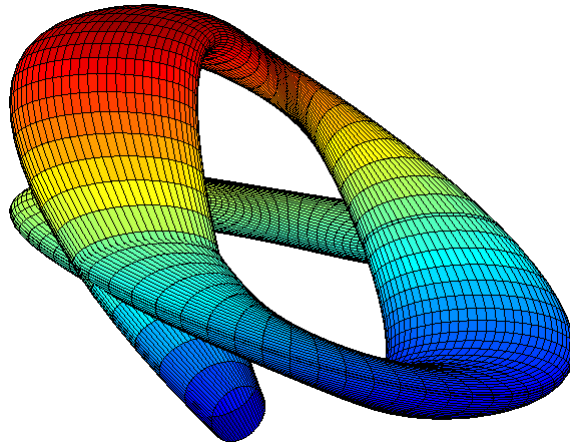


Skript zur Vorlesung

Optimierung



gelesen von

Prof. Dr. S. Volkwein

Konstanz, Sommersemester 2011

Inhaltsverzeichnis

1	Einleitung	3
2	Optimalitätskriterien	4
2.1	Allgemeiner Fall	4
2.2	Konvexe Funktionen	7
3	Allgemeine Abstiegsverfahren und Schrittweitenstrategien	9
3.1	Global konvergente Abstiegsverfahren	10
3.2	Schrittweitenstrategien und -algorithmen	11
3.2.1	Die Armijo-Regel	11
3.2.2	Die Wolfe-Powell-Regel	12
3.2.3	Strenge Wolfe-Powell-Regel	14
3.3	Praktische Aspekte	15
4	Gradientenverfahren	16
5	Das Newton-Verfahren	18
5.1	Das lokale Verfahren	18
5.2	Inexakte Newton-Verfahren	21
5.3	Globale Konvergenz	22
5.3.1	Trust-Region-Methoden	22
5.3.2	Globale Konvergenz des Trust-Region-Verfahrens	24
6	Quasi-Newton-Verfahren	26
	Index	29
	Literatur	30

1 Einleitung

Unter einem *endlichdimensionalen Optimierungsproblem* verstehen wir folgende Aufgabe:

$$\begin{cases} \text{Gegeben seien eine Menge } X \subseteq \mathbb{R}^n \text{ und eine stetige Funktion } f : X \rightarrow \mathbb{R}. \\ \text{Gesucht wird ein } x^* \in X \text{ mit } \forall x \in X : f(x^*) \leq f(x) \end{cases} . \quad (1.1)$$

In Kurznotation:

$$\min f(x) \text{ u.d.N. } x \in X \quad \text{bzw.} \quad \min_{x \in X} f(x). \quad (1.2)$$

Ist $X = \mathbb{R}^n$, so heißt (1.1) bzw. (1.2) *unrestringiert*, andernfalls *restringiert*.

Im Allgemeinen nennt man X den *Zulässigkeitsbereich* und f die *Zielfunktion*.

BEMERKUNG 1.1

Soll f maximiert werden für $x \in X$, so ist dies gleichbedeutend damit, dass $-f$ minimiert wird u.d.N. $x \in X$. \blacklozenge

Die Aufgabenstellung (1.1) erhält ihre Bedeutung dadurch, dass sie ein mathematisches Modell für viele Probleme zum Beispiel in der Physik, Medizin, Ökonomie und den Ingenieurwissenschaften ist.

Für den Fall, dass $X \neq \mathbb{R}^n$, lässt sich der Zulässigkeitsbereich sehr häufig in der Form

$$X = \{x \in \mathbb{R}^n \mid c_i(x) = 0, i \in I_1\} \cap \{x \in \mathbb{R}^n \mid c_i(x) \leq 0, i \in I_2\} \cap \{x \in \mathbb{R}^n \mid x_i \in \mathbb{Z}, i \in I_3\} =: \Omega_1 \cap \Omega_2 \cap \Omega_3$$

schreiben für gewisse Indexmengen I_1, I_2, I_3 und Abbildungen $c_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i \in I_1, I_2$. Die Mengen $\Omega_1, \Omega_2, \Omega_3$ werden als *Gleichungs-, Ungleichungs-* bzw. *Ganzzahligkeitsrestriktionen* bezeichnet.

Ist X eine Menge von diskreten Punkten, so spricht man von einem *diskreten* oder *kombinatorischen* Optimierungsproblem, andernfalls von einem *stetigen* Optimierungsproblem.

Ist f nicht differenzierbar, so spricht man von einem *nicht-differenzierbaren* Optimierungsproblem.

DEFINITION 1.2

Sei $f : X \rightarrow \mathbb{R}$ mit $X \subseteq \mathbb{R}^n$. Ein Punkt $x^* \in X$ heißt ...

- (1) ... *globale Minimalstelle* von f (auf X), wenn $f(x^*) \leq f(x)$ für alle $x \in X$. $f(x^*)$ heißt dann *globales Minimum*.
- (2) ... *strikte globale Minimalstelle* von f (auf X), wenn $f(x^*) < f(x)$ für alle $x \in X$. $f(x^*)$ heißt dann *striktes globales Minimum*.
- (3) ... *lokale Minimalstelle* von f (auf X), wenn es eine Umgebung U von x^* gibt, so dass $f(x^*) \leq f(x)$ für alle $x \in X \cap U$; $f(x^*)$ heißt dann *lokales Minimum*;
- (4) ... *strikte lokale Minimalstelle* von f (auf X), wenn es eine Umgebung U von x^* gibt, so dass $f(x^*) < f(x)$ für alle $x \in U \cap X$. $f(x^*)$ heißt dann *striktes lokales Minimum*.

BEMERKUNG 1.3

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

Im Folgenden bezeichne $\nabla f(x) := \left(\frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right)^T$ den Gradienten von f in x .

DEFINITION 1.4

Seien $X \subseteq \mathbb{R}^n$ eine offene Menge und $f : X \rightarrow \mathbb{R}$ eine stetig differenzierbare Funktion.

Ein Punkt $x^* \in X$ heißt *stationärer Punkt* von f , wenn $\nabla f(x^*) = 0$ gilt.

2 Optimalitätskriterien

2.1 Allgemeiner Fall

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

SATZ 2.1

Seien $X \subseteq \mathbb{R}^n$ offen und $f : X \rightarrow \mathbb{R}$ stetig differenzierbar.

Ist $x^* \in X$ eine lokale Minimalstelle von f (auf X), so gilt

$$\nabla f(x^*) = 0, \quad (2.1)$$

d.h. x^* ist ein stationärer Punkt.

BEWEIS (Analysis II)

Sei $x^* \in X$ lokale Minimalstelle von f , aber $\nabla f(x^*) \neq 0$. Dann existiert $d \in \mathbb{R}^n$ mit $\nabla f(x^*)^T d < 0$ (z.B. $d := -\nabla f(x^*)$). Da nach Voraussetzung f stetig differenzierbar ist, existiert die Richtungsableitung $f'(x^*; d)$ von f in x^* in Richtung d . Es gilt

$$f'(x^*; d) = \lim_{t \searrow 0} \frac{f(x^* + td) - f(x^*)}{t} = \nabla f(x^*)^T d < 0.$$

Folglich existiert $\bar{t} > 0$ mit $x^* + td \in X$ und $\frac{f(x^* + td) - f(x^*)}{t} < 0$ für alle $t \in (0, \bar{t}]$. Somit ist auch $f(x^* + td) < f(x^*)$ für alle $t \in (0, \bar{t}]$, was einen Widerspruch zur Voraussetzung ergibt. ■

BEMERKUNG 2.2

- (a) Da Satz 2.1 nur Ableitungen bis zur ersten Ordnung verwendet, gibt er eine **notwendige Bedingung erster Ordnung** an.
- (b) Die Bedingung $\nabla f(x^*) = 0$ ist nicht hinreichend dafür, dass x^* eine lokale Minimalstelle von f (auf X) ist (betrachte z.B. $f(x) = x^3$). ◆

Wir zitieren folgendes Resultat:

LEMMA 2.3

Sei \mathcal{S}_n der Vektorraum der symmetrischen $(n \times n)$ -Matrizen. Für $A \in \mathcal{S}_n$ sei $\lambda(A)$ der kleinste Eigenwert von A .

Dann gilt

$$|\lambda(A) - \lambda(B)| \leq \|A - B\|_2 \quad \text{für alle } A, B \in \mathcal{S}_n,$$

wobei $\|\cdot\|_2$ hier die Spektralnorm symmetrischer Matrizen bezeichnet, d.h.

$$\|A\|_2 := \max\{|\lambda| \mid \lambda \text{ ist Eigenwert von } A\}.$$

Mit Hilfe von Lemma 2.3 folgt aus der Stetigkeit der Hessematrix $\nabla^2 f$ von f , dass $\nabla^2 f(x)$ positiv definit ist in einer Umgebung von x^* , falls $\nabla^2 f(x^*)$ positiv definit ist.

Eine analoge Folgerung gilt für den Fall, dass $\nabla^2 f(x^*)$ negativ definit ist.

SATZ 2.4

Seien $X \subseteq \mathbb{R}^n$ offen und $f : X \rightarrow \mathbb{R}$ zweimal stetig differenzierbar.

Ist $x^* \in X$ eine lokale Minimalstelle von f auf X , so sind $\nabla f(x^*) = 0$ und $\nabla^2 f(x^*)$ positiv semidefinit.

BEWEIS

Die Bedingung $\nabla f(x^*) = 0$ folgt aus Satz 2.1. Sei x^* eine lokale Minimalstelle von f , jedoch $\nabla^2 f(x^*)$ nicht positiv semidefinit. Dann existiert ein $d \in \mathbb{R}^n$ mit

$$d^T \nabla^2 f(x^*) d < 0.$$

Mit Hilfe des Satzes von Taylor ergibt sich

$$f(x^* + td) = f(x^*) + \underbrace{t \nabla f(x^*)^T d}_{=0} + \frac{1}{2} t^2 d^T \nabla^2 f(\xi_t) d$$

für kleines $t > 0$. Dabei ist $\xi_t = x^* + \vartheta_t td$ für $\vartheta_t \in (0, 1)$.

Aus Lemma 2.3 und der Stetigkeit der zweiten Ableitung von f folgt die Existenz von $\bar{t} > 0$ mit

$$d^T \nabla^2 f(\xi_t) d < 0 \quad \text{für alle } t \in (0, \bar{t}];$$

wegen $\nabla f(x^*) = 0$ also

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

was einen Widerspruch zur Voraussetzung ergibt. ■

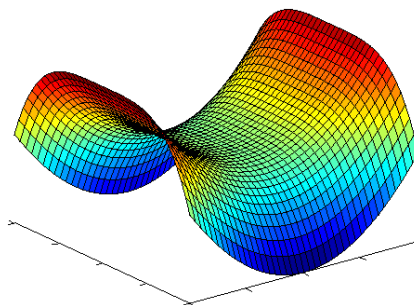
BEMERKUNG 2.5

Die Bedingungen aus Satz 2.4 (und Satz 2.1) sind nicht hinreichend dafür, dass x^* eine lokale Minimalstelle ist.

Betrachte z.B. $f(x) = x_1^2 - x_2^4$ mit $x^* = (0, 0)$, dann

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

Da Satz 2.4 Ableitungen bis zur zweiten Ordnung verwendet, gibt er **notwendige Bedingungen zweiter Ordnung** an.♦



Nun kommen wir zu hinreichenden Bedingungen.

SATZ 2.6

Seien $X \subseteq \mathbb{R}^n$ offen und $f : X \rightarrow \mathbb{R}$ zweimal stetig differenzierbar.

Gelten

- (a) $\nabla f(x^*) = 0$ und
- (b) $\nabla^2 f(x^*)$ ist positiv definit,

dann ist x^* eine strikte lokale Minimalstelle von f auf X .

BEWEIS

Aus (b) folgt die Existenz einer Konstanten $\mu > 0$ mit

$$d^T \nabla^2 f(x^*) d \geq \mu d^T d \quad \text{für alle } d \in \mathbb{R}^n$$

(z.B. $\mu := \min\{\lambda \mid \lambda \text{ ist Eigenwert von } \nabla^2 f(x^*)\}$).

Nach dem Satz von Taylor gilt für alle $d \in \mathbb{R}^n$, die hinreichend nahe bei 0 sind, dass

$$\begin{aligned} f(x^* + d) &= f(x^*) + \nabla f(x^*)^T d + \frac{1}{2} d^T \nabla^2 f(\xi_d) d \\ &\stackrel{(a)}{=} f(x^*) + \frac{1}{2} d^T \nabla^2 f(\xi_d) d \end{aligned}$$

mit $\xi_d := x^* + \vartheta_d d$ für $\vartheta_d \in (0, 1)$.

Man erhält so unter Verwendung der Cauchy-Schwarzschen Ungleichung

$$\begin{aligned} f(x^* + d) &= f(x^*) + \frac{1}{2}d^T \nabla^2 f(x^*)d + \frac{1}{2}d^T (\nabla^2 f(\xi_d) - \nabla^2 f(x^*))d \\ &\geq f(x^*) + \frac{1}{2}(\mu - \|\nabla^2 f(\xi_d) - \nabla^2 f(x^*)\|_2) \|d\|_2^2 \\ &> f(x^*), \end{aligned}$$

falls $\|\nabla^2 f(\xi_d) - \nabla^2 f(x^*)\|_2$ klein ($< \mu$) und $\|d\|_2$ klein ($d \neq 0$).

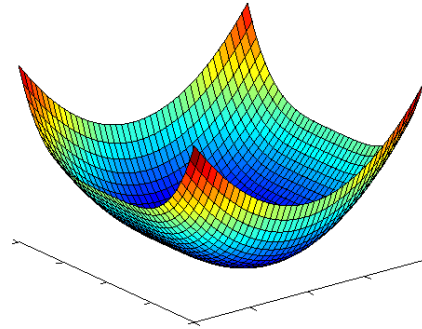
Damit ist x^* strikte lokale Minimalstelle. ■

BEMERKUNG 2.7

Die Bedingungen aus Satz 2.6 sind nicht notwendig dafür, dass x^* strikte lokale Minimalstelle von f (auf X) ist.

Betrachte z.B. $f(x) = x_1^2 + x_2^4$ mit $x^* = (0, 0)$.

Ist $\nabla^2 f(x^*)$ indefinit, so spricht man von einem *Sattelpunkt*. ♦



2.2 Konvexe Funktionen

DEFINITION 2.8

(1) Eine Menge $X \subseteq \mathbb{R}^n$ heißt **konvex**, wenn für alle $x, y \in X$ und alle $\lambda \in (0, 1)$ auch $\lambda x + (1 - \lambda)y$ in X liegt.

(2) Sei $X \subseteq \mathbb{R}^n$ eine konvexe Menge. Eine Funktion $f : X \rightarrow \mathbb{R}$ heißt ...

(a) ... **strikt konvex** bzw. **konvex**, wenn für alle $x, y \in X$, $x \neq y$ und alle $\lambda \in (0, 1)$ gilt

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y) \quad \text{bzw.} \quad f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

(b) ... **gleichmäßig konvex**, falls es ein $\mu > 0$ gibt, so dass für alle $x, y \in X$, $\lambda \in (0, 1)$ gilt:

$$f(\lambda x + (1 - \lambda)y) + \mu\lambda(1 - \lambda)\|x - y\|^2 \leq \lambda f(x) + (1 - \lambda)f(y).$$

Man sagt dazu auch: f ist gleichmäßig konvex mit **Modul** $\mu > 0$.

Aus der Definition folgt, dass jede gleichmäßig konvexe Funktion auch strikt konvex ist und jede strikt konvexe Funktion konvex ist. Die Umkehrungen gelten i.A. nicht.

BEMERKUNG 2.9

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ sei eine quadratische Funktion, d.h.

$$f(x) = \frac{1}{2}x^T Qx + c^T x + \gamma$$

mit $Q \in \mathcal{S}_n$, $c \in \mathbb{R}^n$, $\gamma \in \mathbb{R}$. Dann gelten:

(a) f ist konvex $\Leftrightarrow Q$ ist positiv semidefinit.

(b) f ist strikt konvex $\Leftrightarrow f$ ist gleichmäßig konvex $\Leftrightarrow Q$ ist positiv definit. \blacklozenge

Ohne Beweis geben wir folgende Charakterisierungen zweimal stetig differenzierbarer, strikt gleichmäßig konvexer Funktionen an.

SATZ 2.10

Seien $X \subseteq \mathbb{R}^n$ eine offene und konvexe Menge sowie $f : \mathbb{R}^n \rightarrow \mathbb{R}$ zweimal stetig differenzierbar.

Dann gelten:

(a) f ist konvex auf $X \Leftrightarrow \nabla^2 f(x)$ ist für alle $x \in X$ positiv semidefinit.

(b) Ist $\nabla^2 f(x)$ für alle $x \in X$ positiv definit, so ist f strikt konvex (auf X).

(c) f ist genau dann gleichmäßig konvex (auf X), wenn $\nabla^2 f(x)$ gleichmäßig positiv definit ist, d.h. wenn es ein $\mu > 0$ gibt, so dass für alle $x \in X$ und alle $d \in \mathbb{R}^n$ gilt:

$$d^T \nabla^2 f(x) d \geq \mu \|d\|^2.$$

BEMERKUNG 2.11

Aussage (b) von Satz 2.10 lässt sich i.A. nicht umkehren. Betrachte z.B. $f(x) = x^4$ mit $X = \mathbb{R}$. \blacklozenge

Der folgende Hilfssatz beschäftigt sich mit Niveaumengen gleichmäßig konvexer Funktionen.

LEMMA 2.12

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar, $x^0 \in \mathbb{R}^n$ beliebig gegeben, die Niveaumenge

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

konvex und f gleichmäßig konvex auf $\mathcal{L}(x^0)$. Dann ist $\mathcal{L}(x^0)$ kompakt.

SATZ 2.13

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar und $X \subseteq \mathbb{R}^n$ konvex. Betrachtet man das Optimierungsproblem

$$\min f(x) \text{ u.d.N. } x \in X, \quad (2.2)$$

dann gelten:

- (a) Ist f konvex auf X , so ist die Lösungsmenge von (2.2) konvex (evtl. leer).
- (b) Ist f strikt konvex auf X , so besitzt (2.2) höchstens eine Lösung.
- (c) Sind f gleichmäßig konvex auf X , $X \neq \emptyset$ und abgeschlossen, so besitzt (2.2) genau eine Lösung.

BEWEIS

- (a) Seien x^1, x^2 Lösungen von (2.2), also $f(x^1) = f(x^2) = \min_{x \in X} f(x)$.

Für $\lambda \in (0, 1)$ ist dann auch $\lambda x^1 + (1 - \lambda)x^2 \in X$. Weiter

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

Also nimmt f auch an $\lambda x^1 + (1 - \lambda)x^2$ sein Minimum an.

- (b) Angenommen, (2.2) hat zwei verschiedene Lösungen x^1, x^2 . Für $\lambda \in (0, 1)$ gelten $\lambda x^1 + (1 - \lambda)x^2 \in X$ und

$$f(\lambda x^1 + (1 - \lambda)x^2) < \lambda f(x^1) + (1 - \lambda)f(x^2) = \min_{x \in X} f(x),$$

was einen Widerspruch ergibt.

- (c) Sei $x^0 \in X$ beliebig gewählt. Dann ist wegen

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \leq f(x^0) \quad \text{für alle } x, y \in \mathcal{L}(x^0)$$

die Menge $\mathcal{L}(x^0)$ konvex und nach Lemma 2.11 kompakt. Dann ist $X \cap \mathcal{L}(x^0)$ kompakt und nichtleer. Daher nimmt die stetige Funktion f ihr globales Minimum in $X \cap \mathcal{L}(x^0)$ an, welches natürlich auch ein Minimum von (2.2) ist. ■

BEMERKUNG 2.14

- (a) Das Problem (2.2) muss selbst für strikt konvexes f keine Lösung besitzen. Betrachte dazu z.B. $f(x) = \exp(x)$ auf $X = \mathbb{R}$.
- (b) Auf die Forderung nach Abgeschlossenheit von X in Satz 2.12 (c) kann nicht verzichtet werden, betrachte z.B. $f(x) = x^2$ für $x \in (0, 1]$. ◆

Das zentrale Resultat dieses Abschnitts wird in Satz 2.15 angegeben. Man kann daraus sehen, dass $\nabla f(x^*) = 0$ auch hinreichend ist dafür, dass x^* globale Minimalstelle ist.

SATZ 2.15

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar und konvex und $x^* \in \mathbb{R}^n$ ein stationärer Punkt von f .

Dann ist x^* globale Minimalstelle von f auf \mathbb{R}^n .

BEWEIS

Mit Taylor folgt

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

mit $\xi := x^* + \vartheta(x - x^*)$, $\vartheta \in (0, 1)$. Nach Satz 2.9 (a) ist $\nabla^2 f(\xi)$ positiv semidefinit. Ferner gilt $\nabla f(x^*) = 0$, daher $f(x) \leq f(x^*)$, was zu zeigen war. ■

3 Allgemeine Abstiegsverfahren und Schrittweitenstrategien

Wir betrachten ein *Abstiegsverfahren* zur Lösung von $\min_{x \in \mathbb{R}^n} f(x)$ ($f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar).

Die zentrale Idee der Verfahren in diesem Abschnitt ist wie folgt:

- Ist man in einem Punkt $x \in \mathbb{R}^n$, so sucht man eine Richtung $d \in \mathbb{R}^n$ aus, in welcher der Funktionswert fällt („Abstiegsverfahren“).
- Entlang dieser Richtung d geht man so lange, bis man den Funktionswert von f hinreichend verkleinert hat („Schrittweitenstrategie“).

Diese Schritte wollen wir formalisieren.

DEFINITION 3.1

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ und $x \in \mathbb{R}^n$. Ein Vektor $d \in \mathbb{R}^n$ heißt *Abstiegsrichtung* von f im Punkt x , wenn es ein $\bar{t} > 0$ gibt mit

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

Ist f stetig differenzierbar, dann ist

$$\nabla f(x)^T d < 0 \tag{3.1}$$

hinreichend dafür, dass $d \in \mathbb{R}^n$ eine Abstiegsrichtung von f in x ist.

Um dies einzusehen, definieren wir $\varphi(t) := f(x + td)$. Aus $f \in \mathcal{C}^1$ folgt

$$\varphi(t) = \varphi(0) + t\varphi'(0) + r(t), \tag{3.2}$$

wobei $\frac{r(t)}{t} \rightarrow 0$ für $t \searrow 0$ ($r(t) = o(t)$). Es gelten

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

Umformen von (3.2) und Division durch $t > 0$ liefert

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

Aus $r(t) = o(t)$ und (3.1) folgt, dass ein $\bar{t} > 0$ existiert mit

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

d.h. $f(x + td) < f(x)$ und d ist Abstiegsrichtung von f in x .

BEMERKUNG 3.2

- Die Bedingung (3.1) bedeutet, dass der Winkel zwischen d und dem negativen Gradienten von f in x kleiner als 90° ist.
- Das Kriterium (3.1) ist nicht notwendig. Ist x z.B. eine strikte lokale Maximalstelle, so sind alle $d \in \mathbb{R}^n$ Abstiegsrichtungen, aber (3.1) gilt nicht. \blacklozenge

BEISPIEL 3.3

Mögliche Kandidaten für d sind ...

- ... $d = -\nabla f(x)$, die Richtung des steilsten Abstiegs:

$$\nabla f(x)^T d = -\|\nabla f(x)\|^2;$$

- ... $d = -M\nabla f(x)$, $M \in \mathcal{S}_n$ positiv definit („gradientenähnliche Verfahren“):

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

3.1 Global konvergente Abstiegsverfahren

Wir wollen ein allgemeines Abstiegsverfahren angeben:

ALGORITHMUS 3.4 (Allgemeines Abstiegsverfahren)

Input: $f : \mathbb{R}^n \rightarrow \mathbb{R}$, Startpunkt $x^0 \in \mathbb{R}^n$

Begin

$k := 1$

While Konvergenzkriterium nicht erfüllt Do

bestimme Abstiegsrichtung d^k von f in x^k ;

bestimme eine 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(While)

End

In theoretischen Konvergenzuntersuchungen betrachten wir kein Konvergenzkriterium, d.h. wir nehmen an, dass eine unendliche Folge $(x^k)_{k \in \mathbb{N}}$ erzeugt wird.

SATZ 3.5

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar und $(x^k)_{k \in \mathbb{N}}$ eine durch den Algorithmus 3.4 erzeugte Folge so, dass

(a) ... eine Konstante $\theta_1 > 0$ existiert mit

$$-\nabla f(x^k)^T d^k \geq \theta_1 \|\nabla f(x^k)\| \|d^k\| \quad \text{für alle } k \in \mathbb{N}; \quad \text{(Winkelbedingung)}$$

(b) ... eine von $(x^k)_{k \in \mathbb{N}}$ und $(d^k)_{k \in \mathbb{N}}$ unabhängige Konstante $\theta_2 > 0$ existiert mit

$$f(x^k + t_k d^k) \leq f(x^k) - \theta_2 \left(\frac{\nabla f(x^k)^T d^k}{\|d^k\|} \right)^2 \quad \text{mit } t_k > 0 \text{ für alle } k \in \mathbb{N}.$$

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

BEWEIS

Da jedes t_k die Bedingung (b) erfüllt, folgt mit (a), dass

$$f(x^{k+1}) - f(x^k) \leq -\theta_2 \left(\frac{\nabla f(x^k)^T d^k}{\|d^k\|} \right)^2 \leq -\theta_1^2 \theta_2 \|\nabla f(x^k)\|^2 \leq 0. \quad (3.3)$$

Sei nun x^* ein Häufungspunkt von $(x^k)_{k \in \mathbb{N}}$. Da $(f(x^k))_{k \in \mathbb{N}}$ wegen (3.3) monoton fällt und auf einer Teilfolge $(x^{k_\nu})_{\nu \in \mathbb{N}}$ mit $x^{k_\nu} \xrightarrow{\nu \rightarrow \infty} x^*$ gegen $f(x^*)$ konvergiert, folgt daraus, dass die gesamte Folge der Funktionswerte gegen $f(x^*)$ konvergiert.

Insbesondere: $f(x^{k+1}) - f(x^k) \xrightarrow{k \rightarrow \infty} 0$ und mit (3.3) folgt $\lim_{k \rightarrow \infty} \|\nabla f(x^k)\| = 0$.

Damit ist jeder Häufungspunkt stationärer Punkt:

$$\|\nabla f(x^*)\| = \lim_{\nu \rightarrow \infty} \|\nabla f(x^{k_\nu})\| = 0. \quad \blacksquare$$

BEMERKUNG 3.6

Bezeichne η_k den Winkel zwischen d^k und $-\nabla f(x^k)$, dann bedeutet die Forderung (a) aus Satz 3.5, dass

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

gleichmäßig größer als 0 ist. Ein wichtiges Beispiel ist $d^k := -\nabla f(x^k)$. \blacklozenge

3.2 Schrittweitenstrategien und -algorithmen

Das allgemeine Abstiegsverfahren (Algorithmus 3.4) besitzt in der Wahl der Abstiegsrichtung d^k und der Schrittweite $t_k > 0$ große Freiheitsgrade.

Die nahe liegende Minimierungsregel $t_k := t_k^{\min}$ mit

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

ist unter der Annahme, dass $\mathcal{L}(x^0)$ kompakt ist und ∇f Lipschitz-stetig ist auf $\mathcal{L}(x^0)$, wohldefiniert. Allerdings ist diese Regel i.A. nicht praktikabel (Aufwand!).

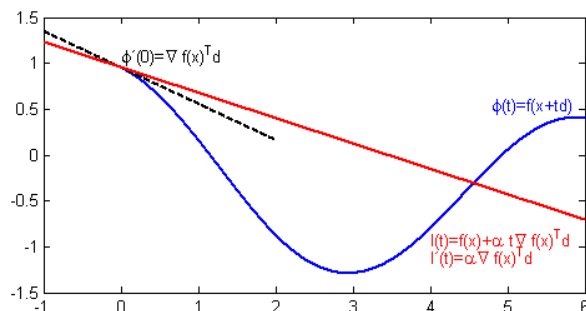
3.2.1 Die Armijo-Regel

Wir behandeln *gradientenähnliche Richtungen* d.h.

$$d := -M \nabla f(x), \quad M \in \mathcal{S}_n \text{ positiv definit.}$$

Sei $\alpha \in [0, 1]$ fest vorgegeben. Die *Armijo-Regel* ist eine Bedingung, die einen hinreichenden Abstieg sichert, und lautet

$$f(x + t d) \leq f(x) + \alpha t \nabla f(x)^T d \quad (3.4)$$



Zur tatsächlichen Berechnung von t überprüft man (3.4) sequenziell z.B. für

$$t = \beta^l, \quad l = 0, 1, 2, \dots \quad (3.5)$$

mit $\beta \in (0, 1)$ fest vorgegeben, z.B. $\beta = \frac{1}{2}$. Bei erstmaliger Gültigkeit von (3.4) bricht man ab. Die Größe t nennt man *Schrittweite*.

Im Folgenden ist (3.5) verallgemeinert, d.h. falls (3.4) für ein $t = t_c$ nicht erfüllt ist, dann wird t_+ so konstruiert, dass

$$t_+ \in [\underline{\nu} t_c, \bar{\nu} t_c] \quad \text{mit } 0 < \underline{\nu} \leq \bar{\nu} < 1 \text{ fest.}$$

ALGORITHMUS 3.7 (Armijo-Schrittweitenalgorithmus)

Input: Abstiegsrichtung d

Begin

$l := 0; t^{(0)} := 1;$

While (3.4) ist nicht erfüllt Do

bestimme $t^{(l+1)} \in [\underline{\nu} t^{(l)}, \bar{\nu} t^{(l)}];$

setze $l := l + 1;$

End(While)

$t_k := t^{(l)};$

End

Im folgenden Satz ist eine Konvergenzaussage für Algorithmus 3.4 mit der Schrittweitenwahl gemäß Algorithmus 3.7 formuliert.

SATZ 3.8

Sei $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar mit ∇f Lipschitz-stetig. Seien $(x^k)_{k \in \mathbb{N}}$ die von Algorithmus 3.4 mit Algorithmus 3.7 erzeugte Iterationsfolge und $(M^k)_{k \in \mathbb{N}}$ eine Folge symmetrischer, positiv definiter Matrizen mit

$$0 < \underline{\lambda} \leq \lambda_s^k \leq \lambda_l^k \leq \bar{\lambda} < \infty \quad \text{für alle } k \in \mathbb{N},$$

wobei λ_s^k und λ_l^k den kleinsten bzw. größten Eigenwert von M^k bezeichnen. $L > 0$ bezeichne die Lipschitz-Konstante von ∇f .

Dann erfüllen die Schritte

$$s^k := x^{k+1} - x^k = t_k d^k = -t_k M^k \nabla f(x^k)$$

die Bedingung

$$t_k \geq \underline{t} = \frac{2\nu \underline{\lambda}(1 - \alpha)}{L\bar{\kappa}} \quad \text{mit } \bar{\kappa} := \frac{\bar{\lambda}}{\underline{\lambda}} \geq \kappa_2(M^k) \geq 1. \quad (\text{vgl. Kap. 4})$$

Ferner ist entweder $(f(x^k))_{k \in \mathbb{N}}$ nach unten unbeschränkt oder $\lim_{k \rightarrow \infty} \nabla f(x^k) = 0$.

Somit ist jeder Häufungspunkt von $(x^k)_{k \in \mathbb{N}}$ stationärer Punkt von f in \mathbb{R}^n .

Insbesondere: Sind $(f(x^k))_{k \in \mathbb{N}}$ nach unten beschränkt und existiert eine Teilfolge $(x^{k(l)})_{l \in \mathbb{N}}$ von $(x^k)_{k \in \mathbb{N}}$ mit $\lim_{l \rightarrow \infty} x^{k(l)} = x^*$, dann ist $\nabla f(x^*) = 0$.

BEMERKUNG 3.9

- (a) Im Allgemeinen gibt es keine Garantie, dass ein (eindeutiger) Häufungspunkt existiert.
 (b) Die folgende Variation der Armijo-Regel führt ebenfalls zu einem im Sinne von Satz 3.8 konvergenten Verfahren:

Seien $r > 0$ ein Skalierungsverfahren und $\beta \in (0, 1)$. Bestimme

$$t = \max\{r\beta^l\}, \quad l = 0, 1, 2, \dots \quad (3.6)$$

so, dass (3.4) erfüllt ist.

Die Bestimmung der Schrittweite t gemäß (3.5) oder (3.6) wird in der englischsprachigen Literatur „Backtracking“ genannt.

- (c) Weitere Strategien basieren auf Polynommodellen, die $\varphi(t) = f(x + td)$ durch ein quadratisches oder kubisches *Modell* ersetzen und dann dieses Modell minimieren. \blacklozenge

3.2.2 Die Wolfe-Powell-Regel

Seien $\alpha \in (0, \frac{1}{2})$ und $\rho \in [\alpha, 1]$ gegeben. Die *Wolfe-Powell-Regel* lautet:

Zu $x, d \in \mathbb{R}^n$ mit $\nabla f(x)^T d < 0$ bestimme eine Schrittweite $t > 0$ mit

$$f(x + td) \leq f(x) + \alpha t \nabla f(x)^T d \quad (\text{oder: } \varphi(t) \leq \varphi(0) + \alpha t \varphi'(0)) \quad (3.7a)$$

$$\nabla f(x + td)^T d \geq \rho \nabla f(x)^T d \quad (\text{oder: } \varphi'(t) \geq \rho \varphi'(0)) \quad (3.7b)$$

Graphik einfügen

Ohne Beweis geben wir den folgenden Satz an.

SATZ 3.10

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar und $\alpha \in (0, \frac{1}{2})$, $\rho \in [\alpha, 1)$, $x^0 \in \mathbb{R}^n$ fest vorgegeben.
Zu $x \in \mathcal{L}(x^0)$ und einer Richtung $d \in \mathbb{R}^n$ mit $\nabla f(x)^T d < 0$ sei

$$\mathcal{T}_{\text{WP}}(x, d) := \{t > 0 \mid (3.7) \text{ ist erfüllt}\}.$$

Dann gelten:

- (a) Ist f nach unten beschränkt, so ist $\mathcal{T}_{\text{WP}}(x, d) \neq \emptyset$, d.h. die Wolfe-Powell-Schrittweitenstrategie ist wohldefiniert.
(b) Ist weiter ∇f auf $\mathcal{L}(x^0)$ Lipschitz-stetig, dann existiert eine Konstante $\theta > 0$ (unabhängig von x und d) mit

$$f(x + td) \leq f(x) - \theta \left(\frac{\nabla f(x)^T d}{\|d\|} \right)^2 \quad \text{für alle } x \in \mathcal{T}_{\text{WP}}(x, d).$$

ALGORITHMUS 3.11 (Wolfe-Powell-Liniensuche)

Input: Abstiegsrichtung $d \in \mathbb{R}^n$

Begin

Wähle $t^{(0)} > 0$, $\gamma > 1$ (z.B. $\gamma = \frac{3}{2}$ oder $\gamma = 2$), $i := 0$

If $\varphi(t^{(i)}) \geq \varphi(0) + \alpha t^{(i)} \varphi'(0)$ (A.1)

$a := 0$; $b := t^{(i)}$; Goto (B.0)

Else

If $\varphi'(t^{(i)}) \geq \rho \varphi'(0)$

$t := t^{(i)}$; Return 1 (Ausgabe der Schrittweite t)

Else

$t^{(i+1)} := \gamma t^{(i)}$; $i := i + 1$; Goto (A.1)

End(If)

End(If)

Wähle $\tau_1, \tau_2 \in (0, \frac{1}{2}]$; $j := 0$; $t_1^{(0)} := a$; $t_2^{(0)} := b$; $\Delta^{(0)} := t_2^{(0)} - t_1^{(0)}$ (B.0)

Wähle $t^{(j)} \in [t_1^{(j)} + \tau_1 \Delta^{(j)}, t_2^{(j)} - \tau_2 \Delta^{(j)}]$ (B.1)

If $\varphi(t^{(j)}) \geq \varphi(0) + \alpha t^{(j)} \varphi'(0)$

$b_1^{(j+1)} := t_1^{(j)}$; $t_2^{(j+1)} := t^{(j)}$; $j := j + 1$; Goto (B.1)

Else

If $\varphi'(t^{(j)}) \geq \rho \varphi'(0)$

$t := t^{(j)}$; Return 2 (Ausgabe der Schrittweite t)

Else

$t_2^{(j+1)} := t_2^{(j)}$; $t_1^{(j+1)} := t^{(j)}$; $j := j + 1$; Goto (B.1)

End(If)

End(If)

End

Das folgende Lemma motiviert teilweise den Algorithmus 3.8.

LEMMA 3.12

Seien $\alpha < \rho$ und $\varphi'(0) < 0$. Ist $[a, b]$ mit $0 \leq a < b$ ein Intervall mit

$$\varphi(a) \leq \varphi(0) + \alpha a \varphi'(0); \quad \varphi(b) \geq \varphi(0) + \alpha b \varphi'(0) \quad \varphi'(a) < \rho \varphi'(0),$$

so enthält $[a, b]$ einen Punkt \bar{t} mit

$$\varphi(\bar{t}) < \varphi(0) + \alpha \bar{t} \varphi'(0); \quad \varphi'(\bar{t}) = \alpha \varphi'(0) > \rho \varphi'(0).$$

Der Punkt \bar{t} ist ein innerer Punkt eines Intervalls I , so dass für alle $t \in I$ gilt

$$\varphi(t) \leq \varphi(0) + \alpha t \varphi'(0); \quad \varphi'(t) \geq \rho \varphi'(0),$$

d.h. $I \subseteq \mathcal{T}_{\text{WP}}(x, d)$.

SATZ 3.13

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar und nach unten beschränkt und $\alpha \in (0, \frac{1}{2})$, $\rho \in (\frac{1}{2}, 1)$.

Dann bricht Algorithmus 3.11 nach endlich vielen Schritten entweder bei **Return 1** oder bei **Return 2** mit einer Schrittweite $t \in \mathcal{T}_{\text{WP}}(x, d)$ ab.

3.2.3 Strenge Wolfe-Powell-Regel

Seien $\alpha \in (0, \frac{1}{2})$ und $\rho \in [\frac{1}{2}, 1)$ fest vorgegeben. Die *Strenge Wolfe-Powell-Regel* lautet:

Zu $x, d \in \mathbb{R}^n$ mit $\nabla f(x)^T d < 0$ bestimme eine Schrittweite $t > 0$ mit

$$f(x + td) \leq f(x) + \alpha t \nabla f(x)^T d \tag{3.8a}$$

$$|\nabla f(x + td)^T d| \leq -\rho \nabla f(x)^T d \tag{3.8b}$$

Graphik einfügen

Der Graph von φ fällt also nicht zu steil ab, steigt aber auch nicht zu steil.

Für sehr kleines ρ (und damit auch kleines α) ist eine Schrittweite, die (3.8) erfüllt, nahe an einem stationären Punkt von φ .

3.3 Praktische Aspekte

Die Algorithmen in Abschnitt 3.2 sind idealisiert. In der Praxis sind f und ∇f maschinen- und/oder problemabhängig genau. Werden diese Ungenauigkeiten nicht berücksichtigt, fährt dies schnell zu Endlosschleifen.

Ideal wäre, wenn mit den Funktionswerten $\varphi(t), \varphi(0)$ und Ableitungen $\varphi'(t), \varphi'(0)$ Fehlerschranken mit geliefert würden: $\varepsilon(t), \varepsilon(0)$ und $\hat{\varepsilon}(t), \hat{\varepsilon}(0)$. Dann wird

$$\varphi(t) \leq \varphi(0) + \alpha t \varphi'(0)$$

ersetzt durch

$$\varphi(t) \leq \varphi(0) + \alpha t (\varphi'(0) + \hat{\varepsilon}(0)) + \varepsilon(0) + \varepsilon(t)$$

und

$$\varphi'(t) \geq \rho \varphi'(0)$$

wird ersetzt durch

$$\varphi'(t) \geq \rho(\varphi'(0) - \hat{\varepsilon}(0)) - \hat{\varepsilon}(t).$$

Weiter ist abzubrechen, wenn $[t_1^{(j)}, t_2^{(j)}]$ „zu klein“ wird, d.h. wenn $t_2^{(j)} - t_1^{(j)}$ klein wird.

Ferner sollte man eine untere Schranke für f mitführen.

4 Gradientenverfahren

Das allgemeine Abstiegsverfahren lässt noch einige Freiheiten in der Wahl der Abstiegsrichtung d^k .

Eine nahe liegende Wahl für d (auch in Hinblick auf die Winkelbedingung aus Satz 3.5 (a)) ergibt sich als Lösung von

$$\min \nabla f(x)^T d \text{ u.d.N. } \|d\| = 1. \quad (4.1)$$

Das Ziel ist also, d als jene Richtung zu bestimmen, in welche f in x am steilsten fällt. Offensichtlich gilt

$$0 \leq |\nabla f(x)^T d| \stackrel{\|d\|=1}{\leq} \|\nabla f(x)\|.$$

Die Wahl

$$d := -\frac{\nabla f(x)}{\|\nabla f(x)\|}$$

liefert

$$\nabla f(x)^T d = -\|\nabla f(x)\|$$

und löst somit (4.1).

Verwenden wir die Wolfe-Powell-Schrittweitenstrategie, dann folgt sofort aus den Sätzen 3.5 und 3.10, dass jeder Häufungspunkt der Folge $(x^k)_{k \in \mathbb{N}}$ ein stationärer Punkt von f ist.

Eine analoge Aussage gilt auch für die strenge Wolfe-Powell-Regel.

Da die Armijo-Bedingung die Forderung (b) aus Satz 3.5 nicht notwendigerweise erfüllt, geben wir folgenden Satz an.

SATZ 4.1

Ist $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar, so ist jeder Häufungspunkt einer durch Algorithmus 3.4 konstruierten Folge $(x^k)_{k \in \mathbb{N}}$ mit

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

und der Armijo-Schrittweitenstrategie ein stationärer Punkt.

Das Konvergenzverhalten des steilsten Abstiegs kann sehr schlecht sein. Für

$$f(x) = \frac{1}{2}x^T Qx + c^T x + \gamma$$

mit $Q \in \mathcal{S}_n$ positiv definit, $c \in \mathbb{R}^n$, $\gamma \in \mathbb{R}$ lässt sich zeigen, dass

$$\|x^k - x^*\| \leq \sqrt{\kappa} \left(\frac{\kappa - 1}{\kappa + 1} \right)^k \|x^0 - x^*\|,$$

wobei $\kappa := \kappa_2(Q)$ die spektrale Konditionszahl von Q bezeichnet, d.h. $\kappa := \frac{\lambda_{\max}}{\lambda_{\min}}$ mit

$$\begin{aligned} \lambda_{\max} &:= \max\{\lambda \mid \lambda \text{ ist Eigenwert von } Q\}; \\ \lambda_{\min} &:= \min\{\lambda \mid \lambda \text{ ist Eigenwert von } Q\}. \end{aligned}$$

Eine mögliche Abhilfe für die langsame Konvergenz des Verfahrens des steilsten Abstiegs besteht darin,

$$d^k := -H^{-1} \nabla f(x^k) \quad \text{mit } H \in \mathcal{S}_n \text{ positiv definit}$$

zu setzen. H soll überdies so sein, dass

$$0 < \frac{\lambda_{\max}(H^{-1}Q)}{\lambda_{\min}(H^{-1}Q)} = \kappa_2(H^{-1}Q) < \frac{\lambda_{\max}(Q)}{\lambda_{\min}(Q)} = \kappa_2(Q).$$

DEFINITION 4.2

Seien $f : \mathbb{R}^n \rightarrow \mathbb{R}$ stetig differenzierbar und $(x^k)_{k \in \mathbb{N}} \subseteq \mathbb{R}^n$. Dann heißt $(d^k)_{k \in \mathbb{N}} \subseteq \mathbb{R}^n$ *gradientenähnlich* bzgl. f und $(x^k)_{k \in \mathbb{N}}$, wenn für jede gegen einen nichtstationären Punkt von f konvergierende Teilfolge Konstanten $c > 0, \varepsilon > 0$ existieren, so dass

- (a) $\forall l \in \mathbb{N} : \|d^{k(l)}\| \leq c$ und
 (b) $\forall l \in \mathbb{N}$ hinreichend groß: $\nabla f(x^{k(l)})^T d^{k(l)} \leq -\varepsilon$.

BEMERKUNG 4.3

- (a) Sei $(H^k)_{k \in \mathbb{N}} \subseteq \mathcal{S}_n$ eine Folge positiv definiter Matrizen, welche

$$\forall x \in \mathbb{R}^n, \forall k \in \mathbb{N} : c_1 \|x\|^2 \leq x^T H^k x \leq c_2 \|x\|^2$$

erfüllen ($c_1, c_2 > 0$ konstant). Dann ist $(d^k)_{k \in \mathbb{N}}$, gegeben durch

$$\forall k \in \mathbb{N} : H^k d^k = -\nabla f(x^k),$$

gradientenähnlich. Denn:

$$\|d^{k(l)}\| = \|-(H^{k(l)})^{-1} \nabla f(x^{k(l)})\| \leq \|(H^{k(l)})^{-1}\| \|\nabla f(x^{k(l)})\| \leq \frac{1}{c_1} \|\nabla f(x^{k(l)})\| \leq C,$$

da $x^{k(l)}$ konvergente Teilfolge, und außerdem

$$-\nabla f(x^{k(l)})^T (H^{k(l)})^{-1} \nabla f(x^{k(l)}) \stackrel{(*)}{\leq} -\frac{1}{c_2} \underbrace{\|\nabla f(x^{k(l)})\|}_{\neq 0} \leq -\varepsilon,$$

wobei (*) erfüllt ist wegen $\frac{1}{c_2} \|x\|^2 \leq x^T (H^{k(l)})^{-1} x \leq \frac{1}{c_1} \|x\|^2$.

- (b) Für Algorithmus 3.4 mit gradientenähnlichen Suchrichtungen und der Armijo-Schrittweitenstrategie gilt eine analoge Aussage zu Satz 4.1.
 (c) Manchmal bringt die Wahl $H^k := \text{diag}(h_{ii}^k)$ mit

$$h_{ii}^k := \frac{\partial^2 f(x^k)}{\partial x_i^2}, \quad 1 \leq i \leq n,$$

eine deutliche Verbesserung. ◆

5 Das Newton-Verfahren

5.1 Das lokale Verfahren

Wir setzen ab nun voraus, dass f und die lokale Minimalstelle x^* (von f) folgende Voraussetzungen erfüllen:

$$\begin{cases} 1. f \text{ ist zweimal stetig differenzierbar mit } \|\nabla^2 f(x) - \nabla^2 f(y)\| \leq \gamma \|x - y\| \text{ für ein } \gamma > 0 \\ 2. \nabla f(x^*) = 0 \\ 3. \nabla^2 f(x^*) \text{ ist positiv definit} \end{cases} \quad (\text{A})$$

Zur Vereinfachung der Schreibweise bezeichne x_a den aktuellen Iterationspunkt und x_+ die neue Iterierte. Wir betrachten ein quadratisches Modell von f um x_a :

$$m_a(x) = f(x_a) + \nabla f(x_a)^T(x - x_a) + \frac{1}{2}(x - x_a)^T \nabla^2 f(x_a)(x - x_a).$$

Wenn $\nabla^2 f(x_a)$ positiv definit ist, dann definieren wir x_+ als die (eindeutige) Minimalstelle unseres Modells. Es gilt

$$0 = \nabla m_a(x_+) = \nabla f(x_a) + \nabla^2 f(x_a)(x_+ - x_a).$$

Umformen liefert die Iterationsvorschrift des *Newton-Verfahrens*, d.h.

$$x_+ = x_a - \nabla^2 f(x_a)^{-1} \nabla f(x_a).$$

Natürlich wird nicht $\nabla^2 f(x_a)^{-1}$ berechnet, sondern es wird

$$\nabla^2 f(x_a)d = -\nabla f(x_a)$$

gelöst und $x_+ := x_a + d$ gesetzt.

Falls x_a weit von einer lokalen Minimalstelle x^* , die (A) erfüllt, entfernt ist, dann kann $\nabla^2 f(x_a)$ negative Eigenwerte haben. Also kann x_+ lokale Minimalstelle oder ein Sattelpunkt sein.

Um dies zu vermeiden, müssen geeignete Modifikationen eingeführt werden. Zunächst sei aber vorausgesetzt, dass x_a hinreichend nahe an x^* ist.

SATZ 5.1

Sei (A) erfüllt. Dann existieren Konstanten $K > 0$ und $\delta > 0$ derart, dass für alle x_a aus der Menge $B(x^*, \delta) = \{x \in \mathbb{R}^n \mid \|x - x^*\| < \delta\}$ der Newton-Schritt

$$x_+ = x_a - \nabla^2 f(x_a)^{-1} \nabla f(x_a)$$

folgende Abschätzung erfüllt:

$$\|x_+ - x^*\| \leq K \|x_a - x^*\|^2.$$

SATZ 5.2

Es sei (A) erfüllt. Dann existiert $\delta > 0$, so dass das Newton-Verfahren

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

für $x^0 \in B(\delta) := B(0, \delta)$ gegen x^* konvergiert mit

$$\|x^{k+1} - x^*\| \leq C \|x^k - x^*\|^2$$

mit $C > 0$ (quadratische Konvergenzordnung).

Ein natürliches Abbruchkriterium für das Newton-Verfahren (wie auch für die Gradientenverfahren aus Abschnitt 4) setzt sich aus einer relativen und einer absoluten Fehlerschranke zusammen.

Sei $\tau_r \in (0, 1)$ eine erwünschte Reduktion in der Gradientennorm und τ_a mit $1 \gg \tau_a > 0$ eine absolute Fehlerschranke, dann stoppt man das Verfahren, wenn

$$\|\nabla f(x^k)\| \leq \tau_r \|\nabla f(x^0)\| + \tau_a$$

gilt. Ist $\|\nabla f(x^0)\|$ groß, so ist $\tau_r \|\nabla f(x^0)\|$ der dominante Term. Ist hingegen $\|\nabla f(x^0)\|$ klein, dann ist τ_a der dominante Term.

Wir nehmen nun an, dass $n = 1$ ist und f nur approximativ ausgewertet werden kann, d.h.

$$\tilde{f}(x) = f(x) + \tilde{\varepsilon}_f(x) \quad \text{mit } \tilde{\varepsilon}_f(x) \geq 0 \text{ und } |\tilde{\varepsilon}_f(x)| \leq \varepsilon_f \text{ für ein } \varepsilon_f > 0.$$

Bestimmen wir nun die Ableitungen numerisch, z.B. durch Vorwärtsdifferenzen, so ergibt sich

$$D_h^+ f(x) = \frac{\tilde{f}(x+h) - \tilde{f}(x)}{h}.$$

Also:

$$\begin{aligned} \|D_h^+ f(x) - f'(x)\| &= \left\| \frac{\tilde{f}(x+h) - \tilde{f}(x)}{h} - f'(x) \right\| \\ &= \left\| \frac{f(x+h) + \tilde{\varepsilon}_f(x+h) - f(x) - \tilde{\varepsilon}_f(x)}{h} - f'(x) \right\| \\ &\leq \left\| \frac{f(x+h) - f(x)}{h} - f'(x) \right\| + \frac{2\varepsilon_f}{h} \\ &\stackrel{\xi \in (x, x+h)}{=} \left\| \frac{f(x) + f'(x)h + \frac{1}{2}f''(\xi)h^2 - f(x)}{h} - f'(x) \right\| + \frac{2\varepsilon_f}{h} \\ &= \frac{h}{2} \|f''(\xi)\| + \frac{2\varepsilon_f}{h} \\ &= \mathcal{O}\left(h + \frac{\varepsilon_f}{h}\right). \end{aligned}$$

Die Minimalstelle h^* der Fehlerfunktion $\text{err}_+(h) = h + \frac{\varepsilon_f}{h}$ erfüllt

$$\text{err}'_+(h^*) = 1 - \frac{\varepsilon_f}{(h^*)^2} = 0,$$

d.h. $h^* = \sqrt{\varepsilon_f}$. Der Fehler im Gradienten ist also $\varepsilon_g = \mathcal{O}(h^*) = \mathcal{O}(\sqrt{\varepsilon_f})$.

Verwenden wir nun abermals Vorwärtsdifferenzen zur Approximation der Hessematrix, so ergibt sich offensichtlich für den Fehler ε_H die Größenordnung

$$\varepsilon_H = \mathcal{O}(\sqrt{\varepsilon_g}) = \mathcal{O}(\varepsilon_f^{\frac{1}{4}}).$$

Dies impliziert, dass Hessematrizen, basierend auf zweifacher numerischer Differenzierung, relativ ungenau sind – selbst wenn $\varepsilon_f = 10^{-16}$ (Maschinen-Epsilon) ist, folgt $\varepsilon_H \approx 10^{-4}$.

Im Falle zentraler Differenzenapproximationen ergibt sich ein besseres Ergebnis:

$$\varepsilon_H = \mathcal{O}(\varepsilon_f^{\frac{4}{9}}).$$

Bei $\varepsilon_f = 10^{-16}$ erhalten wir $\varepsilon_H \approx 10^{-7,1}$.

Konvergenz des Newton-Verfahrens ist nur zu erwarten, wenn $\varepsilon_g \rightarrow 0$ im Laufe der Iteration.

SATZ 5.3

Es sei (A) erfüllt. Dann existieren Konstanten $\bar{K} > 0$, $\delta > 0$ und ein $\varepsilon > 0$, so dass für $x_a \in B(x^*, \delta)$ und $\|\varepsilon_H(x_a)\| < \varepsilon$ gilt:

$$x_+ = x_a - (\nabla^2 f(x_a) + \varepsilon_H(x_a))^{-1}(\nabla f(x_a) + \varepsilon_g(x_a))$$

ist wohldefiniert, d.h. $(\nabla^2 f(x_a) + \varepsilon_H(x_a))$ ist regulär, und erfüllt

$$\|x_+ - x^*\| \leq \bar{K}(\|x_a - x^*\|^2 + \underbrace{\|\varepsilon_H(x_a)\|}_{\text{beeinflusst Konv.geschw.}} \|x_a - x^*\| + \underbrace{\|\varepsilon_g(x_a)\|}_{\text{Genauigkeit!}}).$$

Wir betrachten das Verfahren

$$x^{k+1} = x^k - \nabla^2 f(x^0)^{-1} \nabla f(x^k), \quad x^0 \text{ Startwert, } k = 0, 1, \dots \quad (5.1)$$

Es gilt

$$\begin{cases} \varepsilon_H(x^k) &= \nabla^2 f(x^0) - \nabla^2 f(x^k), \quad \|\varepsilon_H(x^k)\| \leq \varepsilon_H \\ \|\varepsilon_H(x^k)\| &= \|\nabla^2 f(x^0) - \nabla^2 f(x^k)\| \leq \gamma \|x^0 - x^k\| \leq \gamma(\|x^0 - x^*\| + \|x^* - x^k\|) \end{cases} \quad (5.2)$$

Die Konvergenz des Verfahrens (5.1) folgt aus Satz 5.3 mit $\varepsilon_g = 0$ und $\varepsilon_H = \mathcal{O}(\|x^0 - x^*\|)$.

SATZ 5.4

Es sei (A) erfüllt. Dann existieren $K \in (0, 1)$ und $\delta > 0$, so dass für $x^0 \in B(x^*, \delta)$ gilt: Die Folge der Iterierten $(x^k)_{k \in \mathbb{N}}$, erzeugt durch (5.1), konvergiert linear gegen x^* und

$$\|x^{k+1} - x^*\| \leq K \|x^* - x^k\|.$$

BEWEIS

Sei $\delta > 0$ so gewählt, dass Satz 5.2 gilt. Mit (5.2) folgt:

$$\begin{aligned} \|x^{k+1} - x^*\| &\leq \bar{K}(\|x^k - x^*\|^2 + \gamma(\|x^0 - x^*\| + \|x^* - x^k\|))\|x^k - x^*\| \\ &= \bar{K}(\underbrace{\|x^k - x^*\|}_{\leq \delta}(1 + \gamma) + \gamma \underbrace{\|x^0 - x^*\|}_{\leq \delta})\|x^k - x^*\| \\ &\leq \bar{K}(1 + 2\gamma)\delta \|x^k - x^*\|. \end{aligned}$$

Um Konvergenz zu garantieren, verkleinere δ , so dass $\bar{K}(1 + 2\gamma)\delta < 1$. ■

5.2 Inexakte Newton-Verfahren

Betrachte für $k = 0, 1, \dots$:

$$\begin{cases} \nabla^2 f(x^k) d^k &= -\nabla f(x^k) \\ x^{k+1} &= x^k + d^k \end{cases} .$$

Inexakte Newton-Verfahren verwenden einen approximativen Schritt \tilde{d} , welcher für ein $\eta_a > 0$

$$\|\nabla^2 f(x_a) \tilde{d} + \nabla f(x_a)\| \leq \eta_a \|\nabla f(x_a)\| \quad (5.3)$$

erfüllt. Wir wissen, dass $\nabla^2 f(x_a)$ positiv definit ist für x_a nahe x^* . Daher eignet sich z.B. das „CG-Verfahren“ (Verfahren der konjugierten Gradienten) zur iterativen Lösung von

$$\nabla^2 f(x_a) \tilde{d} = -\nabla f(x_a).$$

Man spricht dann vom *Newton-CG-Verfahren*.

SATZ 5.5

Sei (A) erfüllt. Dann existieren $K_I \geq 0$, $\delta > 0$, so dass für $x_a \in B(\delta)$, \tilde{d} aus (5.3) und $x_+ = x_a + \tilde{d}$ gilt:

$$\|x_+ - x^*\| \leq K_I (\|x_a - x^*\| + \eta_a) \|x_a - x^*\|.$$

Für das gesamte Verfahren erhalten wir

SATZ 5.6

Sei (A) erfüllt. Dann existieren $\delta > 0$ und $\bar{\eta} \geq 0$, so dass für $x^0 \in B(x^*, \delta)$ eine Folge $(\eta_k)_{k \in \mathbb{N}} \subseteq [0, \bar{\eta}]$, so dass

$$x^{k+1} := x^k + \tilde{d}^k \quad (\text{inexakte Newton-Iteration})$$

mit

$$\|\nabla^2 f(x^k) \tilde{d}^k + \nabla f(x^k)\| \leq \eta_k \|\nabla f(x^k)\|$$

linear gegen x^* konvergiert.

Ferner: Falls $\eta_k \xrightarrow{k \rightarrow \infty} 0$, dann ist die Konvergenz *superlinear*, d.h.

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

Falls $\eta_k \leq K_\eta \|\nabla f(x^k)\|$ für ein $K_\eta > 0$, so ist die Konvergenz sogar quadratisch.

5.3 Globale Konvergenz

Bisher im Zusammenhang mit dem Newton-Verfahren lokale Konvergenzaussagen, d.h. „ x^0 hinreichend nahe an x^* “.

Jetzt: Globalisierungen, die die Startpunktwahl relaxieren.

Kann man sicherstellen, dass die *Newton-Iterationsmatrix* $\nabla^2 f(x^k)$ oder eine entsprechende Approximation

$$\forall k \in \mathbb{N}, \forall d \in \mathbb{R}^n : c_1 \|d\|^2 \leq d^T \nabla^2 f(x^k) d \leq c_2 \|d\|^2 \quad (0 < c_1 \leq c_2)$$

erfüllt, so ist d^k als Lösung von $\nabla^2 f(x^k) d^k = -\nabla f(x^k)$ eine gradientenähnliche Abstiegsrichtung.

Eingesetzt in das allgemeine Abstiegsverfahren folgt dann aus Abschnitt 4 die globale Konvergenz des *globalisierten Newton-Verfahrens* (x^0 kann beliebig gewählt werden).

5.3.1 Trust-Region-Methoden

Schrittweisen + Newton (Algorithmus 3.4): Problem $(\nabla^2 f(x^k))_{k \in \mathbb{N}}$ positiv definit.

Zu prüfen: Ist $(d^k)^T \nabla^2 f(x^k) d^k \geq \varepsilon \|d^k\|^2$, dann ersetze $\nabla^2 f(x^k)$ durch $H^k \in \mathcal{S}_n$ positiv definit mit $H^k d^k = -\nabla f(x^k)$.

Idee der *Trust-Region-Methoden*:

- (a) Globale Konvergenzeigenschaften von gradientenähnlichen Verfahren ausnützen;
- (b) glatter Übergang zur Newton-Richtung.

Wir verwenden dabei eine Umgebung, in der wir einem Modell von f trauen können.

Sei m_a ein quadratisches Modell von f um x_0 , gegeben durch

$$m_a(x) := f(x_a) + \nabla f(x_a)^T (x - x_a) + \frac{1}{2} (x - x_a)^T \nabla^2 f(x_a) (x - x_a).$$

Weiter sei Δ ein Radius einer Kugel um x_a , in welcher wir dem Modell von f vertrauen. Δ nennt man *Trust-Region-Radius* und die Kugel $\mathcal{T}(\Delta) := \{x \in \mathbb{R}^n \mid \|x - x_a\| \leq \Delta\}$ *Trust-Region*.

Die nächste Iterierte wird als approximative Minimalstelle von m_a in $\mathcal{T}(A)$ gewählt.

Das zugehörige Trust-Region-Hilfsproblem lautet daher

$$\min m_a(x + d) \text{ u.d.N. } \|d\| \leq \Delta. \quad (5.4)$$

Wir bezeichnen die Lösung von (5.4) mit d_V (*Versuchslösung*) und setzen $x_V := x_a + d_V$.

Im Wesentlichen überprüft man, ob das quadratische Modell eine „gute“ Approximation von f in $\mathcal{T}(\Delta)$ ist. Dazu definiere

$$\begin{aligned} \text{ared}_a &:= f(x_a) - f(x_V) && \text{(tatsächliche Reduktion)} \\ \text{pred}_a &:= m_a(x_a) - m_a(x_V) && \text{(erwartete Reduktion)}. \end{aligned}$$

Es gilt (mit $H_a := \nabla^2 f(x_a)$):

$$\begin{aligned} \text{pred}_a &= m_a(x_a) - m_a(x_V) \\ &= f(x_a) - f(x_a) - \nabla f(x_a)^T (x_V - x_a) - \frac{1}{2} (x_V - x_a)^T H_a (x_V - x_a) \\ &= -\nabla f(x_a)^T (x_V - x_a) - \frac{1}{2} (x_V - x_a)^T H_a (x_V - x_a). \end{aligned}$$

Im folgenden Algorithmus benötigen wir die Parameter

$$0 < \mu_0 \leq \underline{\mu} < \bar{\mu},$$

um zu entscheiden, ob der Versuchsschritt verworfen wird ($\frac{\text{ared}_a}{\text{pred}_a} < \mu_0$) und/oder ob der Trust-Region verkleinert ($\frac{\text{ared}_a}{\text{pred}_a} < \underline{\mu}$), vergrößert ($\frac{\text{ared}_a}{\text{pred}_a} > \bar{\mu}$) oder unverändert belassen werden soll.

Die Änderung von Δ wird mit Hilfe von $0 < \underline{\omega} < 1 < \bar{\omega}$ durchgeführt. Weiter sei $C > 1$.

ALGORITHMUS 5.7

Input: $x_a \in \mathbb{R}^n$, $x_V \in \mathbb{R}^n$, $\Delta \in \mathbb{R}^+$

Begin

$z^0 := x_a$; $z_V^0 := x_V$; $\hat{\Delta}^{(0)} := \Delta$; $l := 0$

While $z^l = x_a$

$\text{ared}^{(l)} := f(x_a) - f(z_V^l)$; $d_V^l := z_V^l - x_a$;

$\text{pred}^{(l)} := -\nabla f(x_a)^T d_V^l - \frac{1}{2}(d_V^l)^T H_a d_V^l$

If $\frac{\text{ared}^{(l)}}{\text{pred}^{(l)}} < \mu_0$

$z^{l+1} := x_a$, $\hat{\Delta}^{(l+1)} := \underline{\omega} \hat{\Delta}^{(l)}$

If $l \geq 1$ & $\hat{\Delta}^{(l)} > \hat{\Delta}^{(l-1)}$

$z^{l+1} := z_V^{l-1}$, $\hat{\Delta}^{(l+1)} := \hat{\Delta}^{(l-1)}$

Else

berechne die Lösung d_V^{l+1} des T.-R.-Hilfsproblems mit Radius $\hat{\Delta}^{(l+1)}$

$z_V^{l+1} := x_a + d_V^{l+1}$

End(If)

Elseif $\mu_0 \leq \frac{\text{ared}^{(l)}}{\text{pred}^{(l)}} \leq \underline{\mu}$

$z^{l+1} := z_V^l$; $\hat{\Delta}^{(l+1)} := \underline{\omega} \hat{\Delta}^{(l)}$

Elseif $\underline{\mu} \leq \frac{\text{ared}^{(l)}}{\text{pred}^{(l)}} \leq \bar{\mu}$

$z^{l+1} := z_V^l$, $\hat{\Delta}^{(l+1)} := \hat{\Delta}^{(l)}$

Elseif $\bar{\mu} \leq \frac{\text{ared}^{(l)}}{\text{pred}^{(l)}}$

If $\|d_V^l\| = \hat{\Delta}^{(l)} \leq C \|\nabla f(x_a)\|$

$z^{l+1} := x_a$; $\hat{\Delta}^{(l+1)} := \bar{\omega} \hat{\Delta}^{(l)}$

berechne die Lösung d_V^{l+1} des T.-R.-Hilfsproblems mit Radius $\hat{\Delta}^{(l+1)}$

$z_V^{l+1} := x_a + d_V^{l+1}$

Else

$z^{l+1} := z_V^l$, $\hat{\Delta}^{(l+1)} := \hat{\Delta}^{(l)}$

End(If)

End(If)

$l := l + 1$

End(While)

$x_+ := z^l$; $\Delta_+ := \hat{\Delta}^{(l)}$

End

In Algorithmus 5.7 ist der Trust-Region nach oben durch $C \|\nabla f(x_a)\|$ beschränkt. Die **While**-Schleife sollte nach endlich vielen Schritten terminieren.

ALGORITHMUS 5.8 (Trust-Region-Framework)

Input: $x^0 \in \mathbb{R}^n$; $\Delta_0 \in \mathbb{R}_+$

Begin

$k := 0$; $\tau_0 := \|\nabla f(x^0)\|$

While $\|\nabla f(x^k)\| > \tau_r \tau_0 + \tau_a$

berechne eine Approximation H^k der Hessematrix $\nabla^2 f(x^k)$

berechne d_V^k als Lösung von

$$\min f(x^k) + \nabla f(x^k)d + \frac{1}{2}d^T H^k d \text{ u.d.N. } \|d\| \leq \Delta_k$$

berechne (x^{k+1}, Δ_{k+1}) mit Algorithmus 5.6 mit **Input** x^k ; $x_V^k := x^k + d_V^k$; Δ_k

$k := k + 1$

End(While)

End

5.3.2 Globale Konvergenz des Trust-Region-Verfahrens**ANNAHME 5.9**

(a) Es existiere ein $\sigma > 0$, so dass

$$\text{pred}_a = m_a(x_a) - m_a(x_V) = f(x_a) - m_a(x_V) \geq \sigma \|\nabla f(x_a)\| \min\{\|d_V\|, \|\nabla f(x_a)\|\} \quad (5.5)$$

(b) Es existiere $M > 0$, so dass

$$\|d_V\| \geq \frac{\|\nabla f(x_a)\|}{M} \quad \text{oder} \quad \|d_V\| = \Delta_a. \quad \blacklozenge$$

SATZ 5.10

Sei ∇f Lipschitz-stetig mit Konstante $L > 0$. Sei die Folge $(x_n)_{n \in \mathbb{N}}$ erzeugt von Algorithmus 5.8 und es sei angenommen, dass die Lösungen des Trust-Region-Hilfsproblems Annahme 5.9 erfüllen. Ferner seien die Matrizen $(H^k)_{k \in \mathbb{N}}$ beschränkt.

Dann ist entweder f nach unten unbeschränkt, $\nabla f(x^k) = 0$ für ein k oder $\lim_{k \rightarrow \infty} \|\nabla f(x^k)\| = 0$.

Eine einfache Idee zur Lösung des Hilfsproblems beruht auf Fixieren der Richtung gemäß des Verfahrens des steilsten Abstiegs unter Berücksichtigung des Trust-Regions.

Seien x_a die aktuelle Iterierte und Δ_a der aktuelle Trust-Region-Radius. Der Versuchspunkt $x_V = x_V(t)$ ist dann definiert über die Minimalstelle t_a von

$$\begin{cases} \min_{t \geq 0} \psi_a(t) := m_a(x_a - t\nabla f(x_a)) \\ \text{u.d.N. } x_V(t) - x_a = x_a - t\nabla f(x_a) - x_a = -t\nabla f(x_a) \in \mathcal{T}(\Delta_a) \end{cases}$$

Es gelten

$$\begin{aligned} \psi(t) &= m_a(x_a - t\nabla f(x_a)) \\ &= f(x_a) - t\|\nabla f(x_a)\|^2 + \frac{t^2}{2} \nabla f(x_a)^T H_a \nabla f(x_a) \end{aligned}$$

und

$$\psi'(t) = -\|\nabla f(x_a)\|^2 + t\nabla f(x_a)^T H_a \nabla f(x_a).$$

Fallunterscheidung zur Bestimmung von t_a :

(a) $\nabla f(x_a)^T H_a \nabla f(x_a) \leq 0$: Offensichtlich wird die Trust-Region-Restriktion aktiv, d.h.

$$\|x_V(t_a) - x_a\| = t_a \|\nabla f(x_a)\| = \Delta_a \Rightarrow t_a = \frac{\Delta_a}{\|\nabla f(x_a)\|}.$$

(b) $\nabla f(x_a)^T H_a \nabla f(x_a) > 0$: Dann impliziert $\psi'(\hat{t}_a) \stackrel{!}{=} 0$

$$\hat{t}_a = \frac{\|\nabla f(x_a)\|^2}{\nabla f(x_a)^T H_a \nabla f(x_a)}.$$

Falls $\|x_V(\hat{t}_a) - x_a\| \leq \Delta_a \Rightarrow t_a := \hat{t}_a$. Andernfalls ist die Trust-Region-Restriktion aktiv, d.h.

$$t_a = \frac{\Delta_a}{\|\nabla f(x_a)\|}.$$

Zusammenfassend:

$$t_a = \begin{cases} \frac{\Delta_a}{\|\nabla f(x_a)\|}, & \text{falls } \nabla f(x_a)^T H_a \nabla f(x_a) \leq 0, \\ \min\left\{\frac{\Delta_a}{\|\nabla f(x_a)\|}, \frac{\|\nabla f(x_a)\|^2}{\nabla f(x_a)^T H_a \nabla f(x_a)}\right\}, & \text{sonst.} \end{cases}$$

Die Minimalstelle $x_V(t_a)$ ($:= x_a - t_a \nabla f(x_a)$) des quadratischen Modells m_a in Richtung des negativen Gradienten heißt *Cauchy-Punkt* (Bezeichnung: x_a^{CP}). Man kann zeigen, dass der Cauchy-Punkt Annahme 5.9 erfüllt.

BEMERKUNG 5.11

(a) Die Verwendung des Cauchy-Punktes führt zwar zu globaler Konvergenz, aber unter Umständen ist die Konvergenzgeschwindigkeit lokal nur linear.

(b) „*dogleg-Technik*“ leistet einen „glatten“ Übergang von der Richtung des steilsten Abstiegs zur Newton-Richtung.

Lokal liegt dann quadratische Konvergenz vor, wenn $H^k = \nabla^2 f(x^k)$, $k \in \mathbb{N}$, gilt. ◆

6 Quasi-Newton-Verfahren

Nachteile des Newton-Verfahrens:

- zweite Ableitungen benötigt,
- positive Definitheit,
- $\mathcal{O}(n^3)$ Multiplikationen (Lösung des linearen Systems mit direktem Verfahren).

QN-Verfahren:

- approximieren zweite Ableitungen durch erste Ableitungen,
- erhalten positive Definitheit (bei einigen QN-Verfahren),
- benötigen $\mathcal{O}(n^2)$ Multiplikationen (bei einigen QN-Varianten).

H^k wird von Iteration zu Iteration aufdatiert.

Allgemeine Grundstruktur:

- Setze $d^k := -(H^k)^{-1} \nabla f(x^k)$.
- Bestimme $x^{k+1} := x^k + t_k d^k$ mit Schrittweitenstrategie.
- Verwende x^k, x^{k+1}, H^k , um H^k zu H^{k+1} aufzudatieren.

Der letzte Punkt der QN-Vorteile gilt für Varianten, die direkt $\nabla^2 f(x^k)^{-1}$ approximieren und somit die Lösung des Gleichungssystems ersparen.

Die positive Definitheit der Matrix H^k ist nicht bei allen QN-Verfahren gegeben.

Aufdatierungsformeln für H :

Seien $s_a := t_a d_a (= -t_a H_a^{-1} \nabla f(x_a))$, $y_a := \nabla f(x_+) - \nabla f(x_a)$, $x_+ := x_a + s_a$.

Es gilt

$$\begin{aligned} y_a &= \nabla f(x_+) - \nabla f(x_a) = \nabla f(x_a) + \nabla^2 f(x_a)(x_+ - x_a) + \mathcal{O}(\|x_+ - x_a\|) - \nabla f(x_a) \\ &= \nabla^2 f(x_a) s_a + \mathcal{O}(\|s_a\|). \end{aligned}$$

Naheliegende Forderung für H_+ (Update von H_a):

$$H_+ s_a = y_a \quad (\text{QN-Bedingung oder Sekantenbedingung}). \quad (6.1)$$

Ein einfacher Ansatz für H_+ in (6.1) ist

$$H_+ := H_a + \alpha u u^T \quad (\text{symmetrische Rang 1-Korrektur}).$$

Einsetzen in (6.1) liefert

$$\begin{aligned} H_+ s_a &= H_a s_a + \alpha \underbrace{u^T s_a}_{\in \mathbb{R}} u = y_a \\ \Rightarrow u &\text{ proportional zu } y_a - H_a s_a \\ \Rightarrow u &= y_a - H_a s_a \text{ (die Länge in } \alpha \text{ berücksichtigt)} \\ \alpha(u^T s_a) &= 1, \text{ d.h. } \alpha = \frac{1}{y_a^T s_a - s_a^T H_a s_a}. \end{aligned}$$

Also:

$$H_+ = H_a + \frac{(y_a - H_a s_a)(y_a - H_a s_a)^T}{(y_a - H_a s_a)^T s_a}.$$

Nachteile:

- positive Definitheit geht meist verloren,
- $y_a - H_a s_a$ eventuell nahe bei 0.

Flexibler sind Rang 2-Korrekturen:

$$H_+ := H_a + \alpha u u^T + \beta v v^T. \quad (6.2)$$

Einsetzen in (6.1):

$$H_+ s_a = H_a s_a + \alpha u (u^T s_a) + \beta v (v^T s_a) = y_a.$$

Die Vektoren u und v sind nicht mehr eindeutig bestimmt.

Es bietet sich an,

$$u := y_a \quad \text{und} \quad v := H_a s_a$$

zu wählen.

$$\begin{aligned} H_a s_a + \alpha y_a y_a^T s_a + \beta (H_a s_a) (H_a s_a)^T s_a &= y_a \\ \Leftrightarrow \alpha y_a (y_a^T s_a) + \beta H_a s_a (s_a^T H_a s_a) &= y_a - H_a s_a \\ \Rightarrow \alpha (y_a^T s_a) = 1 \quad \text{und} \quad \beta (s_a^T H_a s_a) &= -1 \\ \Rightarrow \alpha = \frac{1}{y_a^T s_a} \quad \text{und} \quad \beta = \frac{-1}{s_a^T H_a s_a}. \end{aligned}$$

Also (*Broyden/Fletcher/Goldfarb/Shanno*):

$$H_+ = H_a + \frac{y_a y_a^T}{y_a^T s_a} - \frac{H_a s_a (H_a s_a)^T}{s_a^T H_a s_a} \quad (\text{BFGS-Formel}). \quad (6.3)$$

Man kann auch $\nabla^2 f(x^k)^{-1}$ approximieren durch B^k . Die QN-Bedingung lautet dann

$$B_+ y_a = s_a. \quad (6.4)$$

Verwenden wir die symmetrische Rang 2-Korrektur analog zu (6.4) mit $u := s_a$ und $v := B_a y_a$, dann erhalten wir

$$B_+ = B_a + \frac{s_a s_a^T}{s_a^T y_a} - \frac{(B_a y_a) (B_a y_a)^T}{y_a^T B_a y_a} \quad (\text{DFP-Formel}) \quad (6.5)$$

nach *Davidon/Fletcher/Powell*.

LEMMA 6.1

Seien $H_a \in \mathcal{S}_n$ positiv definit, $y_a^T s_a > 0$ und H_+ gemäß (6.3) bestimmt.

Dann ist H_+ symmetrisch und positiv definit.

BEWEIS

H_a positiv definit und $y_a^T s_a \neq 0$ liefern für alle $z \neq 0$

$$\begin{aligned} z^T H_+ z &= z^T H_a z + \frac{z^T y_a y_a^T z}{y_a^T s_a} - \frac{z^T (H_a s_a) (H_a s_a)^T z}{s_a^T H_a s_a} \\ &= \frac{(z^T y_a)^2}{y_a^T s_a} + z^T H_a z - \frac{(z^T H_a s_a)^2}{s_a^T H_a s_a}. \end{aligned}$$

Da H_a symmetrisch und positiv definit ist, existiert $H_a^{1/2}$ mit $H_a = H_a^{1/2} H_a^{1/2}$. Damit gilt

$$z^T H_a s_a = z^T H_a^{1/2} H_a^{1/2} s_a \leq \|H_a^{1/2} z\| \|H_a^{1/2} s_a\|$$

bzw.

$$(z^T H_a s_a)^2 \leq \underbrace{\|H_a^{1/2} z\|^2}_{z^T H_a z} \underbrace{\|H_a^{1/2} s_a\|^2}_{s_a^T H_a s_a},$$

d.h.

$$z^T H_a z - \frac{(z^T H_a s_a)^2}{s_a^T H_a s_a} \geq z^T H_a z - z^T H_a z = 0.$$

Ferner

$$z^T H_a z - \frac{(z^T H_a s_a)^2}{s_a^T H_a s_a} = 0,$$

wenn $H_a^{1/2}z$ und $H_a^{1/2}s_a$ linear abhängig sind. Da $H_a^{1/2}$ regulär ist, gilt zu diesem Fall $z = \lambda s_a \neq 0$, d.h. $\lambda \neq 0$. Also

$$z^T y_a = \lambda s_a^T y_a \neq 0$$

und somit

$$z^T H_+ z \begin{cases} = \frac{(z^T y_a)^2}{y_a^T s_a} > 0 & \text{für } H_a^{1/2} z = \lambda H_a^{1/2} s_a \\ > \frac{(z^T y_a)^2}{y_a^T s_a} \geq 0 & \text{sonst} \end{cases} . \quad \blacksquare$$

Die Bedingung $y_a^T s_a > 0$ ist realistisch. Für quadratische Probleme, d.h. von der Form

$$f(x) = \frac{1}{2} x^T G x + g^T x + b$$

mit symmetrischer, positiv definiter Hessematrix G , gilt die Beziehung

$$y_a^T s_a = (\nabla f(x_+) - \nabla f(x_a))^T (x_+ - x_a) = (x_+ - x_a)^T G (x_+ - x_a) > 0 \quad \text{für } x_+ \neq x_a.$$

Für allgemeine Probleme wird $y_a^T s_a > 0$ durch Verwendung von Schrittweitenstrategien (Wolfe-Powell) sicher gestellt.

SATZ 6.2 (lokale Konvergenz)

Es sei Annahme (A) erfüllt. Dann existiert ein $\delta > 0$, so dass für

$$\|x^0 - x^*\| \leq \delta \quad \text{und} \quad \|H^0 - \nabla^2 f(x^*)\| \leq \delta$$

die BFGS-Methode wohldefiniert ist und superlinear gegen ein x^* konvergiert, d.h.

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

Index

Abstiegsrichtung	9
Abstiegsverfahren	
allgemeines	10
global konvergentes	10
Algorithmus	
Allgemeines Abstiegsverfahren	10
Armijo-Schrittweitenstrategie	11
Trust-Region-Framework	24
Trust-Region-Methode	23
Wolfe-Powell-Algorithmus	13
Armijo-Regel	11
Backtracking	12
Bedingung	
erster Ordnung	4
zweiter Ordnung	5
BFGS-Formel	27
Cauchy-Punkt	25
CG-Verfahren	21
Definitheit der Hessematrix	4
DFP-Regel	27
dogleg-Technik	25
Fehlerfunktion	19
globalisiertes Newton-Verfahren	22
Gradient	3
gradientenähnlich	17
gradientenähnliche Richtungen	11
Gradientenverfahren	16
Hessematrix	4
inexakte Newton-Iteration	21
inexaktes Newton-Verfahren	21
konvex	
gleichmässig	7
strikt	7
konvexe Funktion	7
konvexe Menge	7
Maschinen-Epsilon	19
Maximalstelle	
globale	3
lokale	3
strikte globale	3
strikte lokale	3
Minimalstelle	
globale	3
lokale	3
strikte globale	3
strikte lokale	3
Minimum	
globales	3
lokales	3
striktes globales	3
striktes lokales	3
Modell	12
Modul	7
Nebenbedingung	3
Newton-Iterationsmatrix	22
Newton-Schritt	18
Newton-Verfahren	18
Niveaumenge	7
Optimalitätskriterien	4
Optimierungsproblem	
diskretes	3
endlichdimensionales	3
kombinatorisches	3
nicht-differenzierbares	3
restringiertes	3
stetiges	3
unrestringiertes	3
QN-Bedingung	26
QN-Verfahren	26
Reduktion	
erwartete	22
tatsächliche	22
Restriktion	
Ganzzahligkeits-	3
Gleichungs-	3
Ungleichungs-	3
Sattelpunkt	6
Schrittweite	11
Schrittweitenstrategie	9
spektrale Konditionszahl	16
Spektralnorm	4
stationärer Punkt	3
steilster Abstieg	9
superlinear	21
symmetrische Rang 1-Korrektur	26
TR-Methoden	22
Trust-Region-Radius	22
Trust-Region	22
Versuchslösung	22
Vorwärtsdifferenz	19
Winkelbedingung	10
Wolfe-Powell-Regel, strenge	14
Wolfe-Powell-Regel	12
zentrale Differenz	19
Zielfunktion	3
Zulässigkeitsbereich	3

Literatur

- [1] C. T. Kelley: *Iterative Methods for Optimization*. SIAM Frontiers in Applied Mathematics, Philadelphia, 1999
http://www.siam.org/books/textbooks/fr18_book.pdf
- [2] J. Nocedal und S. J. Wright: *Numerical Optimization*. Springer-Verlag, 2006
- [3] S. Volkwein: *Numerische Verfahren der restringierten Optimierung*. Vorlesungsscript, 2009.
<http://www.math.uni-konstanz.de/numerik/personen/volkwein/teaching/Optimierung2.pdf>