

VORLESUNGSNOTIZEN

# Optimierung

Wintersemester 24/25

*Emma Bach*

Vorlesung gehalten von  
Prof. Rolf BACKOFEN  
und  
Prof. Moritz DIEHL

March 14, 2025

# Inhalt

<b>1</b>	<b>Einführung</b>	<b>2</b>
1.1	Anwendungsbeispiel: Support Vector Machines.....	2
1.2	Typische Problemklassen .....	2
1.3	Konvexität .....	3
<b>2</b>	<b>Gradientenverfahren</b>	<b>4</b>
2.1	Wiederholung: Mehrdimensionale Differentiation.....	4
2.2	Iterative Optimierung .....	5
2.2.1	Welche Richtung ist optimal?.....	5
2.3	Bestimmung der Schrittweite .....	5
2.3.1	Die Armijo-Bedingung.....	6
2.3.2	Die Wolfe-Bedingungen .....	6
2.3.3	Line Search .....	7
2.3.4	Conjugate Gradient-Verfahren.....	8
<b>3</b>	<b>Newton und Quasi-Newton-Verfahren</b>	<b>9</b>
3.1	Das Newton-Verfahren .....	9
3.2	Konvergenzordnungen .....	9
3.2.1	Konvergenzradius des Newton-Verfahren .....	10
3.3	Quasi-Newton Verfahren .....	10
3.3.1	Das Davidon-Fletcher-Powell Verfahren .....	11
3.3.2	Das Broyden-Fletcher-Goldfarb-Shanno Verfahren .....	11
3.4	Least-Squares-Aufgaben.....	11
3.4.1	Lineare Residuen .....	12
3.4.2	Nichtlineare Residuen .....	12
<b>4</b>	<b>Optimierung mit Nebenbedingungen</b>	<b>14</b>
4.1	Optimierung mit Gleichheitsbedingungen.....	14
4.2	Optimierung mit Ungleichheitsbedingungen .....	15
4.3	Dualität.....	16
4.3.1	Beispiel: Optimierung mit Hilfe des dualen Problems.....	17
4.3.2	Wichtige Probleme mit starker Dualität .....	17

# Chapter 1

## Einführung

Ein Optimierungsproblem besteht aus einer **zulässigen Menge**  $G$  und einer **Zielfunktion**  $f : G \rightarrow H$ , wobei wir im Allgemeinen von  $H = \mathbb{R}$  ausgehen werden. Wir schreiben dann zum Beispiel

$$\min_{x \in \mathbb{R}} f(x),$$

um das Problem “finde den kleinsten Wert, den  $f(x)$  bei reellem  $x$  annimmt” zu notieren, und

$$\operatorname{argmin}_{x \in \mathbb{R}} f(x) = 4$$

für das Problem “finde die kleinste reelle Zahl  $x$ , sodass  $f(x) = 4$  ist.

Historisch ist das Feld der mathematischen Optimierung unter Nebenbedingungen stark verankert im Feld der **Operations Research**, welches sich mit der Optimierung von Produktionskosten unter gegebenen Bedingungen beschäftigt. Heutzutage hat die Optimierung zahlreiche Anwendungen, z.B. in der Pfadplanung, in Computer Vision, in der Bioinformatik, im maschinellen Lernen oder im Hardwaredesign.

### 1.1 Anwendungsbeispiel: Support Vector Machines

Ein klassisches Problem des maschinellen Lernens ist die Klassifizierung von Daten durch eine lineare Entscheidungsgrenze - alle Punkte  $(x_i, y_i)$  über einer Ebene werden einer Klasse zugeordnet, und alle Punkte unter der Ebene einer anderen Klasse. Das Problem besteht daraus, eine Optimale Trennungsebene zu finden. Diese wird beschrieben durch eine Gleichung der Form

$$\hat{y}(x) = w^T x + w_0.$$

Sei  $w^*$  der Vektor  $(w_0, w_1, \dots, w_{n-1})$ . Es stellt sich heraus, dass das zu lösende Optimierungsproblem gegeben ist durch:

$$\begin{array}{ll} \operatorname{argmin}_{w^* \in \mathbb{R}^n} & \frac{1}{2} \|w^*\|^2 \\ \text{s.t.} & \forall i \ y_i \hat{y}(x_i) \geq 1 \end{array}$$

Wir wollen die Länge  $\|w^*\|$  des Vektors  $w^*$  minimieren, da dies zu einem größeren Abstand zwischen unserer Ebene und den Datenpunkten führt, wodurch unser Modell besser generalisiert. Die Nebenbedingung entspricht der Anforderung, dass alle Punkte korrekt klassifiziert werden sollen.

### 1.2 Typische Problemklassen

Optimierungsprobleme können gemäß diverser Kriterien klassifiziert werden:

- Probleme mit Nebenbedingungen vs Probleme ohne Nebenbedingungen

- Optimierung mit Variablen aus verschiedenen Mengen, insbesondere kontinuierliche Variablen vs diskrete Variablen
- Lineare vs nichtlineare Funktionen
- Eindimensionale vs mehrdimensionale Funktionen
- Konvexe Funktionen vs nicht konvexe Funktionen
- Konvexe Mengen vs nicht konvexe Mengen

Diese verschiedenen Problemklassen führen zu unterschiedlich schwierigen Problemen. Insbesondere sind konvexe Probleme einfacher zu lösen als nicht konvexe Probleme, kontinuierliche Probleme sind in der Regel einfacher zu lösen als diskrete Probleme, und lineare Probleme sind einfacher als nichtlineare Probleme. Relevante Fragen sind dann:

- Wie schnell konvergiert das Verfahren zu einer Lösung? Wie viele Iterationen sind nötig? Was ist die Komplexität (in  $O$ -Notation) einer einzelnen Iteration?
- Konvergiert das Verfahren immer gegen ein Globales Optimum? Falls nein, gibt es garantierte obere/untere Schranken für die maximale Abweichung vom globalen Optimum?

### 1.3 Konvexität

Eine Menge  $G$  ist **konvex**, wenn für beliebige Punkte  $x, y \in G$  auch beliebige lineare Interpolationen zwischen den Punkten in der Menge enthalten sind:

$$x, y \in G \implies \{(1 - \alpha)x + \alpha y \mid \alpha \in [0, 1]\} \subset G$$

Intuitiv entspricht das der Forderung, dass zwischen für alle Paare von Punkten eine Verbindungsstrecke zwischen den Punkten in der Menge enthalten sein muss.

Die **konvexe Hülle** einer Menge  $G$  ist die kleinste konvexe Menge  $H$  sodass  $G \subset H$ .

Analog zur Definition einer konvexen Menge ist eine **konvexe Funktion** definiert als eine Funktion, für die gilt:

$$x, y \in G \implies f((1 - \alpha)x + \alpha y) \leq (1 - \alpha)f(x) + \alpha f(y) \quad \forall \alpha \in [0, 1]$$

Dies entspricht der Forderung, dass jede Verbindungsstrecke zwischen zwei Punkten auf dem Graphen unterhalb des Graphen liegen muss.

Eine nicht konvexe Funktion, in der jedes lokale Minimum auch ein globales Minimum ist, wird als **quasikonvex** bezeichnet. Da dies dem Hauptvorteil von konvexen Funktionen in der Optimierung entspricht, ist die Optimierung von quasikonvexen Funktionen ebenfalls einfacher als die Optimierung allgemeiner nichtlinearer Funktionen.

## Chapter 2

# Gradientenverfahren

### 2.1 Wiederholung: Mehrdimensionale Differentiation

Zur Erinnerung: der **Gradient**  $\nabla f$  ist ein Vektor, der alle partiellen Ableitungen von  $f$  enthält, also die Ableitung von  $f$  in alle Koordinatenrichtungen:

$$\nabla f = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^T$$

Im eindimensionalen Fall entspricht das genau der gewöhnlichen Ableitung, da partielle Ableitungen letztendlich genau wie ordinäre Ableitungen berechnet werden. Sei  $f_1(x_1, x_2, x_3)$  eine beliebige Funktion  $\mathbb{R}^3 \rightarrow \mathbb{R}$  und sei  $f_2(x_1) = f_2(x_1, x_2, x_3)$  für beliebige Konstante  $x_2, x_3$ . Dann ist

$$\frac{\partial}{\partial x_1} f_1(x_1, x_2, x_3) = \frac{d}{dx_1} f_2(x_1)$$

Zum Beispiel:

$$\begin{aligned} f_1(x_1, x_2, x_3) &= x_1^2 x_2 + x_1 x_3 + 2x_2 \\ \frac{\partial}{\partial x_1} f_1(x_1, x_2, x_3) &= 2x_1 x_2 + x_3 \end{aligned}$$

entspricht der Ordinären Ableitung

$$\begin{aligned} f_2(x_1) &= x_1^2 x_2 + x_1 x_3 + 2x_2 \\ \frac{d}{dx_1} f_2(x_1) &= 2x_1 x_2 + x_3 \end{aligned}$$

Eine mehrdimensionale Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  kann nicht nur in die Koordinatenrichtungen abgeleitet werden, sondern in beliebige Richtungen  $\mathbf{d} \in \mathbb{R}^n$ . Um Formeln verständlicher zu notieren, werde ich so gut es geht Variablen  $\mathbf{x} \in \mathbb{R}^n$  im Fall  $n > 1$  immer durch fette Buchstaben notieren. Diese **Richtungsableitung** wird geschrieben als  $\nabla_{\mathbf{d}} f(\mathbf{x})$ . Analog zur ordinären Definition der Ableitung ist die Richtungsableitung definiert als:

$$\nabla_{\mathbf{d}} f(\mathbf{x}) = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{d}) - f(\mathbf{x})}{h}$$

$f$  wird als differenzierbar bezeichnet, falls die Richtungsableitung in allen Richtungen existiert. Für ein differenzierbares  $f$  kann die Richtungsableitung mit Hilfe des Gradienten sehr einfach als Skalarprodukt berechnet werden:

$$\nabla_{\mathbf{d}} f(\mathbf{x}) = \nabla f(\mathbf{x})^T \mathbf{d}$$

Im Rahmen dieser Vorlesung gehen wir davon aus, dass  $f \in C^1$ , also dass die zu minimierende Funktion mindestens einmal differenzierbar ist. Im Allgemeinen ist es jedoch wichtig, anzumerken, dass  $f$  nicht unbedingt differenzierbar ist, nur weil partielle Ableitungen in alle Richtungen existieren.

## 2.2 Iterative Optimierung

Eine Funktion  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  kann numerisch durch iteratives "Ausprobieren" verschiedener Werte von  $\mathbf{x}$  optimiert werden. Wir starten mit einem beliebigen Startwert  $\mathbf{x}^{(0)}$ <sup>1</sup>, und bewegen uns dann bei jedem Schritt  $k$  in eine Richtung  $\mathbf{d}^{(k)} \in \mathbb{R}^n$ , welche uns hoffentlich näher zum Optimum bringt. Wir nutzen zusätzlich einen Parameter  $\tau^{(k)} \in \mathbb{R}$ , welcher die Schrittweite beschreibt.

$$\begin{aligned}\mathbf{x}^{(0)} &\in \mathbb{R}^n \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}\end{aligned}$$

In der Praxis wird oft einfach  $\mathbf{x}^{(0)} = \mathbf{0}$  gewählt. Ist jedoch  $f$  nicht quasikonvex, so kann eine schlechte Wahl von  $\mathbf{x}^{(0)}$  dazu führen, dass das Verfahren nur gegen ein lokales Minimum konvergiert, welches oft kein globales Minimum ist und somit in vielen Fällen ein schlechtes Ergebnis.

### 2.2.1 Welche Richtung ist optimal?

Das Optimale  $\mathbf{x}^{(k+1)}$  ist das, das  $f(\mathbf{x}^{(k+1)})$  minimiert. Wir können  $f(\mathbf{x}^{(k+1)})$  für kleine Schrittweiten  $\tau^{(k)}$  folgendermaßen approximieren (**Taylor-Approximation**):

$$f(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}) \approx f(\mathbf{x}^{(k)}) + \tau^{(k)} \nabla f(\mathbf{x}^{(k)}) \mathbf{d}^{(k)}$$

Da wir den Wert von  $f$  möglichst stark verringern wollen, sollte  $\tau^{(k)} \nabla f(\mathbf{x}^{(k)}) \mathbf{d}^{(k)}$  ein negativer Term mit möglichst großem Betrag sein. Dabei können wir  $\tau^{(k)}$  aber nicht zu hoch wählen, da sonst die Approximation ungenau wird.

Nach der Definition des Skalarprodukts gilt:

$$\nabla f(\mathbf{x}^{(k)}) \mathbf{d}^{(k)} = \|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{d}^{(k)}\| \cos(\theta) \quad (2.1)$$

Wobei  $\theta$  der Winkel zwischen  $\nabla f(\mathbf{x}^{(k)})$  und  $\mathbf{d}^{(k)}$  ist. Diese Gleichung ist minimal, wenn  $\cos(\theta) = -1$ , also  $\theta = 180$ . Dementsprechend muss  $\mathbf{d}^{(k)}$  in die umgekehrte Richtung von  $\nabla f(\mathbf{x}^{(k)})$  zeigen, also  $\mathbf{d}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$ .<sup>2</sup>

Mit dieser optimalen Wahl der Abstiegsrichtung  $\mathbf{d}^{(k)}$  erhalten wir das **Gradientenverfahren** (auch bekannt als "Gradientenabstieg", auf Englisch "**Gradient Descent**"):

$$\begin{aligned}\mathbf{x}^{(0)} &\in \mathbb{R}^n \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \tau^{(k)} \nabla f(\mathbf{x}^{(k)})\end{aligned}$$

## 2.3 Bestimmung der Schrittweite

Die Wahl der Schrittweite  $\tau^{(k)}$  ist für das Gradientenverfahren von entscheidender Bedeutung. Ist  $\tau^{(k)}$  zu klein, wird die Konvergenz des Algorithmus oft deutlich verlangsamt. Ist  $\tau^{(k)}$  zu groß, wird die Approximation ungenau. Dies kann dazu führen, dass das Verfahren über das Ziel hinausschießt, und im schlimmsten Fall eventuell gar nicht konvergiert.

Die Bestimmung eines optimalen Werts für  $\tau^{(k)} \in \mathbb{R}$  ist ein eigenes Optimierungsproblem. In einigen Fällen ist ein dauerhaft konstantes  $\tau$  genug, um die Konvergenz zu garantieren. In diesen Fällen ist eine konstante Schrittweite oft effizienter.

<sup>1</sup>Ich schreibe hier  $\mathbf{x}^{(k)}$  mit der  $k$  in Klammern, um klarzustellen, dass  $\mathbf{x}^{(k)}$  nicht " $\mathbf{x}$  hoch  $k$ " ist, sondern "das  $k$ -te  $\mathbf{x}$ ". Eigentlich würde ich das am liebsten als  $\mathbf{x}_k$  notieren, aber da  $\mathbf{x} \in \mathbb{R}^n$  ist, ist diese Notation viel zu leicht mit der Notation  $x_k$  für die verschiedenen Komponenten von  $\mathbf{x}$  zu verwechseln ( $\mathbf{x} = (x_1, \dots, x_n)^T$ ). Vermutlich werde ich an einigen Stellen aus Versehen trotzdem  $\mathbf{x}^k$  schreiben, ich entschuldige mich im Voraus. ('-w-')

<sup>2</sup>Technisch gesehen wird die Gleichung durch  $\mathbf{d}^{(k)} = -\alpha \cdot \nabla f(\mathbf{x}^{(k)})$  mit möglichst großen  $\alpha \in \mathbb{R}$  minimiert. Die "Rolle" dieses Parameters  $\alpha$  wird jedoch bereits durch unseren Schrittweiten-Parameter  $\tau^{(k)}$  abgedeckt, und es gilt auch für dieses theoretische  $\alpha$ , dass die Approximation sehr ungenau wird, wenn wir ein hohes  $\alpha$  wählen.

Um ein optimales  $\tau^{(k)}$  für eine gegebene Iteration  $k$  zu finden, definieren wir eine neue Funktion  $h : \mathbb{R} \rightarrow \mathbb{R}$ :

$$h(\tau) = f(\mathbf{x}^{(k)} + \tau \mathbf{d}^{(k)})$$

$$\tau^{(k)} = \underset{\tau \in \mathbb{R}}{\operatorname{argmin}} h(\tau)$$

Um die Ableitung von  $h$  zu finden, beschreiben wir den Term  $\mathbf{x}^{(k)} + \tau \mathbf{d}^{(k)}$  als eine eigene eindimensionale Funktion  $g(\tau)$ . Dann ist  $h(\tau) = f(g(\tau))$ , gemäß der mehrdimensionalen Kettenregel gilt also:

$$\begin{aligned} \frac{d}{d\tau} h(\tau) &= \frac{d}{d\tau} f(g(\tau)) \\ &= \sum_{i=1}^n \left( \frac{d}{d\tau} g(\tau)_i \right) \frac{\partial}{\partial g(\tau)_i} f(g(\tau)) \\ &= \sum_{i=1}^n d_i^{(k)} \frac{\partial}{\partial g(\tau)_i} f(g(\tau)) \\ &= \nabla f(g(\tau)) \cdot \mathbf{d}^{(k)} \\ &= \nabla f(\mathbf{x}^{(k)} + \tau \mathbf{d}^{(k)}) \mathbf{d}^{(k)} \end{aligned}$$

Eine exakte Bestimmung des Optimalwerts ist mit hohem Aufwand Verbunden und ist oft sowieso nur von geringem Vorteil, da in der Regel auch bei wiederholter optimaler Wahl von  $\tau^{(k)}$  für den letztendlichen Gradientenabstieg viele Schritte nötig sind. Stattdessen werden in der Praxis approximative Lösungen gefunden, welche bestimmte **Qualitätsbedingungen** erfüllen.

Eine einfache Suche nach einem beliebigen Wert, welcher den Funktionswert reduziert, also

$$f(\mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}) - f(\tau^{(k)}) \geq 0,$$

ist dabei nicht genug, da bei hinreichend schneller Konvergenz von  $f(\mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}) - f(\tau^{(k)})$  gegen 0 eventuell nie das Optimum erreicht wird.

### 2.3.1 Die Armijo-Bedingung

Die **Armijo-Bedingung** hat das Ziel, sicherzustellen, dass die Zielfunktion bei einem Schritt um die gegebenen Schrittweite  $\tau^{(k)}$  hinreichend reduziert wird. Die Bedingung ist gegeben durch die folgende Ungleichung:

$$f(\mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}) \leq f(\mathbf{x}^{(k)}) + \delta \tau^{(k)} \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)}$$

Die minimal notwendige Reduktion wird dabei durch einen Parameter  $\delta \in (0, 1)$  gesteuert. Ein üblicher Wert ist dabei  $\delta \approx 10^{-4}$ . Die Multiplikation mit dem Gradienten bedeutet, dass bei einem steileren Gradienten auch nach einer höheren Reduktion gesucht wird.

### 2.3.2 Die Wolfe-Bedingungen

Die **schwache Wolfe-Bedingung** fordert, dass der Gradient an der Zielposition entweder weniger steil als an der Startposition ist oder das Vorzeichen ändert:

$$-\nabla f(\mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)})^T \mathbf{d}^{(k)} \leq -\eta \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)}$$

Das Minus auf beiden Seiten kommt davon, dass bei einer Richtung  $\mathbf{d}^{(k)}$ , welche zu einem Abstieg führt, notwendigerweise  $\nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)}$  negativ ist.

Wir wählen hier  $\eta \in (\delta, 1)$ . In der Regel wird  $\eta$  deutlich größer als  $\delta$  gewählt, Wikipedia zitiert einen Richtwert von  $\eta \approx 0.9$ .

Bei der **starken Wolfe-Bedingung** wird der Betrag der Terme genommen, also gefordert, dass auch bei Vorzeichenwechsel die Steigung nach dem Schritt weniger steil als davor sein soll:

$$|\nabla f(\mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)})^T \mathbf{d}^{(k)}| \leq \eta |\nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)}|$$

Intuitiv verhindern die Wolfe-Bedingungen ineffizient kleine Schritte, bei denen man ohne Probleme in die gleiche Richtung weitergehen könnte.

Schrittweiten, welche sowohl die starke Wolfe-Bedingung als auch die Armijo-Bedingung <sup>3</sup> erfüllen, existieren unter den von uns gemachten Annahmen immer. Als Erinnerung: Wir haben folgende Annahmen gemacht:

- $f(x)$  ist kontinuierlich differenzierbar ( $f(x) \in C^1$ )
- $\mathbf{d}^{(k)}$  ist eine Abstiegsrichtung an der Stelle  $\mathbf{x}^{(k)}$  (also  $\nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} < 0$ )
- $0 \leq \delta \leq \eta \leq 1$

Für die schwache Wolfe-Bedingung muss zusätzlich gegeben sein, dass die Menge  $\{f(\mathbf{x}^{(k)} + \tau \mathbf{d}^{(k)}) \mid \tau > 0\}$  von unten beschränkt ist.

**Beweis**, nach Nocedal und Wright <sup>4</sup> (In der Vorlesung nicht gegeben):

- $h(\tau) = f(\mathbf{x}^{(k)} + \tau \mathbf{d}^{(k)})$  ist nach Annahme für alle  $\tau$  nach unten beschränkt.
- Da  $\mathbf{d}^{(k)}$  eine Abstiegsrichtung und  $\delta$  positiv ist, ist die Linie  $l_{Arm}(\tau) := f(\mathbf{x}^{(k)}) + \delta \tau \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)}$  nach unten unbeschränkt.
- Es gilt außerdem  $h(0) = l_{Arm}(0)$  und  $l_{Arm}(0) = \delta h'(0) < h'(0)$ , also muss für kleine  $\tau$  in Abstiegsrichtung gelten, dass  $h(\tau) < l_{Arm}(\tau)$ .
- Somit muss  $l_{Arm}(\tau)$  den Graphen von  $h(\tau)$  für mindestens einen Wert  $\tau > 0$  schneiden. Sei  $\tau'$  der kleinste Wert, bei dem es das tut.
- Es folgt  $\forall \tau < \tau' : h(\tau) < l_{Arm}(\tau)$ , da sonst  $\tau'$  nicht der minimale Schnittpunkt wäre. Also erfüllen alle  $\tau \in (0, \tau')$  die Armijo-Bedingung.
- Nach dem Mittelwertsatz muss es nun ein  $\tau'' \in (0, \tau')$  geben, sodass

$$f(\mathbf{x}^{(k)} + \tau' \mathbf{d}^{(k)}) - f(\mathbf{x}^{(k)}) = \tau' \nabla f(\mathbf{x}^{(k)} \tau'' \mathbf{d}^{(k)})^T \mathbf{d}^{(k)}$$

- Da  $\tau'$  ein Schnittpunkt ist, gilt nun

$$\begin{aligned} f(\mathbf{x}^{(k)} + \tau' \mathbf{d}^{(k)}) &= f(\mathbf{x}^{(k)}) + \tau' \delta \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} \\ \implies f(\mathbf{x}^{(k)} + \tau' \mathbf{d}^{(k)}) - f(\mathbf{x}^{(k)}) &= \tau' \delta \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} \\ \implies \tau' \nabla f(\mathbf{x}^{(k)} \tau'' \mathbf{d}^{(k)})^T \mathbf{d}^{(k)} &= \tau' \delta \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} \\ \implies \nabla f(\mathbf{x}^{(k)} \tau'' \mathbf{d}^{(k)})^T \mathbf{d}^{(k)} &= \delta \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} \\ \xrightarrow{\eta > \delta} \nabla f(\mathbf{x}^{(k)} \tau'' \mathbf{d}^{(k)})^T \mathbf{d}^{(k)} &> \eta \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} \end{aligned}$$

Also erfüllt  $\tau''$  sowohl die Armijo-Bedingung als auch die schwache Wolfe-Bedingung. Desweiteren ist der Term auf der linken Seite der finalen Ungleichung negativ, also gilt sogar die starke Wolfe-Bedingung. Da die Ungleichungen strikt sind existiert desweiteren ein Intervall um  $\tau''$ , welches die beiden Bedingungen ebenfalls erfüllt.

### 2.3.3 Line Search

Ein einfaches Verfahren zur Bestimmung eines geeigneten  $\tau$  ist die sogenannte **Backtracking Line Search**:

- Definiere  $\tau^{(0)} = 1$

<sup>3</sup>In dieser Vorlesung wurden die starke Wolfe-Bedingung und die Armijo-Bedingung als separate Bedingungen vorgestellt. In der Literatur wird die Armijo-Bedingung oft als Teil der starken Wolfe-Bedingung definiert, mit "Erfüllung der starken Wolfe-Bedingungen" ist dann die Erfüllung beider Bedingungen gemeint.

<sup>4</sup>J. Nocedal, S. J. Wright: Numerical Optimization, Springer, 2006, Seite 35ff



- Falls die Armijo-Bedingung nicht erfüllt ist, probiere als nächstes  $\tau^{(k+1)} = \beta \tau^{(k)}$ , wobei  $\beta \in (0, 1)$ .

Unter Nutzung des Gradienten lässt sich auch ein **Interpolationsverfahren** anwenden, welches schneller konvergiert:

- Definiere  $\tau^{(0)} = 1$
- Betrachte die Quadratische Taylor-Approximation von  $h(\tau)$ :

$$h_q(\tau) = \left( \frac{h(\tau^{(0)}) - h(0) - \tau^{(0)} h'(0)}{\tau^{(0)^2} \right) \tau^2 + h'(0) \tau + h(0)$$

Diese ist minimiert durch

$$\tau^{(1)} = \frac{h'(0) \tau^{(0)^2}{2(h(\tau^{(0)}) - h(0) - h'(0) \tau^{(0)})}$$

- Solange  $\tau^{(k)}$  die Armijo-Bedingung noch nicht erfüllt, bestimme  $\tau^{(k+1)}$  durch “kubische Interpolation mit  $h(0), h'(0), h(\tau^{(0)}), h(\tau^{(1)})$ ”<sup>5</sup>

Line Search Verfahren, welche die Wolfe-Bedingungen sicherstellen, sind komplizierter. Solche Verfahren bestehen aus zwei Komponenten:

- **Bracketing** erweitert das Suchintervall, bis darin geeignete Schrittweiten garantiert werden können.
- **Zooming** reduziert das Suchintervall, bis darin per Interpolationsverfahren eine geeignete Schrittweite gefunden wird.

### 2.3.4 Conjugate Gradient-Verfahren

Angenommen, unsere Zielfunktion ist eine quadratische Funktion in folgender Form:

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{b}^T \mathbf{x} \quad (2.2)$$

Das sogenannte **Conjugate Gradient-Verfahren** funktioniert dann folgendermaßen:

- Wähle als erste Richtung den Gradienten an der Startposition:

$$\mathbf{d}^{(0)} = -\nabla f(\mathbf{x}^{(0)}) = \mathbf{b} - A \mathbf{x}^0$$

- Bestimme die optimale Schrittweite:

$$\tau^{(k)} = \frac{|\nabla f(\mathbf{x}^{(k)})|}{\mathbf{d}^{(k)^T} A \mathbf{d}^{(k)}}$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}$$

- Wähle eine neue Richtung so, dass  $\mathbf{d}^{(k+1)}$  und  $\mathbf{d}^{(k)}$  gemäß der von  $A$  induzierten Metrik orthogonal sind, also  $\mathbf{d}^{(k+1)^T} \cdot A \cdot \mathbf{d}^{(k)} = 0$ . Dies wird garantiert durch:

$$\beta^{(k)} = \frac{|\nabla f(\mathbf{x}^{(k+1)})|^2}{|\nabla f(\mathbf{x}^{(k)})|^2}$$

$$\mathbf{d}^{(k+1)} = -\nabla f(\mathbf{x}^{(k+1)}) + \beta^{(k)} \mathbf{d}^{(k)}$$

Das Minimieren von  $f$  entspricht genau dem Lösen des Gleichungssystems  $A \mathbf{x} = \mathbf{b}$ . Das CG-Verfahren garantiert dadurch für die Matrix  $A = \mathbb{R}^{n \times n}$  eine bei exakter Zahlendarstellung exakte Lösung nach nur  $n$  Schritten. Wenn die **Konditionszahl** von  $A$ , definiert als der größte Eigenwert geteilt durch den kleinsten, klein ist, konvergiert es sogar wesentlich schneller.

Für nicht quadratische Funktionen kann das Conjugate-Gradient verfahren nahezu identisch angewandt werden, allerdings existiert im Allgemeinen für die Schrittweitenbedingung keine analytische Lösung, weshalb wieder Line Search nötig wird.

<sup>5</sup>Die Vorlesung war hier extrem vage - ich vermute, wir interpolieren iterativ immer  $h(0), h'(0), h(\tau^{(k-1)}), h(\tau^{(k)})$ ? Wie genau man eine kubische Interpolation durchführt wurde aber nicht erklärt.

## Chapter 3

# Newton und Quasi-Newton-Verfahren

### 3.1 Das Newton-Verfahren

Das **Newton-Verfahren** ähnelt dem Gradientenverfahren, verwendet jedoch auch die Krümmung der Funktion, also die zweite Ableitung, gegeben durch die Hesse-Matrix  $H_f$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau^{(k)} H_f^{-1} \mathbf{x}^{(k)} \nabla f(\mathbf{x}^{(k)})$$

Dadurch ermöglicht das Newton-Verfahren größere Schritte in Richtungen mit schwacher Krümmung, bei denen die Approximation genauer und somit das "Risiko" geringer ist.

Zur Erinnerung: Die Hesse-Matrix  $H_f$  ist gegeben durch:

$$H_f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

Die Hesse-Matrix ist symmetrisch, da  $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$ .

Während das Gradientenverfahren also mit einer linearen Taylor-Approximation arbeitet, nutzt das Newton-Verfahren eine quadratische Taylor-Approximation:

$$\begin{aligned} \mathbf{m}^{(k)}(\mathbf{d}) &:= f(\mathbf{x}^{(k)} + \mathbf{d}) \approx f(\mathbf{x}^{(k)}) + \mathbf{d}^T \nabla f(\mathbf{x}^{(k)}) + \frac{1}{2} \mathbf{d}^T H_f \mathbf{x}^{(k)} \mathbf{d} \\ \nabla \mathbf{m}^{(k)}(\mathbf{d}) &= \nabla f(\mathbf{x}^{(k)}) + H_f(\mathbf{x}^{(k)}) \mathbf{d} \end{aligned}$$

Für das optimale  $\mathbf{d}$  ist  $\nabla \mathbf{m}^{(k)}(\mathbf{d}) = 0$  und wir erhalten:

$$\begin{aligned} H_f(\mathbf{x}^{(k)}) \mathbf{d} &= -\nabla f(\mathbf{x}^{(k)}) \\ \implies \mathbf{d} &= -H_f^{-1} \mathbf{x}^{(k)} \nabla f(\mathbf{x}^{(k)}) \end{aligned}$$

Die Optimale Schrittweite  $\tau^{(k)}$  muss auch beim Newton-Verfahren durch Line Search bestimmt werden.

### 3.2 Konvergenzordnungen

Sei  $s$  eine beliebige Folge und  $\bar{s} := \lim_{k \rightarrow \infty} s_k$ . Wir betrachten nun die Konvergenzordnung der Folge  $s$ , also die Geschwindigkeit, mit der die Folge konvergiert. In dieser Vorlesung unterscheiden wir folgende Konvergenzordnungen:

- $s$  konvergiert **linear**, wenn:

$$\lim_{k \rightarrow \infty} \frac{|s_{k+1} - \bar{s}|}{|s_k - \bar{s}|} < 1$$

ein Beispiel für eine linear konvergierende Folge ist  $s_k = 0.9^k$ . Das Gradientenverfahren konvergiert ebenfalls linear.

- $s$  konvergiert **superlinear**, wenn:

$$\lim_{k \rightarrow \infty} \frac{|s_{k+1} - \bar{s}|}{|s_k - \bar{s}|} = 0$$

ein Beispiel für eine superlinear konvergierende Folge ist  $s_k = \frac{1}{k!}$ .

- $s$  konvergiert **quadratisch**, wenn:

$$\lim_{k \rightarrow \infty} \frac{|s_{k+1} - \bar{s}|}{|s_k - \bar{s}|^2} < 1$$

ein Beispiel für eine superlinear konvergierende Folge ist  $s_k = 0.9^{(2^k)}$ .

Das Newton-Verfahren konvergiert für glatte, konvexe Funktionen nahe der Lösung quadratisch. Nominal ist es also wesentlich effizienter als Gradientenabstieg. Dabei ist jedoch zu berücksichtigen, dass die Berechnung und Invertierung der Hesse-Matrix wesentlich aufwendiger ist als die Berechnung des Gradienten. Demenssprechend ist ein einzelner Iterationsschritt des Newton-Verfahrens auch wesentlich aufwendiger als ein einzelner Schritt des Gradientenverfahrens.

### 3.2.1 Konvergenzradius des Newton-Verfahren

Damit  $\mathbf{d}^{(k)} = -H_f^{-1}(\mathbf{x}^{(k)}) \nabla f(\mathbf{x}^{(k)})$  eine Abstiegsrichtung ist, muss  $H_f$  positiv definit sein, also die Krümmung positiv. Bei negativer Krümmung läuft das Newton-Verfahren in die falsche Richtung! Dies ist kein Problem für konvexe Funktionen, da diese überall positive Krümmung haben. Bei nicht konvexen Funktionen muss die Hesse-Matrix erst manipuliert werden, damit alle ihre Eigenwerte positiv werden.

## 3.3 Quasi-Newton Verfahren

Die Idee hinter Quasi-Newton-Verfahren ist es, die Hesse-Matrix mithilfe der Gradienten sukzessiv zu approximieren. Man ersetzt also in der Iterationsformel des Newton-Verfahrens die Matrix  $H_f$  durch eine Approximation  $B$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \tau^{(k)} (B^k)^{-1} \nabla f(\mathbf{x}^{(k)})$$

Die einfachste Möglichkeit ist eine lineare Taylor-Approximation des Gradienten. Sei  $\mathbf{s} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ . Dann haben wir:

$$\begin{aligned} \nabla f(\mathbf{x}^{(k)} + \mathbf{s}) &\approx \nabla f(\mathbf{x}^{(k)}) + H_f(\mathbf{x}^{(k)}) \mathbf{s} \\ \implies H_f(\mathbf{x}^{(k)}) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) &\approx \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) \end{aligned}$$

Wir suchen also in jedem Schritt eine Matrix  $B^{(k)}$ , sodass:

$$B^{(k)} (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) = \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) \quad (3.1)$$

Abgekürzt schreiben wir auch:

$$B^{(k)} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}$$

Aus  $B^{(k)} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}$  folgt  $\mathbf{s}^{(k)T} B^{(k)} \mathbf{s}^{(k)} = \mathbf{s}^{(k)T} \mathbf{y}^{(k)}$ . Wir gehen von einer konvexen Funktion  $f$  aus, also ist  $B^{(k)}$  positiv definit. Daraus folgt  $\mathbf{s}^{(k)T} B^{(k)} \mathbf{s}^{(k)} > 0$ , damit gilt auch  $\mathbf{s}^{(k)T} \mathbf{y}^{(k)} > 0$ . Diese Ungleichung können wir als Bedingung für die Schrittweitensuche nutzen.

Um die Anzahl an möglichen  $B^{(k)}$  einzuschränken, können wir als weitere Bedingung fordern, dass  $B^{(k)}$  in jedem Schritt möglichst ähnlich zur letzten Iteration sein soll:

$$B^{(k)} = \min_B \|B - B^{(k-1)}\|$$

$\|\cdot\|$  kann eine beliebige Matrixnorm sein. An diesem Punkt führen verschiedene Normen zu verschiedenen Unterverfahren. Ein beliebtes Verfahren ist das **Davidon-Fletcher-Powell Verfahren**, welches eine sog. gewichtete Frobeniusnorm benutzt.

### 3.3.1 Das Davidon-Fletcher-Powell Verfahren

Beim Davidon-Fletcher-Powell (DFPV)-Verfahren werden die Approximation  $B^{k+1}$  und ihr Inverses  $Q^{(k+1)}$  direkt aus den vorherigen Werten berechnet:

$$\begin{aligned}\rho^{(k)} &= \frac{1}{(\mathbf{y}^{(k)})^\top \mathbf{y}^{(k)}} \\ B^{k+1} &= (I - \rho^{(k)} \mathbf{y}^{(k)} (\mathbf{s}^{(k)})^\top) B^{(k)} (I - \rho^{(k)} \mathbf{s}^{(k)} (\mathbf{y}^{(k)})^\top) + \rho^{(k)} \mathbf{y}^{(k)} (\mathbf{y}^{(k)})^\top \\ Q^{(k+1)} &= Q^{(k)} - \frac{Q^{(k)} \mathbf{y}^{(k)} (\mathbf{y}^{(k)})^\top Q^{(k)}}{(\mathbf{y}^{(k)})^\top Q^{(k)} \mathbf{y}^{(k)}} + \frac{\mathbf{s}^{(k)} (\mathbf{s}^{(k)})^\top}{(\mathbf{y}^{(k)})^\top \mathbf{s}^{(k)}}\end{aligned}$$

### 3.3.2 Das Broyden-Fletcher-Goldfarb-Shanno Verfahren

Das Broyden-Fletcher-Goldfarb-Shanno (BFGS)-Verfahren ist eine etwas effizientere Variante des DFPV-Verfahrens, bei dem  $Q^{(k+1)}$  direkt folgendermaßen berechnet wird:

$$Q^{k+1} = (I - \rho^{(k)} \mathbf{s}^{(k)} (\mathbf{y}^{(k)})^\top) Q^{(k)} (I - \rho^{(k)} \mathbf{y}^{(k)} (\mathbf{s}^{(k)})^\top) + \rho^{(k)} \mathbf{y}^{(k)} (\mathbf{y}^{(k)})^\top$$

$Q^{k+1}$  ist immer positiv definit, wenn  $Q^{(k)}$  positiv definit war.

Das volle BFGS-Verfahren verläuft dann also folgendermaßen:

1.  $\mathbf{x}^{(0)} \in \mathbb{R}$ ,  $Q^{(0)} \in \mathbb{R}^{n \times n}$  (z.B.  $Q^{(0)} = I_n$ )
2.  $\mathbf{d}^{(k)} = Q^{(k)} \nabla f(\mathbf{x}^{(k)})$
3. Bestimme optimale Schrittweite  $\tau^{(k)}$  durch Line Search mit Wolfe-Bedingungen. Es sollte immer zuerst  $\tau^{(k)} = 1$  getestet werden, da diese Schrittweite zur schnellsten Konvergenz führt und in den meisten Iterationen in der Praxis akzeptiert wird.
4. Berechne neue Lösung  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \tau^{(k)} \mathbf{d}^{(k)}$
5. Berechne neue Approximation der inversen Hesse Matrix:

$$\begin{aligned}\mathbf{s}^{(k)} &= \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \\ \mathbf{y}^{(k)} &= \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) \\ \rho^{(k)} &= \frac{1}{(\mathbf{y}^{(k)})^\top \mathbf{y}^{(k)}} \\ Q^{k+1} &= (I - \rho^{(k)} \mathbf{s}^{(k)} (\mathbf{y}^{(k)})^\top) Q^{(k)} (I - \rho^{(k)} \mathbf{y}^{(k)} (\mathbf{s}^{(k)})^\top) + \rho^{(k)} \mathbf{y}^{(k)} (\mathbf{y}^{(k)})^\top\end{aligned}$$

Die Vorlesung gibt folgende Iterationszahlen für ein Beispielpproblem aus der Praxis:

- Gradientenabstieg: 5264 Iterationen
- BFGS: 34 Iterationen
- Newton: 21 Iterationen (Aber wie bereits gesagt ist dafür jede Iteration sehr viel teurer als bei BFGS)

## 3.4 Least-Squares-Aufgaben

Eine sehr häufige Klasse von Zielfunktionen, welche insbesondere in der Statistik und somit auch im maschinellen Lernen wichtig sind, ist gegeben durch **Least Squares**, auch bekannt als  $L_2$ -Loss:

$$f(\mathbf{x}) = \frac{1}{2} \sum_{j=1}^m r_j(\mathbf{x})^2$$

Wobei  $r_j$  in der Regel ein Klassifikationsfehler für einen Datenpunkt  $x_j, y_j$  ist, also:

$$r_j(\mathbf{x}) = (\hat{y}_j(\mathbf{x}) - y_j) \in \mathbb{R}$$

Wenn wir die Residuen  $r_j$  als einen **Residuenvektor**  $\mathbf{r}(\mathbf{x}) = (r_1(\mathbf{x}), \dots, r_m(\mathbf{x}))^\top$  schreiben, erhalten wir eine gekürzte Schreibweise des Problems:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{r}(\mathbf{x})\|^2$$

Der Residuenvektor ist eine Funktion  $\mathbf{r} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , also eine Funktion mit mehrdimensionaler Bildmenge. Die Ableitung einer solchen Funktion ist gegeben durch die sog. **Jacobimatrix**:

$$J_{\mathbf{r}}(\mathbf{x}) = \begin{pmatrix} \nabla r_1(\mathbf{x}^\top) \\ \vdots \\ \nabla r_m(\mathbf{x}^\top) \end{pmatrix} \in \mathbb{R}^{m \times n}$$

Im Least-Squares-Problem lässt sich damit der Gradient von  $f(\mathbf{x})$  sehr kompakt darstellen:

$$\nabla f(\mathbf{x}) = \frac{1}{2} \sum_{j=1}^m 2r_j(\mathbf{x}) \nabla r_j(\mathbf{x}) = J_{\mathbf{r}}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})$$

Die Hesse-Matrix ist dann gegeben durch:

$$\begin{aligned} H_f(\mathbf{x}) &= \sum_{j=1}^m \nabla r_j(\mathbf{x}) \nabla r_j(\mathbf{x})^\top + \sum_{j=1}^m r_j(\mathbf{x}) \nabla^2 r_j(\mathbf{x}) \\ &= J(\mathbf{x})^\top J(\mathbf{x}) + \sum_{j=1}^m r_j(\mathbf{x}) \nabla^2 r_j(\mathbf{x}) \end{aligned}$$

### 3.4.1 Lineare Residuen

Ist die Residuenfunktion  $\mathbf{r}(\mathbf{x})$  linear, dann ist  $\mathbf{r}(\mathbf{x}) = J\mathbf{x} - \mathbf{b}$  mit  $J \in \mathbb{R}^{m \times n}$ . In diesem Fall lässt sich das Ergebnis weiter vereinfachen: Wir haben dann  $J_{\mathbf{r}}(\mathbf{x}) = J$ , und der zweite Teil der Hesse-Funktion fällt weg. Dann gilt also ganz einfach

$$H_f(\mathbf{x}) = J^\top J.$$

Das Problem der Optimierung von  $f$  ist dann durch folgendes lineares Gleichungssystem beschrieben:

$$\begin{aligned} \nabla f(\mathbf{x}) &= 0 \\ \implies J^\top \mathbf{r}(\mathbf{x}) &= 0 \\ \implies J^\top (J\mathbf{x} - \mathbf{b}) &= 0 \\ \implies J^\top J\mathbf{x} &= J^\top \mathbf{b} \end{aligned}$$

### 3.4.2 Nichtlineare Residuen

Für nichtlineare Residuen bietet es sich an, die Hesse-Matrix einfach durch der ersten Term zu approximieren:

$$\begin{aligned} H_f(\mathbf{x}) &= J(\mathbf{x})^\top J(\mathbf{x}) + \sum_{j=1}^m r_j(\mathbf{x}) \nabla^2 r_j(\mathbf{x}) \\ &\approx J(\mathbf{x})^\top J(\mathbf{x}) \end{aligned}$$

Dies entspricht der in dieser Vorlesung bereits oft verwendeten linearen Taylor-Approximation einer Funktion. Anders als bei BFGS wird hier die spezifische Problemstruktur von Summen quadrierter

Residuen ausgenutzt. Wir können die Newtonsche Gleichung von ganz am Anfang nutzen, um ein simples lineares Gleichungssystem zu erhalten:

$$\begin{aligned} H_f(\mathbf{x}^{(k)})\mathbf{d} &= -\nabla f(\mathbf{x}^{(k)}) = -J_r(\mathbf{x}^{(k)})^\top \mathbf{r}(\mathbf{x}) \\ \implies J_r(\mathbf{x}^{(k)})^\top J_r(\mathbf{x}^{(k)})\mathbf{d}^{(k)} &= -J_r(\mathbf{x}^{(k)})^\top \mathbf{r}(\mathbf{x}) \end{aligned}$$

Die Lösung ist dann eine optimale Abstiegsrichtung  $\mathbf{d}^{(k)}$ . Dieses Verfahren ist bekannt als **Gauss-Newton-Verfahren**, da es das Newton-Verfahren durch das Lösen einer Sequenz von linearen Gleichungssystem approximiert (Welche dann wiederum durch das Gauss-Verfahren gelöst werden können). Das Verfahren funktioniert genau dann gut, wenn alle Eigenwerte von  $J^\top J$  deutlich positiv sind. Bei Eigenwerten Nahe null sind andere Verfahren notwendig, welche in dieser Vorlesung nicht besprochen werden, z.B. das Verfahren von Levenberg-Marquardt.

## Chapter 4

# Optimierung mit Nebenbedingungen

Bisher war der Suchraum für ein Optimum einer Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  immer der gesamte Urbildraum  $\mathbb{R}^n$ . In vielen Fällen ist der Suchraum aber eine Teilmenge  $G \subset \mathbb{R}^n$ , welche eingeschränkt ist durch Nebenbedingungen wie " $x_1 \leq 5$ ". Die Lösung eines Problems ohne Nebenbedingung gibt offensichtlich immer eine untere Schranke für die tatsächliche Lösung, da natürlich jede Lösung des Problems mit Nebenbedingungen auch eine Lösung des Problems ohne Nebenbedingungen ist. Oft, aber nicht immer, liegt das Optimum auf dem Rand von  $G$ . Wir gehen in der Regel davon aus, dass  $G$  Kompakt ist, da sonst in diesem Fall kein tatsächliches Optimum existiert.

Ein Optimierungsproblem wird im Allgemeinen geschrieben als:

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} & f(x) \\ \text{subject to} & c_i(x) = 0 \quad \forall i \in \mathcal{E} \\ & c_i(x) \geq 0 \quad \forall i \in \mathcal{I} \end{array}$$

### 4.1 Optimierung mit Gleichheitsbedingungen

Die  $m$  Gleichheitsbedingungen sind zusammengefasst durch eine vektorwertige Funktion  $\mathbf{c}(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  gegeben. Es gilt dementsprechend genau dann  $\mathbf{c}(\mathbf{x}) = \mathbf{0}$ , wenn alle Gleichheitsbedingungen erfüllt sind. Ein Schritt in Richtung  $\mathbf{d}$  ist dementsprechend genau dann erlaubt, wenn  $\nabla \mathbf{c}^\top \mathbf{d} = 0$  und  $\nabla f^\top \mathbf{d} \leq 0$ , also wenn  $\mathbf{d}$  eine Abstiegsrichtung ist und wenn sich der Wert von  $\mathbf{c}(\mathbf{x})$  nach dem Schritt nicht ändert (und dementsprechend die Gleichheitsbedingungen danach immer noch erfüllt sind).

Es stellt sich heraus, dass in einem Extrempunkt die Gradienten  $\nabla f$  und  $\nabla c$  immer parallel sein müssen, also

$$\nabla f(\mathbf{x}) = \lambda(\mathbf{x}) \quad \lambda \in \mathbb{R}^m$$

Ist diese Bedingung nicht gegeben, so existiert ein Schritt, welcher  $f$  ändert, aber nicht  $c$ , und somit ist die momentane Position  $\mathbf{x}$  kein Extremum.

Diese Bedingung können wir nun ausnutzen. Angenommen,  $f$  ist eine Funktion  $\mathbb{R}^n \rightarrow \mathbb{R}$  und  $c$  ist eine Funktion  $\mathbb{R}^n \rightarrow \mathbb{R}^m$ . Dann ist die **Lagrange-Funktion**  $\mathcal{L}$  des Problems die folgende Funktion  $\mathbb{R}^{n+m} \rightarrow \mathbb{R}$ :

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \boldsymbol{\lambda} \mathbf{c}(\mathbf{x})$$

Die Schreibweise " $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ " fand ich persönlich sehr gewöhnungsbedürftig und anfangs sehr verwirrend. Es ist wichtig, sich zu merken, dass  $\mathbf{x}$  selbst ein  $n$ -dimensionaler Vektor ist. Letzendlich "erfinden" wir quasi eine neue Menge an Eingabevariablen, gegeben durch den Vektor  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)$ , welchen wir hinten an den "Eingabevektor"  $\mathbf{x} = (x_1, \dots, x_n)^\top$  hängen. So erhalten wir einen neuen Vektor  $\mathbf{x}' = (x_1, \dots, x_n, \lambda_1, \dots, \lambda_m)^\top$ , welchen wir dann in die Lagrange-Funktion einsetzen. Eine intuitivere, und arguably korrektere, Schreibweise dafür wäre  $\mathcal{L}(x_1, \dots, x_n, \lambda_1, \dots, \lambda_m)$ . Da das jedoch den meisten

Mathematikern zu viel Schreibarbeit ist wird in der Praxis nunmal  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$  geschrieben. Die "Original-Variablen"  $x_i$  werden oft als die **primale Variablen** bezeichnet, die neuen Variablen  $\lambda_i$  hingegen als die **Lagrange-Multiplikatoren**.

Die Bedingung für ein Minimum von  $f$  kann nun mithilfe der Lagrange-Funktion ausgedrückt werden als:

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \nabla f(\mathbf{x}) - \boldsymbol{\lambda} \nabla \mathbf{c}(\mathbf{x}) = 0$$

Wir können also nun unser Optimierungsproblem mit Gleichheitsbedingungen in ein äquivalentes Optimierungsproblem ohne Nebenbedingungen umformen. Wenn für eine Komponente  $i$  diese Bedingung bereits durch  $\lambda_i = 0$  erfüllt ist, dann ist  $c_i$  eine **Nebenbedingung ohne Relevanz**, d.h. das Minimum unter Einschränkung durch die Nebenbedingung  $c_i$  ist dasselbe wie das Minimum ohne die Nebenbedingung. Man sagt in diesem Fall auch, dass die Nebenbedingung **inaktiv** ist. Die Menge der aktiven Nebenbedingungen wird auch als die **aktive Menge** bezeichnet.

## 4.2 Optimierung mit Ungleichheitsbedingungen

Auch hier ist die Parallelität der Gradienten eine notwendige Bedingung. Allerdings zählt nun das Vorzeichen von  $\boldsymbol{\lambda}$ , da wir  $\mathbf{c}$  am Rand zwar nicht verringern können, aber beliebig erhöhen. Wir erhalten also die Bedingung

$$f(\mathbf{x}) = \boldsymbol{\lambda} \mathbf{c}(\mathbf{x}) \quad \boldsymbol{\lambda} \in \mathbb{R}_+^m$$

bzw.

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \nabla f(\mathbf{x}) - \boldsymbol{\lambda} \nabla \mathbf{c}(\mathbf{x}) = 0, \quad \forall i \lambda_i > 0$$

Zusätzlich gilt am Optimum immer

$$\boldsymbol{\lambda} \mathbf{c}(\mathbf{x}) = 0,$$

da sich für jedes  $i$  entweder das Optimum auf dem Rand befindet, also  $c_i(\mathbf{x}) = 0$ , oder die Nebenbedingung inaktiv ist, also  $\lambda_i = 0$ . Wir sagen, eine Nebenbedingung ist **streng komplementär**, wenn immer nur genau einer dieser beiden Fälle auftritt.

Zusammenfassend sind die **notwendigen Bedingungen** für ein lokales Optimum die folgenden:

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= 0 \\ c_i(\mathbf{x}) &= 0 \quad \forall i \in \mathcal{E} \\ \lambda_i c_i(\mathbf{x}) &= 0 \quad \forall i \in \mathcal{I} \\ c_i(\mathbf{x}) &\geq 0 \quad \forall i \in \mathcal{I} \\ \lambda_i &\geq 0 \quad \forall i \in \mathcal{I} \end{aligned}$$

Diese Bedingungen sind bekannt als die **Karush-Kuhn-Tucker-Bedingungen (KKT)**. Später werden wir sehen, dass wir durch die Einführung weiterer Lagrange-Multiplikatoren Ungleichheitsbedingungen als weitere Gleichheitsbedingungen ausdrücken können. So werden wir letztendlich bei der Folgenden deutlich übersichtlicheren Liste landen:

$$\begin{aligned} \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= \mathbf{0} \\ c_i(\mathbf{x}) &= 0 \quad \forall i \in (1, \dots, m), \end{aligned}$$

oder noch genauer:

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})}{\partial x_i} &= 0 \quad \forall i \in (1, \dots, n) \\ c_i(\mathbf{x}) &= 0 \quad \forall i \in (1, \dots, m) \end{aligned}$$



### 4.3 Dualität

Mithilfe der Lagrange-Funktion können wir nun aus dem Originalproblem (auch **primales Problem** genannt)

$$\begin{array}{ll} \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} & f(x) \\ \text{subject to} & c_i(\mathbf{x}) = 0 \quad \forall i \in \mathcal{E} \\ & c_i(\mathbf{x}) \geq 0 \quad \forall i \in \mathcal{I} \end{array}$$

ein sogenanntes **duales Problem** erhalten:

$$q(\boldsymbol{\lambda}) := \inf_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}), \forall i : \lambda_i \geq 0$$

$$\max_{\boldsymbol{\lambda} \in \mathbb{R}^m} q(\boldsymbol{\lambda})$$

Intuitiv kann man sich vorstellen, dass für ein festes  $\mathbf{x}$  die Lagrange-Funktion einen höheren Wert hat, je mehr Bedingungen erfüllt sind. Somit führt die Maximierung der dualen Funktion auch zu einer Optimierung der Originalfunktion. Die duale Problemlösung ist hauptsächlich dann sinnvoll, wenn es eine analytische Lösung zur Maximierung von  $q(\boldsymbol{\lambda})$  gibt, und/oder wenn im dualen Problem keine Nebenbedingungen mehr existieren.

**Lemma:** Für beliebige  $\hat{\boldsymbol{\lambda}} \in \mathbb{R}^m$ ,  $\hat{\mathbf{x}} \in \mathbb{R}^n$  gilt  $q(\hat{\boldsymbol{\lambda}}) \leq f(\hat{\mathbf{x}})$ .

**Beweis:**

$$q(\hat{\boldsymbol{\lambda}}) = \inf_{\mathbf{x}} f(\mathbf{x}) - \hat{\boldsymbol{\lambda}}^\top \mathbf{c}(\mathbf{x}) \leq f(\hat{\mathbf{x}}) - \hat{\boldsymbol{\lambda}}^\top \mathbf{c}(\hat{\mathbf{x}}) \leq f(\hat{\mathbf{x}})$$

Also erhalten wir eine untere Schranke für  $f(\mathbf{x})$ , indem wir  $q(\hat{\boldsymbol{\lambda}})$  maximieren.

**Lemma:**  $q(\boldsymbol{\lambda})$  ist konkav.

**Beweis:**

$$\begin{aligned} \mathcal{L}(\mathbf{x}, (1-\alpha)\boldsymbol{\lambda}_0 + \alpha\boldsymbol{\lambda}_1) &= f(\mathbf{x}) - (1-\alpha)\boldsymbol{\lambda}_0^\top \mathbf{c}(\mathbf{x}) - \alpha\boldsymbol{\lambda}_1^\top \mathbf{c}(\mathbf{x}) \\ &= (1-\alpha)f(\mathbf{x}) + \alpha f(\mathbf{x}) - (1-\alpha)\boldsymbol{\lambda}_0^\top \mathbf{c}(\mathbf{x}) - \alpha\boldsymbol{\lambda}_1^\top \mathbf{c}(\mathbf{x}) \\ &= (1-\alpha)\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_0) + \alpha\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_1) \end{aligned}$$

Für beliebige Funktionen  $f, g$  gilt im Allgemeinen  $\inf_{\mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) \leq \inf_{\mathbf{x}} f(\mathbf{x}) + \inf_{\mathbf{x}} g(\mathbf{x})$ , also:

$$\inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, (1-\alpha)\boldsymbol{\lambda}_0 + \alpha\boldsymbol{\lambda}_1) \leq (1-\alpha) \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_0) + \alpha \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_1). \quad \square$$

In vielen Fällen ist das duale Problem einfacher zu lösen als das primale Problem. Das duale Problem liefert eine untere Schranke für die Lösung des Originalproblems:

$$\inf_{\mathbf{x}} f(\mathbf{x}) - \boldsymbol{\lambda}^\top \mathbf{c}(\mathbf{x}) \leq f(\mathbf{x}) - \boldsymbol{\lambda}^\top \mathbf{c}(\mathbf{x}) \leq f(\mathbf{x})$$

Wir sprechen von **starker Dualität**, wenn das Optimum des dualen Problems auch das primale Problem minimiert. Dies ist der Fall unter folgenden Bedingungen:

- $f(\mathbf{x})$  ist konvex
- Die gültige Menge  $G$  ist konvex
- $\mathcal{L}(\mathbf{x}, \hat{\boldsymbol{\lambda}})$  ist **streng** konvex, wobei  $\hat{\boldsymbol{\lambda}}$  die Lösung des dualen Problems ist.

Für ein "gut benommenes" Problem ist die Lösung des dualen Problems äquivalent zum Wert der konvexen Entspannung des primalen Problems, also zu dem Problem, was wir erhalten, wenn wir das feasible

Set  $G$  durch seine konvexe Hülle und die Funktion  $f(\mathbf{x})$  durch ihre konvexe Schließung ersetzen<sup>1</sup>. "Gut benommen" bedeutet in diesem Kontext, dass das Minimum die KKT-Bedingungen erfüllt - Gegenbeispiele sind hier oft Probleme mit seltsamen Feasible Sets, bei denen z.B. Ungleichheitsbedingungen existieren, welche nur einen einzigen Punkt zulassen. Es gibt eine lange Reihe an sogenannten Regularitätsbedingungen, die zu dieser Bedingung äquivalent sind. Einige Beispiele sind:

- **Linearity Constraint Qualification (LCQ)**: Alle  $c_i$  sind affin.
- **Linear Independence Constraint Qualification (LICQ)**: Die Gradienten der Ungleichheitsbedingungen und die Gradienten der Gleichheitsbedingungen sind am Minimum linear unabhängig.
- **Slater's Condition**: In einem Problem mit konvexem  $f$ , konvexen Ungleichheitsbedingungen und affinen Gleichheitsbedingungen existiert ein Punkt  $\mathbf{x}$ , welcher alle Gleichheitsbedingungen erfüllt und alle Ungleichheitsbedingungen strikt erfüllt (also  $c_i(\mathbf{x}) > 0$  anstatt nur  $c_i(\mathbf{x}) \geq 0$ ).
- etc. etc.

### 4.3.1 Beispiel: Optimierung mit Hilfe des dualen Problems

Gegeben sei folgendes konvexes Problem:

$$\begin{array}{ll} \underset{x_1, x_2 \in \mathbb{R}}{\operatorname{argmin}} & \frac{1}{2}(x_1^2 + x_2^2) \\ \text{s.t.} & x_1 - 1 \geq 0 \end{array}$$

Die Lagrange-Funktion ist dann:

$$\mathcal{L}(x_1, x_2, \lambda) = \frac{1}{2}(x_1^2 + x_2^2) - \lambda(x_1 - 1)$$

Nötige Bedingung ist, dass die partiellen Ableitung von  $\mathcal{L}$  Null sind, also:

$$\begin{aligned} \frac{\partial \mathcal{L}(x_1, x_2, \lambda)}{\partial x_1} &= x_1 - \lambda = 0 \\ \frac{\partial \mathcal{L}(x_1, x_2, \lambda)}{\partial x_2} &= x_2 = 0 \end{aligned}$$

Wir setzen diese Werte in die Lagrange-Funktion ein und erhalten:

$$\mathcal{L}(x_1, x_2, \lambda) = \frac{1}{2}(x_1^2 + x_2^2) - \lambda(x_1 - 1) = \frac{1}{2}(\lambda^2 + 0^2) - \lambda(\lambda - 1) = -\frac{1}{2}\lambda^2 + \lambda$$

Das Maximum dieser Funktion ist  $\lambda = 1$ . Unsere Lösung ist also  $x_1 = \lambda = 1, x_2 = 0$ .

### 4.3.2 Wichtige Probleme mit starker Dualität

**Lineare Programme** sind Probleme der Form:

$$\begin{array}{ll} \min_{\mathbf{x}} & \mathbf{w}^\top \mathbf{x} \\ \text{s.t.} & A\mathbf{x} - \mathbf{b} \geq \mathbf{0} \end{array}$$

Das duale Problem eines linearen Programms ist:

$$\max_{\boldsymbol{\lambda}} \mathbf{b}^\top \boldsymbol{\lambda}, A^\top \boldsymbol{\lambda} = \mathbf{w}, \boldsymbol{\lambda} \geq \mathbf{0}$$

---

<sup>1</sup>"Duality Gap", Wikipedia

**Quadratische Programme** sind Probleme der Form:

$$\begin{array}{ll}\min_{\mathbf{x}} & \frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \mathbf{w}^\top \mathbf{x} \\ \text{s.t.} & A\mathbf{x} - \mathbf{b} \geq \mathbf{0}\end{array}$$

Das duale Problem eines quadratischen Programms ist:

$$\max_{\boldsymbol{\lambda}} -\frac{1}{2}(\mathbf{A}^\top \boldsymbol{\lambda} - \mathbf{w})^\top Q^{-1}(\mathbf{A}^\top \boldsymbol{\lambda} - \mathbf{w}) + \mathbf{b}^\top \boldsymbol{\lambda}, \boldsymbol{\lambda} \geq \mathbf{0}$$