrow_{k \rightarrow \infty} f(\bar{z}) - \frac{\epsilon}{2}\rho^2 < f(\bar{z}). \quad \square$$

## II.8.5 Minimization of cubic upper bound

Assumption:  $F' = f''$  Lipschitz continuous

$$\|F'(x) - F'(y)\|_H \leq \omega \|x - y\|_H$$

$$\text{Then } f(x + \delta x) = f(x) + F'(x) \delta x + \frac{1}{2} \delta x^T F'(x) \delta x + \frac{\omega}{6} \|\delta x\|_H^3 =: w(x)$$

(II.8.6)

$$\delta x \in \arg \min w(x)$$

Choice of  $\omega$

$\omega \rightarrow 0$ : ordinary Newton method

$\omega \rightarrow \infty$ : preconditioned gradient method,

small step size

$\delta x(\omega)$  "defines" a path for  $\omega \geq 0$ .

## II.8.7 Theorem

The paths  $\delta x(\eta)$ ,  $\delta x(\rho)$ ,  $\delta x(\omega)$  are identical for sufficiently large  $\eta, \rho$  small  $\rho$  up to parametrization.

Proof: (i) Let  $\delta x = \delta x(\omega)$ . Then (II.8.6) implies

$$\begin{aligned} 0 &= F'(x) + F'(x) \delta x + \frac{\omega}{6} \cdot \frac{3}{2} (\delta x^T \delta x)^{1/2} \cdot 2 \delta x \\ &= F'(x) + (F'(x) \delta x + \underbrace{\frac{\omega}{2} \|\delta x\|_H}_{\eta} \delta x) \delta x \end{aligned}$$

$$\Rightarrow \delta x = \delta x(\eta) \text{ with } \eta = \frac{\omega}{2} \|\delta x\|_H.$$

$$(ii) \quad a) \|Sx\|_Y < \rho \Rightarrow 0 = F(x) + F'(x) Sx$$

$$\Rightarrow Sx = Sx(\eta) \text{ with } \eta = 0 \\ = Sx(\omega) \text{ with } \omega = 0$$

$$b) \|Sx\|_Y = \rho \Rightarrow \exists \lambda \in \mathbb{R}: 0 = F(x) + F'(x) Sx + \lambda \Pi Sx \\ = F(x) + \underbrace{(F'(x) + \lambda \Pi)}_{\eta} Sx$$

$$\Rightarrow Sx(\eta) \text{ with } \eta = \lambda.$$

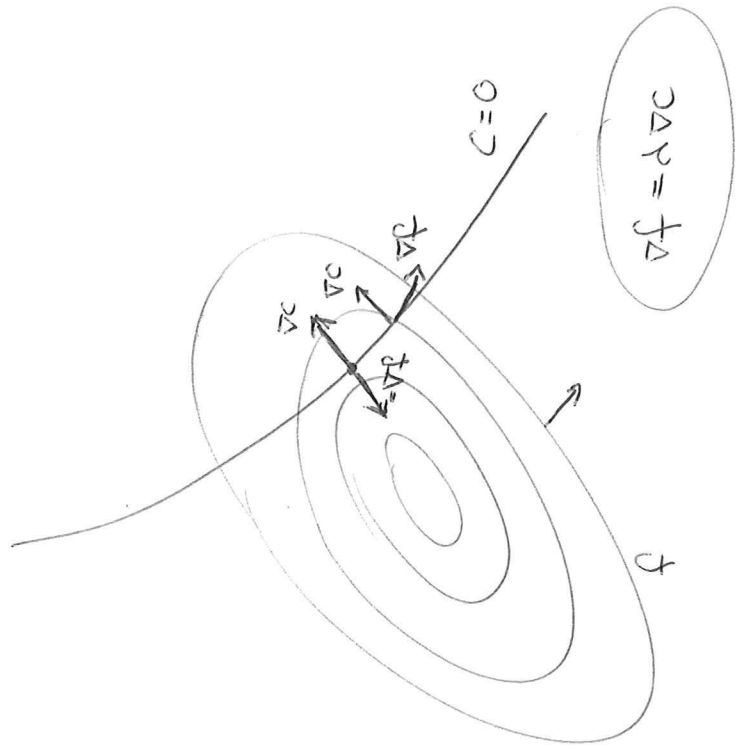
□

### III. Equality constrained optimization

$$(III.1) \quad \min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad c(x) = 0, \quad \begin{array}{l} f \in C^2(\mathbb{R}^n, \mathbb{R}) \\ c \in C^2(\mathbb{R}^n, \mathbb{R}^m) \end{array}$$

$\rightarrow$  admissible set  $U = \{x \in \mathbb{R}^n \mid c(x) = 0\}$

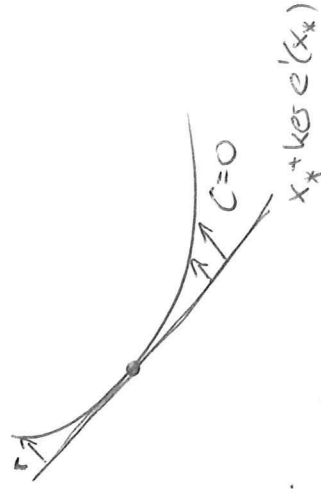
### III.2 Necessary optimality conditions



II. 2.1 Theorem Let  $x_*$  minimizer of (III.1). Then  $f'(x_*) \cdot \delta x = 0$  for all  $\delta x \in \ker c'(x_*)$ .

Proof Let  $A = c'(x_*) \in \mathbb{R}^{m \times n}$  and  $\mathbb{R}^n = \ker A \oplus (\ker A)^\perp$   
 $\rightarrow$  unique representation of  $x = x^0 + x^\perp$ ,  $c(x) = c(x^0, x^\perp)$ .

The implicit function theorem guarantees the existence of  $r \in C^1(\ker A, (\ker A)^\perp)$  with  $c(x^0, x^\perp) = 0 \Leftrightarrow x^\perp = x^0 + r(x^0)$  in a neighborhood of  $x_*$ .



Let  $\delta x \in \ker A$  and define

$$x(\varepsilon) = x_* + \varepsilon \delta x + r(x_* + \varepsilon \delta x).$$

$\rightarrow x(\varepsilon)$  is admissible ( $c(x(\varepsilon)) = 0$ ).

As  $x_*$  is minimizer,  $\varepsilon = 0$  is a minimizer

of  $f(x(\varepsilon))$  and hence

$$0 = \left. \frac{\partial f(x(\varepsilon))}{\partial \varepsilon} \right|_{\varepsilon=0} = f'(x_*) \cdot x'(0) = f'(x_*) [\delta x + r'(x_*) \delta x]$$

$$\text{Imp. fct. thm. } r'(x_*) = \frac{\partial c}{\partial x^\perp}(x_*)^{-1} \underbrace{\frac{\partial c}{\partial x^0}(x_*)}_{=0} = 0.$$

$$0 = f'(x_*) \delta x.$$

□

### III.2.2 Definition

Let  $U \subseteq X$  be a subspace. Then

$U^\circ = \{v \in X^* \mid \forall u \in U : v \cdot u = 0\}$  is called annihilator of  $U$ .

III.2.3 Theorem Let  $A \in \mathbb{R}^{m \times n}$  have full rank (linear independence constraint qualification)  
LICQ

$$\text{Then } (\ker A)^\circ = \text{im } \bar{A}^T \quad m \leq n$$

Proof: (i)  $Sx \in \ker A: (\bar{A}^T \lambda)^T Sx = \bar{A}^T A Sx = 0 \Rightarrow \text{im } \bar{A}^T \subset (\ker A)^\circ$

(ii)  $\dim (\ker A)^\circ = n - \dim \ker A = n - (n-m) = m \quad \left. \begin{array}{l} \dim \text{im } \bar{A}^T = \dim (\ker A)^\circ \\ \dim \text{im } \bar{A}^T = m \end{array} \right\}$

$$\Rightarrow \text{im } \bar{A}^T = (\ker A)^\circ$$

□

### III.2.4

#### Corollary (Karush-Kuhn-Tucker)

Let  $x_*$  be the minimizer of (III.1) and let  $c'(x_*)$  have full rank. Then there is a Lagrange multiplier  $\lambda \in \mathbb{R}^m$  such that

$$f'(x_*) = \bar{A}^T c'(x_*) \Leftrightarrow \nabla f(x_*) = c'(x_*)^T \lambda$$

$$c(x_*) = 0$$



III.2.1 Theorem Let  $x_*$  be minimizer of (III.1). Then  $f'(x_*)\delta x = 0$  for all  $\delta x \in \ker C'(x_*)$ .

Proof Let  $A = C'(x_*) \in \mathbb{R}^{m \times n}$  and  $\mathbb{R}^n = \ker A \oplus (\ker A)^\perp$ .  
 Unique representation of  $x = x^0 + x^\perp$ . Let  $c(x) = c(x^0, x^\perp)$ .  
 The Impl. fct. then guarantees existence of  $r \in C'(\ker A, (\ker A)^\perp)$   
 with  $c(x^0, x^\perp) = 0 \Leftrightarrow x^\perp = x^0 + r(x^0)$  in a neighborhood of  $x_*$ .

Let  $\delta x \in \ker A$  and define  $x(\varepsilon) = x_* + \varepsilon \delta x + r(x_* + \varepsilon \delta x)$   
 $\Rightarrow x(\varepsilon)$  admissible ( $c(x(\varepsilon)) = 0$ ) for small  $|\varepsilon|$ .

As  $x_*$  is minimizer,  $\varepsilon = 0$  is minimizer of  $f(x(\varepsilon))$  and

$$\text{hence } 0 = \frac{\partial f(x(\varepsilon))}{\partial \varepsilon} \bigg|_{\varepsilon=0} = f'(x_*) \cdot x'(0) = f'(x_*) [\delta x + r'(x_*)\delta x].$$

$$\text{Impl. fct. then: } r'(x_*) = \frac{\partial c}{\partial x^\perp}(x_*)^{-1} \cdot \underbrace{\frac{\partial c}{\partial x^0}(x_*)}_{=0} = 0.$$

$$\Rightarrow 0 = f'(x_*)\delta x$$

□

[vice but algorithmically not practical since we'd need a basis of  $\ker A$  in every point]  
 Can we characterize  $f'(x_*)$  better?

III.2.2 Definition Let  $U \subseteq X$  be a subspace. Then  
 $U^\circ = \{v \in X^* \mid \forall u \in U : v \cdot u = 0\}$  is called the annihilator  
 of  $U$ .

### III.2.5 Example

$$\min_{x \in \mathbb{R}^2} \frac{1}{2} x^T x \quad \text{s.t.} \quad x_1 + x_2 = 1$$

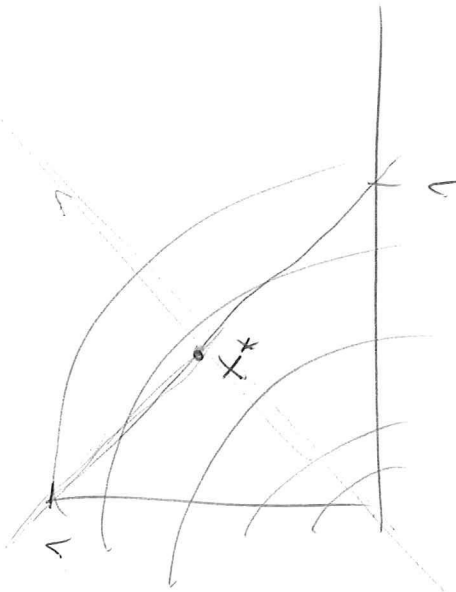
$$\rightarrow c' = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$\begin{cases} x - \begin{bmatrix} 1 \\ 1 \end{bmatrix} \lambda = 0 \\ \begin{bmatrix} 1 & 1 \end{bmatrix} x - 1 = 0 \end{cases}$$

$$x = \lambda \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$\lambda \underbrace{\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix}}_2 = 1$$

$$\lambda = \frac{1}{2}$$



### III.3. Second order conditions

$$\min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t.} \quad c(x) = 0 \in \mathbb{R}^m, \quad m < n$$

$$\exists \lambda \in \mathbb{R}^m: \nabla f(x_*) - c'(x_*)^T \lambda = 0$$

$$-c(x_*) = 0$$

#### III.3.1 Definition

The Lagrange function is

$$L(x, \lambda) = f(x) - c(x)^T \lambda.$$

$$\nabla_x L(x, \lambda) = \nabla f(x) - c'(x)^T \lambda$$

$$\nabla_\lambda L(x, \lambda) = -c(x)$$

$\text{KKT points are}$   
 $\Rightarrow$  stationary points  
of Lagrangian

2nd derivative

$$L_{xx} = f'' - c''(x)^T \lambda = f'' - \sum_i \lambda_i c_i''$$

### III. 3.2 Theorem

Let  $c'(x_*) = 0$  have full rank and  $(x_*, d)$  be a KKT point.

(i)  $x_*$  is local minimizer  $\Rightarrow L_{xx}(x_*, d)$  positive semidefinite on  $\ker c'(x_*)$

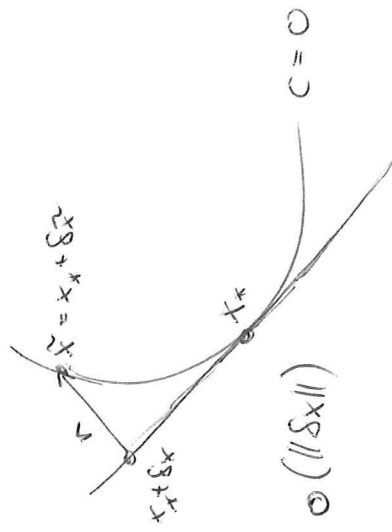
(ii)  $L_{xx}$  positive definite on  $\ker c'(x_*) \Rightarrow x_*$  is local minimizer

Proof Let  $\tilde{x} = \tilde{x} + x_*$  with  $c(\tilde{x}) = 0$ . Then

$$f(\tilde{x}) = f(\tilde{x}) - \underbrace{\tilde{L}c(\tilde{x})}_{=0} = L(\tilde{x}, d) = L(x_*, d) + \underbrace{L_x(x_*, d)\tilde{x}}_{=0} + \frac{1}{2}\tilde{x}^T L_{xx}(x_*, d)\tilde{x} + o(\|\tilde{x}\|^2)$$

$$= f(x_*) + \frac{1}{2}\tilde{x}^T L_{xx}\tilde{x} + o(\|\tilde{x}\|^2)$$

Now let  $\delta x \in \ker c'(x_*)$  and  $\tilde{x} = \delta x + r(x_* + \delta x)$



$$\Rightarrow \|r(x_* + \delta x)\| = o(\|\delta x\|)$$

Then  $S_x^T L_{xx} \tilde{S}_x - S_x^T L_{xx} S_x$

$$= \underbrace{(S_x - \tilde{S}_x)^T L_{xx} (S_x + \tilde{S}_x)}_{\|S_x - \tilde{S}_x\| = o(\|S_x\|)} = o(\|S_x\|^2)$$

Thus,  $f(\tilde{x}) - f(x_*) = \frac{1}{2} S_x^T L_{xx} S_x + o(\|S_x\|^2)$ .

(i)  $x_*$  minimizes  $\Rightarrow f(\tilde{x}) - f(x_*) \geq 0 \quad \forall \tilde{x} = x_* + S_x(\tilde{x}), S_x \in \ker c'(x_*) \cap B_\varepsilon(0)$   
 $\Rightarrow S_x^T L_{xx} S_x \geq -o(\|S_x\|^2) \Rightarrow S_x^T L_{xx} S_x \geq 0$ .

(ii)  $L_{xx}$  pos. def. on  $\ker c'(x_*)$

$$\begin{aligned} \Rightarrow \exists \alpha > 0 : 0 < \alpha &: f(\tilde{x}) - f(x_*) \geq \frac{1}{2} S_x^T L_{xx} S_x - o(\|S_x\|^2) \\ &= \frac{\alpha}{2} \|S_x\|^2 - o(\|S_x\|^2) \\ &= \frac{\alpha}{2} \|S_x\|^2 \text{ for } \|S_x\| \text{ small enough.} \quad \square \end{aligned}$$

### III.4 Newton-KKT

$L'(x, \lambda) = 0 \Rightarrow$  Newton's method

$$F(x) = \begin{bmatrix} \nabla f(x) - C'(x)^T \lambda \\ -c(x) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$F'(x) = \begin{bmatrix} L_{xx} & L_{x\lambda} \\ L_{\lambda x} & L_{\lambda\lambda} \end{bmatrix} = \begin{bmatrix} H & C^T \\ C & -C'(x) \end{bmatrix} \quad \text{with } H = f''(x) - C''(x)^T \lambda$$

III.4.1 Theorem Assume  $f, c \in C^{2,1}(\mathbb{R}^n)$  and  $(x_*, \lambda)$

satisfy sufficient 2nd order conditions.

Then, Newton's method converges locally quadratically to  $(x_*, \lambda)$ .

Proof: Thm 15.1: need to show  $F'$  is Lipschitz continuous and  $F'$  is invertible

$$\text{Let } F'(x_*, \lambda) \begin{bmatrix} \delta x \\ \delta \lambda \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow C \delta x = 0 \Rightarrow \delta x \in \ker C$$

$$0 = \begin{bmatrix} \delta x^T & \delta \lambda^T \end{bmatrix} \underbrace{F'(x_*, \lambda) \begin{bmatrix} \delta x \\ \delta \lambda \end{bmatrix}}_{=0} = \delta x^T H \delta x \geq \alpha \|\delta x\|^2$$

$$\Rightarrow \delta x = 0 \quad | \quad C^T \delta \lambda = 0 \Rightarrow \delta \lambda = 0 \quad \square$$

### III. 5. Penalty method

$$(III. 5.1) \quad \min_{x \in \mathbb{R}^n} f(x) \quad \text{s.t. } c(x) = 0 \in \mathbb{R}^m, \quad c, f \in C^2(\mathbb{R}^n)$$

Idea: transform (approximately) to unconstrained problem  
• apply methods for unconstr. problems

#### III. 5.2 quadratic penalty

$$f_\mu(x) = f(x) + \frac{1}{\mu} \|c(x)\|^2, \quad \mu > 0$$

$$x_\mu \in \arg \min f_\mu(x)$$

$$x_\mu \rightarrow ?$$

III. 5.3 Theorem Let  $\mu_k \rightarrow 0$  and  $x_k$  a global minimizer of  $f_{\mu_k}$ . Then, every accumulation point  $\tilde{x}$  of  $(x_k)_k$  is a global minimizer of (III. 5.1), and  $c(x_k) = O(\sqrt{\mu_k})$ .

Proof Let  $x_0$  be feasible, i.e.  $c(x_0) = 0$ . Wlog

$x_k \rightarrow \tilde{x}$ . Then

$$f(x_0) = f_{\mu_k}(x_0) \geq f_{\mu_k}(x_k) = f(x_k) + \frac{1}{\mu_k} \|c(x_k)\|^2$$

$$(i) \Rightarrow \|c(x_k)\|^2 \leq \mu_k (f(x_0) - f(x_k)) \leq \mu_k \cdot \text{const} \Rightarrow c(x_k) = O(\sqrt{\mu_k})$$

$$\Rightarrow c(\xi) = 0 \text{ feasibility}$$

$$(ii) f(x_0) \geq f(\xi) \text{ optimality}$$

□

III.5.4. Theorem let  $(x_*, d)$  satisfy sufficient 2nd order cond.

$$\text{Then } x_\mu - x_* = O(\sqrt{\mu})$$

$$\frac{2}{\mu} c(x_\mu) - d = O(\sqrt{\mu})$$

$$\text{Proof } L'(x_\mu, \frac{2}{\mu} c(x_\mu)) = \underbrace{L'(x_*, d)}_{=0} + L''(x_*, d) \left[ \frac{x_\mu - x_*}{\frac{2}{\mu} c(x_\mu) - d} \right] + o\left(\|x_\mu - x_*\| + \left\| \frac{2}{\mu} c(x_\mu) - d \right\|\right)$$

$$\begin{bmatrix} \nabla f(x_\mu) - \frac{2}{\mu} c'(x_\mu) c(x_\mu) \\ -c(x_\mu) \end{bmatrix} = \begin{bmatrix} 0 \\ -c(x_\mu) \end{bmatrix} = O(\|c(x_\mu)\|).$$

$$= O(\sqrt{\mu}).$$

□



$$f_{\mu}(x) = f(x) + \frac{1}{\mu} \|c(x)\|^2 \quad \rightarrow x_{\mu} \quad \|x_{\mu} - x_{\mu}^*\| = O(\sqrt{\mu})$$

III. 5.5 Corollary Assumptions as in III. 5.4. Then,  $\|x_{\mu} - x_{\mu}^*\| = O(\sqrt{\mu})$ .

Proof:  $\| \frac{2}{\mu} c(x_{\mu}) - \lambda \| = O(\sqrt{\mu})$   
 $\Rightarrow \|c(x_{\mu})\| = \frac{\mu}{2} (\|\lambda\| + O(\sqrt{\mu})) = O(\mu) \quad \square$   
 $\leq \frac{2}{\mu} \|c(x_{\mu})\| - \|\lambda\|$

Drawbacks of penalty

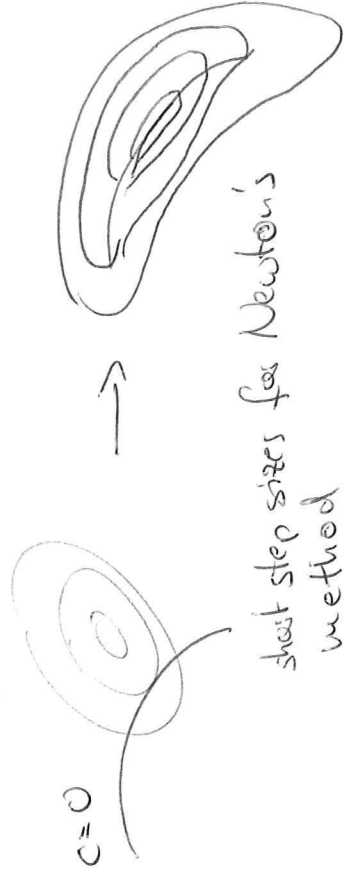
(i) Hessian is ill-conditioned

$$f_{\mu}''(x) = \underbrace{f''(x)}_{=O(1)} + \underbrace{\frac{2}{\mu} c'(x)^T c'(x)}_{O(\frac{1}{\mu})} + \underbrace{\frac{2}{\mu} c(x)^T c''(x)}_{O(1)}$$

$\rightarrow \frac{2}{\mu} c'(x)^T c'(x)$  not regular due to

$\text{rank} \leq m \Rightarrow \ker c'(x)$  of dim  $n-m > 0$

(ii)  $f_{\mu}$  develops steep curved valleys (banana shaped)



partial remedy: continuation



### III. 5.6 Algorithm

$x_0, \mu_0$  given

choose  $\sigma < 1$

for  $k=0, \dots$

$$\mu_{k+1} = \sigma \mu_k$$

$$x_{k+1} = x_k - f''_{\mu_{k+1}}(x_k)^{-1} \nabla f_{\mu_{k+1}}(x_k)$$

if  $\sigma$  sufficiently large  $\rightarrow$  Newton converges quickly

### III.6 Augmented Lagrangian method

observations: (i)  $L(x, d) = f(x)$  for admissible  $x$

(ii)  $x_* = \arg \min_x L(x, d_*)$  if  $f$  convex

$$\text{Idea: } L_\mu(x, d) = L(x, d) + \frac{1}{\mu} \|C(x)\|^2$$

III.6.1 Theorem Let  $x, d$  satisfy necessary & sufficient optimality conditions.

Assume  $L'CC$  holds. Then  $\exists \bar{\mu} > 0: \forall \mu \leq \bar{\mu}$ :

$x$  is locally unique minimizer of  $L_\mu(x, d)$

$$\text{Proof: (i) } L'_\mu(x, d) = \underbrace{f'(x) + C'(x)^T d}_{=0} + \frac{2}{\mu} C'(x)^T \underbrace{C(x)}_{=0} = 0$$

$$\begin{aligned} \text{(ii) } L''_\mu(x, d) &= \underbrace{L''_{xx}(x, d)}_H + \frac{2}{\mu} C'(x)^T C'(x) + \frac{2}{\mu} \underbrace{C(x)^T C''(x)}_{=0} \\ &= H + \frac{2}{\mu} C^T C \end{aligned}$$

Let  $u \in \mathbb{R}^n$  arbitrary. Then  $\exists w \in \mathbb{R}^m, v \in \mathbb{R}^m$ :

$$u = w + C^T v$$

$$u^T L''_\mu u = u^T H u + \frac{2}{\mu} \|C u\|^2$$

$$= w^T H w + 2w^T H C^T v + v^T C H C^T v + \frac{2}{\mu} \|C C^T v\|^2$$

$$\begin{aligned}
&\geq \alpha \|w\|^2 - 2 \|HC^T\| \cdot \|w\| \cdot \|v\| - \|CHC^T\| \cdot \|v\|^2 + \frac{2}{\mu} \|CC^T v\|^2 \\
&= \underbrace{\left( \sqrt{\alpha} \|w\| - \frac{\|HC^T\|}{\sqrt{\alpha}} \|v\| \right)^2}_{\geq 0} + \underbrace{\left( \frac{2}{\mu} \beta - \|CHC^T\| - \frac{\|HC^T\|^2}{\alpha} \right) \|v\|^2}_{> 0 \text{ for suff. small } \mu} \geq \delta \|w\|^2 \quad \square
\end{aligned}$$

III.6.2 Theorem Let  $x_*, d_*$  satisfy necessary & sufficient conditions & LICQ.

Then there is a neighborhood of  $(x_*, d_*)$  such that for  $\mu < \bar{\mu}$ :

$x_\mu := \arg \min_x L_\mu(x, d)$  satisfies

$$\|x_\mu - x_*\| + \|d - \frac{2}{\mu} c(x_\mu) - d_*\| \leq M_\mu \|d - d_*\|$$

Proof  $x_\mu(d)$  exists as implicit function for  $L'_\mu(x, d) = 0$  in  $U$ .

Since  $L''_\mu$  is spd (see proof III.6.1),  $x_\mu$  is a minimizer.

$$\begin{bmatrix} 0 \\ c(x_\mu) \end{bmatrix} = L'(x_\mu, d - \frac{2}{\mu} c(x_\mu)) = \underbrace{L'(x_\mu, d_*)}_{=0} + \underbrace{L''(x_\mu, d_*)}_{\text{invertible}} \begin{bmatrix} x_\mu - x_* \\ d - \frac{2}{\mu} c(x_\mu) - d_* \end{bmatrix} + o(\|x_\mu - x_*\| + \|d - \frac{2}{\mu} c(x_\mu) - d_*\|)$$

$$\Rightarrow \|x_\mu - x_*\| + \|d - \frac{2}{\mu} c(x_\mu) - d_*\| \leq \|L''(x_\mu, d_*)^{-1}\| \cdot \|c(x_\mu)\| + o(\|x_\mu - x_\mu\| + \dots)$$

$$\Rightarrow \|x_\mu - x_*\| + \|d - \frac{2}{\mu} c - d_*\| \leq 2\delta \|c(x_\mu)\|$$

$$\Rightarrow \frac{2}{\mu} \|c(x_\mu)\| - \|d - d_*\| \leq 2\delta \|c(x_\mu)\|$$

$$\Rightarrow \|c(x_\mu)\| \leq \frac{\mu}{2} (1 - \mu\delta)^{-1} \|d - d_*\|$$

$$\Rightarrow \|x_\mu - x_*\| + \|d - \frac{2}{\mu} c - d_*\| \leq \mu\delta \underbrace{(1 - \mu\delta)^{-1}}_M \|d - d_*\| \quad \square$$

Apparently,  $d - \frac{2}{\mu} c(x_n)$  is a better estimate of  $d_*$  than  $d$  itself (by a factor  $\mu/2$ ): first order multiplier update.

### III. 6.3 Algorithm (augmented Lagrangian method)

for  $i = 1, \dots$

$$x_{i+1} = \arg \min_x L_\mu(x, d_i)$$

$$d_{i+1} = d_i - \frac{2}{\mu} c(x_{i+1})$$

III. 6.4 Corollary AL method converges locally linearly if  $\mu$  is sufficiently small.

Proof:  $\|d_{i+1} - d_*\| \leq \mu/2 \|d_i - d_*\|$  converges linearly

$$\|x_{i+1} - x_*\| \leq \mu/2 \|d_i - d_*\|$$

□

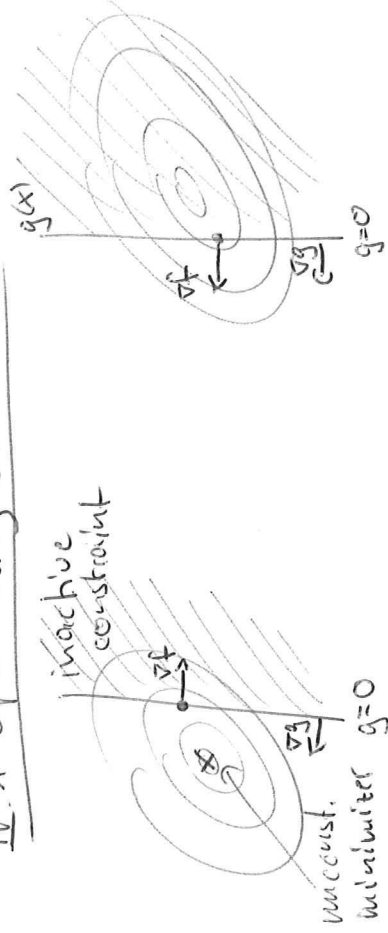
## IV Inequality Constraints

min  $f(x)$  s.t.  $g(x) \geq 0$  component-wise in  $\mathbb{R}^k$   
 $x \in \mathbb{R}^n$

observation:  $k > n$  possible, e.g.

$$\begin{bmatrix} n=2 \\ k=4 \end{bmatrix}$$

### IV.1 Optimality conditions



Necessary conditions:  $\exists \lambda \in \mathbb{R}^k$ :

$$\nabla f(x_*) = g'(x_*)^T \lambda \quad \text{adjoint equation}$$

$$g(x) \geq 0, \lambda \geq 0$$

$$\lambda^T g(x_*) = 0 \quad \text{complementarity}$$

if constraint qualifications are satisfied  
 e.g. LICQ (linear independence of active  
 constraint gradients).

## IV. 1.1 Definition

active set  $A(x) = \{i \in \{1, \dots, k\} \mid g_i(x) \leq 0\}$

inactive set  $I(x) = \{i \in \{1, \dots, k\} \mid g_i(x) > 0\}$

Sometimes, LICQ is too restrictive.

### Example

(i)  $\min x+y$  s.t.  $x \geq 0, y \geq 0, x+y \geq 0$

(ii)  $\min x$  s.t.  $y \geq 0, x^2 - y \geq 0$  "wrong directions" of constraint gradients



Mangasarian - Fromovitz - CQ (MFCQ):

$\exists d \in \mathbb{R}^n : g'_A(x_*) d > 0$  ("all constraints point in similar directions")

(i) MFCQ ✓

(ii) MFCQ ✗



## IV.2 Active set methods

Idea: ignore inactive constraints with  $g_i(x_*) > 0$ .  
then  $\min_x f(x)$  s.t.  $g_A(x) = 0$  has solution  $x_*$ .

strategy: (i) find correct active set  
(ii) solve equality constrained problem

### IV.2.1 Simple algorithm

$$A_0 \leftarrow \emptyset$$

for  $k=0, \dots$

$$x_k \leftarrow \arg \min f(x) \text{ s.t. } g_{A_k}(x) = 0$$

if KKT satisfied

break

$$A_{k+1} \leftarrow A_k(x_k) \setminus \{i \mid \lambda_i = 0\}$$

nasty details:

- solvability of equality constrained subproblems
- cycling: choose an active set  $A_k$  encountered before

efficiency:

- in practice, often good
- in theory, exponential run time possible  
(combinatorial number of possible active sets)

### IV.3. Penalty & barrier methods

Idea: penalize infeasibility

$$\min f(x) + \frac{1}{\mu} \sum_i \min(0, g_i(x))^2$$

consider  $\mu \rightarrow 0$

- same approach as for equality constraints
- similar performance, same issues (infeasible iterates, ill-conditioned Hessians, banana-shaped valleys)

augmented Lagrangian approach possible as well

Idea: penalize approach to infeasibility

$$\min f(x) - \mu \sum_i \log g_i(x)$$

consider  $\mu \rightarrow 0$ .

- feasible iterates, but starting points are difficult to define
- ill-conditioned Hessians
- path-following (homotopy) for  $\mu \rightarrow 0$  (central path) is efficient, both practically and theoretical: polynomial run time for linear problems.