

# Optimization for AI

*Rade Kutil, 14. November 2024*

# 1 Einführung

Optimierung wird in der Mathematik oft auch als *Programm/Programming* bezeichnet. Die generelle Definition von Optimierung einer Funktion  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  ist:

$$\mathbf{x}^* = \operatorname{arg\,min}_{\mathbf{x} \in X} f(\mathbf{x}), \quad X \subseteq \mathbb{R}^n$$

Es gibt mehrere Unterscheidungskriterien, jeweils von einfach bis schwierig:

- Art der Funktion  $f$

**Linear**  $f(\mathbf{x}) = \mathbf{Ax}$ . Optimierung macht nur in Zusammenhang mit Constraints Sinn.

**Least Squares**  $f(\mathbf{x}) = \|\mathbf{Ax} + \mathbf{b}\|_2^2$ .

**Quadratische Programmierung**  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{Ax} + \mathbf{b}^\top \mathbf{x}$ .

**Konvexe Optimierung**  $f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y})$ . Auch die Definitionsmenge  $X$  muss üblicherweise konvex sein.

**Generell nichtlinear** Ein lokales Minimum ist nicht automatisch ein globales Minimum.

- Verfügbarkeit von Gradienten

**Keine Ableitung** Nur Funktionswerte verfügbar.

**Erste Ableitung** Gradient-Descent.

**Zweite Ableitung** Hesse-Matrix. Newton-Methode.

- Constraints (Nebenbedingungen)

**Unconstrained**

**Equality Constraints**  $g_i(\mathbf{x}) = 0$ . Lagrange-Multiplikator.

**Inequality Constraints**  $h_i(\mathbf{x}) \geq 0$ . Karush-Kuhn-Tucker (KKT).

All diese Möglichkeiten gibt es in vielen Kombinationen. Daher ist es schwierig, eine vernünftige Reihenfolge aller Methoden zu finden.

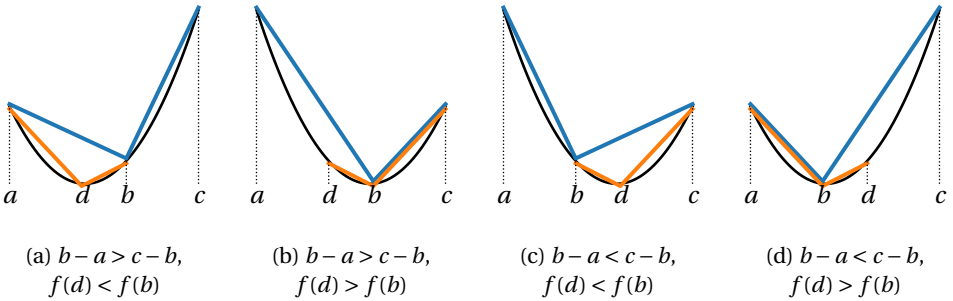


Abbildung 1: Golden-Section-Search. Blau: altes Tripel  $(a, b, c)$ , orange: neues Tripel.

## 2 Golden-Section-Search

Eindimensional, konvex, keine Ableitung, unconstrained.

Dieses Verfahren ist ähnlich zur Nullstellensuche mittels Binary-Search (Bisektionsverfahren). Dabei wird ein Intervall  $[a, b]$  mit entweder  $f(a) < 0$  und  $f(b) > 0$  (oder umgekehrt) in zwei Hälften  $[a, \frac{a+b}{2}]$  und  $[\frac{a+b}{2}, b]$  geteilt. Wenn  $f(\frac{a+b}{2}) > 0$  ist, wird mit ersterem Intervall rekursiv fortgefahren, ansonsten mit zweiterem (oder umgekehrt).

Wenn nicht eine Nullstelle sondern ein Minimum gesucht wird, reicht ein Intervall mit einem Mittelpunkt nicht aus, weil aus  $f(\frac{a+b}{2}) < f(a) < f(b)$  nicht ersichtlich ist, ob das Minimum links oder rechts von  $\frac{a+b}{2}$  liegt. Stattdessen braucht man ein Tripel  $(a, b, c)$  mit  $a < b < c$ ,  $f(b) < f(a)$ ,  $f(b) < f(c)$ .

Es gibt nun zwei Möglichkeiten:  $b - a > c - b$ , also das linke Teilintervall ist größer, oder umgekehrt. Wir betrachten vorerst nur die erste Möglichkeit. Es wird nun ein  $d$  mit  $a < d < b$  gewählt. Ist nun  $f(d) < f(b)$  wird mit dem Tripel  $(a, d, b)$  fortgefahren (Abbildung 1 (a)). Ist  $f(d) > f(b)$  wird mit  $(d, b, c)$  fortgefahren (Abbildung 1 (b)). Dabei wäre von Vorteil, wenn die Verhältnisse der Teilintervalle gleich bleiben. Das Verhältnis von Gesamtintervall zu größeren Teilintervall ist für  $(a, b, c)$  und die möglichen neuen Tripel  $(a, d, b)$  und  $(d, b, c)$  sollte dann gleich sein:

$$\frac{c - a}{b - a} = \varphi \quad \frac{b - a}{d - a} = \varphi \quad \frac{c - d}{c - b} = \varphi$$

$\varphi$  ist dann der goldene Schnitt, für den gilt:

$$\varphi = \frac{1 + \sqrt{5}}{2} \quad \varphi^2 = \frac{(1 + \sqrt{5})^2}{4} = \frac{1 + 2\sqrt{5} + 5}{4} = \frac{3 + \sqrt{5}}{2} = 1 + \varphi$$

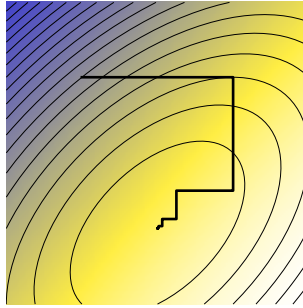


Abbildung 2: Coordinate-Descent

*Beweis:* Wenn  $b$  und  $d$  so gewählt wird, dass die ersten zwei Gleichungen gelten, dann gilt auch die dritte:

$$\frac{c-d}{c-b} = \frac{(c-a)-(d-a)}{(c-a)-(b-a)} \stackrel{1(b-a)}{=} \frac{\varphi-1/\varphi}{\varphi-1} \stackrel{\cdot\varphi}{=} \frac{\varphi^2-1}{\varphi^2-\varphi} = \frac{1+\varphi-1}{1+\varphi-\varphi} = \varphi$$

Wenn  $b-a < c-b$  ist, also das rechte Teilintervall größer, dann gilt das gleiche, nur mit  $a$  und  $c$  vertauscht, sodass sich Abbildung 1 (c) und (d) ergibt.

Also: Entweder ist  $b = a + (c-a)/\varphi$  (wenn  $b-a > c-b$ ) und dann

$$d = a + \frac{b-a}{\varphi} = a + \frac{c-a}{\varphi^2} \stackrel{\varphi^2-\varphi=1}{=} \frac{a\varphi^2 + (c-a)(\varphi^2-\varphi)}{\varphi^2} = c + \frac{a-c}{\varphi}$$

oder  $b = c + (a-c)/\varphi$  (wenn  $b-a < c-b$ ) und dann  $d = a + (c-a)/\varphi$ .

### 3 Coordinate Descent

Konvex/generell, keine (komplette) Ableitung, unconstrained.

Hier wird in jeder Iteration eine Koordinate (oder mehrere) ausgewählt. Danach wird eine eindimensionale Optimierung entlang dieser Koordinate durchgeführt.

Beispiel (siehe Abbildung 2): Für die Funktion

$$f(x, y) = x^2 + y^2 - xy$$

wird abwechselnd die  $y$ - und  $x$ -Koordinate ausgewählt. Dann wird die jeweils andere Koordinate als konstant betrachtet. Wir minimieren die Funktion eindimensional durch Ableiten und Nullsetzen:

$$\frac{df(x, y)}{dx} = 2x - y = 0 \quad \Rightarrow \quad x = \frac{y}{2}$$

Für die andere Koordinate:

$$\frac{df(x,y)}{dy} = 2y - x = 0 \quad \Rightarrow \quad y = \frac{x}{2}$$

Mögliche Varianten:

- Statt Ableitung-Nullsetzen z.B. Golden-Section-Search oder ähnliches.
- Nicht unbedingt volle eindimensionale Iterationen für jede Koordinaten-Iteration. Z.B. nur ein Schritt in die Koordinaten-Richtung, dann schon die andere Koordinate.
- Auswahl mehrerer Koordinaten. Wenn z.B.  $f(x, y)$  nach  $x$  und  $y$  ableitbar und nullsetzbar ist, aber nicht nach  $(x, y)$  gemeinsam. Das nennt man dann üblicherweise *alternierende* Optimierung.
- Optimierung in jeder Koordinate versuchen. Dann die beste Koordinate auswählen.

## 4 Pattern-Search

Konvex/generell, keine Ableitung, unconstrained.

Hier werden um einen Suchpunkt in einem Muster um den Suchpunkt mehrere Punkte evaluiert. Der Suchpunkt wird dann auf den Punkt mit dem kleinsten Funktionswert verschoben. Hat der Suchpunkt selbst den kleinsten Funktionswert, wird das Muster verkleinert.

Siehe Abbildung 3. Der Algorithmus ist:

$$P = \text{pattern}(\mathbf{x}) \quad (\text{Beispiel: } P = \mathbf{x} + s \cdot \{(\pm 1, 0, \dots), (0, \pm 1, \dots), \dots\})$$

$$f(\mathbf{x}) < \min f(P) \quad \Rightarrow \quad \underline{s} = s/2$$

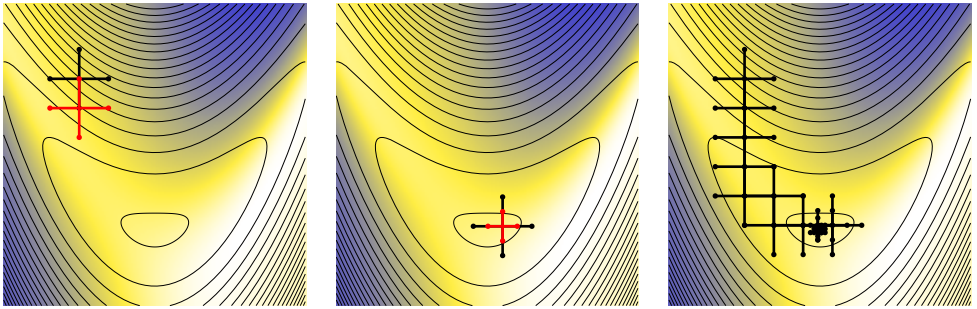
$$\mathbf{u} \in P \wedge f(\mathbf{u}) = \min f(P \cup \{\mathbf{x}\}) \quad \Rightarrow \quad \underline{\mathbf{x}} = \mathbf{u}$$

## 5 Nelder-Mead

Konvex/generell, keine Ableitung, unconstrained.

Der aktuelle Zustand im Nelder-Mead-Algorithmus ist ein Simplex in  $n$  Dimensionen, das sind  $n+1$  Ecken  $X = \{\mathbf{x}_i \mid i \in 1 \dots n+1\}$ . Die Ecken werden nach Funktionswert sortiert:

$$f(\mathbf{x}_1) \leq f(\mathbf{x}_2) \leq \dots \leq f(\mathbf{x}_{n+1})$$

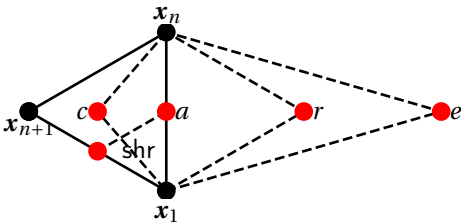


(a) Randpunkt minimal

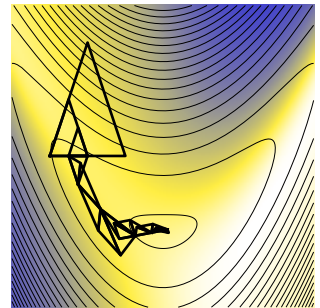
(b) Suchpunkt minimal

(c) Iteration

Abbildung 3: Pattern-Search. Schwarz: Ausgangszustand, rot: Folgezustand.



(a) Diagramm: Zentroid  $a$ , Contraction  $s$ , Reflection  $r$ , Expansion  $e$ , Shrink  $shr$



(b) Plot

Abbildung 4: Nelder-Mead-Algorithmus

Dann wird eines der Operationen Reflection, Expansion, Contract, Shrink durchgeführt, mit Hilfe des Zentroids der Ecken ohne der Ecke mit dem größten Funktionswert (Mittelpunkt von  $\{x_1, \dots, x_n\}$ ).

Der Algorithmus in Abbildung 4 und Abbildung 5 probiert zuerst die Reflection der schlechtesten Ecke  $x_{n+1}$  am Zentroid der anderen  $n$  Ecken aus. Führt diese einem besseren Funktionswert als die beste bisherige Ecke  $x_1$ , wird zusätzlich noch eine Expansion in die selber Richtung probiert. Falls diese noch besser ist, wird sie an Stelle der schlechtesten Ecke übernommen. Falls die Reflection wenigstens besser als die bisher zweitschlechteste Ecke ist, wird sie übernommen. Andernfalls wird eine Contraction zwischen  $x_{n+1}$  und dem Zentroid übernommen, falls sie eine Verbesserung gegenüber der Reflection und der schlechtesten

$\mathbf{a} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$	Zentroid
$\mathbf{r} = \mathbf{a} + q_r(\mathbf{a} - \mathbf{x}_{n+1})$	Reflection, $q_r = 1$
$\mathbf{e} = \mathbf{a} + q_e(\mathbf{a} - \mathbf{x}_{n+1})$	Expansion, $q_e = 2$
$\mathbf{c} = \mathbf{a} - q_c(\mathbf{a} - \mathbf{x}_{n+1})$	Contraction, $q_c = 0.5$
$f(\mathbf{e}) < f(\mathbf{r}) < f(\mathbf{x}_1) \Rightarrow$	
$\underline{X} = \text{sort}_f\{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{e}\}$	
else $f(\mathbf{r}) < f(\mathbf{x}_n) \Rightarrow$	
$\underline{X} = \text{sort}_f\{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{r}\}$	
else $f(\mathbf{c}) < \min(f(\mathbf{r}), f(\mathbf{x}_{n+1})) \Rightarrow$	
$\underline{X} = \text{sort}_f\{\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{c}\}$	
else $\Rightarrow$	
$\underline{X} = \text{sort}_f(\underline{X} + \mathbf{x}_1)/2$	Shrink

Abbildung 5: Nelder-Mead-Algorithmus.  $\text{sort}_f$  sortiert aufsteigend nach dem Funktionswert  $f(\mathbf{x}_i)$ .

Ecke darstellt. Falls so keine bessere Ecke gefunden wird, wird der ganze Simplex in Richtung der besten Ecke geschrumpft.

Es gibt einige modifizierte Varianten des Algorithmus mit anderen Konstanten  $q_r, q_e, q_c$  und leicht anderen Bedingungen/Entscheidungen.

## 6 Line-Search

Generell, keine/erste Ableitung, unconstrained.

Line-Search ist ein Hilfsmittel, um bei gegebener Optimierungs-Richtung  $\mathbf{d}$  eine geeignete Schrittweite  $\alpha$  zu finden.

$$\alpha = \arg \min_{\alpha} f(\underline{\mathbf{x}}) = \arg \min_{\alpha} f(\mathbf{x} + \alpha \mathbf{d})$$

Sehr oft ist die Optimierungsrichtung  $\mathbf{d} = -\nabla f(\mathbf{x})$  (Gradient-Descent, Abschnitt 7). Es reicht aber, wenn in einem (kleinen) Bereich  $X$  um  $\mathbf{x} \in X$  gilt:

$$\forall \mathbf{y} \in X: f(\mathbf{y}) - f(\mathbf{x}) \geq -\mathbf{d}^\top(\mathbf{y} - \mathbf{x})$$

Diese Definition funktioniert auch an Stellen, an denen  $f$  nicht differenzierbar ist. Siehe Abbildung 6

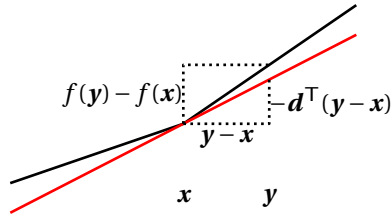


Abbildung 6: Subgradient. Schwarz: Eine Funktion  $f$ , die an  $\mathbf{x}$  nicht differenzierbar ist. Rot: Die Gerade (Ebene, Hyperebene), die durch  $-\mathbf{d}$  aufgespannt wird.

## 6.1 Exact Line-Search

Durch übliches Ableiten (nach  $\alpha$ ) und Nullsetzen kann man die optimale Schrittweite ermitteln:

$$0 = \frac{d f(\mathbf{x} + \alpha \mathbf{d})}{d\alpha} = \nabla f(\mathbf{x} + \alpha \mathbf{d})^\top \frac{d(\mathbf{x} + \alpha \mathbf{d})}{d\alpha} = \nabla f(\mathbf{x} + \alpha \mathbf{d})^\top \mathbf{d} =: g(\alpha)$$

$\mathbf{d}$  und  $\mathbf{x}$  sind Konstanten in der Nullstellensuche von  $g(\alpha)$ .  $g(\alpha) = 0$  ist univariat und oft ein Polynom, und lässt sich daher mit bekannten Methoden lösen. Man beachte, dass  $g(\alpha) = 0 \Leftrightarrow \nabla f(\mathbf{x})^\top \nabla \mathbf{d} = 0$  ist, d.h. man sucht eine Stelle entlang der Geraden  $\mathbf{x} + \alpha \mathbf{d}$ , an der der Gradient orthogonal zur aktuellen Optimierungsrichtung ist.

## 6.2 Numeric Line-Search

Benutze eine eindimensionale Methode (z.B. Golden-Section-Search), um  $f(\mathbf{x} + \alpha \mathbf{d})$  zu minimieren.

## 6.3 Backtracking Line-Search

Es ist nicht unbedingt notwendig, im Line-Search wirklich das (ein) Minimum zu finden. Es reicht, wenn der Funktionswert ausreichend kleiner wird. Das ist für ein genügend kleines  $\alpha$  sicher der Fall. Dafür reicht es,  $\alpha$  sukzessive zu verkleinern.

$$\underline{\alpha} = \rho \alpha \quad 0 < \rho < 1$$

Wiederhole das, bis z.B.  $f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x})$ .



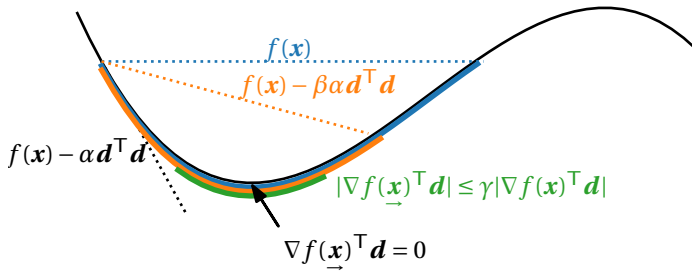


Abbildung 7: Wolfe-Conditions. Blau: Erlaubte Region für normales Backtracking. Orange: Für die erste Wolfe-Condition. Grün: Für die zweite Wolfe-Condition.

## 6.4 Wolfe-Conditions

Die Bedingung  $f(\mathbf{x} + \alpha \mathbf{d}) < f(\mathbf{x})$  garantiert leider keine Konvergenz. Die Funktionsreduktion könnte einfach immer kleiner werden. Daher wird folgende *erste Wolfe-Condition* oder *Armijo-Regel* eingeführt:

$$f(\mathbf{x} + \alpha \mathbf{d}) \leq f(\mathbf{x}) + \beta \alpha \mathbf{d}^T \nabla f(\mathbf{x}) = f(\mathbf{x}) - \beta \alpha \mathbf{d}^T \mathbf{d} \quad 0 < \beta < 1$$

Beachte: Nicht in allen Optimierungsansätzen ist  $\mathbf{d} = -\nabla f(\mathbf{x})$ . Man kann grundsätzlich jede Richtung mit  $\mathbf{d}^T \nabla f(\mathbf{x}) < 0$  wählen.

Die zweite Wolfe-Condition benutzt, dass im Idealfall  $\nabla f(\underline{\mathbf{x}})^T \nabla f(\mathbf{x}) = 0$  sein soll. Annähernd soll das Produkt zumindest kleiner als das aktuelle werden:

$$|\nabla f(\underline{\mathbf{x}})^T \mathbf{d}| \leq \gamma |\nabla f(\mathbf{x})^T \mathbf{d}| \quad 0 < \gamma < 1$$

Siehe Abbildung 7.

## 7 Gradient-Descent

Generell, erste Ableitung, unconstrained.

Gradient-Descent-Methoden ermitteln die Richtung des steilsten Abstiegs der Funktionswertes, was dem negativen Gradienten entspricht, und geht dann einen Schritt in diese Richtung, wobei die Schrittweite mit einem Parameter  $\alpha$  kontrolliert wird.

$$\underline{\mathbf{x}} = \mathbf{x} - \alpha \nabla f(\mathbf{x}) = \mathbf{x} + \alpha \mathbf{d}$$

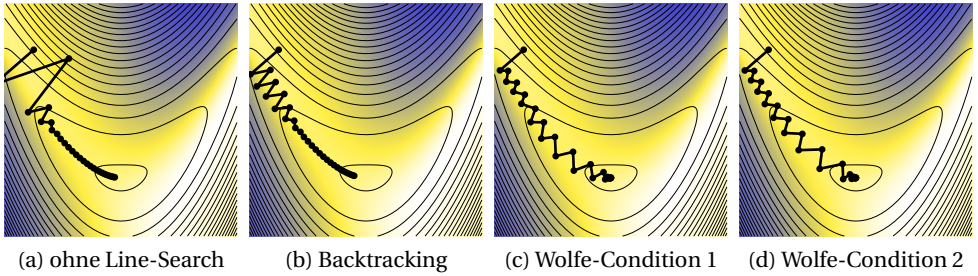


Abbildung 8: Gradient-Descent

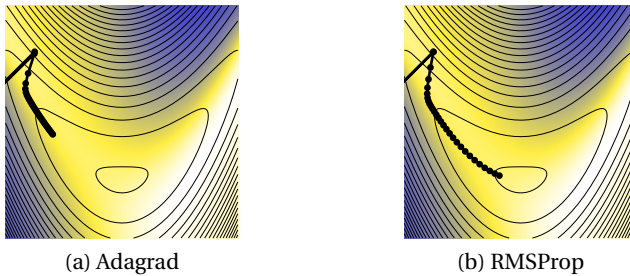


Abbildung 9: Gradient-Descent mit Koordinaten-adaptiven Lernraten

Es gibt viele Varianten und Methoden,  $\alpha$  festzulegen, und die Schritttrichtung zu optimieren.

Abbildung 8 zeigt das Verhalten von Gradient-Descent mit diversen Line-Search-Varianten.

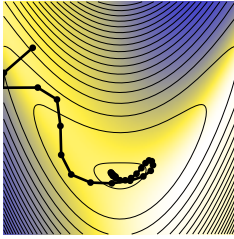
## 7.1 Adagrad

Oft benötigen manche Parameter eine andere Lernrate  $\alpha$  als andere. Die Methode der adaptive Gradients (Adagrad) löst das, indem für jede Komponente  $x$  ein adaptiver Lernfaktor aus der Summe der bisherigen quadrierten Gradienten-Komponenten berechnet wird. Das führt allerdings zu immer kleiner werdenden Lernschritten, oft noch bevor das Minimum erreicht wird. Siehe Abbildung 9 (a).

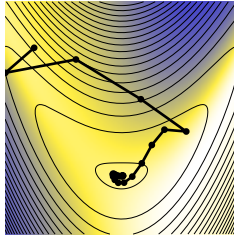
$$x_{\leftarrow i} = x_i + \frac{\alpha}{\epsilon + \sqrt{s_i}} d_i$$

$$s_{\circ i} = 0$$

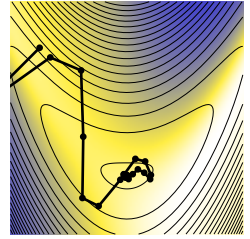
$$s_{\leftarrow i} = s_{\leftarrow i} + d_i^2$$



(a) Momentum



(b) Nesterov



(c) Adam

Abbildung 10: Gradient-Descent mit Momentum-Methoden

$\epsilon$  ist dabei eine sehr kleine Zahl, und dient nur dazu, Divisionen durch 0 zu vermeiden.  $\alpha$  ist wie üblich die generelle Lernrate.

## 7.2 RMSProp

Die RMSProp-Methode versucht daher, einen abklingenden Mittelwert der quadrierten Gradienten-Komponenten zu berechnen. Siehe Abbildung 9 (b).

$$\begin{aligned} \underline{x}_i &= x_i + \frac{\alpha}{\epsilon + \sqrt{s_i}} d_i \\ s_i &= 0 \\ s_i &= \gamma \underline{s}_i + (1 - \gamma) d_i^2 \quad \gamma = 0.9 \end{aligned}$$

## 7.3 Momentum

Da die Konvergenz entlang eines flachen Tales bei den bisherigen Methoden relativ langsam ist, versucht man mit Hilfe eines Bewegungsimpulses (Momentum) eine Beschleunigung zu erzielen, vergleichbar mit einem Ball, der das Tal entlang rollt und dabei beschleunigt. Siehe Abbildung 10 (a).

$$\begin{aligned} \underline{x} &= \underline{x} + \underline{v} \\ \underline{v} &= \beta \underline{v} + \alpha \underline{d} \end{aligned}$$

## 7.4 Nesterov-Momentum

Ein Problem bei der normalen Momentums-Methode ist das Überschießen an der Minimums-Stelle. Die Nesterov-Momentum-Methoden versucht, das in den Griff zu bekommen, indem statt des aktuellen Gradienten den an der potentiell nächsten Iterationsstelle  $\mathbf{x} + \beta \mathbf{v}$  zu verwenden. Siehe Abbildung 10 (b).

$$\begin{aligned}\underline{\mathbf{x}} &= \mathbf{x} + \underline{\mathbf{v}} \\ \underline{\mathbf{v}} &= \beta \mathbf{v} - \alpha \nabla f(\mathbf{x} + \beta \mathbf{v})\end{aligned}$$

## 7.5 Adam

Die Adaptive-Momentum-Estimation-Methode (Adam) vereint die Momentum- und RMSProp-Methode. Außerdem werden korrigierte Varianten von  $\mathbf{s}$  und  $\mathbf{v}$  berechnet ( $\hat{\mathbf{s}}, \hat{\mathbf{v}}$ ), die das Problem der Null-Initialisierung von  $\mathbf{s}$  und  $\mathbf{v}$  ausgleichen sollen. Siehe Abbildung 10 (c).

$$\begin{aligned}\underline{x}_i &= x_i + \frac{\alpha}{\epsilon + \sqrt{\hat{s}_i}} \hat{v}_i \\ \underline{s}_i &= 0 \\ \underline{s}_i &= \gamma \underline{s}_i + (1 - \gamma) d_i^2 \quad \gamma = 0.9 \\ \hat{\mathbf{s}} &= \mathbf{s} / (1 - \gamma^k) \\ \underline{\mathbf{v}} &= \beta \mathbf{v} + (1 - \beta) \mathbf{d} \\ \hat{\mathbf{v}} &= \underline{\mathbf{v}} / (1 - \beta^k) \\ \underline{k} &= 0 \quad k = \underline{k} + 1\end{aligned}$$

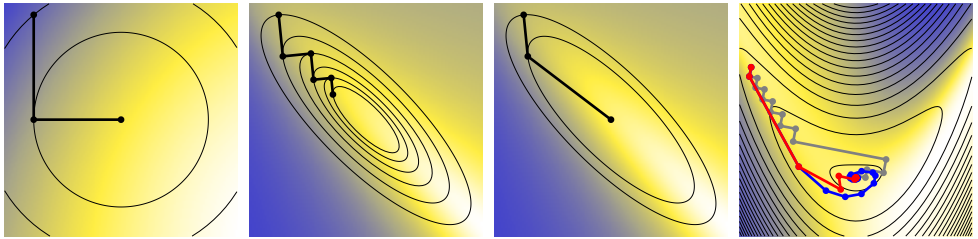
Für konstantes  $d_i^2$  gilt

$$\begin{aligned}s_i &= d_i^2 (1 - \gamma^k) \quad \text{weil} \quad \circ_i s_i = d_i^2 (1 - \gamma^0) = 0 \checkmark \\ \text{und} \quad \underline{s}_i &= \gamma s_i + (1 - \gamma) d_i^2 = \gamma d_i^2 (1 - \gamma^k) + (1 - \gamma) d_i^2 \\ &= d_i^2 (\gamma - \gamma^{k+1} + 1 - \gamma) = d_i^2 (1 - \gamma^{k+1}) \checkmark\end{aligned}$$

Daher ist  $s_i / (1 - \gamma^k)$  eine bessere Schätzung für  $d_i^2$ , vor allem für kleine  $k$  für die ersten Iterationen. Für  $\mathbf{v}$  gilt dasselbe.

## 7.6 Conjugate Gradient

Generell/quadratisch, erste Ableitung, unconstrained.



(a) Orthogonale Richtungen für  $A = I$ :  
 $\mathbf{x}^\top A \mathbf{x} = \mathbf{x}^\top \mathbf{x}$

(b) Orthogonale Richtungen für  $A$  nicht-diagonal

(c) Konjugierte Richtungen

(d) Nicht-quadratische Funktion  $f$ . Grau: normaler Gradient-Descent, blau: Fletcher-Reeves, rot: Polak-Ribière.

Abbildung 11: Conjugate-Gradient

Quadratische Funktionen haben die Form

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top A \mathbf{x} + \mathbf{b}^\top \mathbf{x} + c,$$

wobei  $A$  eine symmetrische positiv definite Matrix ist, für die  $\mathbf{x}^\top A \mathbf{x} > 0$  für alle  $\mathbf{x} \neq 0$  gilt. Für solche Matrizen gibt es eine (invertierbare) Transformationsmatrix  $B$  mit  $A = B^\top B$ . Es gilt:

$$\nabla f(\mathbf{x}) = A \mathbf{x} + \mathbf{b}.$$

Wir verwenden hier *exact Line-Search*. In diesem Fall sind aufeinander folgende Gradienten orthogonal (siehe Abschnitt 6.1). Wenn  $A = I$  (Identitätsmatrix) ist, können wir das Minimum in  $n$  Schritten ( $n$  ist die Dimension, hier 2) finden (siehe Abbildung 11(a)), selbst wenn orthogonale Richtungen unabhängig vom Gradienten vorgegeben sind. Wenn  $A$  nicht-diagonal ist, führt das hingegen zu einem Zick-Zack-Pfad (siehe Abbildung 11(b)). Die Idee ist nun, eine Art Orthogonalität auf einem durch  $B$  transformierten Koordinatensystem einzuführen (siehe Abbildung 11(c)). Das nennt man dann *konjugiert* statt orthogonal.

Zwei Vektoren  $\mathbf{x}$  und  $\mathbf{y}$  sind *konjugiert* bezüglich  $A$ , wenn

$$0 = \mathbf{x}^\top A \mathbf{y} = \mathbf{x}^\top B^\top B \mathbf{y} = (B \mathbf{x})^\top (B \mathbf{y}),$$

also wenn die transformierten Vektoren orthogonal sind.

Wir starten wie üblich mit einer Optimierungsrichtung

$$\underline{\mathbf{d}} = -\nabla f(\mathbf{x}), \quad \underline{\mathbf{x}} = \mathbf{x} + \alpha \underline{\mathbf{d}}.$$

In folgenden Iterationen soll die Richtung  $\underline{\mathbf{d}}$  nun so gewählt werden, dass sie zur vorherigen konjugiert ist. Und zwar wählen wir den Ansatz

$$\underline{\mathbf{d}} = -\nabla f(\mathbf{x}) + \beta \underline{\mathbf{d}}.$$

$\beta$  muss dann wie folgt berechnet werden:

$$\underline{\mathbf{d}}^T A \underline{\mathbf{d}} = 0 \quad \Rightarrow \quad \beta = \frac{\nabla f(\mathbf{x})^T A \underline{\mathbf{d}}}{\underline{\mathbf{d}}^T A \underline{\mathbf{d}}}$$

Für nicht-quadratische Funktionen kann diese Methode auch verwendet werden, weil die meisten Funktionen lokal durch eine quadratische Funktion approximiert werden können, zumindest am Konvergenzpunkt. Leider ist in diesem Fall aber  $A$  nicht bekannt, bzw. wäre nur mittels der zweiten Ableitung erudierbar. Es gibt aber äquivalente Formeln für  $\beta$ , die auch ohne Kenntnis von  $A$  berechnet werden können. Diese dienen als gut funktionierende Annäherungen, und zwar einmal nach *Fletcher-Reeves*:

$$\beta = \frac{\nabla f(\mathbf{x})^T \nabla f(\mathbf{x})}{\nabla f(\underline{\mathbf{x}})^T \nabla f(\underline{\mathbf{x}})},$$

oder nach *Polak-Ribière*:

$$\beta = \frac{\nabla f(\mathbf{x})^T (\nabla f(\mathbf{x}) - \nabla f(\underline{\mathbf{x}}))}{\nabla f(\underline{\mathbf{x}})^T \nabla f(\underline{\mathbf{x}})},$$

wobei bei letzterem Schema  $\beta$  nach unten mit 0 beschränkt werden sollte (Reset). Einen Vergleich zeigt Abbildung 11(d).

## 7.7 Stochastic Gradient-Descent

Oft ist die Zielfunktion  $f(\mathbf{x})$  als Summe (bzw. Mittelwert) von  $n$  einzelnen Teilfunktionen  $g$  gegeben, die von Datenpunkten  $\mathbf{y}_i$  abhängen. Dann ergibt sich:

$$f(\mathbf{x}) = \frac{1}{n} \sum_i g(\mathbf{x}, \mathbf{y}_i) = \frac{1}{n} \sum_i h_i(\mathbf{x})$$

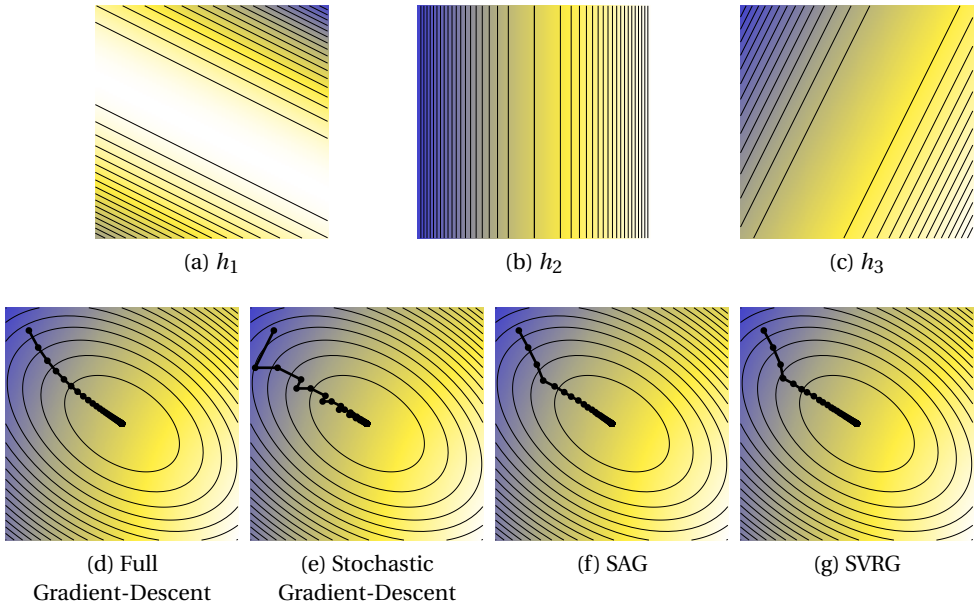


Abbildung 12: Stochastic gradient descent

Zum Beispiel ist beim Trainieren von künstlichen neuronalen Netzen  $\mathbf{y}$  ein Satz von Trainingsdaten (Batch), die sich in Inputs  $\mathbf{y}_{\text{in}}$  und gewünschten Outputs  $\mathbf{y}_{\text{out}}$  unterteilen. Die Zielfunktion ist dann die Summe aller Abweichungen der Prediction-Funktion  $\mathbf{p}$  und den Outputs aus  $\mathbf{y}_{\text{out}}$ .  $\mathbf{p}(\mathbf{x}, \mathbf{y}_{\text{in},i})$  berechnet mit Hilfe der Parameter (Gewichte) des neuronalen Netzes aus den Input-Daten  $\mathbf{y}_{\text{in},i}$  eine Annäherung an den gewünschten Output  $\mathbf{y}_{\text{out},i}$ .

$$g(\mathbf{x}, \mathbf{y}_i) = |\mathbf{p}(\mathbf{x}, \mathbf{y}_{\text{in},i}) - \mathbf{y}_{\text{out},i}|^2$$

Als Beispiel wählen wir  $\mathbf{p}(\mathbf{x}, \mathbf{y}_{\text{in}}) = x_1 \mathbf{y}_{\text{in}} + x_2$ ,  $\mathbf{y}_{\text{in}} = (2, 0, -0.5)$ ,  $\mathbf{y}_{\text{out}} = (5, 1, 0)$ , was zur Lösung  $\mathbf{x} = (2, 1)$  führen sollte. Die einzelnen Fehlerfunktionen  $h_1$ ,  $h_2$ ,  $h_3$  sind in Abbildung 12(a-c) dargestellt.

Üblicherweise müsste man hier den Gradienten von  $f$  berechnen, der sich aus der Summe der Gradienten von  $g(\mathbf{x}, \mathbf{y}_i)$  für alle  $i$  zusammensetzt.

$$\mathbf{d} = -\nabla f(\mathbf{x}) = -\nabla \frac{1}{n} \sum_i g(\mathbf{x}, \mathbf{y}_i) = -\frac{1}{n} \sum_i \nabla_{\mathbf{x}} g(\mathbf{x}, \mathbf{y}_i) = -\frac{1}{n} \sum_i \nabla h_i(\mathbf{x})$$

Stattdessen führt man bei der Stochastic Gradient-Descent-Methode den Update-Schritt für jeden der  $n$  Einträge in den Trainingsdaten durch.

$$\underline{\mathbf{x}} = \mathbf{x} - \alpha \nabla g(\mathbf{x}, \mathbf{y}_i)$$

$$\underline{i} = 0 \quad \underline{i} = (i + 1) \bmod n$$

Den Vergleich zwischen „Full Gradient-Descent“ und Stochastic Gradient-Descent sieht man in Abbildung 12(d) und (e).

Der Nachteil dabei ist, dass je nach (zufälliger) Sortierung der Trainingsdaten  $\mathbf{x}$  in beinahe zufällige Richtungen springen kann, die sich nur durch spätere Iterationen wieder ausgleichen. Der Vorteil ist allerdings, dass die Berechnung der Gradienten viel schneller ist, und daher viel mehr Iterationen pro Zeit ausgeführt werden können.

Um die Zufälligkeit dieser Methode doch etwas einzuschränken, kann man die Trainingsdaten (Batch, Größe  $n$ ) in sogenannte *Minibatches* der Größe  $b$  unterteilen und die Gradienten über die Minibatches summieren.

$$\underline{\mathbf{x}} = \mathbf{x} - \alpha \sum_{k=0}^{b-1} \nabla g(\mathbf{x}, \mathbf{y}_{(i+k) \bmod n})$$

$$\underline{i} = 0 \quad \underline{i} = (i + b) \bmod n$$

## 7.8 Stochastic Variance-Reduction

Eine andere Methode, das zufällige Herumspringen, also die Varianz der Gradienten zu reduzieren, ist *Variance-Reduction*. Dazu gibt es zwei Ansätze. Der erste ist, eine Tabelle von Gradienten für alle  $h_i$  aufzubauen. Immer wenn ein neuer Gradient  $\nabla h_i$  berechnet wird, überschreibt dieser den jeweiligen Eintrag in der Tabelle. Für den Optimierungsschritt wird dann der Mittelwert der Tabelle verwendet. Die Berechnung dieses Mittelwerts kann man übrigens umgehen, indem man bei jedem Update der Tabelle auch den gespeicherten Mittelwert updatet. Folgender Algorithmus nennt sich SAG:

$$\underline{i} = 0 \quad \underline{i} = (i + 1) \bmod n$$

$$\underline{\mathbf{g}}_i = \nabla h_i(\mathbf{x})$$

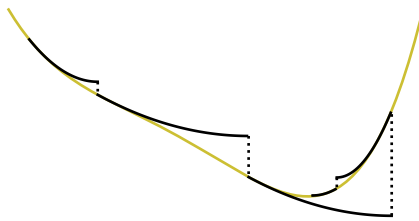
$$\underline{\mathbf{x}} = \mathbf{x} - \alpha \frac{1}{n} \sum_j \underline{\mathbf{g}}_j$$

Das Ergebnis ist eine glattere Konvergenz, wie man in Abbildung 12(f).

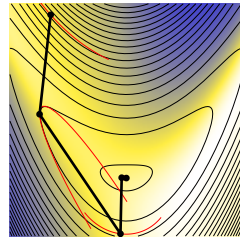
Der zweite Ansatz ist, in regelmäßigen Abständen (z.B. nach  $n$  Schritten) einen vollen Gradienten  $\nabla f(\mathbf{x})$  zu berechnen, und in den folgenden Schritten einen Optimierungsschritt zu machen, der diesen vollen Gradienten ergänzt. Folgender Algorithmus nennt sich SVRG:

$$\underline{i} = 0 \quad \underline{i} = (i + 1) \bmod n$$





(a) eindimensional



(b) zweidimensional

Abbildung 13: Newton-Methode

$$\tilde{\mathbf{x}} = \begin{cases} \mathbf{x} & i = 0 \\ \tilde{\mathbf{x}} & i > 0 \end{cases} \quad \mathbf{g} = \nabla f(\tilde{\mathbf{x}})$$

$$\underline{\mathbf{x}} = \mathbf{x} - \alpha(\nabla h_i(\mathbf{x}) - \nabla h_i(\tilde{\mathbf{x}}) + \mathbf{g})$$

Das Ergebnis ist in Abbildung 12(f) zu sehen.

## 8 Newton-Methode

Generell, zweite Ableitung, unconstrained.

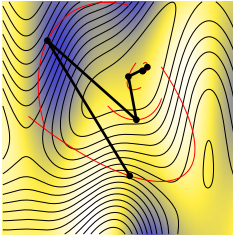
Die Newton-Methode gibt es in zwei Varianten: Einmal zur Nullstellensuche in eindimensionalen Funktionen  $f: \mathbb{R} \rightarrow \mathbb{R}$  (oder  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^n$ ). Dabei approximiert man die Funktion an der Stelle  $\mathbf{x}$  als Gerade mit Steigung  $f'(\mathbf{x})$  und berechnet die Nullstelle  $\underline{\mathbf{x}}$  dieser Gerade, also

$$0 = y = f(\mathbf{x}) + (\underline{\mathbf{x}} - \mathbf{x})f'(\mathbf{x}) \quad \Rightarrow \quad \underline{\mathbf{x}} = \mathbf{x} - \frac{f(\mathbf{x})}{f'(\mathbf{x})},$$

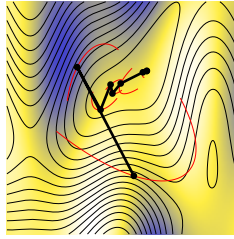
und fährt an dieser Stelle mit neuen Iterationen fort.

Wir beschäftigen uns mit der Variante zur Funktionsoptimierung. Dabei wird die Funktion an der Stelle  $\mathbf{x}$  als Paraboloid oder quadratische Funktion (siehe Abschnitt 7.6) approximiert, und davon die Minimumsstelle  $\underline{\mathbf{x}}$  berechnet. Für die Parameter dieser Funktionen ergibt sich durch zweimaliges Ableiten nach  $\mathbf{d}$  (dem Abstand von  $\mathbf{x}$ ):

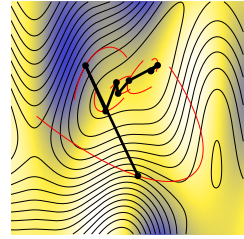
$$f(\mathbf{x} + \mathbf{d}) \approx \frac{1}{2} \mathbf{d}^\top \mathbf{A} \mathbf{d} + \mathbf{b}^\top \mathbf{d} + c \quad \Rightarrow \quad f(\mathbf{x}) = c$$



(a) Gauss-Newton



(b)  
Levenberg-Marquardt



(c) Powell's Dog-Leg

Abbildung 14: Least-Squares-Methoden

$$\begin{aligned}\nabla f(\mathbf{x} + \mathbf{d}) &= A\mathbf{d} + \mathbf{b} & \Rightarrow & \nabla f(\mathbf{x}) = \mathbf{b} \\ \nabla^2 f(\mathbf{x} + \mathbf{d}) &= A & \Rightarrow & \nabla^2 f(\mathbf{x}) = A\end{aligned}$$

Siehe auch Abschnitt 22. Durch Nullsetzen der ersten Ableitung ergibt sich

$$0 = \nabla f(\mathbf{x} + \mathbf{d}) = \nabla^2 f(\mathbf{x})\mathbf{d} + \nabla f(\mathbf{x}) \quad \Rightarrow \quad \nabla^2 f(\mathbf{x})\mathbf{d} = -\nabla f(\mathbf{x}).$$

Das ist ein Gleichungssystem, das mit Hilfe von Gleichungssystemlösern (Cholesky-Faktorisierung) nach  $\mathbf{d}$  gelöst wird.

Abbildung 13 zeigt Ausführungen der Newton-Methode. Man beachte das Überschießen an Stellen mit verringerter Krümmung, was in vielen Fällen Backtracking-Line-Search notwendig macht.

## 9 Gauss-Newton-Methode

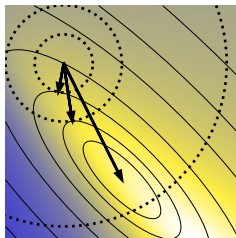
Nichtlineare Least-Squares, erste Ableitung, unconstrained.

Bei Least-Squares-Problemen ist die zu optimierende Funktion als Summe von quadrierten Funktionen gegeben.

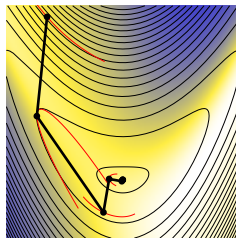
$$f(\mathbf{x}) = \sum_i g_i^2(\mathbf{x}) = \mathbf{g}(\mathbf{x})^T \mathbf{g}(\mathbf{x})$$

Wie in Abschnitt 7.7 ist das häufigste Beispiel die Anpassung einer Prediction-Funktion:

$$g_i^2(\mathbf{x}) = \|\mathbf{p}(\mathbf{x}, \mathbf{y}_{\text{in},i}) - \mathbf{y}_{\text{out},i}\|^2$$



(a) Minima sind abhängig vom Radius der Trust-Region



(b) Newton-Methode mit Trust-Region

Abbildung 15: Trust-Region

$\mathbf{g}$  wird nun als lineare Funktion approximiert:

$$\mathbf{g}(\mathbf{x} + \mathbf{d}) = \nabla \mathbf{g}(\mathbf{x}) \mathbf{d} + \mathbf{g}(\mathbf{x})$$

Hier ist  $\nabla \mathbf{g}$  die *Jacobi-Matrix* (siehe Abschnitt 22). Dann ist:

$$f(\mathbf{x} + \mathbf{d}) = \mathbf{g}(\mathbf{x} + \mathbf{d})^\top \mathbf{g}(\mathbf{x} + \mathbf{d}) \approx \mathbf{d}^\top \nabla \mathbf{g}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{x}) \mathbf{d} + 2\mathbf{d}^\top \nabla \mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x}) + \mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})$$

Und mittels Ableiten und Nullsetzen erhalten wir:

$$\nabla_{\mathbf{d}} f(\mathbf{x} + \mathbf{d}) \approx 2\nabla \mathbf{g}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{x}) \mathbf{d} + 2\nabla \mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x}) = 0$$

$$\nabla \mathbf{g}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{x}) \mathbf{d} = -\nabla \mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})$$

Ähnlich wie bei der Newton-Methode müssen wir also ein Gleichungssystem lösen. Abbildung 14(a) zeigt einen Optimierungslauf.

## 10 Trust-Region

Trust-Region-Methoden sind in gewisser Weise das Gegenstück zu Line-Search. Während bei Line-Search zuerst die Optimierungsrichtung feststeht, und danach eine gute Schrittweite gesucht wird, wird bei Trust-Region zuerst die (maximale) Schrittweite festgelegt (der Radius der Trust-Region), und danach die beste Optimierungsrichtung. Das führt zu einem Optimierungsproblem mit Constraints (Nebenbedingungen), welche wir uns im Detail erst später anschauen.

$$\underline{\mathbf{x}} = \mathbf{x} + \underset{\mathbf{d}}{\operatorname{argmin}} \hat{f}(\mathbf{x} + \mathbf{d}) \text{ mit } |\mathbf{d}|^2 < r^2$$

Dass die Optimierungsrichtung dabei abhängig von der Größe der Trust-Region ist, sieht man in Abbildung 15(a).

## 10.1 Newton-Methode mit Trust-Region

Generell, zweite Ableitung, unconstrained.

Man verwendet dabei üblicherweise nicht die originale Funktion  $f$ , sondern eine lokale Approximation, die sich leichter minimieren lässt. Meist wird dabei eine quadratische Funktion, basierend auf der Hesse-Matrix (also eine Newton-Methode) verwendet.

$$\hat{f}(\mathbf{x} + \mathbf{d}) = \frac{1}{2} \mathbf{d}^T A \mathbf{d} + \mathbf{b}^T \mathbf{d}, \quad \mathbf{b} = \nabla f(\mathbf{x}), \quad A = \nabla^2 f(\mathbf{x})$$

Es gibt dann einen *Lagrange-Multiplikator*  $\lambda$ , sodass man stattdessen die Constraint-lose Funktion

$$\hat{f}_\lambda(\mathbf{x} + \mathbf{d}) = \hat{f}(\mathbf{x} + \mathbf{d}) + \lambda(\mathbf{d}^T \mathbf{d} - r^2)$$

minimieren kann. Durch Ableiten und Nullsetzen erhält man:

$$\nabla_{\mathbf{d}} \hat{f}_\lambda(\mathbf{x} + \mathbf{d}) = A \mathbf{d} + \mathbf{b} + 2\lambda \mathbf{d} = (A + 2\lambda I) \mathbf{d} + \mathbf{b} = 0,$$

wobei  $I$  die Identitätsmatrix (Diagonale 1) ist. Das heißt, dass man das für die Newton-Methode zu lösende Gleichungssystem modifiziert, indem man konstante Werte  $2\lambda$  zur Diagonale von  $A$  addiert. Ist  $\lambda = 0$  erhält man die Newton-Methode. Das kann man verwenden, wenn die Trust-Region das Minimum der quadratischen Approximation enthält. Je größer  $\lambda$  wird, desto kleiner wird der Radius der Trust-Region, und desto mehr nähert sich  $\mathbf{d}$  der Gradienten-Richtung. Abbildung 15(a) zeigt mehrere Lösungen für verschiedene  $\lambda$ .

Die nächste Frage ist nun, welchen Wert  $\lambda$  haben sollte. Die Lösung ist üblicherweise, mit einem bestimmten kleinen  $\lambda$  zu starten, und dann wiederholt zu erhöhen, falls der neue Funktionswert nicht kleiner ist als der aktuelle. Das ist eine Form von Backtracking. Wenn der Funktionswert verkleinert werden kann, wird  $\lambda$  entweder wieder auf den ursprünglichen Wert zurückgesetzt, oder verkleinert. Das gibt folgenden Algorithmus:

Löse  $(\nabla^2 f(\mathbf{x}) + \lambda I) \mathbf{d} = -\nabla f(\mathbf{x})$

if  $f(\mathbf{x} + \mathbf{d}) < f(\mathbf{x})$

$$\vec{\mathbf{x}} = \vec{\mathbf{x}} + \vec{\mathbf{d}}$$

$$\vec{\lambda} = \vec{\lambda} \quad \text{oder} \quad \vec{\lambda} = \lambda / \rho$$

sonst

$$\vec{\lambda} = \rho \lambda \quad \rho = 2$$

Ein Optimierungslauf ist in Abbildung 15(b) zu sehen. Vergleiche dazu Abbildung 13(b), wo im zweiten Schritt der Funktionswert schlechter wird.

## 10.2 Levenberg-Marquard

Nichtlineare Least-Squares, erste Ableitung, unconstrained.

Die Levenberg-Marquard-Methode erweitert die Gauss-Newton-Methode um eine Trust-Region. Es ergibt sich also:

$$(\nabla \mathbf{g}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{x}) + \lambda I) \mathbf{d} = -\nabla \mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})$$

Und die Steuerung von  $\lambda$  wie im Algorithmus in Abschnitt 10. Abbildung 14(b) zeigt einen Testlauf.

## 10.3 Powells Dog-Leg

Nichtlineare Least-Squares, erste Ableitung, unconstrained.

Eine andere Methode, die Gauss-Newton-Methode mit Trust-Region zu vereinen, ist Powells Dog-Leg. Dabei wird zuerst  $\mathbf{d}_{\text{GN}}$  wie in der Gauss-Newton-Methode berechnet. Außerdem der Gradient-Descent-Schritt  $\mathbf{d}_{\text{GD}} = -\nabla \mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x})$ .

Danach wird in Richtung von  $\mathbf{d}_{\text{GD}}$  das Minimum gesucht:

$$\begin{aligned} \frac{d}{dt} f(\mathbf{x} + t\mathbf{d}_{\text{GD}}) &= \nabla f(\mathbf{x} + t\mathbf{d}_{\text{GD}})^\top \mathbf{d}_{\text{GD}} \\ &= 2t \underbrace{\mathbf{d}_{\text{GD}}^\top \nabla \mathbf{g}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{x}) \mathbf{d}_{\text{GD}}}_{\|\nabla \mathbf{g}(\mathbf{x}) \mathbf{d}_{\text{GD}}\|^2} + 2 \underbrace{\mathbf{g}(\mathbf{x})^\top \nabla \mathbf{g}(\mathbf{x}) \mathbf{d}_{\text{GD}}}_{-\mathbf{d}_{\text{GD}}^\top} = 0 \\ t &= \frac{\|\mathbf{d}_{\text{GD}}\|^2}{\|\nabla \mathbf{g}(\mathbf{x}) \mathbf{d}_{\text{GD}}\|^2} \end{aligned}$$

Wenn nun der Gauss-Newton-Schritt innerhalb der Trust-Region liegt, wird dieser durchgeführt. Wenn nicht, und der Gradient-Descent-Schritt liegt innerhalb, dann wird zwischen den zweien interpoliert ( $\mathbf{d} = \beta \mathbf{d}_{\text{GN}} + (1 - \beta) t \mathbf{d}_{\text{GD}}$ ), und zwar so, dass der Interpolationspunkt  $\mathbf{d}$  genau die Distanz  $r$  (Trust-Region-Radius) hat. Das ergibt eine quadratische Gleichung in  $\beta$  ( $\beta^2 \|\mathbf{d}_{\text{GN}} - t \mathbf{d}_{\text{GD}}\|^2 + \beta \cdot 2(\mathbf{d}_{\text{GN}} - t \mathbf{d}_{\text{GD}})^\top t \mathbf{d}_{\text{GD}} + \|t \mathbf{d}_{\text{GD}}\|^2 - r^2 = 0$ ). Andernfalls wird der Gradient-Descent-Schritt auf Länge  $r$  gekürzt/verlängert.

Abbildung 16 zeigt den Algorithmus. Abbildung 14(c) zeigt einen Testlauf.

Berechne  $\mathbf{d}_{\text{GN}}$  (wie in Abschnitt 9)

$$\mathbf{d}_{\text{GD}} = -\nabla \mathbf{g}(\mathbf{x})^T \mathbf{g}(\mathbf{x}), \quad t = \|\mathbf{d}_{\text{GD}}\|^2 / \|\nabla \mathbf{g}(\mathbf{x}) \mathbf{d}_{\text{GD}}\|^2$$

if  $\|\mathbf{d}_{\text{GN}}\|^2 < r^2$

$$\mathbf{d} = \mathbf{d}_{\text{GN}}$$

else if  $\|t\mathbf{d}_{\text{GD}}\|^2 < r^2$

Berechne  $\beta$  (quadratische Gleichung)

$$\mathbf{d} = \beta \mathbf{d}_{\text{GN}} + (1 - \beta) t \mathbf{d}_{\text{GD}} \quad \text{sodass} \quad \|\mathbf{d}\|^2 = r^2$$

else

$$\mathbf{d} = r \mathbf{d}_{\text{GD}} / \|\mathbf{d}_{\text{GD}}\|$$

if  $f(\mathbf{x} + \mathbf{d}) < f(\mathbf{x})$

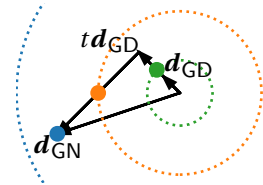
$$\underline{\mathbf{x}} = \mathbf{x} + \mathbf{d}$$

$$\underline{r} = \rho_1 r \quad \rho_1 = 2$$

else

$$\underline{r} = \rho_2 r \quad \rho_2 = 0.8$$

(a) Algorithmus



(b) Schema

Abbildung 16: Powells Dog-Leg

## 11 Quasi-Newton-Methoden

Bei Quasi-Newton-Methoden wird die Hesse-Matrix  $\nabla^2 f(\mathbf{x})$  approximiert, ohne die zweite Ableitung bilden zu müssen. Meist werden dazu aufeinander folgende erste Ableitungen benutzt.

### 11.1 BFGS (Broyden-Fletcher-Goldfarb-Shanno)

Bei dieser Methode wird die inverse Hesse-Matrix approximiert, die wir hier  $G$  nennen. Aus dem sonst benutzten zu lösenden Gleichungssystem wird dann:

$$\nabla^2 f(\mathbf{x}) \mathbf{d} = -\nabla f(\mathbf{x}) \quad \Rightarrow \quad \mathbf{d} = -\underbrace{(\nabla^2 f(\mathbf{x}))^{-1}}_{=:G} \nabla f(\mathbf{x})$$

Wir betrachten nun die Approximation der Funktion  $f$  an der Stelle  $\underline{\mathbf{x}}$  durch die übliche quadratische Funktion  $\hat{f}$ . Der Gradient von  $\hat{f}$  ist bekannterweise:

$$\nabla \hat{f}(\underline{\mathbf{x}} + \mathbf{q}) = \underline{G}^{-1} \mathbf{q} + \nabla f(\underline{\mathbf{x}})$$

Wir wollen nun  $G$  so wählen, dass der Gradient an beiden Stellen  $\underline{x}$  und  $\mathbf{x}$  dem echten Gradienten entspricht:

$$\nabla \hat{f}(\underline{\mathbf{x}}) = \nabla f(\underline{\mathbf{x}}) \quad \nabla \hat{f}(\mathbf{x}) = \nabla f(\mathbf{x})$$

Ersteres ist (mit  $\mathbf{q} = 0$ ) offenbar erfüllt. Für die zweite Bedingung ergibt sich mit  $\mathbf{x} = \underline{\mathbf{x}} - \mathbf{d}$ :

$$\nabla f(\mathbf{x}) = -\underline{G}^{-1} \mathbf{d} + \nabla f(\underline{\mathbf{x}}) \quad \underline{G}^{-1} \mathbf{d} = \nabla f(\underline{\mathbf{x}}) - \nabla f(\mathbf{x}) =: \mathbf{e} \quad \underline{G} \mathbf{e} = \mathbf{d}$$

Wir suchen also ein  $\underline{G}$ , das dem alten  $G$  ähnlich ist, aber diese Bedingung erfüllt.

$$\underline{G} = \underset{\underline{G}}{\operatorname{arg\,min}} \|\underline{G} - G\| \quad \text{mit} \quad \underline{G} = \underline{G}^T, \underline{G} \mathbf{e} = \mathbf{d}$$

Durch geeignete Wahl der Matrix-Norm  $\|\cdot\|$  kann man zeigen, dass folgender Update diese Bedingungen erfüllt.

$$\underline{G} = \left( I - \frac{\mathbf{e} \mathbf{d}^T}{\mathbf{d}^T \mathbf{e}} \right)^T G \left( I - \frac{\mathbf{e} \mathbf{d}^T}{\mathbf{d}^T \mathbf{e}} \right) + \frac{\mathbf{d} \mathbf{d}^T}{\mathbf{d}^T \mathbf{e}}$$

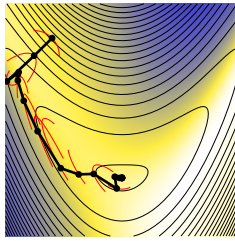
Durch Multiplizieren mit  $\mathbf{e}$  von rechts sieht man schnell, dass  $\underline{G} \mathbf{e} = \mathbf{d}$  ist.  $G$  soll aber auch positiv definit sein. Dafür muss man garantieren, dass:

$$\mathbf{d}^T \mathbf{e} > 0 \quad \Leftrightarrow \quad \mathbf{d}^T \nabla f(\underline{\mathbf{x}}) > \mathbf{d}^T \nabla f(\mathbf{x}) \quad \left( = -\nabla f(\mathbf{x})^T H \nabla f(\mathbf{x}) < 0 \right)$$

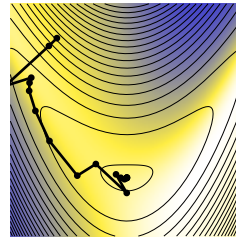
Das erreichen wir durch Line-Search mit den Wolfe-Conditions (beginnend immer mit  $\alpha = 1$ ). Der Algorithmus sieht dann so aus:

$$\begin{aligned} \underline{G} &= \alpha I & \Rightarrow & \mathbf{d} = -\alpha \nabla f(\underline{\mathbf{x}}) \\ \mathbf{d} &= \text{Line-Search in Richtung } -G \nabla f(\underline{\mathbf{x}}) \\ \underline{\mathbf{x}} &= \underline{\mathbf{x}} + \mathbf{d} \\ \mathbf{e} &= \nabla f(\underline{\mathbf{x}}) - \nabla f(\mathbf{x}) \\ \underline{G} &= \text{BFGS-Update}(G, \mathbf{d}, \mathbf{e}) \end{aligned}$$

Abbildung 17(a) zeigt einen Testlauf.



(a) BFGS



(b) L-BFGS ( $m = 5$ )

Abbildung 17: Quasi-Newton-Methoden

## 11.2 L-BFGS (Limited-Memory-BFGS)

Für hoch-dimensionale Probleme wird das Berechnen der approximierten Hesse-Matrix ein Problem in Bezug auf rechnerische Komplexität und Speicherbedarf. Daher wird hier versucht, das Produkt  $G \nabla f(\mathbf{x})$  durch Produkte von Vektoren zu implementieren. Da dieses Produkt wiederum Produkte mit vorangegangenen  $G_{\leftarrow i}$  beinhaltet, ergibt sich ein rekursives Schema, das nach  $m$  Schritten mit einer skalierten Identitätsmatrix statt  $G$  abgebrochen wird. Für die rekursiven Iterationen benötigt man einen Puffer der letzten  $m$  Vektoren  $\mathbf{d}$  und  $\mathbf{e}$ . Es ergibt sich der Algorithmus in Abbildung 18. Abbildung 17(b) zeigt einen Testlauf.

## 12 Regularization

Unter Regularization versteht man Methoden, eine zu optimierende Funktion so zu modifizieren, dass die Optimierungsergebnisse robuster gegen Ungenauigkeiten oder Variabilitäten der Funktionsparameter sind.

### 12.1 Ridge/Tikhonov

Die *Ridge-Regression* oder *Tikhonov-Regularization* ist am besten für lineare Least-Squares-Probleme zu erklären:

$$f(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|^2$$

Durch Ableiten und Nullsetzen bekommt man die sogenannte Normalgleichung:

$$\nabla f(\mathbf{x}) = 2A^T(\mathbf{Ax} - \mathbf{y}) = \mathbf{0} \quad \Rightarrow \quad A^T \mathbf{Ax} = A^T \mathbf{y}$$



$\text{LBFGS}(\mathbf{p}, i) :=$   
 if  $i = m \vee \leftarrow i = \circ$   
     return  $\frac{\mathbf{d}^\top \mathbf{e}}{\mathbf{e}^\top \mathbf{e}} \mathbf{p}$   
 $a = \frac{\mathbf{d}^\top \mathbf{p}}{\mathbf{d}^\top \mathbf{e}}$   
 $\mathbf{q} = \text{LBFGS}(\mathbf{p} - a \mathbf{e}, i + 1)$   
 return  $\mathbf{q} + \left( a - \frac{\mathbf{e}^\top \mathbf{q}}{\mathbf{d}^\top \mathbf{e}} \right) \mathbf{d}$

$\mathbf{d} = -\alpha \nabla f(\mathbf{x})$   
 $\mathbf{d} = \text{Line-Search in Richtung } -\text{LBFGS}(\nabla f(\mathbf{x}), 1)$   
 $\mathbf{x} = \mathbf{x} + \mathbf{d}$   
 $\mathbf{e} = \nabla f(\mathbf{x}) - \nabla f(\mathbf{x})$

Abbildung 18: L-BFGS-Algorithmus

Wie stabil die Lösungen dieser Gleichung sind, hängt von der Matrix  $A^\top A$  ab. Wenn diese knapp singulär (nicht-invertierbar) ist, dann reagiert die Lösung sehr stark auf leicht veränderte Koeffizienten:

$$A = \begin{pmatrix} 1 & 1 \\ 1 & 1.01 \end{pmatrix} \Rightarrow A^\top A \begin{pmatrix} 101 \\ -100 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad A^\top A \begin{pmatrix} 0.5 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1.005 \end{pmatrix}, \quad A^\top A \begin{pmatrix} -100 \\ 100 \end{pmatrix} = \begin{pmatrix} 1 \\ 1.01 \end{pmatrix}$$

Dieses Verhalten wird auch durch die *Konditionszahl* von  $A^\top A$  beschrieben (größer durch kleinster Singulärwert). Sie ist für diese Beispielmatrix  $> 16000$  (die Identitätsmatrix hat 1). So ein Problem nennt man *ill-posed*.

Um das unter Kontrolle zu bekommen, modellieren wir  $A$  als Erwartungswert einer randomisierten Matrix, indem wir eine Matrix von unabhängigen Zufallsvariablen mit Erwartungswert  $E = 0$ , Varianz  $V = \sigma^2$  und Kovarianz  $\text{Cov} = 0$  für verschiedene Zufallsvariablen:

$$\tilde{A} = A + N \quad \text{mit} \quad E(N_{i,j}) = 0, \quad V(N_{i,j}) = E(N_{i,j}^2) = \sigma^2, \quad \text{Cov}(N_{i,j}, N_{k,l \neq i,j}) = 0$$

Dann ergibt sich:

$$E(\tilde{A}) = A, \quad E(\tilde{A}^T \tilde{A})_{i,j} = E\left(\sum_k (A+N)_{i,k} (A+N)_{k,j}\right) = (A^T A)_{i,j} + \begin{cases} \sigma^2 & i = j \\ 0 & i \neq j \end{cases}$$

Das heißt, es wird ein konstanter Wert auf die Diagonale addiert. Dementsprechend ist die Konditionszahl von z.B.  $A^T A + 0.1I$  gleich 41.2. Die Lösungen der Gleichung sind viel stabiler:

$$(A^T A + 0.1I) \begin{pmatrix} 0.25 \\ 0.23 \end{pmatrix} = \begin{pmatrix} 0.987 \\ 0.990 \end{pmatrix}, \quad (A^T A + 0.1I) \begin{pmatrix} 0.25 \\ 0.24 \end{pmatrix} = \begin{pmatrix} 1.01 \\ 1.01 \end{pmatrix}$$

Diese Lösungsgleichung löst auch das modifizierte Optimierungsproblem mit dem Tikhonov-Regularization-Term  $\|\mathbf{x}\|^2$ .

$$\min f(\mathbf{x}) + \lambda \|\mathbf{x}\|^2$$

Denn durch Ableiten und Nullsetzen erhalten wird die gleiche modifizierte Normalgleichung:

$$\nabla(f(\mathbf{x}) + \lambda \|\mathbf{x}\|^2) = 2A^T A \mathbf{x} - 2A^T \mathbf{y} + 2\lambda \mathbf{x} = 2(A^T A + \lambda I) \mathbf{x} - 2A^T \mathbf{y} = \mathbf{0}$$

Für nichtlineare Optimierungen kann die gleiche Art Regularisierung verwendet werden. Man beachte, dass die Newton-Methode mit Trust-Region die gleiche Form hat.

## 12.2 LASSO (least absolute shrinkage and selection operator)

Ähnlich zur Ridge-Regularization, wo die  $\ell_2$ -Norm verwendet wird, wird bei LASSO die  $\ell_1$ -Norm verwendet.

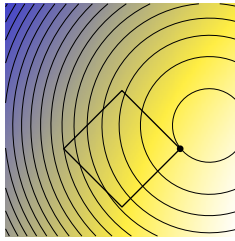
$$\min f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1, \quad \|\mathbf{x}\|_1 = \sum_i |x_i|$$

Das Wichtigste bei dieser Regularisierung ist, dass die Lösung dazu tendiert, *sparse* zu sein, d.h. dass viele Koeffizienten des Lösungsvektors 0 sind. Der Grund dafür ist, dass der Schnitt zwischen Höhenlinien von  $f(\mathbf{x})$  und  $\lambda \|\mathbf{x}\|_1$  aufgrund letzterer geometrischen Form häufig auf Koordinatenachsen liegt. Siehe Abbildung 19.

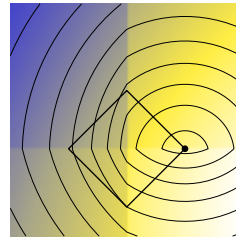
Die starke (und übliche) Version von LASSO ist die Optimierung mit Constraint:

$$\min f(\mathbf{x}) \text{ mit } \|\mathbf{x}\|_1 \leq R$$

Durch Anwendung des Lagrange-Multiplikators  $\lambda$  kommt man dann auf obige Regularisierung. Bei den Optimierungsschritten muss man aber darauf achten, dass  $\mathbf{x}$  tatsächlich die Bedingung  $\|\mathbf{x}\|_1 \leq R$  erfüllt. Siehe dazu Abschnitt 13.



(a) Schnitt zwischen Höhenlinie von  $f(\mathbf{x})$  und  $\|\mathbf{x}\| = 1$



(b) Minimum von  $f(\mathbf{x}) + \|\mathbf{x}\|_1$

Abbildung 19: LASSO

### 13 Projected Gradient-Descent

Die Projected Gradient-Descent-Methode ist eine Variante des Gradient-Descent, bei der garantiert wird, dass die  $\ell_1$ -Bedingung  $\|\mathbf{x}\|_1 \leq R$  erfüllt ist. Es wird also jener Punkt  $\underline{\mathbf{x}}$  gesucht, der dem normalen Gradient-Descent-Schritt  $\mathbf{x} + \alpha \mathbf{d}$  am nächsten liegt, aber die Bedingung erfüllt. Sie löst also das LASSO-Problem:

$$\underline{\mathbf{x}} = \arg \min_{\|\mathbf{u}\|_1 \leq R} \|\mathbf{u} - (\mathbf{x} + \alpha \mathbf{d})\|^2$$

Wenn  $\mathbf{x} + \alpha \mathbf{d} \leq R$  ist, dann wird natürlich einfach  $\underline{\mathbf{x}} = \mathbf{u} = \mathbf{x} + \alpha \mathbf{d}$  gewählt. Falls es aber außerhalb liegt, muss ein Punkt am Rand  $\|\mathbf{x}\|_1 = R$  gewählt werden. Man sieht nun, dass manche Koordinaten von  $\mathbf{x}$  gleich 0 werden, während die anderen Koordinaten alle den gleichen Abstand  $c$  von  $x_i + \alpha d_i$  haben, weil der Vektor-Abstand normal auf den Rand steht und  $R = \sum |x_i| = \sum q_i x_i$  mit  $|q_i| = 1$ , und  $\mathbf{q}$  ist ein Normalvektor. Außerdem sind die Koordinaten, die nicht 0 werden, jene, die betragsmäßig größer sind. Wenn  $b_i$  jene absoluten Koordinaten von  $\mathbf{x} + \alpha \mathbf{d}$  sind, die nicht 0 werden, dann ist  $\sum |x_i| = R$  und  $|x_i| = b_i - c$ . Daraus ergibt sich  $c = \frac{1}{k} ((\sum_{i=1}^k b_i) - R)$ . Weil noch  $c < b_k$  sein muss, sortiert der Algorithmus die  $b_i$  und findet das größte  $k$  sodass dies noch erfüllt ist:

$$\mathbf{a} = \mathbf{x} + \alpha \mathbf{d}$$

$$\mathbf{b} = \text{sort}_{\text{absteigend}} |\mathbf{a}|$$

$$k^* = \max \left\{ k \mid \frac{1}{k} ((\sum_{i=1}^k b_i) - R) < b_k \right\}$$

$$c = \frac{1}{k^*} ((\sum_{i=1}^{k^*} b_i) - R)$$

$$\underline{\mathbf{x}}_i = \begin{cases} a_i - \text{sign}(a_i)c & |a_i| > c \\ 0 & |a_i| \leq c \end{cases}$$

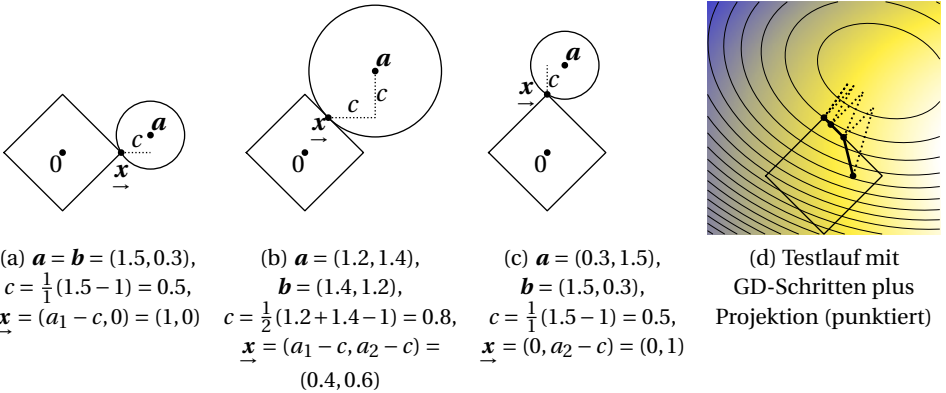


Abbildung 20: Projected Gradient-Descent

Abbildung 20(a-c) zeigt ein 2D-Beispiel. Ein Beispiel höherer Dimension für  $\mathbf{a} = (0.6, -1.5, 1.8, -1.1)$  und  $R = 1$  wäre:

$i$	$b_i$	$c$	$\underline{x}_i$
1	1.8	$(1.8 - 1)/1 = 1$	0 ( $0.6 < 1.15$ )
2	1.5	$(\mathbf{1.8} + \mathbf{1.5} - \mathbf{1})/2 = \mathbf{1.15}$	$-1.5 + 1.15 = -0.35$
3	1.1	$(1.8 + 1.5 + 1.1 - 1)/3 = 1.133 > \mathbf{1.1}$	$1.8 - 1.15 = 0.65$
4	0.6		0 ( $ -1.1  < 1.15$ )

Abbildung 20(d) zeigt einen Testlauf mit Gradient-Descent-Schritt ( $\mathbf{a}$ ) und Projektion ( $\underline{\mathbf{x}}$ ) in jedem Schritt.

## 14 Proximal Gradient-Descent

Das schwächere  $\ell_1$ -Problem

$$\min f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$$

wird von der *Proximal Gradient-Descent*-Methode gelöst. Dabei wird der Gradient-Descent-Schritt mit einem *proximal operator*  $\text{prox}_\varphi$  (für eine Regularization-Funktion  $\varphi$ , meist  $\varphi(\mathbf{u}) = \lambda \|\mathbf{u}\|_1$ ) korrigiert.

$$\underline{\mathbf{x}} = \text{prox}_{\alpha\lambda\|\cdot\|_1}(\underline{\mathbf{x}}) = \arg \min_{\mathbf{u}} \frac{1}{2} \|\mathbf{u} - \underline{\mathbf{x}}\|^2 + \alpha\lambda \|\mathbf{u}\|_1 \quad \underline{\mathbf{x}} := \mathbf{x} + \alpha \mathbf{d}$$

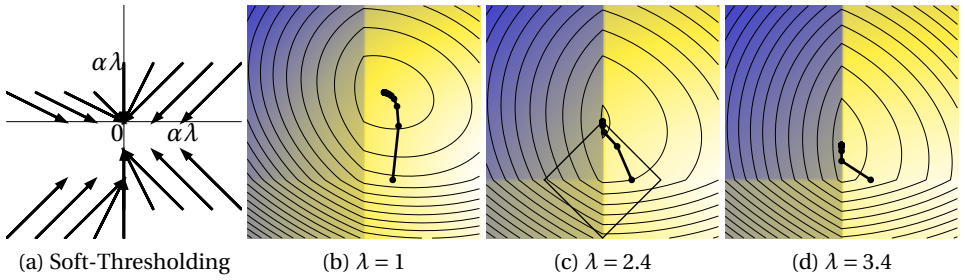


Abbildung 21: Proximal Gradient-Descent

Durch Ableiten und Nullsetzen (nach  $\mathbf{u}$ ) bekommen wir folgende Gleichung:

$$\mathbf{u} - \underline{\mathbf{x}} + \alpha \lambda \nabla \|\mathbf{u}\|_1 = 0, \quad (\nabla \|\mathbf{u}\|_1)_i = \begin{cases} 1 & u_i > 0 \\ [-1, 1] & u_i = 0 \\ -1 & u_i < 0 \end{cases}$$

wobei  $\|\mathbf{u}\|_1$  nicht überall differenzierbar ist. Es gibt an den Stellen mit  $u_i = 0$  eine Subgradienten (siehe Abbildung 6), wodurch ein gewisser Freiheitsgrad entsteht. Dieser ermöglicht die Lösung:

$$\underline{x}_i = u_i = \underline{x}_i + \begin{cases} -\alpha \lambda & \underline{x}_i > \alpha \lambda \\ -\underline{x}_i & -\alpha \lambda \leq \underline{x}_i \leq \alpha \lambda \\ +\alpha \lambda & \underline{x}_i < -\alpha \lambda \end{cases} = \begin{cases} \underline{x}_i - \alpha \lambda & \underline{x}_i > \alpha \lambda \\ 0 & -\alpha \lambda \leq \underline{x}_i \leq \alpha \lambda \\ \underline{x}_i + \alpha \lambda & \underline{x}_i < -\alpha \lambda \end{cases}$$

Eine andere Wahl würde nicht funktionieren, weil wenn z.B.  $0 < \underline{x}_i < \alpha \lambda$  und  $u_i > 0$  wäre, dann müsste  $u_i = \underline{x}_i - \alpha \lambda < 0$  sein (Widerspruch), und umgekehrt. Diesen Update-Schritt nennt man auch *soft threshold* (Abbildung 21(a)).

Der Update-Schritt ist sinnvoll, denn wir können zeigen, dass ein Optimum  $\mathbf{x}^*$  durch den Update-Schritt nicht verändert wird (Fixpunkt), und umgekehrt ein Fixpunkt auch ein Optimum ist. Für ein Optimum  $\mathbf{x}^*$  gilt nämlich, dass die Ableitung von  $f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$  gleich 0 ist:

$$\nabla f(\mathbf{x}^*) + \lambda \nabla \|\mathbf{x}^*\|_1 = \mathbf{0} \Rightarrow \frac{df(\mathbf{x}^*)}{dx_i} = -\lambda \frac{d\|\mathbf{x}^*\|_1}{dx_i^*} = \begin{cases} -\lambda & x_i^* > 0 \\ \lambda[-1, 1] & x_i^* = 0 \\ +\lambda & x_i^* < 0 \end{cases}$$

Für den Fall  $x_i^* > 0$  bleibt  $\underline{x}_i = x_i^*$ , weil:

$$x_i^* > 0 \Leftrightarrow \underline{x}_i = x_i^* - \alpha \frac{df(\mathbf{x}^*)}{dx_i} = x_i^* + \alpha \lambda > \alpha \lambda \Leftrightarrow \underline{x}_i = \underline{x}_i - \alpha \lambda = x_i^*$$

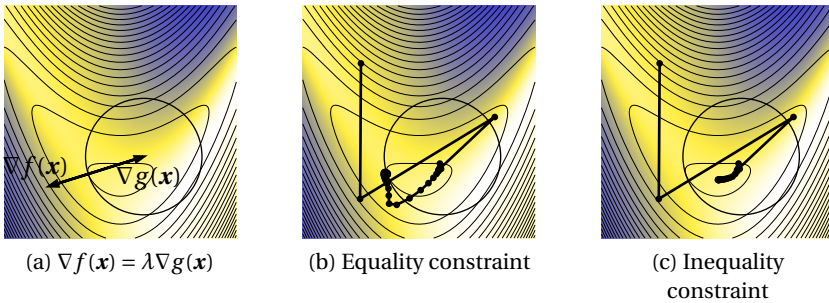


Abbildung 22: Gradient-Descent mit Constraints

Für  $x_i^* < 0$  gilt das analog. Und für  $x_i^* = 0$  gilt ebenfalls:

$$x_i^* = 0 \Leftrightarrow \underline{x}_i \in x_i^* + \alpha\lambda[-1, 1] \Leftrightarrow -\alpha\lambda \leq \underline{x}_i \leq \alpha\lambda \Leftrightarrow \underline{x}_i = 0 = x_i^*$$

Abbildung 21(b-d) zeigt drei Testläufe für die gleiche Funktion  $f$  wie in Abbildung 20, nur eben mit  $\lambda\|\mathbf{x}\|_1$  addiert. (c) findet die gleiche Lösung wie die Projected-Gradient-Descent-Methode, was zeigt, dass es immer ein bestimmtes  $\lambda$  gibt, so dass das schwache LASSO-Problem die Lösung des starken findet. In diesem Fall ist übrigens  $\lambda = -\frac{df((0,1))}{dx_1}$ .

## 15 Constraints

Constraints (Nebenbedingungen) schränken der Wertebereich der erlaubten Lösungsvektoren ein. Es gibt Gleichheits- und Ungleichheits-Constraints. Für ein einzelnes Gleichheits-Constraint müssen wir folgendes Problem lösen:

$$\min f(\mathbf{x}) \text{ mit } g(\mathbf{x}) = 0$$

wobei ein Constraint einer anderen Form (z.B.  $\|\mathbf{x}\|^2 = 5$ ) immer auf diese Form gebracht werden kann (z.B.  $\|\mathbf{x}\|^2 - 5 = 0$ ). Es kann dann gezeigt werden, dass die Gradienten von  $f$  und  $g$  an einer Optimumsstelle  $\mathbf{x}$  parallel sind, d.h. dass es ein  $\lambda \in \mathbb{R}$  (*Lagrange-Multiplikator*) gibt sodass:

$$\nabla_{\mathbf{x}}(f(\mathbf{x}) - \lambda g(\mathbf{x})) = 0, \quad g(\mathbf{x}) = 0$$

Siehe Abbildung 22(a). Die Bedingung  $g(\mathbf{x}) = 0$  muss nach wie vor berücksichtigt werden, weil man diese zusätzliche Gleichung benötigt, um die zusätzliche Variable  $\lambda$  zu fixieren.

Wenn es mehr als ein Constraint gibt, also  $g_i(\mathbf{x}) = 0$ , dann kann man die  $g_i$  als Vektor  $\mathbf{g}$  schreiben:

$$\min f(\mathbf{x}) \text{ mit } \mathbf{g}(\mathbf{x}) = \mathbf{0}$$

Für jedes Constraint wird ein eigener Lagrange-Multiplikator  $\lambda_i$  eingeführt. Wenn man diese  $\lambda_i$  ebenfalls als Vektor  $\boldsymbol{\lambda}$  schreibt, dann gilt für ein Optimum  $\mathbf{x}$ :

$$\nabla_{\mathbf{x}}(f(\mathbf{x}) - \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x})) = \mathbf{0}, \quad \mathbf{g}(\mathbf{x}) = \mathbf{0}$$

Als Beispiel verwenden wir das Constraint  $\|\mathbf{x} - (0.4, 0.4)\|^2 = 1$ , und damit  $g(\mathbf{x}) = 1 - \|\mathbf{x} - (0.4, 0.4)\|^2$ . Um per Gradient-Descent nicht nur  $f(\mathbf{x})$  minimal werden zu lassen, sondern auch  $\lambda$  so zu updaten, dass  $g(\mathbf{x}) = 0$  wird, wollen wir zusätzlich  $\frac{1}{2}g(\mathbf{x})^2$  minimieren, was den Update-Schritt  $-\alpha g(\mathbf{x})\nabla g(\mathbf{x})$  in den Update-Schritt  $-\alpha(\nabla f(\mathbf{x}) - \lambda\nabla g(\mathbf{x}))$  integrieren, was mit einem Update von  $\lambda$  durch  $-g(\mathbf{x})$  gemacht wird. Das sich dadurch ergebende Schema ist:

$$\begin{aligned} \lambda &= \lambda - g(\mathbf{x}) \\ \underline{\mathbf{x}} &= \mathbf{x} - \alpha(\nabla f(\mathbf{x}) - \lambda\nabla g(\mathbf{x})) \end{aligned}$$

Ein Testlauf ( $\alpha = 0.1$ ) ist in Abbildung 22(b) zu sehen. Das Schema ist aber sehr instabil und tendiert zu Schwingungsverhalten um  $g(\mathbf{x}) = 0$ . Bessere Verfahren kommen später.

Wenn zusätzlich Ungleichheitsconstraints nötig sind,

$$\min f(\mathbf{x}) \text{ mit } \mathbf{g}(\mathbf{x}) = \mathbf{0}, \mathbf{h}(\mathbf{x}) \geq \mathbf{0}$$

dann wird das Schema nach Karush-Kuhn-Tucker (KKT) angewendet. Es gibt für die Ungleichheitsconstraints zusätzliche Lagrange-Multiplikatoren  $\mu_i$ . An einer Optimumsstelle  $\mathbf{x}$  gilt:

$$\nabla_{\mathbf{x}}(f(\mathbf{x}) - \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}^T \mathbf{h}(\mathbf{x})) = \mathbf{0}, \quad \mathbf{g}(\mathbf{x}) = \mathbf{0}, \mathbf{h}(\mathbf{x}) \geq \mathbf{0}, \quad \boldsymbol{\mu} \geq \mathbf{0}, \quad \boldsymbol{\mu}^T \mathbf{h}(\mathbf{x}) = 0$$

Auch hier braucht man die ursprünglichen Constraints, um die zusätzlichen Variablen zu fixieren. Außerdem müssen noch die  $\mu_i \geq 0$  sein (vierte Bedingung). Und die fünfte Bedingung bedeutet praktisch nur, dass entweder  $\mu_i = 0$  ist oder  $h_i(\mathbf{x}) = 0$  oder beide. Aus  $h_i(\mathbf{x}) \geq 0$  und  $\mu_i \geq 0$  folgt nämlich:

$$\boldsymbol{\mu}^T \mathbf{h}(\mathbf{x}) = 0 \Leftrightarrow \forall i : \mu_i h_i(\mathbf{x}) = 0 \Leftrightarrow \forall i : \mu_i = 0 \vee h_i(\mathbf{x}) = 0$$

Die KKT-Bedingungen gelten genau genommen aber nur, wenn für die Optimumsstelle die LICQ-Bedingung gilt (linear independence constraint qualification). Dafür müssen die  $g_i$ -Gradienten sowie die *aktiven*  $h_i$ -Gradienten (wo  $h_i(\mathbf{x}) = 0$ , also Minimum am Rand des erlaubten Bereichs

$$\{\nabla g_i(\mathbf{x})\} \cup \{\nabla h_i(\mathbf{x}) \mid h_i(\mathbf{x}) = 0\} \quad \text{linear unabhängig sind.}$$

Lineare Unabhängigkeit bedeutet, dass eine Linearkombination von Vektoren  $\mathbf{v}_i$ :  $\sum_i a_i \mathbf{v}_i$  nur 0 werden kann, wenn alle  $a_i = 0$  sind. Oder anders gesagt: Kein Vektor liegt in der Hyperebene, die durch die jeweils anderen Vektoren aufgespannt wird.

Abbildung 22(c) zeigt einen Testlauf für ein ähnliches Gradient-Descent-Schema wie bei den Equality-Constraints, hier für  $\|\mathbf{x} - (0.4, 0.4)\|^2 \leq 1$ , also  $h(\mathbf{x}) = 1 - \|\mathbf{x} - (0.4, 0.4)\|^2$ , wobei wir  $\mu \geq 0$  im Update-Schritt garantieren:

$$\begin{aligned} \mu &= \max(\mu - g(\mathbf{x}), 0) \\ \underline{\mathbf{x}} &= \mathbf{x} - \alpha(\nabla f(\mathbf{x}) - \mu \nabla g(\mathbf{x})) \end{aligned}$$

## 16 Simplex-Methode

Das Optimieren einer linearen Funktion in  $\mathbf{x}$  macht nur mit Constraints Sinn, da sonst ein Funktionswert  $-\infty$  erreicht werden kann. Die *Standard-Form* für lineare Optimierung ist:

$$\min \mathbf{c}^\top \mathbf{x} \text{ mit } A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}$$

Hier hat  $A$   $m$  Zeilen (Anzahl der Constraints) und  $n$  Spalten (Anzahl der Variablen), wobei üblicherweise  $n > m$  ist.

Um Constraints in Nicht-Standard-Form auf Standard-Form zu bringen, verwendet man zwei Methoden. Erstens werden generelle lineare Ungleichungen über sogenannte Slack-Variablen  $\mathbf{y}$  in absolute Ungleichungen (Vergleich mit 0) plus Gleichungen verwandelt:

$$A\mathbf{x} \geq \mathbf{b} \Leftrightarrow A\mathbf{x} - \mathbf{y} = \mathbf{b}, \mathbf{y} \geq \mathbf{0} \quad \text{oder} \quad A\mathbf{x} \leq \mathbf{b} \Leftrightarrow A\mathbf{x} + \mathbf{y} = \mathbf{b}, \mathbf{y} \geq \mathbf{0}$$

Zweitens werden nicht-ingeschränkte Variablen (also  $\mathbf{x}$ ) in einen positiven  $\mathbf{x}^+$  und negativen Teil  $\mathbf{x}^-$  aufgespalten. Mit  $\mathbf{x} = \mathbf{x}^+ - \mathbf{x}^-$ , wobei  $\mathbf{x}^+ = \max(\mathbf{x}, \mathbf{0})$ ,  $\mathbf{x}^- = \max(-\mathbf{x}, \mathbf{0})$  ist, ergibt sich:

$$\min \begin{pmatrix} \mathbf{c} \\ -\mathbf{c} \\ \mathbf{0} \end{pmatrix}^\top \begin{pmatrix} \mathbf{x}^+ \\ \mathbf{x}^- \\ \mathbf{y} \end{pmatrix} \text{ mit } (A \quad -A \quad I) \begin{pmatrix} \mathbf{x}^+ \\ \mathbf{x}^- \\ \mathbf{y} \end{pmatrix} = \mathbf{b}, \quad \begin{pmatrix} \mathbf{x}^+ \\ \mathbf{x}^- \\ \mathbf{y} \end{pmatrix} \geq \mathbf{0}$$

was die Standard-Form hat. Falls die Bedingungen eine gemischte Form haben, können diese Transformationen auch zeilenweise angewendet werden.

Die Karush-Kuhn-Tucker-Bedingungen (KKT) für ein Optimum der Standardform sind:

$$\nabla_{\mathbf{x}}(\mathbf{c}^\top \mathbf{x} - \boldsymbol{\lambda}^\top (A\mathbf{x} - \mathbf{b}) - \boldsymbol{\mu}^\top \mathbf{x}) = \mathbf{0} \quad \Leftrightarrow \quad A^\top \boldsymbol{\lambda} + \boldsymbol{\mu} = \mathbf{c}$$



$$Ax = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \boldsymbol{\mu} \geq \mathbf{0}, \boldsymbol{\mu}^\top \mathbf{x} = 0$$

Die Menge der erlaubten Punkte, die die Gleichungen und Ungleichungen erfüllen ergeben ein Polytop, das von Hyperebenen begrenzt wird. Es lässt sich nun zeigen (fundamentaler Satz der linearen Optimierung), dass alle Eckpunkte dieses Polytops durch eine Auswahl  $\boldsymbol{\kappa}$  von  $m$  Variablen-Indizes bestimmt ist. Es seien  $\bar{\boldsymbol{\kappa}}$  die Indizes, die nicht in  $\boldsymbol{\kappa}$  sind:  $\bar{\boldsymbol{\kappa}} = \{k \mid k \notin \boldsymbol{\kappa}\}$ . Es seien  $A_{\boldsymbol{\kappa}}, \mathbf{x}_{\boldsymbol{\kappa}}$ , etc. die Matrizen und Vektoren, die auf die Indizes  $\boldsymbol{\kappa}$  beschränkt sind. Dann gilt:

$$\mathbf{x}_{\bar{\boldsymbol{\kappa}}} = \mathbf{0}, \quad A\mathbf{x} = A_{\boldsymbol{\kappa}}\mathbf{x}_{\boldsymbol{\kappa}} = \mathbf{b}, \quad \mathbf{x}_{\boldsymbol{\kappa}} = A_{\boldsymbol{\kappa}}^{-1}\mathbf{b}$$

$$\boldsymbol{\mu}_{\boldsymbol{\kappa}} = \mathbf{0}, \quad A_{\boldsymbol{\kappa}}^\top \boldsymbol{\lambda} = \mathbf{c}_{\boldsymbol{\kappa}}, \quad A_{\bar{\boldsymbol{\kappa}}}^\top \boldsymbol{\lambda} + \boldsymbol{\mu}_{\bar{\boldsymbol{\kappa}}} = \mathbf{c}_{\bar{\boldsymbol{\kappa}}} \Rightarrow \boldsymbol{\lambda} = \left(A_{\boldsymbol{\kappa}}^\top\right)^{-1} \mathbf{c}_{\boldsymbol{\kappa}}, \quad \boldsymbol{\mu}_{\bar{\boldsymbol{\kappa}}} = \mathbf{c}_{\bar{\boldsymbol{\kappa}}} - A_{\bar{\boldsymbol{\kappa}}}^\top \boldsymbol{\lambda}$$

Wenn nun eine Menge  $\boldsymbol{\kappa}$  gefunden wird, wo  $\boldsymbol{\mu}_{\bar{\boldsymbol{\kappa}}} \geq \mathbf{0}$  ist, dann hat man das Minimum gefunden. Es wird daher ein *Entering-Index*  $q \in \bar{\boldsymbol{\kappa}}$  mit  $\mu_q < 0$  ausgewählt. Die Veränderung der Gleichung durch Hinzunahme des  $x_q > 0$  ist dann

$$A_{\boldsymbol{\kappa}}(\mathbf{x}_{\boldsymbol{\kappa}} - \boldsymbol{\Delta}) + A_q x_q = \mathbf{b}, \quad \boldsymbol{\Delta} = A_{\boldsymbol{\kappa}}^{-1} A_q x_q, \quad \mathbf{x}_{\boldsymbol{\kappa}} - \boldsymbol{\Delta} \geq \mathbf{0},$$

wobei  $\mathbf{x}_{\boldsymbol{\kappa}}$  für einen Index  $p$  gleich 0 werden soll, also  $x_p - \Delta_p = 0$ .  $p$  ist der *Leaving-Index* und ergibt sich aus:

$$p = \operatorname{argmin}_{\substack{k \in \boldsymbol{\kappa} \\ d_k > 0}} \frac{x_k}{d_k}, \quad \mathbf{d} = A_{\boldsymbol{\kappa}}^{-1} A_q$$

Für die nächste Iteration gilt:

$$\underline{\boldsymbol{\kappa}} = \boldsymbol{\kappa} \cup \{q\} \setminus \{p\}$$

Für die initiale Index-Auswahl werden für  $\boldsymbol{\kappa}$  meist die Indizes der Slack-Variablen ( $\mathbf{y}$ ) verwendet. Wegen  $A\mathbf{x} = I\mathbf{y} = \mathbf{b}$  gilt  $\mathbf{x} \geq \mathbf{0}$ , da  $\mathbf{b} \geq \mathbf{0}$  ist bzw. gemacht werden kann. Abbildung 23 zeigt einen Testlauf für

$$\min \begin{pmatrix} 1 \\ 1 \\ 0.5 \end{pmatrix}^\top \mathbf{x} \quad \text{mit} \quad \begin{pmatrix} 1 & 1 & 3 \\ 1 & 6 & 1 \\ 6 & 1 & 1 \end{pmatrix} \mathbf{x} \geq \begin{pmatrix} -2.5 \\ -2.5 \\ -2.5 \end{pmatrix}$$

## 17 Penalty-Methode

$$\min f(\mathbf{x}) \quad \text{mit} \quad \mathbf{g}(\mathbf{x}) = \mathbf{0}, \mathbf{h}(\mathbf{x}) \geq \mathbf{0}$$

$$\min f(\mathbf{x}) + l \|\mathbf{g}(\mathbf{x})\|^2 + m \|\max(-\mathbf{h}(\mathbf{x}), \mathbf{0})\|^2$$

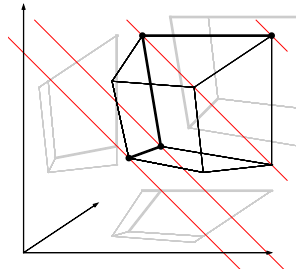


Abbildung 23: Simplex-Algorithmus. Beginnend rechts oben:  $\kappa = (6, 7, 8) \mapsto (3, 6, 7) \mapsto (3, 4, 6) \mapsto (3, 4, 5)$ .

$l, m$  gegen  $\infty$  wachsen lassen.

Andere Penalty-Funktionen als  $\|\cdot\|^2$  möglich. Z.B.  $\ell_1$ -Norm  $\|\cdot\|_1$ . Diese ist oft effizienter.

## 18 Augmented Lagrange-Methode

Verbesserung der Penalty-Methode um Lagrange-Multiplikatoren.

$$\min f(\mathbf{x}) \text{ mit } \mathbf{g}(\mathbf{x}) = \mathbf{0}$$

$$\mathbf{0} = \nabla_{\mathbf{x}} \left( f(\mathbf{x}) - \lambda^T \mathbf{g}(\mathbf{x}) + \frac{1}{2} \mathbf{l}^T \mathbf{g}(\mathbf{x})^2 \right) = \nabla f(\mathbf{x}) - (\lambda - \mathbf{l} \mathbf{g}(\mathbf{x}))^T \nabla \mathbf{g}(\mathbf{x})$$

$\lambda - \mathbf{l} \mathbf{g}(\mathbf{x})$  wäre optimaler Lagrange-Multiplikator. Daher Update-Schema:

$$\underline{\lambda} = \lambda - \mathbf{l} \mathbf{g}(\mathbf{x})$$

## 19 Alternating Direction Multiplier-Methode

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x}) \longrightarrow \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) \text{ mit } \mathbf{x} = \mathbf{y}$$

Lösen mit Augmented Lagrange-Methode abwechselnd

$$\underline{\mathbf{x}} = \arg \min_{\mathbf{x}} f(\mathbf{x}) \text{ mit } \mathbf{x} \approx \underline{\mathbf{y}}$$

$$\underline{\mathbf{y}} = \arg \min_{\mathbf{y}} g(\mathbf{y}) \text{ mit } \mathbf{y} \approx \underline{\mathbf{x}}$$

## 20 Active Set-Methode

$$\min \frac{1}{2} \mathbf{x}^\top G \mathbf{x} + \mathbf{c}^\top \mathbf{x} \quad \text{mit } A \mathbf{x} \geq \mathbf{b}$$

Active set  $\alpha$ :

$$A_\alpha = \mathbf{b}_\alpha$$

$$\min_{\mathbf{d}} \frac{1}{2} (\mathbf{x} + \mathbf{d})^\top G (\mathbf{x} + \mathbf{d}) + \mathbf{c}^\top (\mathbf{x} + \mathbf{d}) \quad \text{mit } A_\alpha (\mathbf{x} + \mathbf{d}) = \mathbf{b}_\alpha$$

Lagrange-Multiplikator:

$$G(\mathbf{x} + \mathbf{d}) + \mathbf{c} - A_\alpha^\top \boldsymbol{\lambda} = \mathbf{0}, \quad A_\alpha (\mathbf{x} + \mathbf{d}) = \mathbf{b}_\alpha$$

$$\begin{pmatrix} G & -A_\alpha^\top \\ A_\alpha & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{d} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} -G\mathbf{x} - \mathbf{c} \\ -A_\alpha \mathbf{x} + \mathbf{b}_\alpha \end{pmatrix}$$

Entering index  $q$  und Verkürzung von  $\mathbf{d}$ , sodass Ungleichungen in  $\bar{\alpha}$  erfüllt sind:

$$A_i(\mathbf{x} + \alpha_i \mathbf{d}) = b_i \quad \Rightarrow \quad \alpha_i = \frac{b_i - A_i \mathbf{x}}{A_i \mathbf{d}}, \quad q = \arg \min_{\substack{i \in \bar{\alpha} \\ 0 \leq \alpha_i < 1}} \alpha_i$$

Leaving index  $p$ :

$$p = \arg \min_{\substack{i \in \alpha \\ \lambda_i < 0}} \lambda_i$$

$\underline{\alpha} = \alpha \cup \{q\} \setminus \{p\}$  sofern  $p$  und  $q$  gefunden wurden

$\underline{\mathbf{x}} = \mathbf{x} + \alpha_q \mathbf{d}$  oder  $\underline{\mathbf{x}} = \mathbf{x} + \mathbf{d}$ , falls kein  $q$  gefunden wurde

Anwendung auf nicht-quadratische Funktionen: Lokale Approximation durch quadratische Funktion ( $G = \nabla^2 f(\mathbf{x})$ ,  $\mathbf{c} = \nabla f(\mathbf{x})$ ). D.h. Newton-Methode mit Active-Set, auch Sequential-Quadratic-Programming (SQP) genannt. Approximation durch lineare Funktion führt zu Simplex-Methode mit Neuberechnung des  $\mathbf{c}$ -Vektors in jedem Schritt, auch Sequential-Linear-Programming (SLP) genannt.

## 21 Interior Point-Methode

Lineare Optimierung:

$$\min \mathbf{c}^\top \mathbf{x} \quad \text{mit } A \mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}$$

KKT:

$$A^\top \boldsymbol{\lambda} + \boldsymbol{\mu} = \mathbf{c}, \quad A \mathbf{x} = \mathbf{b}, \quad \mathbf{x} \geq \mathbf{0}, \quad \boldsymbol{\mu} \geq \mathbf{0}, \quad \boldsymbol{\mu}^\top \mathbf{x} = 0$$

Duales Problem:

$$\min -\mathbf{b}^\top \boldsymbol{\lambda} \quad \text{mit } \mathbf{c} - A^\top \boldsymbol{\lambda} \geq \mathbf{0}$$

KKT mit  $\mathbf{x}$  als Lagrange-Multiplikator:

$$\nabla_{\boldsymbol{\lambda}}(-\mathbf{b}^\top \boldsymbol{\lambda} - \mathbf{x}^\top(\mathbf{c} - A^\top \boldsymbol{\lambda})) = \mathbf{0} \quad \Leftrightarrow \quad A\mathbf{x} = \mathbf{b}$$

Und wenn  $\boldsymbol{\mu} = \mathbf{c} - A^\top \boldsymbol{\lambda}$  gesetzt wird, sind die KKT-Bedingungen genau gleich.

$$\nabla g = \begin{pmatrix} A^\top \boldsymbol{\lambda} + \boldsymbol{\mu} - \mathbf{c} \\ A\mathbf{x} - \mathbf{b} \\ \boldsymbol{\mu}^\top \mathbf{x} \end{pmatrix} \Rightarrow \nabla^2 g = \begin{pmatrix} \mathbf{0} & A^\top & I \\ A & \mathbf{0} & \mathbf{0} \\ \text{diag}(\boldsymbol{\mu}) & \mathbf{0} & \text{diag}(\mathbf{x}) \end{pmatrix}$$

Suchen  $\nabla g = \mathbf{0}$  mittels Newton:

$$(\nabla^2 g)\mathbf{d} = -\nabla g$$

## 22 Differenzieren von Vektoren und Matrizen

Um einen Vektor oder eine Matrix zu konstruieren, verwenden wir hier die Notation

$$\mathbf{a} = (a_i \mid i = 1 \dots n) = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix},$$

$$A = (A_{i,j} \mid i = 1 \dots m, j = 1 \dots n) = \begin{pmatrix} A_{1,1} & \cdots & A_{1,n} \\ \vdots & \ddots & \vdots \\ A_{m,1} & \cdots & A_{m,n} \end{pmatrix}.$$

Der Gradient einer Funktion ist dann

$$\nabla f(\mathbf{x}) = \left( \frac{df}{dx_i}(\mathbf{x}) \mid i = 1 \dots n \right).$$

Der Gradient von inneren Produkten:

$$\nabla \mathbf{a}^\top \mathbf{x} = \left( \frac{d}{dx_j} \sum a_i x_i \mid j \right) = (a_j \mid j) = \mathbf{a}$$

$$\nabla (|\mathbf{x}|^2) = \nabla (\mathbf{x}^\top \mathbf{x}) = \left( \frac{d}{dx_j} \sum x_i^2 \mid j \right) = (2x_j \mid j) = 2\mathbf{x}$$

Gradient einer quadratischen Form:

$$\begin{aligned}\nabla \mathbf{x}^\top A \mathbf{x} &= \left( \frac{d}{dx_k} \sum_{i,j} x_i A_{i,j} x_j \middle| k \right) = \left( \sum_i x_i A_{i,k} + \sum_j A_{k,j} x_j \middle| k \right) \\ &= \left( \sum_i (A_{i,k} + A_{k,i}) x_i \middle| k \right) = (A^\top + A) \mathbf{x} \stackrel{A \text{ symmetrisch}}{=} 2A \mathbf{x}\end{aligned}$$

Jacobi-Matrizen:

$$\nabla \mathbf{f}(\mathbf{x}) = \left( \frac{df_i}{dx_j} \middle| i = 1 \dots m, j = 1 \dots n \right) = \left( \frac{d\mathbf{f}(\mathbf{x})}{dx_1} \dots \frac{d\mathbf{f}(\mathbf{x})}{dx_n} \right) = \begin{pmatrix} \nabla f_1(\mathbf{x})^\top \\ \vdots \\ \nabla f_m(\mathbf{x})^\top \end{pmatrix}$$

$$\nabla A \mathbf{x} = \left( \frac{d}{dx_k} \sum A_{i,j} x_j \middle| i, k \right) = (A_{i,k} \mid i, k) = A$$

Hesse-Matrix:

$$\nabla^2 f(\mathbf{x}) = \nabla(\nabla f(\mathbf{x})) = \left( \frac{d^2 f}{dx_i dx_j}(\mathbf{x}) \middle| i, j \right)$$

$$\nabla^2 \mathbf{x}^\top A \mathbf{x} = \nabla(A^\top + A) \mathbf{x} = A^\top + A \stackrel{A \text{ symmetrisch}}{=} 2A$$