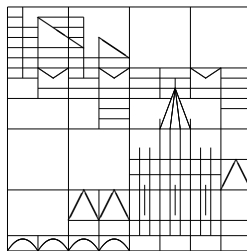


# Skript zur Vorlesung

## Optimierung

Sommersemester 2010

Johannes Schropp



Universität Konstanz

Fachbereich Mathematik und Statistik

Stand: 16. Juli 2010

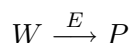


# Kapitel I: Optimierungsaufgaben und Optimierungskriterien

## 1 Einführende Beispiele und Problemstellung

### Beispiel 1.1

Betrachte den Prozess



d.h. das Substrat  $W$  wird mit Hilfe des Enzyms  $E$  in ein Produkt  $P$  umgewandelt. Die Geschwindigkeit  $v$  dieser Umwandlung hängt gemäß der Michaelis-Menten-Theorie aus der Biologie wie folgt von  $w$  ab

$$v = \frac{\mu w}{K + w} = g(\mu, K, w)$$

mit  $w \hat{=} [W]$  Konzentration von  $W$ . Es ist

$$g(\mu, K, K) = \frac{\mu}{2}, \quad \lim_{w \rightarrow \infty} g(\mu, K, w) = \mu.$$

$\mu$  bezeichnet die maximal mögliche Geschwindigkeit und  $K$  die halbmaximale Rate. Die Geschwindig-

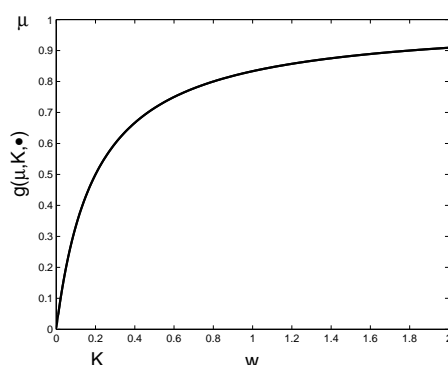


Abbildung 1:  $g(\mu, K, w)$  für  $K = 0.2$  und  $\mu = 1$ .

keit  $v$  kann nun in Abhängigkeit von  $w$  gemessen werden, d.h. wir erhalten die Messwerte  $(v_i, w_i), i = 1, \dots, m$ .

**Ziel:** Bestimme die Parameter  $(\mu, K)$  aus den Messdaten. Setze dazu  $x = (x_1, x_2) = (\mu, K)$

$$f(x_1, x_2) = \sum_{i=1}^m (g(w_i, x_1, x_2) - v_i)^2$$

und bestimme ein  $x^*$  mit

$$f(x^*) \leq f(x) \text{ für } x \in \mathbb{R}^2, x \geq 0.$$

### Beispiel 1.2 (Gleichgewichtslagen eines Ringelwurms)

Verwende als Modell eines Ringelwurms eine Folge von  $N$  aneinandergereihten Segmenten. Als Segmentform verwenden wir symmetrische Pyramidenstümpfe. Aufgrund der Symmetrie wird das Aussehen des Ringelwurms mit  $N$  Segmenten vollständig durch den Vektor

$$x = (b_1, c_1, b_2, c_2, \dots, b_N, c_N, b_{n+1}) \in \mathbb{R}^{2N+1}$$

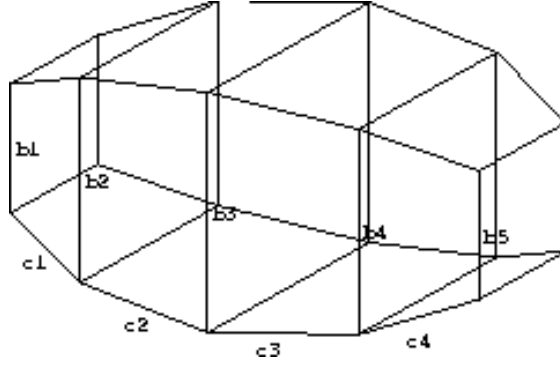


Abbildung 2: Symmetrisches Modell für  $N = 4$ .

beschrieben. Dabei bezeichnet  $c_i, i = 1, \dots, N$  die Länge des Längsmuskels des  $i$ -ten Segmentes und  $b_i, i = 1, \dots, N$  die Länge des Ringmuskels des  $i$ -ten Segmentes.  $b_{N+1}$  bezeichnet die Länge des abschließenden Ringmuskels.

Das Volumen des Ringelwurms ergibt sich zu

$$V(x) = \frac{1}{3} \sum_{i=1}^N \left[ c_i^2 - \frac{(b_i - b_{i+1})^2}{2} \right]^{1/2} \cdot (b_i^2 + b_i b_{i+1} + b_{i+1}^2).$$

Das  $N$ -segmentige Modell hat  $8N + 4$  Kanten, die Längs- bzw. Ringmuskeln entsprechen. Wir modellieren alle Muskeln als Federn. Die in einer Feder mit Ruhelage  $L_0$  in Abhängigkeit der Federkonstante  $\alpha$  enthaltene Energie ist

$$P_\varepsilon(L, \alpha) = \int_{L_0}^L F_\varepsilon(l, \alpha) dl$$

mit dem Kraftgesetz

$$F_\varepsilon(L, \alpha) = \frac{\alpha}{\varepsilon} \tan(\varepsilon(L - L_0)), \quad L_0 - \frac{\pi}{2\varepsilon} < L < L_0 + \frac{\pi}{2\varepsilon}.$$

Mit dem Superpositionsprinzip erhält man für den Ringelwurm im Zustand  $x$  mit den Federkonstanten  $\beta$  und  $\gamma$  die Energie

$$E(x) = 4 \left( \sum_{j=1}^{N+1} P_\varepsilon(b_j, \beta) + \sum_{j=1}^N P_\varepsilon(c_j, \gamma) \right).$$

Die Gleichgewichtslagen eines Ringelwurms mit vorgegebenen Volumen  $V_0$  ergeben sich nun als Zustände  $x^*$  mit minimaler Energie  $E(x^*)$  unter der Nebenbedingung  $V(x^*) = V_0$ .  $x^*$  löst also die Aufgabe:

$$\text{Minimiere } E(x) \text{ unter der Nebenbedingung } V(x) - V_0 = 0.$$

Unter einem **endlichdimensionalen stetigen Optimierungsproblem** verstehen wir die Aufgabe:

$$(I) = \left\{ \begin{array}{l} \text{Sei } D \subset \mathbb{R}^N \text{ offen, } f : D \rightarrow \mathbb{R}, g : D \rightarrow \mathbb{R}^l \text{ und } k : D \rightarrow \mathbb{R}^q. \text{ Setze} \\ Z = \{x \in D \mid g_i(x) = 0, i = 1, \dots, l \text{ und } k_i(x) \geq 0, i = 1, \dots, q\}. \\ \text{Gesucht ist } x^* \in Z \text{ mit } f(x^*) \leq f(x), x \in Z. \end{array} \right.$$

### Bemerkung 1.3

- a.) Ist  $l = q = 0$ , d.h. treten  $g$  und  $k$  nicht auf, so heißt das Problem **unrestringiert**, andernfalls **restringiert**.
- b.)  $Z$  heißt **Zulässigkeitsbereich** und  $f$  nennt man die *Zielfunktion*.
- c.) Soll  $f$  maximiert werden, so ist dies gleichbedeutend damit, dass  $-f$  minimiert wird.
- d.) Ist  $l, q > 0$ , so nennt man  $g(x) = 0$  die **Gleichungs-** und  $k(x) \geq 0$  die **Ungleichheitsrestriktionen**.

### Definition 1.4

Vorgelegt sei (I).  $x^* \in D$  heißt

- a.) **globale Minimalstelle** von  $f$ , falls  $f(x^*) \leq f(x)$  für alle  $x \in Z$ .  $f(x^*)$  heißt dann **globales Minimum**.
- b.) **strikte globale Minimalstelle** von  $f$ , falls  $f(x^*) < f(x)$  für alle  $x \in Z \setminus \{x^*\}$ .  $f(x^*)$  heißt dann **striktes globales Minimum**.
- c.) **lokale Minimalstelle** von  $f$ , wenn es eine Umgebung  $U$  von  $x^* \in \mathbb{R}^N$  gibt mit  $f(x^*) \leq f(x)$  für alle  $x \in U \cap Z$ .  $f(x^*)$  heißt dann **lokales Minimum**.
- d.) **strikte lokale Minimalstelle** von  $f$ , wenn es eine Umgebung  $U$  von  $x^* \in \mathbb{R}^N$  gibt mit  $f(x^*) < f(x)$  für alle  $x \in (U \cap Z) \setminus \{x^*\}$ .  $f(x^*)$  heißt dann **striktes lokales Minimum**.

### Bemerkung 1.5

Ein Punkt  $x^*$  ist genau dann (globale, strikte globale, lokale, strikte lokale) Maximalstelle von  $f$ , wenn  $x^*$  (globale, strikte globale, lokale, strikte lokale) Minimalstelle von  $-f$  ist.

### Definition 1.6

Sei  $D \subset \mathbb{R}^N$  offen und sei  $h \in C^1(D, \mathbb{R})$ . Ein Punkt  $x^* \in D$  heißt **stationärer Punkt** von  $h$ , falls  $\nabla h(x^*) = 0$ .

## 2 Die Kuhn-Tucker Theorie

### 2.1 Optimalitätskriterien für lokale Minimalstellen

Wir behandeln unter geeigneten Differenzierbarkeitsannahmen notwendige und hinreichende Bedingungen für lokale Minimalstellen.

Betrachte zunächst das **freie Minimierungsproblem**:

$$(II) = \left\{ \begin{array}{l} \text{Sei } D \subset \mathbb{R}^N \text{ offen, } f : D \rightarrow \mathbb{R}. \text{ Gesucht ist } x^* \in D \text{ mit } f(x^*) \leq f(x) \text{ für } x \in U(x^*). \end{array} \right.$$

Aus der Analysis II sind die folgenden Sätze bekannt.

#### Satz 2.1 (Notwendige Bedingungen für lokale Minima)

Vorgelegt sei (II).

- a.) Ist  $f \in C^1(D, \mathbb{R})$  und ist  $x^*$  ein lokales Minimum von  $f$ , so gilt  $\nabla f(x^*) = 0$ , d.h.  $x^*$  ist stationärer Punkt von  $f$ .
- b.) Ist  $f \in C^2(D, \mathbb{R})$  und ist  $x^*$  lokales Minimum von  $f$ , so gilt  $\nabla f(x^*) = 0$  und  $\nabla^2 f(x^*)$  positiv semidefinit.

### Bemerkung 2.2

Die Bedingungen aus Satz 2.1 sind nicht hinreichend dafür, dass  $x^*$  lokales Minimum von  $f$  ist. Betrachte z.B.  $f(x_1, x_2) = x_1^2 - x_2^4$ ,  $x^* = (0, 0)$ . Dann gilt:

$$\nabla f(x^*) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \nabla^2 f(x^*) = \begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix}$$

positiv semidefinit, aber  $x^*$  ist kein lokales Minimum.

### Satz 2.3 (Hinreichende Bedingungen für lokale Minima)

Vorgelegt sei (II) mit  $f \in C^2(D, \mathbb{R})$ . Gelten

a.)  $\nabla f(x^*) = 0$

b.)  $\nabla^2 f(x^*)$  positiv definit

für ein  $x^* \in D$ , so ist  $x^*$  strikte lokale Minimalstelle von  $f$ .

### Bemerkung 2.4

Die Bedingungen aus Satz 2.3 sind nicht notwendig dafür, dass  $x^*$  striktes lokales Minimum von  $f$  ist.

Betrachte z.B.

$$f(x_1, x_2) = x_1^2 + x_2^4, \quad x^* = (0, 0).$$

Dann gilt

$$\nabla f(x^*) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \nabla^2 f(x^*) = \begin{pmatrix} 2 & 0 \\ 0 & 0 \end{pmatrix},$$

d.h.  $\nabla^2 f(x^*)$  ist nicht positiv definit, aber  $x^*$  ist striktes lokales Minimum von  $f$ .

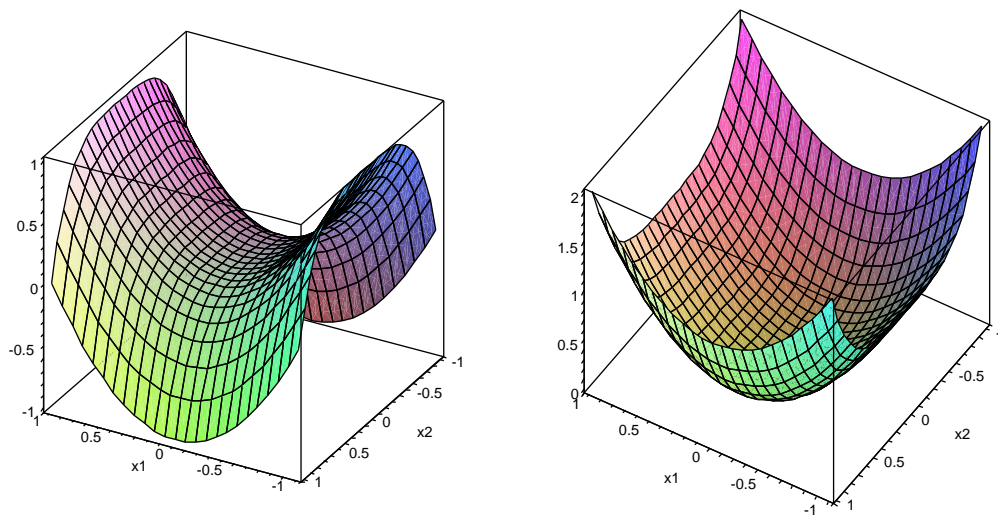


Abbildung 3: Links:  $f(x_1, x_2) = x_1^2 - x_2^4$ . Rechts:  $f(x_1, x_2) = x_1^2 + x_2^4$ .

## 2.2 Tangenten und linearisierte Richtungen

Im Folgenden wollen wir notwendige und hinreichende Bedingungen für lokale Minima für restriktierte Optimierungsprobleme herleiten. Dazu betrachten wir wieder die Aufgabe:

$$(III) = \left\{ \begin{array}{l} \text{Sei } D \subset \mathbb{R}^N \text{ offen und seien } f \in C^2(D, \mathbb{R}), g \in C^2(D, \mathbb{R}^l) \text{ und } k \in C^2(D, \mathbb{R}^q). \text{ Es sei} \\ Z = \{x \in D \mid g(x) = 0 \text{ und } k(x) \geq 0\} \\ \text{abgeschlossen und nicht-leer. Gesucht wird ein } x^* \in Z \text{ mit} \\ f(x^*) \leq f(x) \quad \text{für alle } x \in U(x^*) \cap Z. \end{array} \right.$$

$x^*$  heißt **lokales Minimum** von  $f$  unter den Nebenbedingungen  $g(x) = 0$  und  $k(x) \geq 0$ .  $Z$  heißt Menge der zulässigen Punkte. Ziel dieses Abschnittes ist es, lokale Minima  $x^*$  durch analytische Bedingungen zu charakterisieren.

### Definition 2.5

a.) Sei  $\tilde{x} \in Z$ . Dann heißt

$$\mathcal{A}_{\tilde{x}} = \{i \in \{1, \dots, q\} \mid k_i(\tilde{x}) = 0\}$$

Menge der bei  $\tilde{x} \in Z$  **aktiven Ungleichheitsrestriktionen**.

b.)  $\tilde{x} \in Z$  heißt **regulär**, falls die Vektoren

$$\nabla g_1(\tilde{x}), \dots, \nabla g_l(\tilde{x}), \nabla k_i(\tilde{x}), i \in \mathcal{A}_{\tilde{x}}$$

linear unabhängig sind.

### Definition 2.6

Sei wieder  $x \in Z$ . Eine Folge  $(y_k)_{k \in \mathbb{N}} \in Z^{\mathbb{N}}$  mit  $\lim_{k \rightarrow \infty} y_k = x$  heißt **zulässige Folge**.  $d$  heißt **Tangente** in  $x$ , falls eine zulässige Folge und eine Folge  $(t_k)_{k \in \mathbb{N}}, t_k > 0, \lim_{k \rightarrow \infty} t_k = 0$  existieren mit

$$\lim_{k \rightarrow \infty} \frac{y_k - x}{t_k} = d.$$

Die Menge aller Tangenten zu  $Z$  in  $x$  heißt **Tangentenkegel** und wird mit  $T_Z(x)$  bezeichnet.

### Bemerkung 2.7

Eine Menge  $X$  heißt **Kegel**, falls  $0 \in X$  und falls

$$x \in X \Rightarrow \alpha x \in X, \alpha > 0$$

gilt.

### Definition 2.8

Sei  $x \in Z$  und sei  $\mathcal{A}_x$  die Menge der bei  $x$  aktiven Ungleichheitsrestriktionen. Dann heißt

$$F(x) = \{d \in \mathbb{R}^N \mid d^T \nabla g_i(x) = 0, i = 1, \dots, l \text{ und } d^T \nabla k_i(x) \geq 0, i \in \mathcal{A}_x\}$$

**Menge der bei  $x$  erlaubten linearisierten Richtungen.**

Erläuterungen zum Tangentenkegel  $T_Z(x)$  und der erlaubten linearisierten Richtungen  $F(x)$ :

### Beispiel 2.9

Sei  $f(x_1, x_2) = x_1 + x_2$  und  $g(x_1, x_2) = x_1^2 + x_2^2 - 2 = 0$ . Der Zulässigkeitsbereich lautet

$$Z = \{(x_1, x_2) \in \mathbb{R}^2 \mid x_1^2 + x_2^2 - 2 = 0\}$$

Betrachte  $\bar{x} = (-\sqrt{2}, 0)$ , dann lautet  $\nabla g(-\sqrt{2}, 0) = \begin{pmatrix} -2\sqrt{2} \\ 0 \end{pmatrix}$ .

Eine zulässige Folge  $y_k$  ist  $y_k = \begin{pmatrix} -\sqrt{2 - (\alpha_k)^2} \\ -\alpha_k \end{pmatrix}$ ,  $\alpha_k \rightarrow 0$  für  $k \rightarrow \infty$  mit  $\alpha_k > 0$ . Setze  $t_k =$

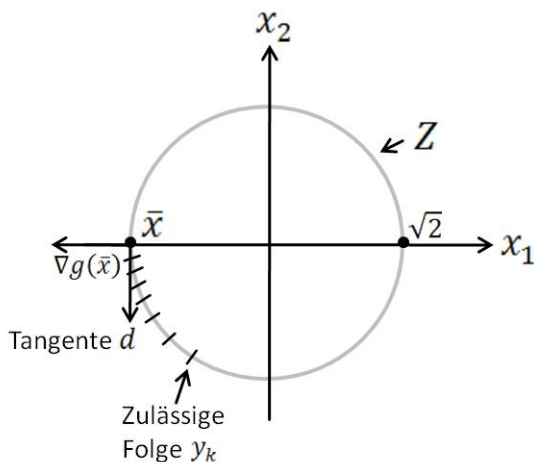


Abbildung 4:

$\alpha_k > 0, t_k \rightarrow 0$  für  $k \rightarrow \infty$  und finde

$$\frac{y_k - \bar{x}}{t_k} = \frac{\begin{pmatrix} y_1^k \\ y_2^k \end{pmatrix} - \begin{pmatrix} -\sqrt{2} \\ 0 \end{pmatrix}}{\alpha_k} = \frac{\begin{pmatrix} -\sqrt{2 - \alpha_k^2} + \sqrt{2} \\ -\alpha_k \end{pmatrix}}{\alpha_k} = \begin{pmatrix} \frac{1}{\alpha_k}(\sqrt{2} - \sqrt{2 - \alpha_k^2}) \\ -1 \end{pmatrix}.$$

Sei  $h(x) = \frac{\sqrt{2} - \sqrt{2 - x^2}}{x}$ .

$$\lim_{x \rightarrow 0} h(x) = \lim_{x \rightarrow 0} \frac{-\frac{1}{2\sqrt{2-x^2}} \cdot 2x}{1} = \lim_{x \rightarrow 0} \frac{-2x}{2\sqrt{2-x^2}} = 0 \Rightarrow \frac{y_k - \bar{x}}{t_k} = \begin{pmatrix} 0 \\ -1 \end{pmatrix} =: d.$$

Ferner folgt  $\nabla g(\bar{x})^T d = (-2\sqrt{2}, 0) \begin{pmatrix} 0 \\ -1 \end{pmatrix} = 0$ , d.h.  $d \in F(x)$ .

### Beispiel 2.10

Sei  $f(x_1, x_2) = x_1 + x_2$  und  $k(x_1, x_2) = 2 - x_1^2 - x_2^2 \geq 0$ . Der Zulässigkeitsbereich lautet

$$Z = \{(x_1, x_2) \in \mathbb{R}^2 \mid 2 - x_1^2 - x_2^2 \geq 0\}$$

Betrachte  $\bar{x} = (-\sqrt{2}, 0)$ . Wähle eine zulässige Folge auf einer Geraden nach  $\bar{x}$ :

$$y_k = \begin{pmatrix} -\sqrt{2} \\ 0 \end{pmatrix} + \frac{1}{k} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix},$$

$w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \in \mathbb{R}^2$  mit  $w_1 > 0$ . Alle derartigen Folgen sind zulässig. Somit folgt



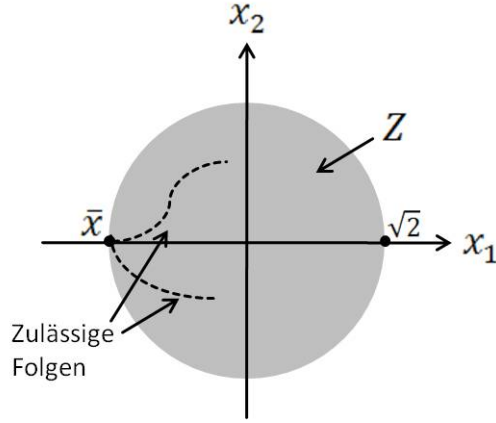


Abbildung 5:

$$y_k = \bar{x} + \frac{1}{k}w = \bar{x} + t_k s_k$$

mit  $t_k = \frac{1}{k}$ ,  $\lim_{k \rightarrow \infty} t_k = 0$ ,  $s_k = w$ ,  $\lim_{k \rightarrow \infty} s_k = w$  und man findet

$$T_Z(-\sqrt{2}, 0) = \{w \in \mathbb{R}^2 \mid w_1 \geq 0\}.$$

Für die erlaubten linearen Richtungen in  $x^*$  gilt nach Definition

$$F(x^*) = \{d \in \mathbb{R}^2 \mid d^T \nabla k(\bar{x}) \geq 0\}.$$

Mit  $\nabla k(-\sqrt{2}, 0) = \begin{pmatrix} 2\sqrt{2} \\ 0 \end{pmatrix}$  folgt

$$F(-\sqrt{2}, 0) = \{d \in \mathbb{R}^2 \mid d^T \begin{pmatrix} 2\sqrt{2} \\ 0 \end{pmatrix} \geq 0\} = \{d \in \mathbb{R}^2 \mid d_1 \geq 0\} = T_Z(-\sqrt{2}, 0).$$

In diesem Beispiel stimmen  $F(x)$  und  $T_Z(x)$  offensichtlich überein.

**Lemma 2.11**

a.) Sei  $x \in Z$ , so gilt

$$T_Z(x) \subset F(x).$$

b.) Ist  $x \in Z$  regulär, so folgt

$$T_Z(x) = F(x).$$

**Beweis:**

a.) Sei  $(y_k)_{k \in \mathbb{N}}$  eine zulässige Folge, d.h.  $\lim_{k \rightarrow \infty} y_k = x$  und sei  $(t_k)_{k \in \mathbb{N}}$ ,  $t_k > 0$ ,  $\lim_{k \rightarrow \infty} t_k = 0$  mit  $s_k = \frac{y_k - x}{t_k}$ ,  $\lim_{k \rightarrow \infty} s_k = d$ , d.h.  $d \in T_Z(x)$ . Dann gilt mit dem Mittelwertsatz

$$0 = g_i(y_k) - g_i(x) = t_k s_k^T \nabla g_i(\eta_i^k), \quad \eta_i^k = x + \theta_i^k (y_k - x), \quad \theta_i^k \in [0, 1], \quad i = 1, \dots, l \quad (1)$$

$$0 \leq \underbrace{k_i(y_k)}_{\geq 0 \text{ da } y_k \in Z} - \underbrace{k_i(x)}_{=0 \text{ da } i \in \mathcal{A}_x} = t_k s_k^T \nabla k_i(\eta_{l+i}^k), \quad \eta_{l+i}^k = x + \theta_{l+i}^k (y_k - x), \quad \theta_{l+i}^k \in [0, 1], \quad i \in \mathcal{A}_x. \quad (2)$$

Dividiert man in (1), (2) durch  $t_k$  und geht dann zum Grenzwert  $k \rightarrow \infty$  über, so erhält man

$$0 = d^T \nabla g_i(x), i = 1, \dots, l, \quad 0 \leq d^T \nabla k_i(x), i \in \mathcal{A}_x$$

d.h.  $d \in F(x)$ .

b.) Sei  $x \in Z$  und sei  $H(z) = \begin{pmatrix} g(z) \\ k_i(z), i \in \mathcal{A}_x \end{pmatrix} \in \mathbb{R}^{l+m}$ ,  $m = \#\mathcal{A}_x$ . Nach Voraussetzung hat dann  $DH(x) \in \mathbb{R}^{l+m, N}$  den Rang  $l + m$ , d.h.

$$\dim(N(DH(x))) = N - \text{rg}(DH(x)) = N - l - m,$$

wobei  $N(DH)$  der Kern (Nullraum) der Abbildung  $DH(x)$  ist. Sei nun  $Z \in \mathbb{R}^{N, N-l-m}$  eine Matrix mit  $R(Z) = N(DH(x))$ , wobei  $R(Z)$  das Bild von  $Z$  ist, d.h. die Spalten von  $Z$  bilden eine Basis des Nullraums von  $DH(x)$ . Sei nun  $d \in F(x)$  beliebig. Zu zeigen:  $d \in T_Z(x)$ . Setze

$$R(z, t) = \begin{pmatrix} H(z) - tDH(x)d \\ Z^T(z - x - td) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

$R : D \times \mathbb{R} \rightarrow \mathbb{R}^{l+m} \times \mathbb{R}^{N-l-m} = \mathbb{R}^N$  mit  $D \subset \mathbb{R}^N$  offen und  $R$  zweimal stetig differenzierbar. Nach Konstruktion gilt

$$R(x, 0) = \begin{pmatrix} H(x) \\ Z^T(x - x) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \frac{\partial R}{\partial z}(z, t) = \begin{pmatrix} DH(z) \\ Z^T \end{pmatrix} \quad \frac{\partial R}{\partial z}(x, 0) = \begin{pmatrix} DH(x) \\ Z^T \end{pmatrix} \in \mathbb{R}^{N, N}$$

invertierbar nach Voraussetzung. Also existiert nach dem Satz über impliziten Funktionen eine Funktion  $\varphi \in C^2[-\varepsilon, \varepsilon[, \mathbb{R}^N)$  mit  $\varphi(0) = x$  und

$$R(\varphi(t), t) = 0, \quad |t| < \varepsilon. \quad (3)$$

Sei nun  $(t_k)_{k \in \mathbb{N}}$ ,  $t_k \in ]0, \varepsilon[$ ,  $t_k \rightarrow 0$ . Dann ist  $(y_k)_{k \in \mathbb{N}}$ ,  $y_k = \varphi(t_k)$  eine zulässige Folge, denn

$$\lim_{k \rightarrow \infty} y_k = \lim_{k \rightarrow \infty} \varphi(t_k) = \varphi(0) = x$$

und nach Konstruktion gilt

$$\begin{pmatrix} g(y_k) \\ k_i(y_k), i \in \mathcal{A}_x \end{pmatrix} = H(y_k) = t_k DH(x)d = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \alpha_1 \\ \vdots \\ \alpha_m \end{pmatrix}$$

mit  $\alpha_i \geq 0$  für  $i = 1, \dots, m$  da

$$DH(x) = \begin{pmatrix} \nabla g_1(x)^T \\ \vdots \\ \nabla g_l(x)^T \\ \nabla k_i(x)^T, i \in \mathcal{A}_x \end{pmatrix}$$

und  $d \in F(x)$ , d.h.  $y_k \in Z$ . Man beachte, dass gilt  $k_i(y_k) > 0$  für  $i \in \{1, \dots, q\} \setminus \mathcal{A}_x$  da  $k_i(x) > 0$  für  $i \in \{1, \dots, q\} \setminus \mathcal{A}_x$ .

Außerdem folgt

$$\frac{y^k - x}{t_k} = \frac{\varphi(t_k) - \varphi(0)}{t_k} \rightarrow \varphi'(0)$$

für  $k \rightarrow \infty$ . Noch zu zeigen bleibt  $\varphi'(0) = d$ , denn dann folgt  $d \in T_Z(x)$ . Differentiation von (3) liefert

$$\frac{\partial R}{\partial z}(\varphi(t), t) \cdot \varphi'(t) + \frac{\partial R}{\partial t}(\varphi(t), t) = 0, \quad |t| < \varepsilon$$

und somit folgt mit  $\varphi(0) = x$  sofort

$$\varphi'(0) = -\frac{\partial R}{\partial z}(x, 0)^{-1} \frac{\partial R}{\partial t}(x, 0).$$

Finde

$$\frac{\partial R}{\partial t}(z, t) = \begin{pmatrix} -DH(x)d \\ -Z^T d \end{pmatrix} = - \begin{pmatrix} DH(x) \\ -Z^T \end{pmatrix} d = -\frac{\partial R}{\partial z}(x, 0)d$$

und erhalte

$$\varphi'(0) = -\frac{\partial R}{\partial z}(x, 0)^{-1} \left( -\frac{\partial R}{\partial z}(x, 0) \right) d = d.$$

□

Desweiteren benötigen wir noch ein fundamentales Resultat über Abstiegsrichtungen von Lösungen von Optimierungsproblemen.

### Lemma 2.12

Sei  $x^* \in Z$  eine Lösung des Optimierungsproblems. Dann gilt

$$\nabla f(x^*)^T d \geq 0 \quad \forall d \in T_Z(x^*).$$

**Beweis:** Sei  $d \in T_Z(x^*)$  beliebig und sei  $(y_k) \in Z^{\mathbb{N}}$  eine zulässige Folge,  $(t_k)_{k \in \mathbb{N}}$ ,  $t_k > 0$ ,  $\lim_{k \rightarrow \infty} t_k = 0$  mit

$$d_k = \frac{y_k - x^*}{t_k}, \quad \lim_{k \rightarrow \infty} d_k = d.$$

Nach Taylor und da  $x^*$  lokale Lösung gilt

$$0 \leq f(y_k) - f(x^*) = t_k d_k^T \nabla f(\eta_0^k), \quad \eta_0^k = x^* + \theta_0^k (y_k - x^*), \quad \theta_0^k \in [0, 1], \quad k \in \mathbb{N}. \quad (4)$$

Division von (4) durch  $t_k$  und Grenzübergang  $k \rightarrow \infty$  liefert

$$0 \leq d^T \nabla f(x^*).$$

□

Der wichtigste Schritt auf dem Weg zur Charakterisierung notwendiger Bedingungen für eine lokale Lösung ist das Lemma von Farkas.

### Lemma 2.13 (Farkas)

Sei der Kegel  $K$  definiert durch

$$K = \{By + Cw \mid y \geq 0\}$$

mit  $B \in \mathbb{R}^{n,m}$ ,  $C \in \mathbb{R}^{n,p}$ ,  $y \in \mathbb{R}^m$ ,  $w \in \mathbb{R}^p$  und sei  $g \in \mathbb{R}^n$  beliebig. Dann gilt

- *entweder*

$$g \in K$$

- *oder*

es existiert ein  $d \in \mathbb{R}^n$  mit  $g^T d < 0$ ,  $B^T d \geq 0$  und  $C^T d = 0$ .

**Beweis:** Wir zeigen zuerst, dass beide Bedingungen zugleich nicht wahr sein können. Sei  $g \in K$ , d.h.  $g = By + Cw$  mit  $y \geq 0$ . Gilt überdies  $d^T g < 0$ , so folgt

$$0 > d^T g = d^T By + d^T Cw = \underbrace{(B^T d)^T}_{\geq 0} \underbrace{y}_{\geq 0} + \underbrace{(C^T d)^T}_{=0} w \geq 0$$

und damit der Widerspruch.

Wir zeigen nun, dass eine der beiden Alternativen eintritt. Konkret konstruieren wir  $d$  mit  $g^T d < 0$ ,  $B^T d \geq 0$ ,  $C^T d = 0$  falls  $g \notin K$ .

Nach Konstruktion ist  $K$  abgeschlossen. Sei jetzt  $\hat{s}$  die Lösung von

$$f(s) = \min\{\|s - g\|_2^2 \mid s \in K\}.$$

Da  $K$  abgeschlossen gilt  $\hat{s} \in K$ .  $K$  ist ein Kegel, da  $0 \in K$  und  $\hat{s} \in K \Rightarrow \alpha \hat{s} \in K$ ,  $\alpha > 0$ .

Betrachte jetzt

$$h(\alpha) = \min\{\|\alpha \hat{s} - g\|_2^2 \mid \alpha \geq 0\}. \quad (5)$$

Wegen  $\alpha \hat{s} \in K$  für  $\alpha \geq 0$  wird (5) minimiert für  $\alpha = 1$  und es folgt

$$\begin{aligned} \frac{d}{d\alpha} \|\alpha \hat{s} - g\|_2^2 \Big|_{\alpha=1} &= \frac{d}{d\alpha} (\alpha^2 \hat{s}^T \hat{s} - 2\alpha \hat{s}^T g + g^T g) \Big|_{\alpha=1} = 0 = (-2\hat{s}^T g + 2\alpha \hat{s}^T \hat{s}) \Big|_{\alpha=1} = -2\hat{s}^T g + 2\hat{s}^T \hat{s} \\ &\iff \hat{s}^T (\hat{s} - g) = 0. \end{aligned}$$

Setze  $d := \hat{s} - g$  und zeige die 'oder'-Alternative: Es gilt:  $d \neq 0$ , da  $\hat{s} \in K$  und  $g \notin K$ .

$$d^T g = d^T (\hat{s} - d) = \underbrace{(\hat{s} - g)^T \hat{s}}_{=0} - d^T d = -\|d\|_2^2 < 0.$$

Zeige nun  $d^T s \geq 0$  für alle  $s \in K$ . Sei  $s \in K$  beliebig. Da  $K$  konvex und  $\hat{s} \in K$  gilt

$$\|\hat{s} + \theta(s - \hat{s}) - g\|_2^2 \geq \|\hat{s} - g\|_2^2 = (\hat{s} - g)^T (\hat{s} - g) \quad \text{für } 0 \leq \theta \leq 1.$$

Also

$$\begin{aligned} (\hat{s} - g)^T (\hat{s} - g) + 2\theta(s - \hat{s})^T (\hat{s} - g) + \theta^2 (s - \hat{s})^T (s - \hat{s}) &\geq (\hat{s} - g)^T (\hat{s} - g) \\ 2\theta(s - \hat{s})^T (\hat{s} - g) + \theta^2 (s - \hat{s})^T (s - \hat{s}) &\geq 0 \\ (s - \hat{s})^T (\hat{s} - g) + \frac{\theta}{2} (s - \hat{s})^T (s - \hat{s}) &\geq 0. \end{aligned} \quad (6)$$

Der Grenzübergang  $\theta \rightarrow 0$  in (6) liefert:

$$(s - \hat{s})^T (\hat{s} - g) = s^T (\hat{s} - g) - \underbrace{\hat{s}^T (\hat{s} - g)}_{=0} = s^T d \geq 0 \quad \forall s \in K.$$

Somit folgt

$$0 \leq d^T s = d^T (By + Cw) \quad \forall y \geq 0 \quad \forall w. \quad (7)$$

Fixiere  $y = 0$  in (7) und erhalte  $(C^T d)^T w \geq 0$  für alle  $w$ , d.h.  $C^T d = 0$ . Fixiere  $w = 0$  in (7) und erhalte  $(B^T d)^T y \geq 0$  für alle  $y \geq 0$  d.h.  $B^T d \geq 0$ .

□

Wir betrachten jetzt wieder unser Originalproblem:

$$(IV) = \begin{cases} \text{Minimiere} \\ f(x), x \in D \\ \text{unter den Nebenbedingungen} \\ g_i(x) = 0, i = 1, \dots, l \text{ und} \\ k_i(x) \geq 0, i = 1, \dots, q. \end{cases}$$

Zu (IV) heißt  $L \in C^2(D \times \mathbb{R}^l \times \mathbb{R}^q, \mathbb{R})$ ,  $L(x, \lambda, \mu) = f(x) - \sum_{i=1}^l \lambda_i g_i(x) - \sum_{i=1}^q \mu_i k_i(x)$  die zugehörige **Lagrange-Funktion** und  $\lambda_i, i = 1, \dots, l, \mu_i, i = 1, \dots, q$  die **Lagrange-Multiplikatoren**.

**Satz 2.14 (Notwendige Bedingungen 1-ter Ordnung)**

Es sei  $x^* \in D$  eine lokale Lösung des Optimierungsproblems. Ferner sei  $x^*$  regulär, d.h.

$$\nabla g_1(x^*), \dots, \nabla g_l(x^*), \nabla k_i(x^*), i \in \mathcal{A}_{x^*}$$

sind linear unabhängig. Dann existiert ein Vektor  $\lambda^* \in \mathbb{R}^l, \mu^* \in \mathbb{R}^q$  mit

$$\frac{\partial L}{\partial x}(x^*, \lambda^*, \mu^*) = \nabla f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla g_i(x^*) - \sum_{i=1}^q \mu_i^* \nabla k_i(x^*) = 0$$

$$g(x^*) = 0, \quad k(x^*) \geq 0, \quad \mu^* \geq 0, \quad \mu_i^* k_i(x^*) = 0, i = 1, \dots, q.$$

**“Karush-Kuhn-Tucker (KKT) - Bedingungen”**

**Beweis:** Wende das Lemma von Farkas an auf

$$g = \nabla f(x^*) \in \mathbb{R}^N \quad B = (\nabla k_i(x^*), i \in \mathcal{A}_{x^*}) \in \mathbb{R}^{N,m}, m = \#\mathcal{A}_{x^*}, C = (\nabla g_i(x^*), i = 1, \dots, l) \in \mathbb{R}^{N,l}$$

und den Kegel

$$K = \{B\hat{\mu} + C\hat{\lambda}, \hat{\mu} \geq 0\} = \left\{ \sum_{i \in \mathcal{A}_x} \hat{\mu}_i \nabla k_i(x^*) + \sum_{i=1}^l \hat{\lambda}_i \nabla g_i(x^*) \mid \hat{\mu} \geq 0 \right\}.$$

Entweder gilt

$$\nabla f(x^*) = \sum_{i=1}^l \hat{\lambda}_i \nabla g_i(x^*) + \sum_{i \in \mathcal{A}_{x^*}} \hat{\mu}_i \nabla k_i(x^*) \quad (8)$$

oder es existiert ein  $d$  mit  $d^T \nabla f(x^*) < 0$  und

$$\nabla g_i(x^*)^T d = 0, i = 1, \dots, l \quad \nabla k_i(x^*)^T d \geq 0, i \in \mathcal{A}_x \quad (9)$$

(Dies bedeutet  $d^T \nabla f(x^*) < 0$  und  $d \in F(x^*) = T_Z(x^*)$ ). Nach Lemma 2.12 gilt aber

$$\nabla f(x^*)^T d \geq 0 \quad \forall d \in T_Z(x^*)$$

d.h. es tritt (8) ein. Setze nun  $\lambda_i^* = \lambda_i$ ,  $i = 1, \dots, l$

$$\mu_i^* = \begin{cases} \mu_i & \text{falls } i \in \mathcal{A}_{x^*} \\ 0 & \text{für } i \in \{1, \dots, q\} \setminus \mathcal{A}_{x^*} \end{cases} \quad (10)$$

und finde

$$\nabla f(x^*) = \sum_{i=1}^l \lambda_i^* \nabla g_i(x^*) + \sum_{i=1}^q \mu_i^* \nabla k_i(x^*)$$

sowie  $\mu_i^* \geq 0$ ,  $i = 1, \dots, q$ . Die Bedingungen  $k(x^*) \geq 0$ ,  $g(x^*) = 0$  folgen direkt aus der Zulässigkeit von  $x^*$ . Aus der Definition (10) folgt direkt

$$\mu_i^* k_i(x^*) = 0$$

denn für  $i \in \mathcal{A}_{x^*}$  ist  $k_i(x^*) = 0$  und für  $i \in \{1, \dots, q\} \setminus \mathcal{A}_{x^*}$  ist  $\mu_i^* = 0$ .

□

### Beispiel 2.15 (Teil 1, notwendige Bedingungen)

a.) Betrachte

$$f(x_1, x_2, x_3) = -x_1 x_2 - x_2 x_3 - x_1 x_3 \stackrel{!}{=} \min$$

unter der Nebenbedingung

$$g(x_1, x_2, x_3) = x_1 + x_2 + x_3 - 3 = 0.$$

Mit

$$x = (x_1, x_2, x_3), \quad L(x, \lambda) = f(x) - \lambda g(x), \quad \nabla f(x) = \begin{pmatrix} -x_2 - x_3 \\ -x_1 - x_3 \\ -x_2 - x_1 \end{pmatrix}, \quad \nabla g(x) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

erhalten wir

$$\begin{aligned} \frac{\partial L}{\partial x}(x, \lambda) &= \begin{pmatrix} -x_2 - x_3 \\ -x_1 - x_3 \\ -x_2 - x_1 \end{pmatrix} - \lambda \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = 0 \\ & \quad x_1 + x_2 + x_3 - 3 = 0 \\ \Leftrightarrow & \begin{pmatrix} 0 & -1 & -1 & -1 \\ -1 & 0 & -1 & -1 \\ -1 & -1 & 0 & -1 \\ 1 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \lambda \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 3 \end{pmatrix}. \end{aligned}$$

Dieses System hat die Lösung  $x^* = (1, 1, 1)$ ,  $\lambda^* = -2$ .

b.) Betrachte

$$f(x_1, x_2, x_3) = x_2 - \frac{1}{2} x_1 \stackrel{!}{=} \min$$

unter den Nebenbedingungen

$$\begin{aligned} k_1(x_1, x_2, x_3) &= -x_1 - \exp(-x_1) - x_3^2 + x_2 \geq 0, \\ k_2(x_1, x_2, x_3) &= x_1 \geq 0. \end{aligned}$$

Wir finden

$$\nabla f(x) = \begin{pmatrix} -1/2 \\ 1 \\ 0 \end{pmatrix}, \quad \nabla k_1(x) = \begin{pmatrix} -1 + \exp(-x_1) \\ 1 \\ -2x_3 \end{pmatrix}, \quad \nabla k_2(x) = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

und somit die KKT-Gleichungen

$$-\frac{1}{2} - \mu_1(-1 + \exp(-x_1)) - \mu_2 = 0 \quad (11)$$

$$1 - \mu_1 = 0 \quad (12)$$

$$2\mu_1 x_3 = 0 \quad (13)$$

sowie

$$\mu_1(-x_1 - \exp(-x_1) - x_3^2 + x_2) = 0, \quad \mu_1 \geq 0 \quad (14)$$

$$\mu_2 x_1 = 0, \quad \mu_2 \geq 0 \quad (15)$$

mit  $x = (x_1, x_2, x_3) \in Z$ .

(12) liefert  $\mu_1 = 1$  und mit (13) folgt  $x_3 = 0$ . Ferner liefert (14)

$$-x_1 - \exp(-x_1) + x_2 = 0 \quad \text{d.h.} \quad x_2 = x_1 + \exp(-x_1). \quad (16)$$

Annahme: Sei  $x_1 = 0$ , d.h. die Ungleichung  $k_2$  ist aktiv. Dann folgt mit (16)  $x_2 = 1$  und mit (11)

$$-\frac{1}{2} - 1 \cdot (-1 + \exp(0)) - \mu_2 = 0 \quad \text{d.h.} \quad \mu_2 = -\frac{1}{2}.$$

Dies ist ein Widerspruch! Also folgt  $x_1 \neq 0$ , d.h.  $k_2(x) = x_1 > 0$  und mit (15)  $\mu_2 = 0$ . Wir erhalten dann

$$-\frac{1}{2} - 1(-1 + \exp(-x_1)) - 0 = 0 \quad \text{d.h.} \quad x_1 = \ln(2)$$

und mit (16)

$$x_2 = \ln(2) + \exp(-\ln(2)) = \frac{1}{2} + \ln(2).$$

Lösung:  $x^* = (\ln(2), \frac{1}{2} + \ln(2), 0)$ ,  $\mu^* = (1, 0)$ .

Wir analysieren jetzt die Rolle der Ableitung 2-ter Ordnung. Es ist wieder

$$F(x) = \{d \in \mathbb{R}^N \mid \nabla g_i(x)^T d = 0, i = 1, \dots, l \text{ und } \nabla k_i(x)^T d \geq 0, i \in \mathcal{A}_x\}$$

der Kegel der linearisierten Richtungen in  $x$ . Sei  $x^*$  ein KKT-Punkt und seien  $\lambda^* \in \mathbb{R}^l$  und  $\mu^* \in \mathbb{R}^q$ ,  $\mu^* \geq 0$  die zu  $x^*$  gehörigen Lagrange-Multiplikatoren. Dann heißt

$$C(x^*, \mu^*) = \{d \in F(x^*) \mid \nabla k_i(x^*)^T d = 0 \text{ für } i \in \mathcal{A}_{x^*} \text{ mit } \mu_i^* > 0\}$$

der **kritische Kegel**, d.h.

$$\begin{aligned} d \in C(x^*, \mu^*) &\iff \nabla g_i(x^*)^T d = 0, \quad i = 1, \dots, l \\ &\quad \nabla k_i(x^*)^T d = 0, \quad i \in \mathcal{A}_{x^*}, \mu_i^* > 0 \\ &\quad \nabla k_i(x^*)^T d > 0, \quad i \in \mathcal{A}_{x^*}, \mu_i^* = 0. \end{aligned}$$

Aus den notwendigen Bedingungen 1-ter Ordnung folgt nun sofort

$$w \in C(x^*, \mu^*) \implies w^T \nabla f(x^*) = \sum_{i=1}^l \lambda_i^* w^T \nabla g_i(x^*) + \sum_{i=1}^q \mu_i^* w^T \nabla k_i(x^*) = 0.$$

**Satz 2.16 (Notwendige Bedingungen 2-ter Ordnung)**

Sei  $x^*$  eine lokale Lösung des Optimierungsproblems und sei  $x^*$  regulär. Ferner seien  $\lambda^* \in \mathbb{R}^l$  und  $\mu^* \in \mathbb{R}^q$  die zugehörigen Lagrange-Multiplikatoren für welche die KKT-Bedingungen erfüllt seien. Dann gilt

$$w^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*) w = w^T \left( \nabla^2 f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(x^*) - \sum_{i=1}^q \mu_i^* \nabla^2 k_i(x^*) \right) w \geq 0$$

für alle  $w \in C(x^*, \mu^*)$ .

**Beweis:** Sei  $w \in C(x^*, \mu^*) \subset F(x^*) = T_Z(x^*)$ . Damit gibt es eine zulässige Folge  $(y_k)_{k \in \mathbb{N}} \in Z^{\mathbb{N}}$ ,  $\lim_{k \rightarrow \infty} y_k = x^*$ , eine Folge  $(t_k)_{k \in \mathbb{N}}, t_k > 0, \lim_{k \rightarrow \infty} t_k = 0$  mit

$$\lim_{k \rightarrow \infty} \underbrace{\frac{y_k - x^*}{t_k}}_{w_k} = w$$

d.h.  $y_k = x^* + t_k w_k, w_k \rightarrow w$  für  $k \rightarrow \infty$ . Durch die Konstruktion mittels

$$R(z, t) = \begin{pmatrix} H(z) - tDH(x^*)w \\ Z^T(z - x^* - tw) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad H(z) = \begin{pmatrix} g(z) \\ k_i(z), i \in \mathcal{A}_{x^*} \end{pmatrix}$$

(vgl. Beweis von Lemma 2.11) lässt sich für die Folge  $(y_k)_{k \in \mathbb{N}}$  überdies

$$g_i(y_k) = t_k \nabla g_i(x^*)^T w, \quad i = 1, \dots, l \tag{17}$$

$$k_i(y_k) = t_k \nabla k_i(x^*)^T w, \quad i \in \mathcal{A}_{x^*} \tag{18}$$

erreichen. (17)-(18) sichert dann

$$\begin{aligned} L(y_k, \lambda^*, \mu^*) &= f(y_k) - \sum_{i=1}^l \lambda_i^* g_i(y_k) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* k_i(y_k) \\ &= f(y_k) - \sum_{i=1}^l \lambda_i^* t_k \underbrace{\nabla g_i(x^*)^T w}_{=0 \text{ für } w \in C(x^*, \lambda^*, \mu^*)} - \underbrace{\sum_{i \in \mathcal{A}_{x^*}} \mu_i^* t_k \nabla k_i(x^*)^T w}_{=0 \text{ für } i \in \mathcal{A}_{x^*}, \mu_i^* > 0} \\ &= f(y_k). \end{aligned}$$



Nach Taylor gilt überdies für  $k \in \mathbb{N}$

$$\begin{aligned} f(y_k) - f(x^*) &= t_k w_k^T \nabla f(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 f(\eta_0^k) t_k w_k, \quad \eta_0^k = x^* + \theta_0^k (y_k - x^*) \\ g_i(y_k) - g_i(x^*) &= t_k w_k^T \nabla g_i(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 g_i(\eta_i^k) t_k w_k, \quad \eta_i^k = x^* + \theta_i^k (y_k - x^*) \\ k_i(y_k) - k_i(x^*) &= t_k w_k^T \nabla k_i(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 k_i(\eta_{l+i}^k) t_k w_k, \quad \eta_{l+i}^k = x^* + \theta_{l+i}^k (y_k - x^*) \end{aligned}$$

mit  $\theta_0^k \in [0, 1]$ ,  $\theta_i^k \in [0, 1]$ ,  $i = 1, \dots, l$  und  $\theta_{l+i}^k \in [0, 1]$ ,  $i \in \mathcal{A}_{x^*}$ . Damit erhalten wir

$$\begin{aligned} f(y_k) &= f(y_k) - \sum_{i=1}^l \lambda_i^* g_i(y_k) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* k_i(y_k) \\ &= f(x^*) + t_k w_k^T \nabla f(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 f(\eta_0^k) t_k w_k \\ &\quad - \sum_{i=1}^l \lambda_i^* \left( \underbrace{g_i(x^*)}_{=0} + t_k w_k^T \nabla g_i(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 g_i(\eta_i^k) t_k w_k \right) \\ &\quad - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \left( \underbrace{k_i(x^*)}_{=0} + t_k w_k^T \nabla k_i(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 k_i(\eta_{l+i}^k) t_k w_k \right) \\ &= f(x^*) + t_k w_k^T \underbrace{\left( \nabla f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla g_i(x^*) - \underbrace{\sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla k_i(x^*)}_{=\sum_{i=1}^q \mu_i^* \nabla k_i(x^*)} \right)}_{=0} \\ &\quad + \frac{1}{2} t_k w_k^T \left( \nabla^2 f(\eta_0^k) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(\eta_i^k) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla^2 k_i(\eta_{l+i}^k) \right) t_k w_k \quad k \in \mathbb{N}. \end{aligned}$$

Somit folgt mit  $f(y_k) \geq f(x^*)$  sofort

$$0 \leq \frac{f(y_k) - f(x^*)}{2t_k^2} = w_k^T \left( \nabla^2 f(\eta_0^k) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(\eta_i^k) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla^2 k_i(\eta_{l+i}^k) \right) w_k$$

und Grenzübergang  $k \rightarrow \infty$  liefert

$$\begin{aligned} 0 &\leq w^T \left( \nabla^2 f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(x^*) - \underbrace{\sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla^2 k_i(x^*)}_{=\sum_{i=1}^q \mu_i^* \nabla^2 k_i(x^*)} \right) w \\ &= w^T \left( \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*) \right) w, \quad w \in C(x^*, \mu^*). \end{aligned}$$

□

Es fehlt noch eine hinreichende Charakterisierung für lokale Lösungen  $x^*$  von nichtlinearen Optimierungsproblemen mit Nebenbedingungen.

**Satz 2.17 (Hinreichende Bedingungen 2-ter Ordnung)**

Vorgelegt sei das Optimierungsproblem  $f(x) \stackrel{!}{=} \min$  unter den Nebenbedingungen  $g(x) = 0$  und  $k(x) \geq 0$ . Es sei  $x^* \in Z$  regulär und es existieren  $\lambda^* \in \mathbb{R}^l, \mu^* \in \mathbb{R}^q, \mu^* \geq 0$  mit

$$\frac{\partial L}{\partial x}(x^*, \lambda^*, \mu^*) = \nabla f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla g_i(x^*) - \sum_{i=1}^q \mu_i^* \nabla k_i(x^*) = 0, \quad \mu_i^* k_i(x^*) = 0, \quad i = 1, \dots, q.$$

Gilt dann

$$w^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*) w > 0 \quad \forall w \in C(x^*, \mu^*) \setminus \{0\}$$

so ist  $x^*$  eine strikte lokale Lösung des Optimierungsproblems.

**Beweis:** Indirekt: Angenommen  $x^*$  sei keine lokale Lösung des Optimierungsproblems. Dann gibt es eine Folge  $(y_k)_{k \in \mathbb{N}} \in Z^{\mathbb{N}}$  mit  $\lim_{k \rightarrow \infty} y_k = x^*$  und  $f(y_k) \leq f(x^*)$ . Schreibe  $y_k = x^* + t_k w_k$  mit  $w_k \in \mathbb{R}^N, \|w_k\|_2 = 1, t_k > 0, \lim_{k \rightarrow \infty} t_k = 0$ . Ferner gelte OE  $\lim_{k \rightarrow \infty} w_k = w$ , da  $\{z \in \mathbb{R}^N \mid \|z\|_2 = 1\}$  kompakt. Wir zeigen nun  $w \in C(x^*, \mu^*)$ . Nach Taylor gilt

$$\begin{aligned} 0 &= g_i(y_k) - g_i(x^*) = t_k w_k^T \nabla g_i(\eta_i^k), \quad \eta_i^k = x^* + \theta_i^k (y_k - x^*), \theta_i^k \in [0, 1], i = 1, \dots, l, k \in \mathbb{N} \\ &\Rightarrow 0 = w_k^T \nabla g_i(\eta_i^k) \xrightarrow{k \rightarrow \infty} 0 = w^T \nabla g_i(x^*), \quad i = 1, \dots, l. \end{aligned}$$

Für  $i \in \mathcal{A}_{x^*}$  betrachten wir

$$0 \leq \underbrace{k_i(y_k)}_{\geq 0} - \underbrace{k_i(x^*)}_{=0} = t_k w_k^T \nabla k_i(\eta_{l+i}^k), \quad \eta_{l+i}^k = x^* + \theta_{l+i}^k (y_k - x^*), \theta_{l+i}^k \in [0, 1], i \in \mathcal{A}_{x^*}, k \in \mathbb{N}.$$

Analog folgt durch Grenzübergang

$$0 \leq w^T \nabla k_i(x^*), \quad i \in \mathcal{A}_{x^*}.$$

Ferner erhalten wir aus

$$0 \geq f(y_k) - f(x^*) = t_k w_k^T \nabla f(\eta_0^k), \quad \eta_0^k = x^* + \theta_0^k (y_k - x^*), \quad \theta_0^k \in [0, 1], \quad k \in \mathbb{N}$$

sofort

$$0 \geq w^T \nabla f(x^*).$$

Wegen der Multiplikator-Regel gilt nun

$$0 \geq w^T \nabla f(x^*) = \sum_{i=1}^l \lambda_i^* \underbrace{w^T \nabla g_i(x^*)}_{=0} + \sum_{i=1}^q \mu_i^* w^T \nabla k_i(x^*) = \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \underbrace{w^T \nabla k_i(x^*)}_{\geq 0}.$$

Ist nun  $\mu_i^* > 0$ , so muss  $w^T \nabla k_i(x^*) = 0$  sein, d.h.

$$w \in \{d \in \mathbb{R}^N \mid \nabla g_i(x^*)^T d = 0, i = 1, \dots, l \quad \nabla k_i(x^*)^T d \geq 0, i \in \mathcal{A}_{x^*}, \mu_i^* = 0, \\ \nabla k_i(x^*)^T d = 0, i \in \mathcal{A}_{x^*}, \mu_i^* > 0\} = C(x^*, \mu^*).$$

Wir folgen nun dem Beweis von Satz (2.16) und entwickeln nach Taylor für  $k \in \mathbb{N}$

$$\begin{aligned} 0 &\geq f(y_k) - f(x^*) = t_k w_k^T \nabla f(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 f(\tilde{\eta}_0^k) t_k w_k, & \tilde{\eta}_0^k &= x^* + \tilde{\theta}_0^k (y_k - x^*) \\ 0 &= g_i(y_k) - g_i(x^*) = t_k w_k^T \nabla g_i(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 g_i(\tilde{\eta}_i^k) t_k w_k, & \tilde{\eta}_i^k &= x^* + \tilde{\theta}_i^k (y_k - x^*) \\ 0 &\leq k_i(y_k) - k_i(x^*) = t_k w_k^T \nabla k_i(x^*) + \frac{1}{2} t_k w_k^T \nabla^2 k_i(\tilde{\eta}_{l+i}^k) t_k w_k, & \tilde{\eta}_{l+i}^k &= x^* + \tilde{\theta}_{l+i}^k (y_k - x^*) \end{aligned}$$

mit  $\tilde{\theta}_0^k \in [0, 1]$ ,  $\tilde{\theta}_i^k \in [0, 1]$ ,  $i = 1, \dots, l$  und  $\tilde{\theta}_{l+i}^k \in [0, 1]$ ,  $i \in \mathcal{A}_{x^*}$ . Somit folgt

$$\begin{aligned} 0 &\geq f(y_k) - f(x^*) - \sum_{i=1}^l \lambda_i^* (g_i(y_k) - g_i(x^*)) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* (k_i(y_k) - k_i(x^*)) \\ &= t_k w_k^T \left( \underbrace{\nabla f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla g_i(x^*) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla k_i(x^*)}_{=0 \text{ nach Multiplikatorregel}} \right) \\ &\quad + \frac{1}{2} t_k w_k^T \left( \nabla^2 f(\tilde{\eta}_0^k) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(\tilde{\eta}_i^k) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla^2 k_i(\tilde{\eta}_{l+i}^k) \right) t_k w_k. \end{aligned}$$

Division durch  $2t_k^2$  liefert

$$\begin{aligned} w_k^T \left( \nabla^2 f(\tilde{\eta}_0^k) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(\tilde{\eta}_i^k) - \sum_{i \in \mathcal{A}_{x^*}} \mu_i^* \nabla^2 k_i(\tilde{\eta}_{l+i}^k) \right) w_k &\leq 0 \\ \downarrow k \rightarrow \infty \\ w^T \left( \underbrace{\nabla^2 f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla^2 g_i(x^*) - \sum_{i=1}^q \mu_i^* \nabla^2 k_i(x^*)}_{= \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*)} \right) w &\leq 0 \end{aligned}$$

für ein  $w \in C(x^*, \mu^*)$  mit  $\|w\|_2 = 1$ . Dies ist ein Widerspruch zur positiven Definitheit von  $\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*)$  auf  $C(x^*, \mu^*)$ .

□

Die Struktur des Kegels  $C(x^*, \mu^*)$  wird entscheidend von den aktiven Ungleichheitsrestriktionen beeinflusst. Ein wichtiger Spezialfall ist die strikte Komplementarität.

### Definition 2.18

Sei  $x^*$  eine lokale Lösung des Optimierungsproblem und seien  $\lambda^* \in \mathbb{R}^l$  und  $\mu^* \in \mathbb{R}^q$  die zugehörigen Lagrange-Multiplikatoren. Ferner bezeichne  $\mathcal{A}_{x^*}$  die Menge der aktiven Ungleichheitsrestriktionen. Dann liegt **strikte Komplementarität** vor, falls  $\mu_i^* > 0$  für alle  $i \in \mathcal{A}_{x^*}$ .

### Bemerkung 2.19

Bei strikter Komplementarität erhalten wir

$$C(x^*, \mu^*) = \{d \in \mathbb{R}^N \mid \nabla g_i(x^*)^T d = 0, i = 1, \dots, l, \nabla k_i(x^*)^T d = 0, i \in \mathcal{A}_{x^*}\} = N(DH(x^*))$$

mit  $H(z) = \begin{pmatrix} g(z) \\ k_i(z), i \in \mathcal{A}_{x^*} \end{pmatrix} \in \mathbb{R}^{l+m}$ ,  $m = \#\mathcal{A}_{x^*}$ . Die hinreichende Bedingung

$$w^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*) w > 0, \quad \forall w \in C(x^*, \mu^*) \setminus \{0\}$$

lautet also, dass  $\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*)$  positiv definit auf dem Nullraum der aktiven Ungleichheits- und den Gleichheitsrestriktionen ist. Ist nun  $Z \in \mathbb{R}^{N-l-m, N}$  mit  $R(Z) = N(DH(x^*))$ , d.h. die Spalten von  $Z$  bilden eine Basis von  $N(DH(x^*))$  so lautet die hinreichende Bedingung von Satz (2.17)  $Z^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*, \mu^*) Z \in \mathbb{R}^{N-l-m, N-l-m}$  ist positiv definit.

### Beispiel 2.20 (Aktive und inaktive Ungleichungen)

$$\begin{aligned} f(x_1, x_2) &= (x_1 - 1)^2 + (x_2 - 2)^2 \stackrel{!}{=} \min \\ k_1(x_1, x_2) &= x_2 - c \geq 0, \quad c \in \mathbb{R}. \end{aligned}$$

Man findet

$$\nabla f(x_1, x_2) = \begin{pmatrix} 2(x_1 - 1) \\ 2(x_2 - 2) \end{pmatrix}, \quad \nabla k_1(x_1, x_2) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

und somit die KKT-Bedingungen

$$\begin{aligned} 2(x_1 - 1) &= 0 \\ 2(x_2 - 2) - \mu_1 &= 0 \\ \mu_1(x_2 - c) &= 0 \end{aligned}$$

mit  $\mu_1 \geq 0, x_2 - c \geq 0$ . Dies liefert

$$x_1 = 1, \quad \mu_1 = 2(x_2 - 2), \quad 2(x_2 - 2)(x_2 - c) = 0 \quad \implies x_2 = 2 \text{ oder } x_2 = c.$$

Unterscheide jetzt 3 Fälle.

a.)  $c > 2$ :  $x_2 = 2$  nicht zulässig, da  $2 - c > 0 \Rightarrow x_2 = c$

$$x^* = (1, c), \quad \mu_1^* = 2(c - 2) > 0, \quad k_1(x^*) = c - c = 0$$

Ungleichung  $k_1$  ist stark aktiv.

b.)  $c = 2$ :  $\Rightarrow x_2 = 2$

$$x^* = (1, 2) \quad \mu_1^* = 2(2 - 2) = 0 \quad k_1(x^*) = 2 - 2 = 0$$

Ungleichung  $k_1$  ist schwach aktiv.

c.)  $c < 2$ :  $x_2 = c$  nicht zulässig, da  $\mu_1 = 2(c - 2) < 0$ .

$$\Rightarrow x_2 = 2, \quad x^* = (1, 2), \quad \mu_1^* = 0, \quad k_1(x^*) = 2 - c > 0.$$

Die Ungleichung  $k_1$  ist jetzt nicht aktiv.

### Beispiel 2.21 (Teil 2, hinreichende Bedingungen)

a.) Betrachte wieder

$$f(x_1, x_2, x_3) = -x_1 x_2 - x_2 x_3 - x_1 x_3 \stackrel{!}{=} \min \tag{19}$$

unter der Nebenbedingung

$$g(x_1, x_2, x_3) = x_1 + x_2 + x_3 - 3 = 0. \quad (20)$$

Lösung  $x^* = (1, 1, 1)$ ,  $\lambda^* = -2$  (stationärer Punkt der Lagrange- Funktion).

Hinreichende Bedingung:

$$\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) = \nabla^2 f(x^*) - \lambda^* \nabla^2 g(x^*)$$

ist positiv definit auf  $N(Dg(x^*))$ . Mit

$$\nabla^2 f(x^*) = \begin{pmatrix} 0 & -1 & -1 \\ -1 & 0 & -1 \\ -1 & -1 & 0 \end{pmatrix} \quad \text{und} \quad \nabla^2 g(x^*) = 0$$

erhalten wir  $\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) = \nabla^2 f(x^*)$ . Sei nun

$$y \in N(Dg(x^*)) = N((1, 1, 1)) = \{z \in \mathbb{R}^3 \mid z_1 + z_2 + z_3 = 0\}, \quad y \neq 0.$$

Dann gilt

$$\begin{aligned} y^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) y &= y^T \begin{pmatrix} 0 & -1 & -1 \\ -1 & 0 & -1 \\ -1 & -1 & 0 \end{pmatrix} y = (y_1, y_2, y_3) \begin{pmatrix} -y_2 - y_3 \\ -y_1 - y_3 \\ -y_1 - y_2 \end{pmatrix} \\ &= y_1 \underbrace{(-y_2 - y_3)}_{=y_1} + y_2 \underbrace{(-y_1 - y_3)}_{=y_2} + y_3 \underbrace{(-y_1 - y_2)}_{=y_3} = y_1^2 + y_2^2 + y_3^2 > 0. \end{aligned}$$

Also ist  $x^* = (1, 1, 1)$  ein lokales Minimum für die Aufgabe (19)-(20).

b.) Betrachte jetzt

$$f(x_1, x_2, x_3) = x_2 - \frac{1}{2}x_1 \stackrel{!}{=} \min \quad (21)$$

unter den Nebenbedingungen

$$k_1(x_1, x_2, x_3) = -x_1 - \exp(-x_1) - x_3^2 + x_2 \geq 0 \quad (22)$$

$$k_2(x_1, x_2, x_3) = x_1 \geq 0. \quad (24)$$

Einzigste Lösung der KKT-Gleichungen

$$x^* = \left( \ln(2), \frac{1}{2} + \ln(2), 0 \right), \quad \mu^* = (1, 0).$$

Es liegt strikte Komplementarität vor, d.h.

$$k_1(x^*) = 0, \mu_1^* > 0 \quad k_2(x^*) > 0, \mu_2^* = 0$$

d.h. der kritische Kegel lautet

$$C(x^*, \mu^*) = \{d \in \mathbb{R}^2 \mid \nabla k_1(x^*)^T d = 0\}.$$

Finde nun

$$\nabla^2 f(x^*) = 0 \in \mathbb{R}^{3,3} \quad \text{und} \quad \nabla^2 k_1(x^*) = \begin{pmatrix} -1/2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & -2 \end{pmatrix}.$$

und somit

$$\frac{\partial^2 L}{\partial x^2}(x^*, \mu^*) = \nabla^2 f(x^*) - \mu_1^* \nabla^2 k_1(x^*) - \mu_2^* \nabla^2 k_2(x^*) = -\nabla^2 k_1(x^*) = \begin{pmatrix} 1/2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{pmatrix}.$$

Ferner gilt  $\nabla k_1(x^*) = \begin{pmatrix} -1/2 \\ 1 \\ 0 \end{pmatrix}$ . Hinreichend ist nun

$$y^T \frac{\partial^2 L}{\partial x^2}(x^*, \mu^*) y > 0 \quad \forall y \in N(\nabla k_1(x^*)^T) \setminus \{0\}.$$

Wir erhalten

$$y^T \begin{pmatrix} 1/2 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{pmatrix} y = \frac{1}{2} y_1^2 + y_3^2 \geq 0 \quad \text{für } y \in N((-1/2, 1, 0)) = \{z \in \mathbb{R}^3 \mid -1/2 z_1 + z_2 = 0\}.$$

Sei nun  $\frac{1}{2} y_1^2 + y_3^2 = 0 \implies y_1 = y_3 = 0$ . Nach Definition von  $N((-1/2, 1, 0))$  folgt damit auch  $y_2 = 1/2 y_1 = 0$  und somit  $y \equiv 0$ . Also ist  $x^* = (\ln(2), 1/2 + \ln(2), 0)$  ein lokales Minimum von (21)-(24).

### Bemerkung 2.22 (Spezialfälle der Kuhn-Tucker Theorie)

a.) **Freie Minimierung:** Dann erhalten wir gemäß Satz 2.14 und 2.16 mit  $L(x) = f(x)$  die notwendigen Bedingungen

$$\nabla f(x^*) = 0, \quad w^T \nabla^2 f(x^*) w \geq 0 \quad \forall w \in \mathbb{R}^N$$

sowie die hinreichenden Bedingungen

$$\nabla f(x^*) = 0, \quad w^T \nabla^2 f(x^*) w > 0 \quad \forall w \in \mathbb{R}^N \setminus \{0\}$$

(vgl. Sätze 2.1 und 2.3).

b.) **Minimierung mit Gleichheitsrestriktionen:** D.h.  $l > 0, q = 0$ . Das Problem lautet

$$f(x) \stackrel{!}{=} \min$$

unter der Nebenbedingung

(25)

$$g(x) = 0.$$

Die Lagrange-Funktion ist

$$L(x, \lambda) = f(x) - \sum_{i=1}^l \lambda_i g_i(x).$$

Anwendung der Sätze 2.14, 2.16 und 2.17 liefert für eine reguläre Lösung  $x^*$  (d.h.  $\text{rg}(Dg(x^*)) = l$ ) die notwendigen Bedingungen

$$\frac{\partial L}{\partial x}(x^*, \lambda^*) = \nabla f(x^*) - \sum_{i=1}^l \lambda_i^* \nabla g_i(x^*) = 0, \quad g(x^*) = 0$$

$$w^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) w \geq 0, \quad \forall w \in C(x^*) = N(Dg(x^*))$$

bzw. die hinreichenden Bedingungen

$$\frac{\partial L}{\partial x}(x^*, \lambda^*) = 0, \quad g(x^*) = 0, \quad \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) \text{ positiv definit auf } C(x^*).$$

# Kapitel II: Numerik freier Optimierungsprobleme

## 1 Grundprinzipien der Abstiegsverfahren

Vorgelegt sei  $f \in C^1(\mathbb{R}^N, \mathbb{R})$ . Gesucht ist ein lokales Minimum von  $f$ .

Idee der Abstiegsverfahren:

- Ist man in einem Punkt  $x \in \mathbb{R}^N$ , so wähle eine Richtung  $d \in \mathbb{R}^N$  aus, in welcher der Funktionswert von  $f$  fällt. (*Abstiegsrichtung*)
- Entlang dieser Richtung  $d$  geht man solange, bis man den Funktionswert von  $f$  hinreichend verkleinert hat. (*Schrittweitensteuerung*)

### Definition 1.1

Sei  $f : \mathbb{R}^N \rightarrow \mathbb{R}$ , und sei  $x \in \mathbb{R}^N$  fest.  $d \in \mathbb{R}^N$  heißt **Abstiegsrichtung** von  $f$  in  $x$ , wenn es ein  $\bar{t} > 0$  gibt mit

$$f(x + td) < f(x) \quad \text{für } t \in ]0, \bar{t}].$$

### Lemma 1.2

Ist  $f \in C^1(D, \mathbb{R})$ ,  $D \subset \mathbb{R}^N$  offen, so ist

$$\nabla f(x)^T d < 0 \tag{26}$$

hinreichend dafür, dass  $d \in \mathbb{R}^N$  eine Abstiegsrichtung ist.

Setze  $\varphi(t) := f(x + td)$ ,  $|t| < \varepsilon$ . Aus  $f \in C^1(D, \mathbb{R})$  folgt  $\varphi \in C^1(]-\varepsilon, \varepsilon[, \mathbb{R})$  und

$$\varphi(t) = \varphi(0) + t\varphi'(0) + r(t)$$

mit  $\frac{r(t)}{t} \rightarrow 0$  für  $t \rightarrow 0$ . Es gilt:

$$\varphi(0) = f(x), \quad \varphi'(t) = \nabla f(x + td)^T d, \quad \varphi'(0) = \nabla f(x)^T d.$$

Somit folgt

$$\frac{\varphi(t) - \varphi(0)}{t} = \nabla f(x)^T d + \frac{r(t)}{t}$$

und mit  $\frac{r(t)}{t} = o(1)$  erhalten wir

$$\frac{\varphi(t) - \varphi(0)}{t} < 0 \quad \text{für } t \in ]0, \bar{t}], \bar{t} < \varepsilon \text{ geeignet.}$$

Dies ist äquivalent zu

$$f(x + td) - f(x) < 0, \quad t \in ]0, \bar{t}].$$

### Bemerkung 1.3

- $\nabla f(x)^T d < 0$  bedeutet, dass der Winkel zwischen  $d$  und  $-\nabla f(x)$  kleiner als  $90^\circ$  ist.
- Mögliche Kandidaten für  $d$  sind

a.)  $d = -\nabla f(x), \quad \nabla f(x)^T d = -\|\nabla f(x)\|_2^2 < 0$

b.)  $d = -M\nabla f(x), \quad M \in \mathbb{R}^{N,N}$  symmetrisch und positiv definit

$$\nabla f(x)^T d = -\nabla f(x)^T M \nabla f(x) < 0$$



## Struktur eines Algorithmus:

**Input:**  $f : \mathbb{R}^N \rightarrow \mathbb{R}$ ,  $x_0 \in \mathbb{R}^N$

$k = 0$

**while** *Konvergenzkriterium nicht erfüllt* **do**

- bestimme Abstiegsrichtung  $d^k$  von  $f$  in  $x^k$ ;
- bestimme Schrittweite  $t_k > 0$  mit

$$f(x^k + t_k d^k) < f(x^k);$$

- setze  $x^{k+1} = x^k + t_k d^k$ ;
- $k = k + 1$ ;

**end**

Algorithmus: ABSTIEG

### 1.1 Konvergenzuntersuchungen

In theoretischen Konvergenzuntersuchungen betrachten wir kein Konvergenzkriterium, d.h. wir nehmen an, dass eine Folge  $(x_k)_{k \in \mathbb{N}}$  erzeugt wird. Wichtig ist nun, dass mögliche Häufungspunkte stationäre Punkte von  $f$  sind.

#### Lemma 1.4

Sei  $f \in C^1(\mathbb{R}^N, \mathbb{R})$ , und sei  $(x^k)_{k \in \mathbb{N}}$  eine durch den Algorithmus ABSTIEG erzeugte Folge. Ferner gelte:

a.) Es existiere ein  $\theta_1 > 0$  mit

$$-\nabla f(x^k)^T d^k \geq \theta_1 \|\nabla f(x^k)\|_2 \cdot \|d^k\|_2 \quad (\geq 0). \quad (\text{Winkelbedingung})$$

b.) Es existiere ein  $\theta_2 > 0$  mit

$$f(x^k + t_k d^k) \leq f(x^k) - \theta_2 \left( \frac{\nabla f(x^k)^T d^k}{\|d^k\|_2} \right)^2.$$

Dann ist jeder Häufungspunkt der Folge  $(x^k)_{k \in \mathbb{N}}$  ein stationärer Punkt von  $f$ .

**Beweis:** Mit den beiden Bedingungen a.) und b.) folgt

$$f(x^{k+1}) - f(x^k) = f(x^k + t_k d^k) - f(x^k) \leq -\theta_2 \left( \frac{\nabla f(x^k)^T d^k}{\|d^k\|_2} \right)^2 \leq -\theta_1^2 \theta_2 \|\nabla f(x^k)\|_2^2 \leq 0. \quad (27)$$

Sei nun  $x^*$  ein Häufungspunkt von  $f$ , d.h. es existiert eine Teilfolge  $(x^{k_n})_{n \in \mathbb{N}}$  mit  $\lim_{n \rightarrow \infty} x^{k_n} = x^*$ . Dies impliziert  $\lim_{n \rightarrow \infty} f(x^{k_n}) = f(x^*)$ . Ferner folgt wegen der Monotonie von  $(f(x^{k_n}))_{n \in \mathbb{N}}$  (vgl. (27)) sofort

$$\lim_{k \rightarrow \infty} f(x^k) = f(x^*).$$

Einsetzen in (27) liefert mit Grenzübergang  $k \rightarrow \infty$

$$0 = f(x^*) - f(x^*) \leq -\theta_1^2 \theta_2 \|\nabla f(x^*)\|_2^2 \leq 0.$$

Damit ist jeder Häufungspunkt stationärer Punkt, d.h.

$$0 = \|\nabla f(x^*)\|_2 = \lim_{n \rightarrow \infty} \|\nabla f(x^{k_n})\|_2.$$

□

### Bemerkung 1.5

- Forderung a.) in Lemma 1.4 bedeutet

$$\cos \left( \angle(d^k, -\nabla f(x^k)) \right) = \frac{-\nabla f(x^k)^T d^k}{\|d^k\|_2 \cdot \|\nabla f(x^k)\|_2} \geq \theta_1 > 0$$

gleichmäßig in  $k$ .

- Schrittweiten  $t_k > 0$ ,  $k \in \mathbb{N}$  mit b.) heißen **effizient**.

## 2 Schrittweitenstrategien

Das allgemeine Abstiegsverfahren

$$x^{k+1} = x^k + t_k d^k, \quad k = 0, 1, 2, \dots$$

besitzt in der Wahl der Abstiegsrichtung  $d^k$  und der Schrittweite  $t_k > 0$  große Freiheitsgrade. Ist

$$L(x^0) = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^0)\} \text{ kompakt,}$$

so ist die Regel  $t_k = t_k^{min}$  mit

$$f(x^k + t_k d^k) = \min_{t \geq 0} f(x^k + t d^k)$$

wohldefiniert und naheliegend. Allerdings ist diese Regel aus Aufwandsgründen meist nicht praktikabel.

### 2.1 Schrittweiten Regeln

Armijo-Regel:

Betrachte gradientenähnliche Richtungen

$$d := -M \nabla f(x), \quad M \in \mathbb{R}^{N,N} \text{ symmetrisch und positiv definit.}$$

Sei  $\alpha \in ]0, 1[$  beliebig, aber fest. Wähle  $t > 0$  mit

$$\varphi(t) := f(x + td) \leq \underbrace{f(x)}_{=\varphi(0)} + \alpha t \underbrace{\nabla f(x)^T d}_{=\varphi'(0)} =: l(t) \tag{28}$$

$l(t)$  ist linear in  $t$ .

...Bild...

Zur tatsächlichen Berechnung von  $t$  überprüft man (28) sequentiell z.B. für

$$t^{(l)} = s \beta^l, \quad l = 0, 1, 2, \dots \quad \text{mit } \beta \in ]0, 1[ \text{ fest, } s > 0 \text{ fest.}$$

Bei erstmaliger Gültigkeit von (28) bricht man ab. Etwas strenger als die Armijo-Regel ist die

Wolfe-Powell Regel:

Sei  $\alpha \in ]0, 1/2[$  und  $\rho \in ]\alpha, 1[$  gegeben. Zu  $x, d \in \mathbb{R}^N$  mit  $\nabla f(x)^T d < 0$  bestimme man eine Schrittweite  $t > 0$  mit

$$f(x + td) \leq f(x) + \alpha t \nabla f(x)^T d \quad (\varphi(t) \leq \varphi(0) + \alpha t \varphi'(0)), \quad (29)$$

$$\nabla f(x + td)^T d \geq \rho \nabla f(x)^T d \quad (\varphi'(t) \geq \rho \varphi'(0) \text{ mit } \varphi(t) = f(x + td)). \quad (30)$$

**Bemerkung 2.1**

Es lässt sich zeigen, dass die Wolfe-Powell Regel wohldefiniert ist für  $f \in C^1$  und  $f$  nach unten beschränkt.

**Satz 2.2 (Theorem von Zoutendijk)**

Vorgelegt sei ein Abstiegsverfahren der Form

$$x^{k+1} = x^k + t_k d^k, \quad k = 0, 1, 2, \dots$$

wobei  $t_k > 0$  die Wolfe Powell Bedingungen (29)-(30) erfülle. Es sei  $f \in C^1(\mathbb{R}^N, \mathbb{R})$  nach unten beschränkt und

$$D \supset L(x^0) = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^0)\}, \quad D \text{ offen}.$$

Außerdem gelte

$$\|\nabla f(x) - \nabla f(\tilde{x})\|_2 \leq L \|x - \tilde{x}\|_2 \quad \forall x, \tilde{x} \in D.$$

Dann gilt

a.) Es existiert ein  $c > 0$  mit

$$f(x^k + t_k d^k) - f(x^k) \leq -c \left( \frac{\nabla f(x^k)^T d^k}{\|d^k\|_2} \right)^2, \quad k \in \mathbb{N}.$$

b.) Mit  $\theta_k := \angle(d^k, -\nabla f(x^k))$  folgt

$$\sum_{k=0}^{\infty} \cos^2(\theta_k) \cdot \|\nabla f(x^k)\|_2^2 < \infty.$$

**Beweis:**

a.) Mit (30) folgt

$$\begin{aligned} (\nabla f(x^{k+1}) - \nabla f(x^k))^T d^k &= (\nabla f(x^k + t_k d^k) - \nabla f(x^k))^T d^k \\ &\geq \rho \nabla f(x^k)^T d^k - \nabla f(x^k)^T d^k = (\rho - 1) \nabla f(x^k)^T d^k. \end{aligned}$$

Andererseits gilt mit der Lipschitzbedingung für  $\nabla f$ :

$$(\nabla f(x^{k+1}) - \nabla f(x^k))^T d^k \leq \|\nabla f(x^{k+1}) - \nabla f(x^k)\|_2 \cdot \|d^k\|_2 \leq L \|t_k d^k\|_2 \cdot \|d^k\|_2 = t_k L \|d^k\|_2^2.$$

Kombiniert man beide Formeln, so erhält man

$$t_k \geq \frac{(\nabla f(x^{k+1}) - \nabla f(x^k))^T d^k}{L \cdot \|d^k\|_2^2} \geq \frac{(\rho - 1) \nabla f(x^k)^T d^k}{L \cdot \|d^k\|_2^2}.$$

Setze dies in (29) ein und finde

$$f(x^{k+1}) \leq f(x^k) + \alpha t_k \underbrace{\nabla f(x^k)^T d^k}_{<0} \leq f(x^k) + \underbrace{\frac{\alpha(\rho - 1)}{L}}_{=:-c, c>0} \cdot \frac{(\nabla f(x^k)^T d^k)^2}{\|d^k\|_2^2}, \quad k \in \mathbb{N}.$$

b.) Über  $a.)$  hinaus folgt

$$f(x^{k+1}) \leq f(x^k) + \frac{\alpha(\rho-1)}{L} \cos^2(\theta_k) \cdot \|\nabla f(x^k)\|_2^2 \quad \text{mit } \theta^k = \angle(-\nabla f(x^k), d^k). \quad (31)$$

Induktiv folgt mit (31) und  $c = \frac{\alpha(\rho-1)}{L} > 0$

$$f(x^{k+1}) - f(x^0) = \sum_{j=0}^k f(x^{j+1}) - f(x^j) \leq c \sum_{j=0}^k \cos^2(\theta_j) \|\nabla f(x^j)\|_2^2, \quad k \in \mathbb{N}.$$

Da  $f$  nach unten beschränkt ist, erhalten wir sofort durch Grenzübergang

$$\sum_{j=0}^{\infty} \cos^2(\theta_j) \cdot \|\nabla f(x^j)\|_2^2 < \infty.$$

□

### Bemerkung 2.3

- Die Voraussetzungen von Satz 2.2 sind nicht zu restriktiv. Ist  $f$  z.B. nicht nach unten beschränkt, so existiert kein globales Minimum und es muss auch kein lokales Minimum existieren.
- Das Zoutendijk-Theorem impliziert insbesondere

$$\lim_{k \rightarrow \infty} \cos^2(\theta_k) \cdot \|\nabla f(x^k)\|_2^2 = 0.$$

- Die Wolfe-Powell Schrittweiten sind nach Satz 2.2 a) effizient.

### Korollar 2.4

Seien die Voraussetzungen von Satz 2.2 erfüllt, und es gelte überdies

$$\cos(\theta_k) \geq \delta > 0, \quad k \in \mathbb{N}$$

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

Dies ist z.B. bei gradientenähnlichen Richtungen

$$d^k = -M_k \nabla f(x^k), \quad M_k \in \mathbb{R}^{N,N} \text{ symmetrisch und positiv definit,}$$

erfüllt, falls gilt

$$0 < \underline{\lambda}_{min} < \lambda_{min}^k < \lambda_{max}^k < \bar{\lambda}_{max}, \quad k \in \mathbb{N}.$$

Dabei bezeichnet  $\lambda_{min}^k, \lambda_{max}^k$  den kleinsten bzw. größten Eigenwert der Matrix  $M_k$ . Es folgt nämlich

$$\cos(\theta_k) = \frac{\nabla f(x^k)^T M_k \nabla f(x^k)}{\|\nabla f(x^k)\|_2 \|M_k \nabla f(x^k)\|_2} \geq \frac{\lambda_{min}^k \nabla f(x^k)^T \nabla f(x^k)}{\|M_k\|_2 \|\nabla f(x^k)\|_2^2} = \frac{\lambda_{min}^k}{\|M_k\|_2} = \frac{\lambda_{min}^k}{\lambda_{max}^k} \geq \frac{\underline{\lambda}_{min}}{\bar{\lambda}_{max}} \quad \text{für } k \in \mathbb{N}.$$

## 2.2 Abschwächung der Winkelbedingung

Es sei  $\theta_k = \angle(d_k, -\nabla f(x^k))$  bzw.

$$\cos(\theta_k) = \frac{-\nabla f(x^k)^T d^k}{\|\nabla f(x^k)\|_2 \|d^k\|_2}$$

Ist  $\theta_k < 90^\circ$ , so ist  $d^k$  eine Abstiegsrichtung. Nach Korollar 2.4 erhält man ein Konvergenzresultat, falls  $\cos(\theta_k) \geq \delta > 0$ ,  $k \in \mathbb{N}$ . Ist  $f$  gleichmäßig konvex, d.h. es gilt  $z^T \nabla^2 f(x) z \geq \mu z^T z$ ,  $\mu > 0$ , so lässt sich die Winkelbedingung abschwächen.

### Satz 2.5

Sei  $f \in C^2(\mathbb{R}^N, \mathbb{R})$ , und sei  $L(x^0) = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^0)\}$  konvex und  $f$  gleichmäßig konvex auf  $L(x^0)$  mit  $\mu$ , d.h.  $z^T \nabla^2 f(x) z \geq \mu z^T z$  für  $x \in L(x^0)$ ,  $z \in \mathbb{R}^N$ . Sei  $(x^k)_{k \in \mathbb{N}}$  erzeugt durch das Abstiegsverfahren d.h.  $x^{k+1} = x^k + t_k d^k$ . Ferner gelte:

$$a.) \sum_{k=0}^{\infty} \delta_k = \infty \text{ mit } \delta_k = \cos^2(\theta_k) = \left( \frac{\nabla f(x^k)^T d^k}{\|\nabla f(x^k)\|_2 \|d^k\|_2} \right)^2.$$

b.) Die Schrittweiten  $t_k > 0$ ,  $k \in \mathbb{N}$  sind effizient.

Dann konvergiert die Folge  $(x^k)_{k \in \mathbb{N}}$  gegen das eindeutig bestimmte globale Minimum von  $f$  in  $L(x^0)$ .

**Beweis:** Sei  $x^*$  dieses Minimum. Aus der Effizienz von  $t_k$  folgt

$$\begin{aligned} f(x^{k+1}) &= f(x^k + t_k d^k) \leq f(x^k) - c \left( \frac{\nabla f(x^k)^T d^k}{\|d^k\|_2} \right)^2 \\ &= f(x^k) - c \delta_k \|\nabla f(x^k)\|_2^2, \quad k = 0, 1, 2, \dots, \text{ mit } c > 0. \end{aligned} \quad (32)$$

Nun gilt trivialerweise

$$\left\| \sqrt{\frac{\mu}{2}}(x^* - x^k) + \sqrt{\frac{1}{2\mu}} \nabla f(x^k) \right\|_2^2 \geq 0.$$

Ausmultiplizieren ergibt

$$\begin{aligned} \frac{\mu}{2} \|x^* - x^k\|_2^2 + 2 \sqrt{\frac{\mu}{2}} \sqrt{\frac{1}{2\mu}} (x^* - x^k)^T \nabla f(x^k) + \frac{1}{2\mu} \|\nabla f(x^k)\|_2^2 &\geq 0 \\ \Rightarrow -\frac{1}{2\mu} \|\nabla f(x^k)\|_2^2 &\leq \frac{\mu}{2} \|x^* - x^k\|_2^2 + \nabla f(x^k)^T (x^* - x^k). \end{aligned}$$

Ferner liefert Taylor-Entwicklung

$$f(x^*) - f(x^k) = \nabla f(x^k)^T (x^* - x^k) + \frac{1}{2} (x^* - x^k)^T \nabla^2 f(\eta_k) (x^* - x^k)$$

mit  $\eta_k = x^k + \theta_k(x^* - x^k)$ ,  $\theta_k \in [0, 1]$ ,  $\eta_k \in L(x^0)$ . Die gleichmäßige Konvexität liefert dann

$$f(x^*) - f(x^k) \geq \nabla f(x^k)^T (x^* - x^k) + \frac{\mu}{2} \|x^* - x^k\|_2^2 \geq -\frac{1}{2\mu} \|\nabla f(x^k)\|_2^2 \quad (33)$$

$$\iff 2\mu c \delta_k (f(x^*) - f(x^k)) \geq -c \delta_k \|\nabla f(x^k)\|_2^2. \quad (34)$$

Mit (32) gilt

$$f(x^{k+1}) \leq f(x^k) - 2\mu c \delta_k (f(x^k) - f(x^*))$$

und damit

$$0 \leq f(x^{k+1}) - f(x^*) = f(x^{k+1}) - f(x^k) + f(x^k) - f(x^*) \leq (1 - 2\mu c \delta_k)(f(x^k) - f(x^*)).$$

Induktiv folgt hieraus

$$\begin{aligned} 0 \leq f(x^{k+1}) - f(x^*) &\leq \prod_{j=0}^k (1 - 2\mu c \delta_j)(f(x^0) - f(x^*)) \\ &\leq \prod_{j=0}^k \exp(-2\mu c \delta_j)(f(x^0) - f(x^*)) = \exp\left(-2c\mu \sum_{j=0}^k \delta_j\right) (f(x^0) - f(x^*)). \end{aligned}$$

Wegen  $\sum_{j=0}^k \delta_j \rightarrow \infty$  für  $k \rightarrow \infty$  und  $(f(x^k))_{k \in \mathbb{N}}$  monoton gilt

$$\lim_{k \rightarrow \infty} f(x^k) = f(x^*).$$

Ferner gilt mit der gleichmäßigen Konvexität

$$f(x^k) - f(x^*) = \underbrace{\nabla f(x^*)}_{=0}(x^k - x^*) + \frac{1}{2}(x^k - x^*) \nabla^2 f(\eta_k)(x^k - x^*) \geq \frac{\mu}{2} \|x^k - x^*\|_2^2, \text{ d.h.}$$

$$\|x^k - x^*\|_2 \leq \sqrt{\frac{2}{\mu}(f(x^k) - f(x^*))} \rightarrow 0 \text{ für } k \rightarrow \infty, \text{ d.h. } \lim_{k \rightarrow \infty} x^k = x^*.$$

□

### Lemma 2.6

Sei  $f \in C^2(\mathbb{R}^N, \mathbb{R})$ , und sei  $L(x^0) = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^0)\}$  konvex, und  $f$  gleichmäßig konvex auf  $L(x^0)$  mit  $\mu$ , d.h.  $z^T \nabla^2 f(x) z \geq \mu z^T z$  für  $x \in L(x^0)$ ,  $z \in \mathbb{R}^N$ . Dann ist  $L(x^0)$  kompakt und die Aufgabe

$$f(x) \stackrel{!}{=} \min$$

hat genau eine globale Lösung  $x^*$  und  $x^*$  liegt in  $L(x^0)$ .

**Beweis:** Zeige zunächst:  $L(x^0)$  ist kompakt. Aus der gleichmäßigen Konvexität von  $f$  folgt (vgl. Formelzeile (33), Beweis Satz 2.5)

$$f(x) - f(x^0) \geq \nabla f(x^0)^T (x - x^0) + \frac{\mu}{2} \|x - x^0\|_2^2.$$

Sei nun  $x \in L(x^0)$ , d.h.  $f(x) \leq f(x^0)$ . Dann gilt

$$0 \geq f(x) - f(x^0) \geq \nabla f(x^0)^T (x - x^0) + \frac{\mu}{2} \|x - x^0\|_2^2$$

$$\iff \frac{\mu}{2} \|x - x^0\|_2^2 \leq -\nabla f(x^0)^T (x - x^0) \leq \|\nabla f(x^0)\|_2 \cdot \|x - x^0\|_2$$

$$\|x - x^0\|_2 \leq \frac{2}{\mu} \|\nabla f(x^0)\|_2 =: c$$

d.h.  $\|x - x^0\|_2 \leq c$  und somit  $L(x^0) \subset K_c(x^0)$  beschränkt. Nach Definition ist  $L(x^0)$  abgeschlossen und somit kompakt. Die Aussagen von Lemma 2.6 folgen dann aus dem folgenden Lemma.

□

**Lemma 2.7**

Sei  $f \in C^2(\mathbb{R}^N, \mathbb{R})$ , und  $X \subset \mathbb{R}^N$  konvex. Betrachte das Problem

$$f(x) \stackrel{!}{=} \min$$

unter der Nebenbedingung  $x \in X$ . Dann gilt:

- a.) Ist  $f$  konvex auf  $X$ , so ist die Lösungsmenge von  $f(x) \stackrel{!}{=} \min, x \in X$  konvex.
- b.) Ist  $f$  strikt konvex (d.h.  $f(\lambda x^1 + (1 - \lambda)x^2) < \lambda f(x^1) + (1 - \lambda)f(x^2)$ ,  $0 < \lambda < 1$ ,  $x_1, x_2 \in X$ ) auf  $X$ , so existiert höchstens eine Lösung.
- c.) Ist  $f$  gleichmäßig konvex auf  $X$ ,  $X$  nicht leer und abgeschlossen, so gibt es genau eine Lösung  $x^* \in X$ .

**Beweis:**

- a.) Seien  $x^1, x^2$  zwei Lösungen von  $f(x) \stackrel{!}{=} \min, x \in X$ , also  $f(x^1) = f(x^2) = \min\{f(x) \mid x \in X\}$ . Für  $\lambda \in ]0, 1[$  ist dann auch  $\lambda x^1 + (1 - \lambda)x^2 \in X$ , da  $X$  konvex. Da  $f$  konvex ist, folgt

$$f(\lambda x^1 + (1 - \lambda)x^2) \leq \lambda f(x^1) + (1 - \lambda)f(x^2) = f(x^1) = \min\{f(x) \mid x \in X\}$$

d.h. auch  $\lambda x^1 + (1 - \lambda)x^2$  ist ein Minimum.

- b.) Angenommen, die Aufgabe

$$f(x) \stackrel{!}{=} \min \quad \text{unter der Nebenbedingung } x \in X$$

besitzt 2 verschiedene Lösungen  $x^1, x^2$ ,  $x^1 \neq x^2$ . Für  $\lambda \in ]0, 1[$  ist dann wieder  $\lambda x^1 + (1 - \lambda)x^2 \in X$  und

$$f(\lambda x^1 + (1 - \lambda)x^2) < \lambda f(x^1) + (1 - \lambda) \underbrace{f(x^2)}_{=f(x^1)} = f(x^1) = \min\{f(x) \mid x \in X\}, \quad 0 < \lambda < 1.$$

Somit haben wir einen Widerspruch! Also folgt  $x^1 = x^2$ .

- c.) Nach b.) existiert höchstens eine Lösung von

$$f(x) \stackrel{!}{=} \min \quad \text{unter der Nebenbedingung } x \in X.$$

Sei nun  $x^0 \in X$ ,  $X \neq \emptyset$  beliebig gewählt. Für  $x, y \in L(x^0)$  gilt dann

$$f(\lambda x + (1 - \lambda)y) \leq \lambda \underbrace{f(x)}_{\leq f(x^0)} + (1 - \lambda) \underbrace{f(y)}_{\leq f(x^0)} \leq f(x^0), \quad 0 \leq \lambda \leq 1$$

d.h.  $L(x^0)$  ist eine konvexe Menge und somit nach Lemma 2.6 kompakt. Folglich ist  $X \cap L(x^0)$  kompakt und nicht leer. Daher nimmt  $f$  als stetige Funktion sein globales Minimum  $x^*$  in  $L(x^0) \cap X$  an.  $x^*$  ist dann natürlich auch das globale Minimum von  $f$  auf  $X$ .

□

### 3 Newton und Quasi-Newton Verfahren

#### 3.1 Das lokale Verfahren

Betrachte das Problem

$$f(x) \stackrel{!}{=} \min$$

für  $f \in C^2(D, \mathbb{R})$ ,  $D \subset \mathbb{R}^N$ . Ferner existiere ein  $x^* \in D$  mit  $\nabla f(x^*) = 0$  und  $\nabla^2 f(x^*)$  positiv definit. Verwende wieder ein Verfahren der Form

$$x^{k+1} = x^k + t_k d^k, \quad k = 0, 1, 2, \dots$$

An den Iterierten  $x^k$  betrachten wir ein quadratisches Modell für  $f$ , d.h.

$$m_k(d) = f(x^k) + \nabla f(x^k)^T d + \frac{1}{2} d^T \nabla^2 f(x^k) d.$$

Für die Lösung  $d$  gilt

$$0 = \nabla m_k(d) = \nabla f(x^k) + \nabla^2 f(x^k) d \iff \nabla^2 f(x^k) d = -\nabla f(x^k).$$

Man setzt dann  $x^{k+1} = x^k + d^k$  (d.h.  $t_k = 1$ ). Dies entspricht gerade dem Newton-Verfahren

$$x^{k+1} = x^k - \nabla^2 f(x^k)^{-1} \nabla f(x^k), \quad k = 0, 1, 2, \dots$$

für die Gleichung  $\nabla f(x) = 0$ . Damit erhalten wir sofort einen lokalen Konvergenzsatz.

#### Satz 3.1

Vorgelegt sei das Problem  $f(x) \stackrel{!}{=} \min$  für  $f \in C^2(D, \mathbb{R})$ ,  $D \subset \mathbb{R}^N$  offen. Ferner existiere ein  $x^* \in D$  mit  $\nabla f(x^*) = 0$  und  $\nabla^2 f(x^*)$  positiv definit. Dann gibt es eine Kugel  $K_\varphi(x^*) = \{x \in D \mid \|x - x^*\|_\infty \leq \varphi\} \subset D$ ,  $\varphi > 0$ , so dass für jedes  $x^0 \in K_\varphi(x^*)$  die Newtonfolge

$$x^{k+1} = x^k - (\nabla^2 f(x^k))^{-1} \nabla f(x^k)$$

in  $K_\varphi(x^*)$  liegt und gegen  $x^*$  konvergiert. Überdies gibt es ein  $C > 0$  mit

$$\|x^{k+1} - x^*\|_\infty \leq C \|x^k - x^*\|_\infty^2 \quad \forall k \geq 0, \quad x^0 \in K_\varphi(x^*).$$

“Lokal quadratische Konvergenzordnung”

#### 3.2 Fehlerentwicklung von Differenzenquotienten

Da die Durchführung des Newton-Verfahrens die Auswertung der 2-ten Ableitung des Originalproblems erfordert, bietet sich eine numerische Approximation von Ableitungen an.

Analysiere jetzt die Fehlerentwicklung bei approximativer Auswertung einer Funktion  $h$  bei  $x$  und Verwendung von Differenzenquotienten. Sei

$$\tilde{h}(x) = h(x) + \tilde{\varepsilon}(x) \quad \text{mit } |\tilde{\varepsilon}(x)| \leq \varepsilon, \quad \varepsilon > 0.$$

Bestimme die Ableitungen von  $h$  numerisch, z.B. durch den Vorwärtsdifferenzenquotient

$$D_{\Delta x}^+ h(x) = \frac{\tilde{h}(x + \Delta x) - \tilde{h}(x)}{\Delta x}.$$



Dann gilt

$$\begin{aligned}
\|D_{\Delta x}^+ h(x) - h'(x)\| &= \left\| \frac{\tilde{h}(x + \Delta x) - \tilde{h}(x)}{\Delta x} - h'(x) \right\| \\
&= \left\| \frac{h(x + \Delta x) + \tilde{\varepsilon}(x + \Delta x) - h(x) - \tilde{\varepsilon}(x)}{\Delta x} - h'(x) \right\| \\
&\leq \left\| \frac{h(x + \Delta x) - h(x)}{\Delta x} - h'(x) \right\| + \frac{2\varepsilon}{\Delta x} \\
&= \left\| \frac{h(x) + h'(x)\Delta x + \frac{1}{2}h''(\xi)\Delta x^2 - h(x)}{\Delta x} - h'(x) \right\| + \frac{2\varepsilon}{\Delta x}, \quad \xi \in ]x, x + \Delta x[ \\
&= \frac{\Delta x}{2} \|h''(\xi)\| + \frac{2\varepsilon}{\Delta x} \\
&= O\left(\Delta x + \frac{\varepsilon}{\Delta x}\right).
\end{aligned}$$

Die Minimalstelle  $\Delta x^*$  der Fehlerfunktion

$$err(\Delta x) = \Delta x + \frac{\varepsilon}{\Delta x}.$$

erfüllt

$$err'(\Delta x) = 1 - \frac{\varepsilon}{\Delta x^2} = 0, \quad \text{d.h. } \nabla x^* = \sqrt{\varepsilon}.$$

Der Fehler in der 1-ten Ableitung ist also

$$err(\Delta x^*) = \Delta x^* + \frac{\varepsilon}{\Delta x^*} = \sqrt{\varepsilon} + \frac{\varepsilon}{\sqrt{\varepsilon}} = O(\sqrt{\varepsilon})$$

Ist nun z.B.  $\varepsilon = 10^{-16}$  die Maschinengenauigkeit, so wendet man deshalb den Vorwärtsdifferenzenquotienten  $D_{\Delta x}^+$  mit der Schrittweite  $\Delta x^* = \sqrt{\varepsilon} = 10^{-8}$  an.

Für die Approximation der 2-ten Ableitung basierend auf zweifacher numerischer Differentiation hat man den Fehler  $O(\sqrt[4]{\varepsilon})$ , d.h. die 2-te Ableitung wird basierend auf zweimalig numerischer Differentiation relativ ungenau approximiert. Deshalb wertet man normalerweise bei  $f(x) \stackrel{!}{=} \min$  die Funktion  $f$  und  $\nabla f$  explizit aus und approximiert höchstens  $\nabla^2 f(x)$ .

### 3.3 Quasi-Newton Verfahren

Die Nachteile des Newton-Verfahrens für die Gleichung  $\nabla f(x) = 0$  lauten:

- Es werden 2-te Ableitungen von  $f$  benötigt.
- $\nabla^2 f(x)$  muss positiv definit sein.
- Bei Lösung des Systems  $\nabla^2 f(x^k)d = -\nabla f(x^k)$  mit dem Cholesky-Verfahren werden  $O(N^3)$ -Multiplikationen benötigt.

Quasi-Newton Verfahren haben folgende Grundstruktur:

- $x^{k+1} = x^k + t_k d^k, \quad k = 0, 1, 2, \dots$
- $d^k = -(B_k)^{-1} \nabla f(x^k), \quad (B_k)^{-1} \in \mathbb{R}^{N,N}$ .
- Verwende  $B_k$  um  $B_{k+1}$  zu erhalten.

## Quasi-Newton Verfahren

- Approximiere 2-te Ableitungen durch 1-te Ableitungen.
- Erhalten die positive Definitheit von  $B^k$ .

Zur Motivation der Abstiegsrichtung  $d^{k+1}$  in  $x^{k+1}$  betrachten wir das quadratische Modell

$$m_{k+1}(d) = f(x^{k+1}) + \nabla f(x^{k+1})^T d + \frac{1}{2} d^T B_{k+1} d$$

mit  $\nabla m_{k+1}(d) = \nabla f(x^{k+1}) + B_{k+1} \cdot d$ .

### Bedingungen an $B_{k+1}$ :

Der Gradient von  $m_{k+1}$  und  $f$  soll bei  $x^k$  und  $x^{k+1}$  übereinstimmen, d.h.

$$\begin{aligned} \nabla m_{k+1}(0) &= \nabla f(x^{k+1}), \\ \nabla m_{k+1}(-t_k d^k) &= \nabla f(x^k). \end{aligned}$$

Die erste Bedingung ist nach Konstruktion immer erfüllt, und die zweite Bedingung erfordert

$$\nabla m_{k+1}(-t_k d^k) = \nabla f(x^{k+1}) + B_{k+1}(-t_k d^k) \stackrel{!}{=} \nabla f(x^k) \iff B_{k+1} t_k d^k = \nabla f(x^{k+1}) - \nabla f(x^k).$$

Mit  $s^k = t_k d^k$ ,  $y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$  schreibt sich dies zu

$$B_{k+1} s^k = y^k. \quad \text{"Sekantengleichung"} \quad (35)$$

Die Forderung  $B_{k+1}$  positiv definit erzwingt dann insbesondere

$$0 < (s^k)^T B_{k+1} s^k = (s^k)^T y^k = t_k (d^k)^T \left( \nabla f(x^{k+1}) - \nabla f(x^k) \right), \quad k \in \mathbb{N}. \quad (36)$$

(36) ist z.B. erfüllt wenn  $t_k > 0$  der Wolfe-Powell Bedingung

$$f(x^{k+1}) \leq f(x^k) + c_1 t_k \nabla f(x^k)^T d^k, \quad c_1 \in ]0, \frac{1}{2}[ \quad (37)$$

$$\nabla f(x^{k+1})^T d^k \geq c_2 \nabla f(x^k)^T d^k, \quad c_2 \in ]c_1, 1[ \quad (38)$$

genügt. Dann gilt nämlich

$$(y^k)^T s^k = \left( \nabla f(x^{k+1}) - \nabla f(x^k) \right)^T t_k d^k \geq \underbrace{(c_2 - 1)}_{<0} t_k \underbrace{\nabla f(x^k)^T d^k}_{<0} > 0$$

d.h. (36) ist erfüllt.

Unter der Voraussetzung  $(s^k)^T y^k > 0$  ist die Sekantengleichung lösbar. Eine eindeutige Bestimmung von  $B_{k+1}$  ist möglich als Lösung der Aufgabe

$$\min \{ \|B - B_k\| \mid B \in \mathbb{R}^{N,N} \text{ symmetrisch und positiv definit, } B s^k = y^k \}.$$

Wählt man als Matrixnorm eine geeignete gewichtete Frobenius-Norm, so erhält man als eindeutige Lösung

$$B_{k+1} = (I - \rho_k y^k (s^k)^T) B_k (I - \rho_k s^k (y^k)^T) + \rho_k y^k (y^k)^T \quad \text{mit } \rho_k = \frac{1}{(y^k)^T s^k}. \quad (39)$$

(39) heißt die **DFP-Aufdatierung** und geht auf Davidon, Fletcher und Powell zurück.

### 3.4 Alternative Aufdatierung

Eine Alternative zur DFP-Aufdatierung erhält man, indem man nicht Bedingungen an die Matrix  $B_k$  sondern an ihre Inverse  $H_k = B_k^{-1}$  stellt. Die Sekantengleichung lautet dann

$$s^k = (B_{k+1})^{-1}y^k = H_{k+1}y^k. \quad (40)$$

Man bestimmt dann  $H_{k+1}$  eindeutig als Lösung der Aufgabe

$$\min \{ \|H - H_k\| \mid H \in \mathbb{R}^{N,N} \text{ symmetrisch und positiv definit, } Hy^k = s^k \}$$

in einer gewichteten Frobenius-Norm. Die eindeutige Lösung lautet dann

$$H_{k+1} = (I - \rho_k s^k (y^k)^T) H_k (I - \rho_k y^k (s^k)^T) + \rho_k s^k (s^k)^T \quad \text{mit } \rho_k = \frac{1}{(y^k)^T s^k}. \quad (41)$$

(41) **BFGS-Update-Formel** (BFGS steht für Broyden, Fletcher, Goldfarb und Shanno).

#### Bemerkung 3.2

Die zu (41) äquivalente Update Formel für  $B_{k+1} = H_{k+1}^{-1}$  lautet

$$B_{k+1} = B_k - \frac{B_k s^k (s^k)^T B_k}{(s^k)^T B_k s^k} + \frac{y^k (y^k)^T}{(y^k)^T s^k}. \quad (42)$$

#### Lemma 3.3

$H_{k+1}$  aus (41) ist symmetrisch und positiv definit, falls  $H_k$  dies war.

**Beweis:** Die Symmetrie  $H_{k+1}^T = H_{k+1}$  ist nach Definition klar. Bzgl. der Definitheit berechnet man

$$w^T H_{k+1} w = w^T (I - \rho_k s^k (y^k)^T) H_k (I - \rho_k y^k (s^k)^T) w + \underbrace{\rho_k}_{>0} w^T s^k (s^k)^T w = z^T H_k z + \rho_k (w^T s^k)^2 \geq 0$$

mit  $z = (I - \rho_k y^k (s^k)^T) w$ . Weiter gilt: Ist  $w^T H_{k+1} w = 0$ , so folgt  $w^T s^k = 0$  und damit

$$w^T (I - \rho_k s^k (y^k)^T) H_k (I - \rho_k y^k (s^k)^T) w = w^T H_k w = 0 \quad \Rightarrow \quad w = 0$$

da  $H_k$  symmetrisch und positiv definit.

□

Dies führt nun zum BFGS-Algorithmus.

**Input:**  $x^0 \in \mathbb{R}^N, \varepsilon > 0$  und  $H_0 \in \mathbb{R}^{N,N}$  (Approximation für  $(\nabla^2 f(x^0))^{-1}$ ) positiv definit

$k = 0$

**while**  $\|\nabla f(x^k)\|_2 \geq \varepsilon$  **do**

    berechne  $d^k = -H_k \nabla f(x^k)$ ;

    setze  $x^{k+1} = x^k + t_k d^k$  mit  $t_k$  gemäß Wolfe-Powell;

    setze  $s^k = x^{k+1} - x^k, y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$  und berechne  $H_{k+1}$  gemäß

$$H_{k+1} = (I - \rho_k s^k (y^k)^T) H_k (I - \rho_k y^k (s^k)^T) + \rho_k s^k (s^k)^T;$$

    setze  $k = k + 1$ ;

**end**

BFGS-Algorithmus

### Bemerkung 3.4

- Jede Iteration des BFGS-Algorithmus kann in  $O(N^2)$ -Operationen durchgeführt werden.
- Der BFGS-Algorithmus ist numerisch robust.
- Man benutze (41) statt (42) zur Aufdatierung wegen

$$d^k = -H_k \nabla f(x^k) = -B_k^{-1} \nabla f(x^k). \quad (\text{Aufwand!})$$

Einige Details zur Implementierung des BFGS-Algorithmus:

- Bei der Liniensuche gemäß Wolfe-Powell sollte man (aus theoretischen Gründen) die Wahl  $t_k = 1$  probieren. Praktische Werte für Wolfe-Powell sind  $\alpha = 10^{-4}$  und  $\rho = 0.9$ .
- Für  $H_0$  wählt man oft  $H_0 = \beta I$  mit  $\beta > 0$ .

### 3.5 Konvergenz-Analyse der BFGS-Methode

#### Satz 3.5

Sei  $f \in C^2(\mathbb{R}^N, \mathbb{R})$ , und sei

$$L(x^0) = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^0)\}$$

konvex. Ferner existieren  $m, M > 0$  mit  $mz^T z \leq z^T \nabla^2 f(x) z \leq Mz^T z$  für  $x \in L(x^0)$ ,  $z \in \mathbb{R}^N$ . Sei nun  $B_0 \in \mathbb{R}^{N,N}$  symmetrisch und positiv definit, und sei  $x^0 \in \mathbb{R}^N$ .

Dann konvergiert die Folge  $(x^k)_{k \in \mathbb{N}}$ ,  $x^{k+1} = x^k + t_k d^k$ ,  $t_k$  gemäß Wolfe-Powell und  $d^k = -B_k^{-1} \nabla f(x^k)$ ,  $B_k^{-1} = H_k$  aus Algorithmus (BFGS) gegen die eindeutige Lösung  $x^* \in L(x^0)$ .

#### Bemerkung 3.6

Gemäß Lemma 2.6 und 2.7 hat die Aufgabe  $f(x) \stackrel{!}{=} \min$  genau eine Lösung  $x^*$  in  $L(x^0)$ .

**Beweis:** Wir bemerken, dass das Verfahren in  $L(x^0)$  wohldefiniert ist. Es gilt nämlich

$$(s^k)^T y^k > 0 \quad \text{für } k \in \mathbb{N}$$

gemäß (35) und der Wolfe-Powell Schrittweitensteuerung, gemäß Lemma 3.3 sind die Matrizen  $H_k$ ,  $k \in \mathbb{N}$  symmetrisch und positiv definit. Ferner ist  $L(x^0)$  gemäß Lemma 2.6 kompakt und somit  $f$  nach unten beschränkt, d.h. nach Bemerkung 2.1 existiert eine Wolfe-Powell Schrittweite  $t_k$ ,  $k \in \mathbb{N}$ . Sei nun  $B_k = H_k^{-1}$ ,  $k \in \mathbb{N}$  und sei

$$m_k = \frac{(y^k)^T s^k}{(s^k)^T s^k}, \quad M_k = \frac{(y^k)^T y^k}{(y^k)^T s^k}.$$

Nach Definition gilt

$$y^k = \nabla f(x^{k+1}) - \nabla f(x^k) = \nabla^2 f(\eta^k)(x^{k+1} - x^k) = \nabla^2 f(\eta^k) s^k$$

mit  $\eta^k = x^{k+1} + \theta^k(x^{k+1} - x^k) \in L(x^0)$ ,  $\theta^k \in [0, 1]$ . Somit folgt

$$m_k = \frac{(s^k)^T \nabla^2 f(\eta^k) s^k}{(s^k)^T s^k} \geq m, \quad k \in \mathbb{N}$$
$$M_k = \frac{(s^k)^T (\nabla^2 f(\eta^k))^2 s^k}{(s^k)^T \nabla f(\eta^k) s^k} = \frac{(w^k)^T \nabla^2 f(\eta^k) w^k}{(w^k)^T w^k} \leq M$$

für  $k \in \mathbb{N}$  mit  $w^k = (\nabla^2 f(\eta^k))^{1/2} s^k$ ,  $\nabla^2 f(\eta^k)$  positiv definit. Die Analyse des BFGS-Algorithmus erfolgt mit der Aufdatierungsformel

$$B_{k+1} = B_k - \frac{B_k s^k (s^k)^T B_k}{(s^k)^T B_k s^k} + \frac{y^k (y^k)^T}{(y^k)^T s^k}.$$

Mit einigen Hilfslemmata aus der lineare Algebra findet man

$$\text{tr}(B_{k+1}) = \text{tr}(B_k) - \frac{\|B_k s^k\|_2^2}{(s^k)^T B_k s^k} + \frac{\|y^k\|_2^2}{(y^k)^T s^k}, \quad (43)$$

$$\det(B_{k+1}) = \det(B_k) \cdot \frac{(y^k)^T s^k}{(s^k)^T B_k s^k}. \quad (44)$$

Setze nun  $\theta_k := \angle(s^k, B_k s^k)$ , d.h.

$$\cos(\theta_k) = \frac{(s^k)^T B_k s^k}{\|s^k\|_2 \cdot \|B_k s^k\|_2} \quad \text{und} \quad q_k := \frac{(s^k)^T B_k s^k}{(s^k)^T s^k}, \quad k \in \mathbb{N}.$$

Einsetzen in (44) liefert

$$\det(B_{k+1}) = \det(B_k) \cdot \frac{(y^k)^T s^k}{(s^k)^T s^k} \cdot \frac{(s^k)^T s^k}{(s^k)^T B_k s^k} = \det(B_k) \cdot \frac{m_k}{q_k}.$$

Definiere jetzt die Funktion

$$\psi : SP_n \rightarrow \mathbb{R}, \quad \psi(B) = \text{tr}(B) - \ln(\det(B))$$

mit  $SP_n = \{C \in \mathbb{R}^{N,N} \mid C \text{ symmetrisch und positiv definit}\}$ . Dann gilt

$$\psi(B) = \text{tr}(B) - \underbrace{\ln(\det(B))}_{>0} = \sum_{i=1}^N \lambda_i - \ln\left(\prod_{i=1}^N \lambda_i\right) = \sum_{i=1}^N \lambda_i - \sum_{i=1}^N \ln(\lambda_i) = \sum_{i=1}^N \underbrace{\lambda_i - \ln(\lambda_i)}_{>0} > 0 \quad \forall B$$

mit  $\lambda_1, \lambda_2, \dots, \lambda_N > 0$  Eigenwerte von  $B \in SP_n$ . Zusammen mit (43)-(44) folgt dann

$$\begin{aligned} \psi(B_{k+1}) &= \text{tr}(B_{k+1}) - \ln(\det(B_{k+1})) \\ &= \text{tr}(B_k) - \frac{\|B_k s^k\|_2^2}{(s^k)^T B_k s^k} + \underbrace{\frac{\|y^k\|_2^2}{(y^k)^T s^k}}_{=M_k} - \ln\left(\det(B_k) \cdot \frac{m_k}{q_k}\right) \\ &= \text{tr}(B_k) - \underbrace{\frac{(B_k s^k)^T B_k s^k}{(s^k)^T B_k s^k}}_{=\frac{q_k}{\cos^2(\theta_k)}} + M_k - \ln(\det(B_k)) - \ln(m_k) + \ln(q_k) \\ &= \text{tr}(B_k) - \frac{q_k}{\cos^2(\theta_k)} + M_k - \ln(\det(B_k)) - \ln(m_k) + \ln(q_k) \\ &= \psi(B_k) + M_k - \ln(m_k) - 1 + \left[1 - \frac{q_k}{\cos^2(\theta_k)} + \ln\left(\frac{q_k}{\cos^2(\theta_k)}\right)\right] + \ln(\cos^2(\theta_k)). \end{aligned}$$

Mit  $h(t) = 1 - t + \ln(t) \leq 0$  für  $t \geq 0$  folgt

$$h\left(\frac{q_k}{\cos^2(\theta_k)}\right) = 1 - \frac{q_k}{\cos^2(\theta_k)} + \ln\left(\frac{q_k}{\cos^2(\theta_k)}\right) \leq 0$$

und somit

$$\begin{aligned} 0 < \psi(B_{k+1}) &\leq \psi(B_k) + (M_k - \ln(m_k) - 1) + \ln(\cos^2(\theta_k)) \\ &\leq \psi(B_k) + (M - \ln(m) - 1) + \ln(\cos^2(\theta_k)). \end{aligned} \quad (45)$$

Induktiv folgt aus (45)

$$\begin{aligned} 0 < \psi(B_{k+1}) &\leq \psi(B_0) + \underbrace{(M - \ln(m) - 1)}_{=c}(k+1) + \sum_{j=0}^k \ln(\cos^2(\theta_j)) \\ &= \varphi(B_0) + c(k+1) + \sum_{j=0}^k \ln(\cos^2(\theta_j)). \end{aligned} \quad (46)$$

Für die Quasi-Newton BFGS-Algorithmen gilt

$$s^k = -t_k B_k^{-1} \nabla f(x^k) \quad (\text{Quasi-Newton Richtung})$$

und somit

$$\theta_k = \angle(s^k, B_k s^k) = \angle(s^k, -t_k \nabla f(x^k)) = \angle(s^k, -\nabla f(x^k))$$

d.h.  $\theta_k$  ist der Winkel zwischen der Richtung des steilsten Abstiegs ( $-\nabla f(x^k)$ ) und der Suchrichtung  $s^k$ . Somit gilt gemäß Satz 2.2 (Theorem von Zoutendijk) sofort

$$\sum_{k=0}^{\infty} \cos^2(\theta_k) \|\nabla f(x^k)\|_2^2 < \infty.$$

Annahme:

$$\cos(\theta_j) \rightarrow 0 \text{ für } j \rightarrow \infty.$$

Dann gilt  $\ln(\cos^2(\theta_j)) < -2c$  für  $j \geq \bar{k}$ . Für  $k > \bar{k}$  folgt mit (46)

$$0 < \psi(B_0) + c(k+1) + \sum_{l=0}^{\bar{k}} \ln(\cos^2(\theta_l)) + \underbrace{\sum_{l=\bar{k}+1}^k (-2c)}_{=(-2c)(k-\bar{k})} = \psi(B_0) + \sum_{l=0}^{\bar{k}} \ln(\cos^2(\theta_l)) + 2c\bar{k} + c - ck < 0$$

für  $k \in \mathbb{N}$  hinreichend groß und somit Widerspruch!

Also existiert eine Teilfolge  $(\theta_{j_l})_{l \in \mathbb{N}}$  mit

$$\cos^2(\theta_{j_l}) \geq \delta > 0, \quad l \in \mathbb{N}$$

und somit erhalten wir

$$\sum_{j=0}^{\infty} \cos^2(\theta_j) \geq \sum_{l=0}^{\infty} \cos^2(\theta_{j_l}) \geq \sum_{l=0}^{\infty} \delta = \infty, \quad \text{d.h. } \sum_{j=0}^{\infty} \cos^2(\theta_j) = \infty.$$

Da die Wolfe-Powell Schrittweiten  $t_k > 0$  effizient sind, ist Satz 2.5 anwendbar und liefert  $\lim_{k \rightarrow \infty} x^k = x^*$ .

□

**Bemerkung 3.7**

- Man kann nun sogar zeigen, dass der BFGS-Algorithmus unter der Voraussetzung von Satz 3.5 nicht nur global, sondern auch superlinear konvergiert und zwar weiterhin mit beliebigen Startvektoren  $x^0$  und symmetrisch und positiv definiten Startmatrix  $H_0 \in \mathbb{R}^{N,N}$ .
- Superlineare Konvergenz heißt

$$\lim_{k \rightarrow \infty} \frac{\|x^{k+1} - x^*\|_2}{\|x^k - x^*\|_2} = 0.$$

# Kapitel III: Numerische Verfahren für allgemeine Optimierungsprobleme

## 1 Verfahren für quadratische Optimierungsaufgaben

Vorgelegt sei die Aufgabe

$$(I) = \begin{cases} f(x) = \frac{1}{2}x^T Qx + c^T x + \gamma \stackrel{!}{=} \min \\ \text{unter den Nebenbedingungen} \\ g_i(x) = a_i^T x - b_i = 0, \quad i = 1, \dots, l \\ k_i(x) = a_{l+i}^T x - b_{l+i} \geq 0, \quad i = 1, \dots, q \\ \text{mit } Q \in \mathbb{R}^{N,N} \text{ symmetrisch und } c, a_1, \dots, a_{l+q} \in \mathbb{R}^N \text{ sowie } \gamma, b_1, \dots, b_{l+q} \in \mathbb{R}. \end{cases}$$

Die Karush-Kuhn-Tucker Bedingungen für (I) lauten

$$\begin{aligned} Qx + c - \sum_{i=1}^l \lambda_i a_i - \sum_{i=1}^q \mu_i a_{l+i} &= 0 \\ \mu_i (a_{l+i}^T x - b_{l+i}) &= 0, \quad i = 1, \dots, q \\ a_i^T x - b_i &= 0 \\ a_{l+i}^T x - b_{l+i} &\geq 0, \quad \mu \geq 0. \end{aligned}$$

Ist nun  $Q \in \mathbb{R}^{N,N}$  positiv semidefinit, so wird (I) zu einem konvexen Optimierungsproblem (d.h.  $f$  konvex,  $k_i$  konkav,  $i = 1, \dots, q$ ) und es gilt das folgende Lemma.

### Lemma 1.1

Vorgelegt sei ein konvexes Optimierungsproblem der Form (I). Erfüllt  $x^*$  die Karush-Kuhn-Tucker Bedingungen, so ist  $x^*$  globales Minimum von (I).

**Beweis:** Siehe Aufgabe 1b), Blatt 3.

### Bemerkung 1.2

- Im Fall der Konvexität folgt also aus den KKT-Bedingungen automatisch die Optimalität.
- Ein konvexes quadratisches Problem (I) kann aber unter Umständen keine zulässigen Punkte besitzen, wie z.B.

$$f(x) = x^2 \stackrel{!}{=} \min$$

unter der Nebenbedingung  $g(x) = 0^T x - 1 = 0$ .

Ist die Menge der zulässigen Punkte

$$Z = \{x \in \mathbb{R}^N \mid a_i^T x - b_i = 0, i = 1, \dots, l, a_{l+i}^T x - b_{l+i} \geq 0, i = 1, \dots, q\} \neq \emptyset,$$

so ist diese konvex und die Existenz einer Lösung folgt aus Lemma 2.6, Kapitel II.

- Ist  $Q \in \mathbb{R}^{N,N}$  nicht positiv definit, so kann das quadratische Optimierungsproblem (I) mehrere lokale Lösungen haben.



## 1.1 Innere-Punkte-Verfahren für quadratische Probleme

Vorgelegt sei die Aufgabe

$$(II) = \begin{cases} f(x) = \frac{1}{2}x^T Qx + c^T x + \gamma \stackrel{!}{=} \min, \\ Q \in \mathbb{R}^{N,N} \text{ symmetrisch, positiv definit, } c \in \mathbb{R}^N \text{ und } \gamma \in \mathbb{R} \text{ unter den Nebenbedingungen} \\ k(x) = Ax - b \geq 0, \quad A \in \mathbb{R}^{q,N}, \quad b \in \mathbb{R}^q. \end{cases}$$

Lineare Gleichungsrestriktionen lassen wir aus technischen Gründen beiseite, da man sie durch Reduktion der Variablenanzahl beseitigen kann.

Ist  $x^*$  eine lokale Lösung von (II), so lauten die notwendigen Optimalitätsbedingungen 1-ter Ordnung

$$\begin{aligned} Qx + c - A^T \mu &= 0 \\ Ax - b &\geq 0 \\ (Ax - b)_i \mu_i &= 0, \quad i = 1, \dots, q \text{ und } \mu \geq 0. \end{aligned}$$

Führe nun die Slack-Variablen  $y = Ax - b \in \mathbb{R}^q$  ein und erhalte die modifizierten KKT-Gleichungen

$$Qx + c - A^T \mu = 0 \tag{47}$$

$$Ax - b - y = 0 \tag{48}$$

$$y_i \mu_i = 0, \quad i = 1, \dots, q \text{ und } y, \mu \geq 0. \tag{49}$$

Mit  $Y = \text{diag}(y_1, \dots, y_q)$ ,  $M = \text{diag}(\mu_1, \dots, \mu_q)$  schreibt sich (47)-(49) in der Form

$$F(x, y, \mu) = \begin{pmatrix} Qx - A^T \mu + c \\ b - Ax + y \\ YM\mathbb{I} \end{pmatrix} = 0 \quad \text{mit } y, \mu \geq 0 \text{ und } \mathbb{I} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^q.$$

Sei  $(x, y, \mu)$  eine aktuelle Iterierte. Dann nennt man

$$\eta := \frac{y^T \mu}{q}$$

die **gewichtete Dualitätslücke**.

Sei  $\tau > 0$ . Die Lösungen  $(x_\tau, y_\tau, \mu_\tau)$  von

$$F(x, y, \mu) = \begin{pmatrix} Qx - A^T \mu + c \\ b - Ax + y \\ YM\mathbb{I} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \tau \mathbb{I} \end{pmatrix}, \quad (y, \mu) > 0 \tag{50}$$

heißen **zentraler Pfad**  $\mathcal{C}$ .  $\mathcal{C}$  ist eine Kurve in  $\mathbb{R}^{N+2q}$  welche für  $\tau \rightarrow 0$  gegen die Lösung des quadratischen Problems (II) konvergiert.

Zur Lösung von (50) geht man wie folgt vor. Wähle  $\sigma \in [0, 1]$ , setze  $\tau = \sigma \eta$  und wende ein gedämpftes Newton-Verfahren zur Lösung von (50) an. Ein Newton Schritt für (50) lautet

$$\begin{pmatrix} Q & 0 & -A^T \\ -A & I & 0 \\ 0 & M & Y \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta y \\ \Delta \mu \end{pmatrix} = \begin{pmatrix} -Qx + A^T \mu - c \\ Ax - y - b \\ -YM\mathbb{I} + \sigma \eta \mathbb{I} \end{pmatrix}. \tag{51}$$

Die nächste Iterierte ist dann für  $\alpha \in ]0, 1]$  gegeben durch

$$(x^+, y^+, \mu^+) = (x, y, \mu) + \alpha(\Delta x, \Delta y, \Delta \mu).$$

Dabei ist  $\alpha \in ]0, 1]$  so zu wählen, dass  $(y^+, \mu^+) > 0$  erfüllt ist.

Der Hauptaufwand eines Inneren-Punkte Verfahrens besteht in der Regel in der Lösung des Systems (51). Deshalb ist es wichtig, die spezielle Struktur von (51) auszunutzen.

Aus der 3-ten Zeile von (51) folgt

$$\Delta y = M^{-1}(-MY\mathbb{I} + \sigma\eta\mathbb{I} - Y\Delta\mu) = -Y\mathbb{I} + \sigma\eta M^{-1}\mathbb{I} - M^{-1}Y\Delta\mu = -y + \sigma\eta M^{-1}\mathbb{I} - M^{-1}Y\Delta\mu.$$

Setze dies in die 2-te Zeile von (51) ein und finde

$$\begin{aligned} b + y - Ax &= A\Delta x - \Delta y = A\Delta x + y - \sigma\eta M^{-1}\mathbb{I} + M^{-1}Y\Delta\mu \\ &= (A, M^{-1}Y) \begin{pmatrix} \Delta x \\ \Delta\mu \end{pmatrix} - (-y + \sigma\eta M^{-1}\mathbb{I}), \text{ d.h.} \end{aligned}$$

$$\begin{pmatrix} Q & -A^T \\ A & M^{-1}Y \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta\mu \end{pmatrix} = \begin{pmatrix} -Qx + A^T\mu - c \\ -Ax + y + b + (-y + \sigma\eta M^{-1}\mathbb{I}) \end{pmatrix}. \quad (52)$$

Aus der 2-ten Zeile von (52) folgt

$$\Delta\mu = Y^{-1}M(-Ax + b + \sigma\eta M^{-1}\mathbb{I} - A\Delta x)$$

d.h. die erste Zeile von (52) liefert das System

$$(Q + A^T Y^{-1} M A) \Delta x = -Qx + A^T \mu - c - A^T Y^{-1} M (-Ax + b + \sigma\eta M^{-1} \mathbb{I}). \quad (53)$$

Das System (53) kann z.B. mit dem Cholesky-Verfahren gelöst werden.

### Schrittweiten-Bestimmung des gedämpften Newton-Verfahrens

Die am meisten verwendete Variante des Innere-Punkte Algorithmus basiert auf dem Prädiktor-Korrektor-Algorithmus von Mehrotra.

Zuerst wird ein affiner Skalierungsschritt  $(\Delta x^{aff}, \Delta y^{aff}, \Delta \mu^{aff})$  bestimmt, indem (51) mit  $\sigma = 0$  gelöst wird. Die erhaltene Richtung wird dann in einem Korrekturschritt verbessert, wobei die Lösung von (51) mit

$$\begin{aligned} \sigma &= \left( \frac{\eta^{aff}}{\eta} \right)^3, \quad \eta^{aff} = \frac{1}{q} (y + \alpha^{aff} \Delta y^{aff})^T (\mu + \alpha^{aff} \Delta \mu^{aff}) \\ \text{und} \quad \alpha^{aff} &= \max\{\alpha \in ]0, 1] \mid (y, \mu) + \alpha(\Delta y^{aff}, \Delta \mu^{aff}) \geq 0\} \end{aligned}$$

berechnet wird. Im Korrekturschritt löst man das System

$$\begin{pmatrix} Q & 0 & -A^T \\ -A & I & 0 \\ 0 & M & Y \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta y \\ \Delta \mu \end{pmatrix} = \begin{pmatrix} -Qx + A^T \mu - c \\ Ax - y - b \\ -MY\mathbb{I} - \Delta M^{aff} \Delta Y^{aff} \mathbb{I} + \sigma\eta\mathbb{I} \end{pmatrix} \quad (54)$$

mit  $\Delta M^{aff} = \text{diag}(\Delta \mu^{aff}), \Delta Y^{aff} = \text{diag}(\Delta y^{aff})$ .

Insgesamt erhalten wir den folgenden Algorithmus:

Wähle  $(x^0, y^0, \mu^0)$  mit  $y^0, \mu^0 > 0$  und wähle  $\tau \in ]0, 1[$ .

**for**  $k = 0, 1, 2, \dots$  **do**

- mit  $(x^k, y^k, \mu^k)$  löse (51) mit  $\sigma = 0$  für  $(\Delta x^{aff}, \Delta y^{aff}, \Delta \mu^{aff})$ ;
- berechne  $\eta = \frac{1}{q}(y^k)^T \mu^k$ ;
- setze  $\hat{\alpha}^{aff} = \max\{\alpha \in ]0, 1[ \mid (y^k, \mu^k) + \alpha(\Delta y^{aff}, \Delta \mu^{aff}) \geq 0\}$ ;
- bestimme  $\eta^{aff} = (y^k + \hat{\alpha}^{aff} \Delta y^{aff})^T (\mu^k + \hat{\alpha}^{aff} \Delta \mu^{aff})$  und setze  $\sigma = \left(\frac{\eta^{aff}}{\eta}\right)^3$ ;
- löse (54) für  $(\Delta x, \Delta y, \Delta \mu)$ ;
- setze  $\hat{\alpha} := \min(\alpha_y, \alpha_\mu)$  mit
  - $\alpha_y = \max\{\alpha \in ]0, 1[ \mid y^k + \alpha \Delta y \geq (1 - \tau)y^k\}$ ;
  - $\alpha_\mu = \max\{\alpha \in ]0, 1[ \mid \mu^k + \alpha \Delta \mu \geq (1 - \tau)\mu^k\}$ ;
- setze  $(x^{k+1}, y^{k+1}, \mu^{k+1}) = (x^k, y^k, \mu^k) + \hat{\alpha}(\Delta x, \Delta y, \Delta \mu)$ ;

**end**

PRAED.-KORR Algorithmus für Quadratische Probleme

Abbruchkriterium:  $\|F(x^k, y^k, \mu^k)\|_\infty \leq \text{EPS}$ .

## 2 Lokale SQP (sequential quadratic programming) Methoden

Wir betrachten zunächst das Problem

$$(III) = \begin{cases} f(x) \stackrel{!}{=} \min \\ \text{unter der Nebenbedingung} \\ g(x) = 0 \\ \text{für } f \in C^2(D, \mathbb{R}), g \in C^2(D, \mathbb{R}^l), D \subset \mathbb{R}^N \text{ offen.} \end{cases}$$

Die Lagrange-Funktion zu (III) lautet

$$L(x, \lambda) = f(x) - \sum_{i=1}^l \lambda_i g_i(x)$$

und die KKT-Bedingungen ergeben sich zu

$$F(x, \lambda) = \begin{pmatrix} \frac{\partial L}{\partial x}(x, \lambda) \\ \frac{\partial L}{\partial \lambda}(x, \lambda) \end{pmatrix} = \begin{pmatrix} \nabla f(x) - Dg(x)^T \lambda \\ -g(x) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}. \quad (55)$$

Ist  $x^*$  eine reguläre Lösung von (III), so existiert zu  $x^*$  ein Lagrange-Multiplikator  $\lambda^*$  mit  $F(x^*, \lambda^*) = 0$ . Wir lösen (55) mit dem Newton-Verfahren. Die Jacobi-Matrix lautet

$$DF(x, \lambda) = \begin{pmatrix} \frac{\partial^2 L}{\partial x^2}(x, \lambda) & -Dg(x)^T \\ -Dg(x) & 0 \end{pmatrix}.$$

Damit ist ein Newton-Schritt gegeben als

$$\begin{pmatrix} x^{k+1} \\ \lambda^{k+1} \end{pmatrix} = \begin{pmatrix} x^k \\ y^k \end{pmatrix} + \begin{pmatrix} \Delta x^k \\ \Delta \lambda^k \end{pmatrix} \quad (56)$$

mit

$$\frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k) \Delta x^k - Dg(x^k)^T \Delta \lambda^k = -\nabla f(x^k) + Dg(x^k)^T \lambda^k, \quad -Dg(x^k) \Delta x^k = g(x^k). \quad (57)$$

(56)-(57) heißt **Lagrange-Newton-SQP Verfahren**.

**Bemerkung 2.1**

Ist  $x^*$  eine reguläre Lösung von (III) und ist  $\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*)$  positiv definit auf dem kritischen Kegel  $C(x^*) = N(Dg(x^*))$ , so konvergiert das Lagrange-Newton-SQP-Verfahren lokal quadratisch, d.h. es gilt

$$\|(x^{k+1}, \lambda^{k+1}) - (x^*, \lambda^*)\| \leq C \|(x^k, \lambda^k) - (x^*, \lambda^*)\|^2,$$

falls  $\|(x^0, \lambda^0) - (x^*, \lambda^*)\|$  hinreichend klein.

Nach dem lokalen Konvergenzsatz des Newton-Verfahrens ist hierzu zu zeigen, dass

$$DF(x^*, \lambda^*) \text{ invertierbar ist.}$$

Sei  $h = \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} \in \mathbb{R}^{N+l}$  in  $N(DF(x^*, \lambda^*))$ , d.h.

$$\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) h_1 - Dg(x^*)^T h_2 = 0 \quad (58)$$

$$Dg(x^*) h_1 = 0. \quad (59)$$

Gemäß (59) gilt  $h_1 \in N(Dg(x^*))$ . Multipliziere (58) mit  $h_1^T$  und finde

$$h_1^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) h_1 - \underbrace{h_1^T Dg(x^*)^T}_{=(Dg(x^*) h_1)^T = 0^T} h_2 = h_1^T \frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*) h_1 = 0 \implies h_1 = 0$$

da  $\frac{\partial^2 L}{\partial x^2}(x^*, \lambda^*)$  positiv definit auf  $N(Dg(x^*))$ . (58) liefert dann

$$Dg(x^*)^T h_2 = 0$$

und mit  $\text{rg}(Dg(x^*)) = l$  folgt  $h_2 = 0$ .

**2.1 Eine neue Motivation für die Iteration (56) – (57)**

Betrachte das quadratische Problem

$$(IV) = \begin{cases} \tilde{f}(\Delta x) = \nabla f(x^k)^T \Delta x + \frac{1}{2} \Delta x^T \frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k) \Delta x \stackrel{!}{=} \min \\ \text{unter der Nebenbedingung} \\ \tilde{g}(\Delta x) = g(x^k) + Dg(x^k) \Delta x = 0. \end{cases}$$

Die KKT-Bedingungen für (IV) lauten

$$\nabla \tilde{f}(\Delta x) - D\tilde{g}(\Delta x)^T \tilde{\lambda} = 0, \quad \tilde{g}(\Delta x) = 0.$$

Mit

$$\nabla \tilde{f}(\Delta x) = \nabla f(x^k) + \frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k) \Delta x \quad \text{und} \quad D\tilde{g}(\Delta x) = Dg(x^k)$$

erhalten wir

$$\begin{aligned}\nabla f(x^k) + \frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k) \Delta x - Dg(x^k)^T \tilde{\lambda} &= 0 \\ g(x^k) + Dg(x^k) \Delta x &= 0.\end{aligned}$$

Dies ist äquivalent zu

$$\begin{aligned}\frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k) \Delta x - Dg(x^k)^T (\tilde{\lambda} - \lambda^k) &= -\nabla f(x^k) + Dg(x^k)^T \lambda^k \\ -Dg(x^k) \Delta x &= g(x^k).\end{aligned}\tag{60}$$

$$\tag{61}$$

Der Vergleich von (56)-(57) mit (60)-(61) liefert

$$\Delta x^k = \Delta x \text{ und } \Delta \lambda^k = \tilde{\lambda} - \lambda^k, \quad \text{d.h.} \quad \Delta x^k = \Delta x \text{ und } \lambda^k + \Delta \lambda^k = \tilde{\lambda} = \lambda^{k+1}, \text{ d.h.}$$

das quadratische Problem (IV) hat die Lösung  $\Delta x = \Delta x^k$  mit Lagrange-Multiplikator  $\lambda^{k+1}$ .

Damit sind ein Schritt eines Newton-Verfahrens und das Lösen des quadratischen Problems äquivalent.

Wähle Startwerte  $(x^0, y^0)$ .

**for**  $k = 0, 1, 2, \dots$  **do**

- löse das quadratische Problem (IV) für  $\Delta x, \tilde{\lambda}$ ;
- setze  $x^{k+1} = x^k + \Delta x, \lambda^{k+1} = \tilde{\lambda}$ ;
- prüfe die Abbruchkriterien;

**end**

#### LOKALE SQP-METHODE

#### Bemerkung 2.2

- Die lokal quadratische Konvergenz des obigen Algorithmus folgt aus der lokal quadratischen Konvergenz des Newton-Verfahren an einer regulären Lösung  $(x^*, \lambda^*)$ .
- Unter Umständen ist die Matrix

$$\frac{\partial^2 L}{\partial x^2} L(x^k, \lambda^k) = \nabla^2 f(x^k) - \sum_{i=1}^l \lambda_i \nabla^2 g_i(x^k)$$

schwer zu berechnen oder nicht positiv definit auf  $N(Dg(x^k))$ . Eine Alternative ist daher, die Matrix  $\frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k)$  durch ein Quasi-Newton-BFGS-Update  $B_k$  zu ersetzen.

- Der SQP-Rahmen kann leicht auf allgemeine Optimierungsprobleme  $f(x) \stackrel{!}{=} \min$  unter der Nebenbedingung  $g(x) = 0$  und  $k(x) \geq 0$  erweitert werden. Ersetze hierzu an der Stelle  $(x^k, \lambda^k, \mu^k)$  (IV) durch das quadratische Problem

$$(V) = \begin{cases} \tilde{f}(\Delta x) = \nabla f(x^k)^T \Delta x + \frac{1}{2} \Delta x^T \frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k, \mu^k) \Delta x \stackrel{!}{=} \min \\ \text{unter den Nebenbedingungen} \\ \tilde{g}(\Delta x) = g(x^k) + Dg(x^k) \Delta x = 0, \\ \tilde{k}(\Delta x) = k(x^k) + Dk(x^k) \Delta x \geq 0. \end{cases}$$

Eventuell wird in (V)  $\frac{\partial^2 L}{\partial x^2}(x^k, \lambda^k, \mu^k)$  durch eine positiv definite Matrix  $H^k$  ersetzt.

### 3 Globale SQP-Methoden

#### 3.1 Penalty-Funktionen

Betrachte die Aufgabe

$$(VI) = \begin{cases} f(x) \stackrel{!}{=} \min \\ f \in C^2(\mathbb{R}^N, \mathbb{R}) \text{ mit einem zulässigen Bereich} \\ Z = \{x \in \mathbb{R}^N \mid g(x) = 0, k(x) \geq 0\} \\ \text{mit } g \in C^2(\mathbb{R}^N, \mathbb{R}^l), k \in C^2(\mathbb{R}^N, \mathbb{R}^q). \end{cases}$$

Eine Klasse von **Penalty-Funktionen** erhält man aus dem Ansatz

$$P(x, \alpha) = f(x) + \alpha r(x) \tag{62}$$

mit  $r \in C(\mathbb{R}^N, \mathbb{R})$ ,  $r(x) \geq 0 \forall x \in \mathbb{R}^N$  und  $r(x) = 0 \Leftrightarrow x \in Z$ .

Die Funktion  $r$  bestraft also gerade das Verlassen des zulässigen Bereiches  $Z$ .

##### Definition 3.1

Eine Penalty-Funktion der Form (62) heißt **exakt** in einem lokalen Minimum  $x^*$  von (VI), falls ein  $\bar{\alpha} > 0$  existiert, so dass  $x^*$  für alle  $\alpha \geq \bar{\alpha}$  auch ein lokales Minimum von  $P(\cdot, \alpha)$  ist.

##### Bemerkung 3.2

Es lässt sich zeigen, dass bei exakten Penalty-Funktionen in  $x^*$  die zugehörige Funktion  $r$  in  $x^*$  nicht differenzierbar sein kann.

Wir betrachten hier die  $l^1$ -Penalty-Funktion

$$P_1(x, \alpha) = f(x) + \alpha \left( \sum_{j=1}^l |g_j(x)| - \sum_{i=1}^q \min(0, k_i(x)) \right).$$

Der Nachweis der Exaktheit gelingt am einfachsten bei konvexen Optimierungsproblemen.

##### Lemma 3.3

Sei  $(x^*, \lambda^*, \mu^*)$  KKT-Punkt des Optimierungsproblems (VI) mit  $g_i(x) = \alpha_i^T x - \beta_i = 0$ ,  $i = 1, \dots, l$ ,  $k_i(x) \geq 0$ ,  $i = 1, \dots, q$  mit  $f$  konvex und  $k_i$  konkav,  $i = 1, \dots, q$ . Dann existiert ein Penalty-Parameter  $\bar{\alpha} > 0$ , so dass  $x^*$  für  $\alpha \geq \bar{\alpha}$  auch ein Minimum von  $P_1(\cdot, \alpha)$  ist.

##### Bemerkung 3.4

Nach Aufgabe 1b), Blatt 3, ist  $x^*$  auch globales Minimum von (VI).

##### Lemma 3.5

Voraussetzungen wie in Lemma 3.3. Es sei

$$L(x, \lambda, \mu) = f(x) - \sum_{i=1}^l \lambda_i g_i(x) - \sum_{i=1}^q \mu_i k_i(x).$$

$(x^*, \lambda^*, \mu^*)$ ,  $\mu^* \geq 0$  sei Sattelpunkt d.h. es gilt

$$L(x^*, \lambda, \mu) \leq L(x^*, \lambda^*, \mu^*) \leq L(x, \lambda^*, \mu^*) \quad \forall x, \lambda, \mu, \quad \mu \geq 0.$$

Erfüllt  $(x^*, \lambda^*, \mu^*)$  die KKT-Bedingungen, so ist  $(x^*, \lambda^*, \mu^*)$  ein Sattelpunkt.

**Beweis:**  $(x^*, \lambda^*, \mu^*)$  ist KKT-Punkt, d.h.  $\frac{\partial L}{\partial x}(x^*, \lambda^*, \mu^*) = 0$ , d.h.  $x^*$  ist stationärer Punkt von  $L(\cdot, \lambda^*, \mu^*)$ . Ferner ist  $L(\cdot, \lambda^*, \mu^*)$  konvex, da  $f$  konvex und  $k_i$ ,  $i = 1, \dots, q$  konvex,  $\mu^* \geq 0$ . Somit gilt

$$L(x, \lambda^*, \mu^*) - L(y, \lambda^*, \mu^*) \geq \frac{\partial L}{\partial x}(y, \lambda^*, \mu^*)(x - y) \quad \forall x, y \in \mathbb{R}^N.$$

Wende dies an mit  $y = x^*$  und finde

$$L(x, \lambda^*, \mu^*) - L(x^*, \lambda^*, \mu^*) \geq 0 \quad \forall x \in \mathbb{R}^N.$$

Desweiteren folgt mit  $k_i(x^*) \geq 0, g_i(x^*) = 0, \mu_i^* k_i(x^*) = 0$  sofort

$$\begin{aligned} L(x^*, \lambda^*, \mu^*) &= f(x^*) - \sum_{i=1}^l \lambda_i^* g_i(x^*) - \sum_{i=1}^q \mu_i^* k_i(x^*) = f(x^*) \\ &\geq f(x^*) - \sum_{i=1}^l \lambda_i \underbrace{g_i(x^*)}_{=0} - \sum_{i=1}^q \underbrace{\mu_i}_{\geq 0} \underbrace{k_i(x^*)}_{\geq 0} = L(x^*, \lambda, \mu) \quad \forall \lambda, \mu \text{ mit } \mu \geq 0. \end{aligned}$$

□

**Beweis von Lemma 3.3:** Sei  $(x^*, \lambda^*, \mu^*)$  KKT-Punkt, so ist nach Lemma 3.5  $(x^*, \lambda^*, \mu^*)$  Sattelpunkt, d.h.

$$L(x^*, \lambda^*, \mu^*) \leq L(x, \lambda^*, \mu^*) \quad \forall x \in \mathbb{R}^N.$$

Setze jetzt  $\bar{\alpha} := \|(\lambda^*, \mu^*)\|_\infty = \max\{|\lambda_1^*|, \dots, |\lambda_l^*|, \mu_1, \dots, \mu_q\}$  und wähle  $\alpha \geq \bar{\alpha}$ . Es gilt

$$\begin{aligned} P_1(x^*, \alpha) &= f(x^*) + \alpha \underbrace{\left( \sum_{j=1}^l \overbrace{|g_j(x^*)|}^{=0} - \sum_{j=1}^q \overbrace{\min(0, k_j(x^*))}^{=0} \right)}_{=r(x^*)} \\ &= f(x^*) \\ &= f(x^*) - \sum_{i=1}^l \lambda_i^* g_i(x^*) - \sum_{i=1}^q \mu_i^* k_i(x^*) \\ &= L(x^*, \lambda^*, \mu^*) \\ &\leq L(x, \lambda^*, \mu^*) \\ &= f(x) - \sum_{i=1}^l \lambda_i^* g_i(x) - \sum_{i=1}^q \mu_i^* k_i(x) \\ &\leq f(x) + \sum_{i=1}^l |\lambda_i^*| |g_i(x)| - \sum_{i=1}^q \mu_i^* \min(0, k_i(x)) \\ &\leq f(x) + \bar{\alpha} \underbrace{\left( \sum_{i=1}^l |g_i(x)| - \sum_{i=1}^q \min(0, k_i(x)) \right)}_{=r(x)} \\ &= f(x) + \bar{\alpha} r(x) \leq f(x) + \alpha r(x) = P_1(x, \alpha) \quad \forall x \in \mathbb{R}^N. \end{aligned}$$

□

### 3.2 Globalisierung von SQP-Verfahren

Es wird sich herausstellen, dass eine Lösung  $\Delta x^k$  eines bei  $x^k$  linearisierten quadratischen Teilproblems unter gewissen Voraussetzungen eine Abstiegsrichtung der  $l_1$ -Penalty-Funktion

$$P_1(x, \alpha) = f(x) + \alpha \left( \sum_{j=1}^l |g_j(x)| - \sum_{i=1}^q \min(0, k_i(x)) \right)$$

ist, so dass sich ein SQP-Verfahren mit einer geeigneten Schrittweitenstrategie für  $P_1(\cdot, \alpha)$  globalisieren lässt. Allerdings ist  $P_1(\cdot, \alpha)$  nicht überall differenzierbar. Wir berechnen stattdessen die Richtungsableitung im Punkt  $x$  in Richtung  $d$  gemäß

$$h'(x; d) = \lim_{t \rightarrow 0} \frac{h(x + td) - h(x)}{t}, \quad h : \mathbb{R}^N \rightarrow \mathbb{R}.$$

#### Lemma 3.6

a.) Sei  $N = 1$ ,  $h(x) = |x|$ , so gilt

$$h'(x; d) = \begin{cases} d & \text{für } x > 0 \\ |d| & \text{für } x = 0 \\ -d & \text{für } x < 0 \end{cases}.$$

b.) Sei  $N = 1$ ,  $h(x) = \min(0, x)$ , so folgt

$$h'(x; d) = \begin{cases} 0 & \text{für } x > 0 \\ \min(0, d) & \text{für } x = 0 \\ d & \text{für } x < 0 \end{cases}.$$

#### Beweis:

a.) Ist  $x \neq 0$ , so ist  $h$  differenzierbar in  $x$  und daher folgt

$$h'(x; d) = h'(x)d = \begin{cases} d & \text{für } x > 0 \\ -d & \text{für } x < 0 \end{cases}.$$

Für  $x = 0$  gilt

$$h'(0; d) = \lim_{t \searrow 0} \frac{h(0 + td) - h(0)}{t} = \lim_{t \searrow 0} \frac{t|d|}{t} = |d|.$$

b.) Ist  $x \neq 0$ , so ist  $h$  differenzierbar in  $x$  mit

$$h'(x; d) = h'(x)d = \begin{cases} d & \text{für } x > 0 \\ 0 & \text{für } x < 0 \end{cases}.$$

Für  $x = 0$  erhält man

$$h'(0; d) = \lim_{t \searrow 0} \frac{h(0 + td) - h(0)}{t} = \lim_{t \searrow 0} \frac{\min(0, td)}{t} = \min(0, d).$$

#### Bemerkung 3.7

Ferner gilt für richtungsdifferenzierbare Funktionen die Kettenregel.



### Korollar 3.8

Die Richtungsableitung der  $l^1$ -Penalty-Funktion

$$P_1(x, \alpha) = f(x) + \alpha \left[ \sum_{j=1}^l |g_j(x)| - \sum_{i=1}^q \min(0, k_i(x)) \right]$$

im Punkt  $x$  in Richtung  $d$  ist gegeben durch

$$P'_1(x, \alpha; d) = \nabla f(x)^T d + \alpha \left[ \sum_{g_j(x) > 0} \nabla g_j(x)^T d + \sum_{g_j(x) = 0} |\nabla g_j(x)^T d| + \sum_{g_j(x) < 0} (-\nabla g_j(x)^T d) \right. \\ \left. - \sum_{k_j(x) < 0} \nabla k_j(x)^T d - \sum_{k_j(x) = 0} \min(0, \nabla k_j(x)^T d) - 0 \right].$$

**Beweis:** Anwendung von Lemma 3.6 und Bemerkung 3.7.

□

An dieser Stelle kehren wir zurück zum SQP-Verfahren und betrachten zu einem Punkt  $x^k$  und einer gegebenen symmetrisch und positiv definiten Matrix  $H_k$  das quadratische Teilproblem

$$(VII) = \begin{cases} \nabla f(x^k)^T \Delta x + \frac{1}{2} \Delta x^T H_k \Delta x \stackrel{!}{=} \min \\ \text{unter den Nebenbedingungen} \\ g(x^k) + Dg(x^k) \Delta x = 0 \\ k(x^k) + Dk(x^k) \Delta x \geq 0. \end{cases}$$

### Satz 3.9

Sei  $\Delta x^k \neq 0$  Lösung des quadratischen Teilproblems (VII) mit  $H_k \in \mathbb{R}^{N,N}$  symmetrisch und positiv definit, und es sei

$$\alpha \geq \max\{|\lambda_1^{k+1}|, \dots, |\lambda_l^{k+1}|, \mu_1^{k+1}, \dots, \mu_q^{k+1}\},$$

wobei  $\lambda^{k+1}, \mu^{k+1}$  die zu  $\Delta x^k$  gehörigen Lagrange-Multiplikatoren seien. Dann hat die  $l^1$ -Penalty-Funktion

$$P_1(x, \alpha) = f(x) + \alpha \left[ \sum_{j=1}^l |g_j(x)| - \sum_{j=1}^q \min(0, k_j(x)) \right]$$

die Richtungsableitung  $P'_1(x^k, d; \Delta x^k)$  mit

$$P'_1(x^k, \alpha; \Delta x^k) \leq -(\Delta x^k)^T H_k \Delta x^k < 0,$$

d.h.  $\Delta x^k$  ist eine Abstiegsrichtung für  $P_1(\cdot, \alpha)$  im Punkte  $x^k$ .

**Beweis:** Gemäß Korollar 3.8 erhalten wir

$$P'_1(x, \alpha; \Delta x^k) = \nabla f(x)^T \Delta x^k + \alpha \left[ \sum_{g_j(x^k) > 0} \nabla g_j(x^k)^T \Delta x^k + \sum_{g_j(x^k) = 0} |\nabla g_j(x^k)^T \Delta x^k| \right. \\ \left. + \sum_{g_j(x^k) < 0} (-\nabla g_j(x^k)^T \Delta x^k) - \sum_{k_j(x^k) < 0} \nabla k_j(x^k)^T \Delta x^k - \sum_{k_j(x^k) = 0} \min(0, \nabla k_j(x^k)^T \Delta x^k) \right]. \quad (63)$$

Ferner erfülle  $\Delta x^k, \lambda^{k+1}, \mu^{k+1}$  die zu (VII) gehörige KKT-Bedingung, d.h. es gilt

$$\nabla f(x^k) + H_k \Delta x^k - \sum_{j=1}^l \lambda_j^{k+1} \nabla g_j(x^k) - \sum_{j=1}^q \mu_j^{k+1} \nabla k_j(x^k) = 0$$

und somit folgt

$$\nabla f(x^k)^T \Delta x^k = -\Delta x^k H_k \Delta x^k + \sum_{j=1}^l \lambda_j^{k+1} (\Delta x^k)^T \nabla g_j(x^k) + \sum_{j=1}^q \mu_j^{k+1} (\Delta x^k)^T \nabla k_j(x^k). \quad (64)$$

Einsetzen von (64) in (63) liefert

$$\begin{aligned} P_1'(x^k, \alpha; \Delta x^k) &= -(\Delta x^k)^T H_k \Delta x^k + \sum_{j=1}^l \lambda_j^{k+1} (\Delta x^k)^T \nabla g_j(x^k) + \sum_{j=1}^q \mu_j^{k+1} (\Delta x^k)^T \nabla k_j(x^k) \\ &+ \alpha \left[ \sum_{g_j(x^k) > 0} \nabla g_j(x^k)^T \Delta x^k + \sum_{g_j(x^k) = 0} |\nabla g_j(x^k)^T \Delta x^k| + \sum_{g_j(x^k) < 0} (-\nabla g_j(x^k)^T \Delta x^k) \right. \\ &\quad \left. - \sum_{k_j(x^k) < 0} \nabla k_j(x^k)^T \Delta x^k - \sum_{k_j(x^k) = 0} \min(0, \nabla k_j(x^k)^T \Delta x^k) \right]. \end{aligned} \quad (65)$$

Aus der Zulässigkeit von  $\Delta x^k$  für (VII) folgt

$$g_j(x^k) + \nabla g_j(x^k)^T \Delta x^k = 0 \iff g_j(x^k) = -\nabla g_j(x^k)^T \Delta x^k, \quad j = 1, \dots, l$$

sowie

$$k_j(x^k) + \nabla k_j(x^k)^T \Delta x^k \geq 0 \iff k_j(x^k) \geq -\nabla k_j(x^k)^T \Delta x^k, \quad j = 1, \dots, q.$$

Wir erhalten damit

$$\sum_{g_j(x^k) < 0} -\nabla g_j(x^k)^T \Delta x^k = \sum_{g_j(x^k) < 0} g_j(x^k) < 0 \quad (66)$$

$$\sum_{g_j(x^k) > 0} \nabla g_j(x^k)^T \Delta x^k = \sum_{g_j(x^k) > 0} -g_j(x^k) \quad (67)$$

$$\sum_{g_j(x^k) = 0} |\nabla g_j(x^k)^T \Delta x^k| = \sum_{g_j(x^k) = 0} |-g_j(x^k)| = 0 \quad (68)$$

$$- \sum_{k_j(x^k) = 0} \min(0, \nabla k_j(x^k)^T \Delta x^k) = \sum_{k_j(x^k) = 0} \min(0, -\nabla k_j(x^k)^T \Delta x^k) \leq \sum_{k_j(x^k) = 0} \min(0, k_j(x^k)) = 0 \quad (69)$$

$$- \sum_{k_j(x^k) < 0} \nabla k_j(x^k)^T \Delta x^k = \sum_{k_j(x^k) < 0} -\nabla k_j(x^k)^T \Delta x^k \leq \sum_{k_j(x^k) < 0} k_j(x^k) \leq 0 \quad (70)$$

$$\sum_{j=1}^l \lambda_j^{k+1} (\Delta x^k)^T \nabla g_j(x^k) = - \sum_{j=1}^l \lambda_j^{k+1} g_j(x^k) \quad (71)$$

$$\begin{aligned} \sum_{j=1}^q \mu_j^{k+1} (\Delta x^k)^T \nabla k_j(x^k) &= \underbrace{\sum_{j=1}^q \mu_j^{k+1} \left( (\Delta x^k)^T \nabla k_j(x^k) + k_j(x^k) \right)}_{=0 \text{ nach KKT-Bed.}} - \sum_{j=1}^q \mu_j^{k+1} k_j(x^k) \\ &= - \sum_{j=1}^q \mu_j^{k+1} k_j(x^k). \end{aligned} \quad (72)$$

Setze (66)-(72) in (65) ein und finde

$$\begin{aligned}
P_1'(x^k, \alpha; \Delta x^k) &\leq -(\Delta x)^T H_k \Delta x^k - \sum_{j=1, g_j(x^k) \neq 0}^l \lambda_j^{k+1} g_j(x^k) - \sum_{j=1, k_j(x^k) \neq 0}^q \mu_j^{k+1} k_j(x^k) \\
&\quad + \alpha \left[ \sum_{g_j(x^k) > 0} -g_j(x^k) + \sum_{g_j(x^k) < 0} g_j(x^k) + \sum_{k_j(x^k) < 0} k_j(x^k) \right] \\
&= (-\Delta x^k)^T H_k \Delta x^k + \sum_{g_j(x^k) > 0} (-\lambda_j^{k+1} - \alpha) g_j(x^k) + \sum_{g_j(x^k) < 0} (-\lambda_j^{k+1} + \alpha) g_j(x^k) \\
&\quad + \sum_{k_j(x^k) < 0} (-\mu_j^{k+1} + \alpha) k_j(x^k) - \sum_{k_j(x^k) > 0} \mu_j^{k+1} k_j(x^k) \\
&\leq -(\Delta x)^T H_k \Delta x^k, \quad \text{da } \alpha \geq \bar{\alpha} = \{|\lambda_1^{k+1}|, \dots, |\lambda_l^{k+1}|, \mu_1^{k+1}, \dots, \mu_q^{k+1}\}.
\end{aligned}$$

□

Aus Satz 3.9 folgt nun zusammen mit den Abstiegsstechniken der freien Minimierung, dass eine Armijo-artige Schrittweitensteuerung wohldefiniert ist, d.h. es existiert eine Schrittweite  $t_k = \beta^{l_k}$ ,  $\beta \in ]0, 1[$ ,  $l_k \in \mathbb{N}$  und  $\sigma \in ]0, 1[$  mit

$$P_1(x^k + t_k \Delta x^k, \alpha) \leq P_1(x^k, \alpha) + \sigma t_k P_1'(x^k, \alpha; \Delta x^k), \quad k \in \mathbb{N}.$$

Damit erhalten wir den Algorithmus:

Wähle  $(x^0, y^0, \mu^0)$ ,  $H_0 \in \mathbb{R}^{N,N}$  symmetrisch, positiv definit,  $\alpha > 0$ ,  $\sigma, \beta \in ]0, 1[$ , und setze  $k = 0$ .

**while**  $(x^k, \lambda^k, \mu^k)$  ist kein KKT-Punkt **do**

- berechne eine Lösung  $\Delta x \in \mathbb{R}^N$  des quadratischen Problems

$$\nabla f(x^k)^T \Delta x + \frac{1}{2} \Delta x^T H_k \Delta x \stackrel{!}{=} \min$$

unter den Nebenbedingungen

$$g(x^k) + Dg(x^k) \Delta x = 0$$

$$k(x^k) + Dk(x^k) \Delta x \geq 0$$

mit zugehörigen Multiplikatoren  $\lambda^{k+1}, \mu^{k+1}$ ;

- bestimme Schrittweite  $t_k = \max\{\beta^l \mid l = 0, 1, 2, \dots\}$  mit

$$P_1(x^k + t_k \Delta x^k, \alpha) \leq P_1(x^k, \alpha) + \sigma t_k P_1'(x^k, \alpha; \Delta x^k);$$

- setze  $x^{k+1} = x^k + t_k \Delta x^k$ , wähle  $H_{k+1}$  symmetrisch und positiv definit;
- setze  $k = k + 1$ ;

**end**

#### Globalisierte SQP-Methode

Für den obigen Algorithmus lassen sich gute globale Konvergenzresultate beweisen. Trotzdem bleiben noch einige Probleme für eine praktische Realisierung zu lösen, wie z.B. die Aufdatierung von  $\alpha$  und

die Wahl der Matrizen  $H_k$ . Der heikelste Punkt bleibt allerdings die Lösbarkeit der quadratischen Teilprobleme (VII). Für symmetrisch und positiv definite Matrizen  $H_k$  existiert eine Lösung, sofern der zulässige Bereich

$$\tilde{Z} = \{\Delta x \in \mathbb{R}^N \mid g(x^k) + Dg(x^k)\Delta x = 0, k(x^k) + Dk(x^k)\Delta x \geq 0\} \neq \emptyset$$

ist. Dies kann aber vorkommen, wie das folgende Beispiel zeigt: Vorgelegt sei

$$f(x) = x^2 \stackrel{!}{=} \min$$

unter der Nebenbedingung  $k(x) = x^2 - 1 \geq 0$  mit  $Z = \{x \in \mathbb{R} \mid |x| \leq 1\} \neq \emptyset$ .

Quadratisches Teilproblem bei  $x^k = 0$ :

$$2x\Delta x + \frac{1}{2}\Delta x^2 H_k \stackrel{!}{=} \min$$

unter der Nebenbedingung

$$k(x^k) + Dk(x^k)\Delta x = -1 + 2 \cdot 0\Delta x = -1$$

d.h.  $\tilde{Z} = \{\Delta x \mid k(0) + Dk(0)\Delta x \geq 0\} = \emptyset$ .

Im konvexen Fall sind die quadratischen Teilprobleme allerdings wohldefiniert.

**Lemma 3.10**

*Vorgelegt sei das Problem  $f(x) \stackrel{!}{=} \min$  unter den Nebenbedingungen  $g(x) = Ax - b = 0$ ,  $k_i(x) \geq 0$ ,  $i = 1, \dots, q$  mit  $f$  konvex und  $k_i$  konkav,  $i = 1, \dots, q$ . Dann besitzen auch die quadratischen Teilprobleme (VII) zulässige Punkte, falls das Ausgangsproblem zulässige Punkte besitzt.*

**Beweis:** Sei  $\tilde{x}$  ein zulässiger Punkt des Originalproblems. Setze  $\Delta x^k = \tilde{x} - x^k$ . Dann gilt

$$g(x^k) + Dg(x^k)\Delta x^k = g(x^k) + Dg(x^k)(\tilde{x} - x^k) = Ax^k - b + A(\tilde{x} - x^k) = A\tilde{x} - b = 0.$$

Da  $k_j$  konkav ist, folgt

$$k_j(x^k) + \nabla k_j(x^k)^T \Delta x^k = k_j(x^k) + \nabla k_j(x^k)(\tilde{x} - x^k) \geq k_j(\tilde{x}) \geq 0, \quad j = 1, \dots, q.$$

□

Im konvexen Fall ist damit die globalisierte SQP-Methode wohldefiniert.