

Mehrdimensionale Optimierung ohne Nebenbedingungen

Sonja Biedermann, Felix Freynschlag, Bernhard Hayden,
Stepan Kharin, Christoph Pressler

14. Januar 2016

Inhaltsverzeichnis

1	Das Nelder-Mead-Verfahren	1
1.1	Einführung	1
1.2	Funktionsweise	2
1.3	Erläuterung der einzelnen Transformationsschritte	3
1.3.1	Reflektion	3
1.3.2	Expansion	3
1.3.3	Kontraktion	4
1.3.4	Komprimierung	5
1.4	Ein Beispiel: die Himmelblau-Funktion	6
1.5	Analyse	10
2	Das Gradientenverfahren	11
2.1	Die Idee und der Zweck des Gradientenverfahren	11
2.2	Der Ablauf des Gradientenverfahrens	11
2.3	Beispiele	13
2.3.1	Beispiel 1	13
3	Die Koordinatenabstiegsmethode	16
3.1	Einführung	16
3.2	Beispiel	17
3.3	Problemfälle	18
	Literatur	19

1 Das Nelder-Mead-Verfahren

1.1 Einführung

Das Nelder-Mead-Verfahren (auch Downhill-Simplex genannt) verfolgt einen geometrischen Ansatz zur iterativen Optimierung: ein Simplex bestehend aus $N + 1$ Punkten (im Falle $N = 2$, dem zweidimensionalen Raum, also ein Dreieck) wird fortwährend so transformiert, sodass sich der Simplex dem Minimum nähert, um welches er sich dann zusammenzieht. [1]

1.2 Funktionsweise

Der Algorithmus besteht im wesentlichen aus drei Phasen [2], die in jeder Iteration durchlaufen werden:

1. Ordering
2. Centroid
3. Transformation

Das Verhalten des Algorithmus' kann mittels den Parametern α , β , γ und δ beeinflusst werden, wobei α den Reflektionsfaktor, γ den Expansionsfaktor, β den Kontraktionsfaktor und δ den Komprimierungsfaktor bezeichnet. Eine typische Belegung wäre

$$\begin{aligned}\alpha &= 1 \\ \gamma &= 2 \\ \beta &= 0.5 \\ \delta &= 0.5\end{aligned}$$

... es sind aber auch andere, beliebige Belegungen möglich, solange folgende Invarianten gelten: $0 < \alpha$, $0 < \beta < 1$, $1 < \gamma$ und $\alpha < \gamma$ sowie $0 < \delta < 1$. [2]

Ordering Zu Beginn besitzen wir $N + 1$ Punkte, die anhand irgendeiner Kriterien (durchaus auch zufällig) gewählt wurden. Diese müssen wir nun anhand ihrer Güte ordnen, also sortieren. Im Falle des zweidimensionalen Raumes erhalten wir also die drei mnemonisch benannten Punkte B (*best*), G (*good*) und W (*worst*). [3] B ist im Falle eine Minimierung der Punkt, dessen Funktionswert am geringsten ist, G der mit dem zweit geringsten Wert, und so weiter.

Centroid Als nächstes berechnen wir den Mittelpunkt der beiden besseren Punkte

$$M = \frac{B + G}{2}$$

und nehmen diesen als Basis für unsere weiteren Berechnungen.

Transformation Nun reflektieren wir den schlechtesten Punkt W am Mittelpunkt

$$R = M + \alpha(M - W)$$

und prüfen, ob dieser Punkt R die Relation $f(B) \leq f(R) \leq f(W)$ erfüllt.

Wenn ja, ersetzen wir den schlechtesten Punkt mit R und brechen die jetzige Iteration ab. Diesen Vorgang nennt man auch das „Akzeptieren“ eines Punktes.

Wenn nein, prüfen wir, ob der reflektierte Punkt R besser als unser bisher bester Punkt ist, also ob $f(R) < f(B)$ gilt. Wenn dem so ist, expandieren wir den Punkt R testweise weiter:

$$E = M + \gamma(M - W)$$

und prüfen wiederum, ob dieser Punkt E besser als unser reflektierte Punkt R ist. Wenn ja, akzeptieren wir E . Wenn nein, akzeptieren wir R .

Gilt $f(R) \geq f(G)$, so ist der reflektierte Punkt kein guter Kandidat. Wir müssen nun die zwei möglichen kontrahierten Punkte C_1, C_2 berechnen. [2], [3] Gilt $f(R) < f(W)$, berechnen wir den äußeren kontrahierten Punkt C_1 , ansonsten den inneren kontrahierten Punkt C_2 .

$$\begin{aligned} C_1 &= M + \beta(R - M) \\ C_2 &= M + \beta(W - M) \end{aligned}$$

Sind die berechneten kontrahierten Punkte besser als die zur Berechnung verwendeten Punkte (R und W respektive), so werden diese akzeptiert. Ist dem nicht so, müssen wir zum nächsten Schritt, der Komprimierung, fortfahren. Bei der Komprimierung werden lediglich die zwei schlechteren Punkte G und W in Richtung des besten Punktes B gezogen. [2] Dieser Schritt wird bei den meisten Funktionen nie feuern — die Kontraktion liefert zumeist einen Punkt, der akzeptiert werden kann. In seltenen Fällen, nämlich wenn sich der Simplex über ein Tal spannt und sich ein Punkt des Simplex viel weiter vom Tal entfernt befindet als die anderen zwei besseren, ist dieser Schritt aber durchaus notwendig. Bei der Reflektion würde sich der reflektierte Punkt weiter vom Tal weg bewegen statt in das Tal hinein. Die Komprimierung wirkt dem entgegen und zieht alle Punkte in Richtung des besten Punktes, der in diesem Falle dem Tal am nächsten wäre. Ohne Komprimierung würde der Simplex in diesem Fall möglicherweise in einem falschen Punkt konvergieren.

$$\begin{aligned} G &= B + \delta(B - G) \\ W &= B + \delta(B - W) \end{aligned}$$

Diese Schritte werden fortgeführt, bis eine Abbruchbedingung eintritt, wie etwa dass der Unterschied zwischen dem besten Funktionswert $f(B)$ und dem schlechtesten Funktionswert $f(W)$ oder der Umfang des Simplex' ein gegebenes ε unterschreitet.

1.3 Erläuterung der einzelnen Transformationsschritte

Folgend werden die einzelnen Transformationsschritte anhand eines Dreiecks als Simplex vorgeführt und erläutert. Für unsere Koeffizienten wählen wir die Werte $\alpha = 1, \gamma = 2, \beta = \frac{1}{2}, \delta = \frac{1}{2}$. Unser Dreieck $\triangle BGW$ ist definiert mit $B = (1, 3), G = (4, 6), W = (6, 2)$.

1.3.1 Reflektion

Wir reflektieren den schlechtesten Punkt an dem Mittelpunkt in die vermeintlich bessere Richtung.

$$R = M + M - W = 2M - W = \begin{pmatrix} 5 \\ 9 \end{pmatrix} - \begin{pmatrix} 6 \\ 2 \end{pmatrix} = \begin{pmatrix} -1 \\ 7 \end{pmatrix}$$

Dies entspricht dem „Heruntertaumeln“ des Simplex, welches die Optimierung in den ersten Iterationsschritten unter anderem antreibt.

1.3.2 Expansion

Nehmen wir an, dass $f(R)$ besser als unser bis jetzt bester Wert ist. Wir expandieren also weiter:

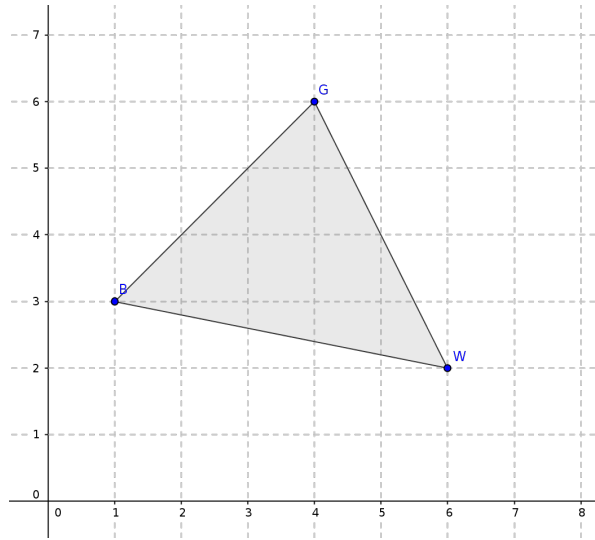


Abbildung 1: Unser Ursprungssimplex

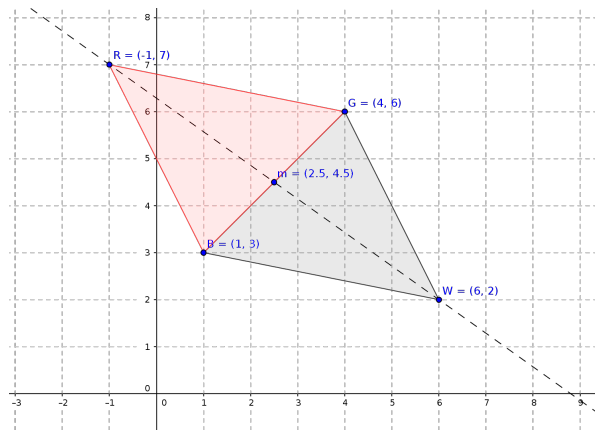


Abbildung 2: Nach der Reflektion

$$E = M + 2(M - W) = \begin{pmatrix} 2.5 \\ 4.5 \end{pmatrix} + \begin{pmatrix} -7 \\ 5 \end{pmatrix} = \begin{pmatrix} -4.5 \\ 9.5 \end{pmatrix}$$

Dies entspricht einem *greedy approach* — dieser Schritt ist nicht notwendig per se, minimiert aber zusätzlich die Anzahl der nötigen Iterationen. In schwer zu optimierenden Funktionen, wie etwa Rosenbrock's Bananenfunktion, kann der Expansionsfaktor γ dazu genutzt werden, den Simplex in flachem Gelände schneller wandern zu lassen.

1.3.3 Kontraktion

Gehen wir jetzt vom Gegenteil aus — Punkt R war also kein geeigneter Punkt und Punkt E wurde gar nicht berechnet — so müssen wir die beiden Kontraktionspunkte C_1, C_2 berechnen. C_1 ist hierbei der *äußere* Kontraktionspunkt, also derjenige, der näher an R liegt; dieser Punkt wird gewählt, wenn der reflektierte Punkt besser als der schlechteste Punkt W ist, also eine Verbesserung darstellt. Ist der

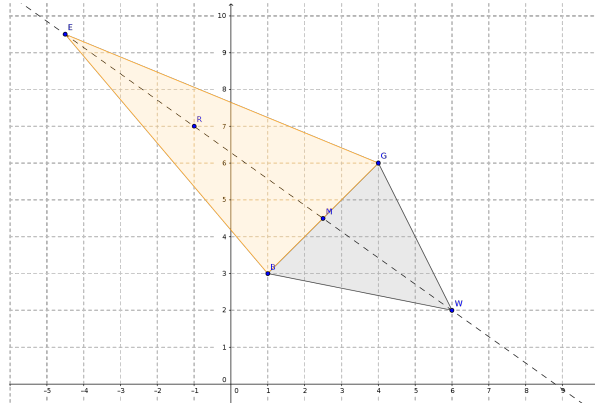


Abbildung 3: Nach der Expansion

reflektierte Punkt schlechter als der schlechteste Punkt, wird der *innere* Kontraktionspunkt gewählt, der näher an W liegt.

$$C_1 = M + \frac{R - M}{2} = \begin{pmatrix} 2.5 \\ 4.5 \end{pmatrix} + \frac{\begin{pmatrix} -1 \\ 7 \end{pmatrix} - \begin{pmatrix} 2.5 \\ 4.5 \end{pmatrix}}{2} = \begin{pmatrix} -1.75 \\ 1.25 \end{pmatrix}$$

$$C_2 = M + \frac{W - M}{2} = \begin{pmatrix} 2.5 \\ 4.5 \end{pmatrix} + \frac{\begin{pmatrix} 6 \\ 2 \end{pmatrix} - \begin{pmatrix} 2.5 \\ 4.5 \end{pmatrix}}{2} = \begin{pmatrix} 4.25 \\ 1.25 \end{pmatrix}$$

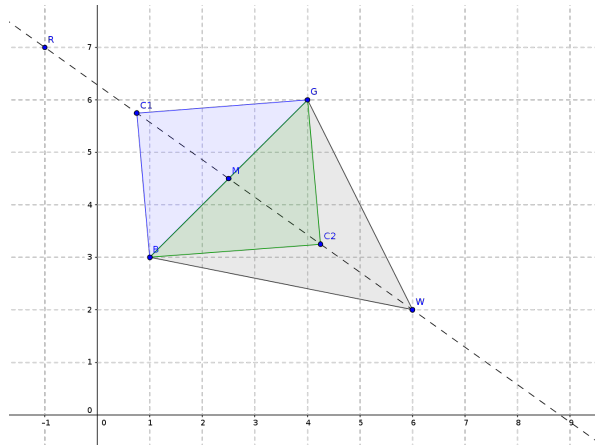


Abbildung 4: Nach der Kontraktion

1.3.4 Komprimierung

Wir komprimieren die Punkte G und W in Richtung B .

$$G' = B + \frac{G - B}{2} = \begin{pmatrix} 1 \\ 3 \end{pmatrix} + \frac{\begin{pmatrix} 4 \\ 6 \end{pmatrix} - \begin{pmatrix} 1 \\ 3 \end{pmatrix}}{2} = \begin{pmatrix} 2.5 \\ 4.5 \end{pmatrix}$$

$$W' = B + \frac{W - B}{2} = \begin{pmatrix} 1 \\ 3 \end{pmatrix} + \frac{\begin{pmatrix} 6 \\ 2 \end{pmatrix} - \begin{pmatrix} 1 \\ 3 \end{pmatrix}}{2} = \begin{pmatrix} 3.5 \\ 2.5 \end{pmatrix}$$

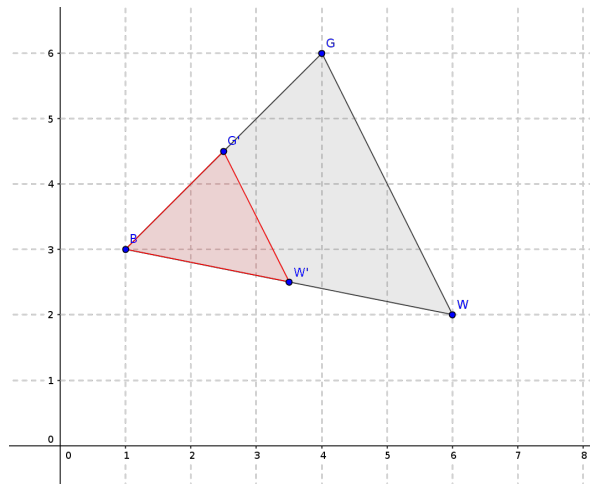


Abbildung 5: Nach der Komprimierung Richtung B

Dies entspricht dem Zusammenziehen des Simplex' in Richtung eines im Tal liegenden Punktes — lässt man diesen Schritt weg, schießt der Simplex in einigen Fällen über das Tal hinweg und konvergiert gegebenenfalls in einem falschen Punkt.

1.4 Ein Beispiel: die Himmelblau-Funktion

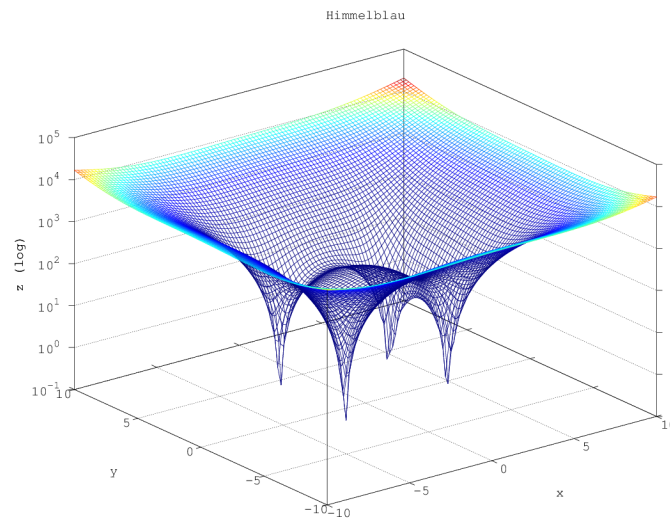


Abbildung 6: Die Himmelblau Funktion auf logarithmischer z -Skala

Um unseren Algorithmus exemplarisch vorzuführen, wollen wir die Himmelblau- Funktion minimieren. Dazu verwenden wir die beiliegende C++-Implementierung, um die hier aufgelisteten Abbildungen zu generieren.

Die Himmelblau-Funktion ist definiert als

$$f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$$

und hat 4 gleiche lokale Minima an den Stellen

$$\begin{aligned} f(3.2, 2.0) &= 0.0 \\ f(-2.805118, 3.131312) &= 0.0 \\ f(-3.779310, -3.283186) &= 0.0 \\ f(3.584428, -1.848126) &= 0.0 \end{aligned}$$

Wir wollen unseren Algorithmus auf sie loslassen, um eines davon zu finden. Als Abbruchbedingung wählen wir

$$\Delta f_{W,B} \leq 0.1,$$

also dass der Unterschied des besten und schlechtesten Funktionswert kleiner gleich 0.1 ist. Wir wählen diesen so groß, weil wir lieber schnell konvergieren als ein möglichst genaues Resultat wollen.

Unsere Startpunkte sind folgende:

$$\begin{aligned} P_1 &= (-1, -5) \\ P_2 &= (3, -8) \\ P_3 &= (8, 8) \end{aligned}$$

Wie bereits besprochen, sind die Startpunkte relativ egal und wurden auch hier recht zufällig gewählt. Die einzige Einschränkung des Startsimplex ist, dass es sich hierbei auch um einen tatsächlichen Simplex handeln muss (ein degenerierter Simplex wäre nutzlos). Weiters ist ein großer Simplex wünschenswert, da gegebenenfalls zu Beginn einer Optimierung große Schritte in Richtung Minimum gemacht werden müssen, um schnell konvergieren zu können. Ein kleiner Simplex benötigt weit mehr Zeit um zu konvergieren — dieses Problem kann man beispielsweise bei Rosenbrock's Bananenfunktion beobachten, die in der beiliegenden Implementierung ebenfalls vorhanden ist. Für die Parameter α , β , γ und δ wurde die Standardbelegung aus Unterabschnitt 1.2 Funktionsweise gewählt. (Abbildung 7)

Der erste Schritt ist das Ordnen der Punkte nach ihrem Funktionswert. Wir berechnen:

$$\begin{aligned} f(P_1) &= 512 \\ f(P_2) &= 3700 \\ f(P_3) &= 7964 \end{aligned}$$

Netterweise sind unsere Punkte also schon geordnet und P_1 , P_2 und P_3 entsprechen jeweils B , G , und W . Wir fahren fort und berechnen den Mittelpunkt

$$M = \frac{B + G}{2} = \begin{pmatrix} 1 \\ -6.5 \end{pmatrix}.$$

Als nächstes berechnen wir den reflektierten Punkt

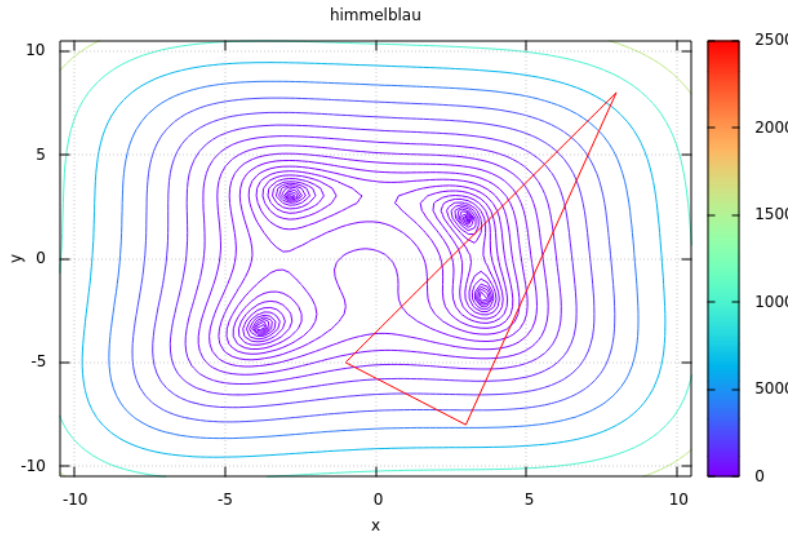


Abbildung 7: Unser Anfangssimplex aufgetragen auf einen Contourplot der Himmelblau-Funktion.

$$R = 2M - W = \begin{pmatrix} -6 \\ -21 \end{pmatrix}.$$

Berechnen wir nun $f(R) = 183200$, sehen wir, dass der reflektierte Punkt keine gute Idee ist. Da R bei weitem schlechter als unser schlechtester Punkt W ist, berechnen wir den inneren kontrahierten Punkt

$$C_2 = M + \frac{W - M}{2} = \begin{pmatrix} 4.5 \\ 0.75 \end{pmatrix}.$$

Berechnen wir nun $f(C_2) = 103.75391$, sehen wir, dass wir einen guten Fortschritt gemacht haben. Wir ersetzen W und ordnen unsere Punkte neu. (Abbildung 8)

$$\begin{aligned} B &= (4.5, 0.75) \\ G &= (-1, -5) \\ W &= (3, -8) \end{aligned}$$

Nun berechnen wir den neuen Mittelpunkt und reflektierten Punkt

$$\begin{aligned} M &= \frac{B + G}{2} = \begin{pmatrix} 1.75 \\ -2.125 \end{pmatrix} \\ R &= 2M - W = \begin{pmatrix} 0.5 \\ 3.75 \end{pmatrix} \end{aligned}$$

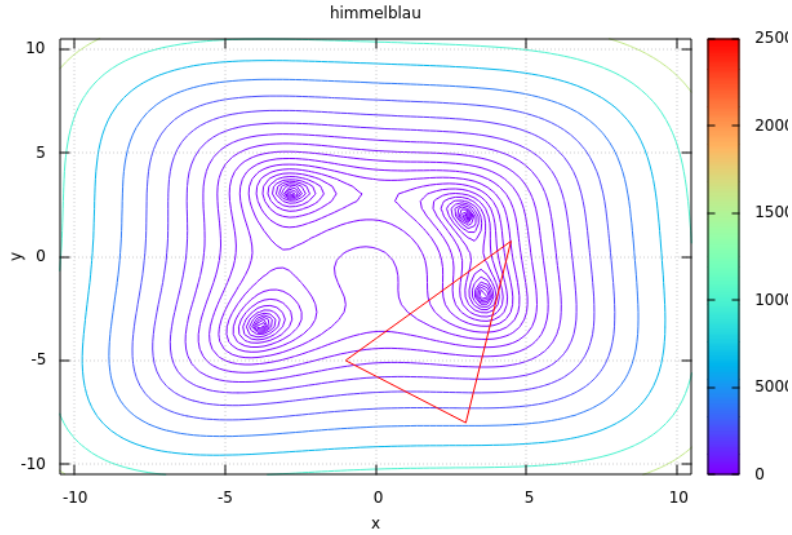


Abbildung 8: Unser Simplex nach dem ersten Transformationsschritt

Nun berechnen wir $f(R) = 106.19141$ und sehen, dass die Relation $f(B) \leq f(R) \leq f(W)$ damit erfüllt ist. Wir akzeptieren also den reflektierten Punkt und ordnen unsere Werte neu. (Abbildung 9)

$$\begin{aligned} B &= (4.5, 0.75) \\ G &= (0.5, 3.75) \\ W &= (-1, -5) \end{aligned}$$

Wir berechnen den neuen Mittelpunkt $M = \begin{pmatrix} 2.5 \\ 2.25 \end{pmatrix}$ und den neuen reflektierten Punkt $R = \begin{pmatrix} 6 \\ 9.5 \end{pmatrix}$. $f(R) = 9155.815$, wir erreichen also keine Verbesserung. Wir berechnen wiederum den inneren kontrahierten Punkt $C_2 = \begin{pmatrix} 0.75 \\ -1.375 \end{pmatrix}$ und sehen, dass $f(C_2) = 158.53931$, also eine Verbesserung gegenüber $f(W) = 514$. Wir akzeptieren also C_2 und ordnen unsere Punkte neu. (Abbildung 10)

$$\begin{aligned} B &= (4.5, 0.75) \\ G &= (0.5, 3.75) \\ W &= (0.75, -1.375) \end{aligned}$$

Die beiden besten Punkte haben sich nicht geändert, also bleibt auch der Mittelpunkt gleich. Wir berechnen den reflektierten Punkt als $R = \begin{pmatrix} 4.25 \\ 5.875 \end{pmatrix}$. Da $f(R) = 1176.43384$ haben wir keine Verbesserung erreicht und sind sogar schlechter als der schlechteste Punkt unseres Simplex', demnach berechnen wir wiederum den inneren kontrahierten Punkt $C_2 = \begin{pmatrix} 1.625 \\ 0.4375 \end{pmatrix}$. Der Funktionswert an der Stelle C_2 ist 89.62574 und damit besser als B . Wir akzeptieren diesen Punkt daher und ordnen unsere Punkte neu. (Abbildung 11)

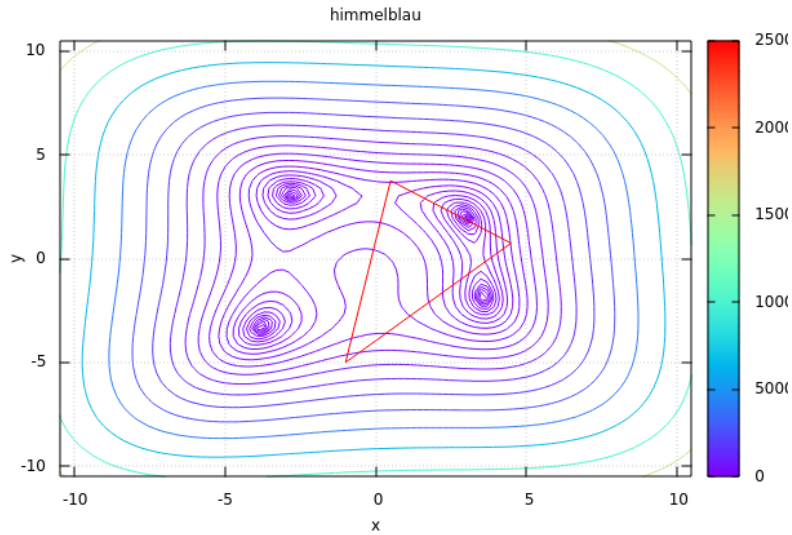


Abbildung 9: Unser Simplex nach dem zweiten Transformationsschritt

$$\begin{aligned} B &= (1.625, 0.4375) \\ G &= (4.5, 0.75) \\ W &= (0.5, 3.75) \end{aligned}$$

Es ist mittlerweile klar, dass wir im Punkt $\binom{3}{2}$ konvergieren werden, da zwei Punkte ihm am nächsten sind. Das beliegende C++-Programm konvergierte mit den angegebenen Parametern in 23 Iterationsschritten. Die weiteren Schritte werden hier mit dem Verweis auf die Implementierung und die .gif-Datei auf der Projektwebseite ausgelassen.

1.5 Analyse

Die formale mathematische Analyse des Nelder-Mead-Algorithmus' ist nicht möglich. Es kann nicht eindeutig bewiesen werden, dass der Algorithmus beim Bearbeiten einer gegebenen Funktion f_x konvergiert. Die Annahme, dass der Simplex konvergiert, beruht auf zwei Gegebenheiten:

- Der Simplex degeneriert nicht während der Laufzeit.
- Irgendeine Art eines Abstiegskriteriums wird verwendet.

Beides ist *nicht* gegeben — der Simplex kann und wird in einigen Fällen beinahe degenerieren, sofern die Landschaft der Funktion dies zulässt.¹ Weiters ist das Abstiegskriterium mehr zufällig — es ist nicht mathematisch fundiert wie etwa die Verwendung des negativen Gradienten, der mathematisch beweisbar immer in die Richtung des steilsten Abstiegs zeigt. [2]

Nichtsdestotrotz ist der Nelder-Mead-Algorithmus ein guter ableitungsfreier Optimierungsalgorithmus, der zumeist zuverlässig korrekte Resultate liefert. Tatsächlich gibt es aber eine ganze Familie von

¹Man betrachte Rosenbrock's Bananenfunktion, in deren Tal der Simplex rapide schrumpft und sich somit nur mehr sehr langsam bewegen kann.

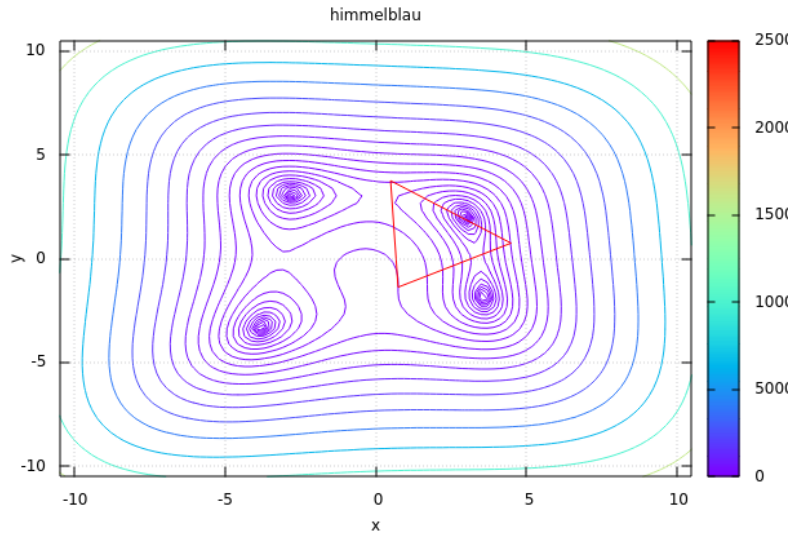


Abbildung 10: Nach dem dritten Transformationsschritt

Funktionen, für die das Nelder-Mead-Verfahren immer zu einem Punkt konvergiert, der nicht das Minimum ist. [4]

In vielen naturwissenschaftlichen Bereichen sind Funktionswerte fehlerbehaftet und damit genaue Lösungen unnötig. In diesen Fällen ist der Nelder-Mead-Algorithmus von großem Vorteil, da er schnell Ergebnisse liefert die oft bereits „gut genug“ sind. Dazu kommt selbstverständlich die Tatsache, dass der vorliegende Algorithmus anschaulich und relativ leicht verständlich ist, und viele natürlich mehr inkliniert sind, Methoden zu verwenden, die sie auch verstehen.

2 Das Gradientenverfahren

2.1 Die Idee und der Zweck des Gradientenverfahren

Grundsätzlich ist das Gradientenverfahren ein Optimierungsverfahren, welches das Minimum einer mehrdimensionalen Funktion errechnen. Der Gradient ($\nabla f(x, y)$) einer Funktion $f(x)$ zeigt, aus später gezeigten ² Gründen, immer in die Richtung des steilsten Anstiegs, folglich zeigt $-\nabla f(x, y)$ in die Richtung des stärksten Abstiegs. Das Gradientenverfahren macht sich dies zunutze und man „wandert“ eine gewisse Dauer in diese Richtung, wie "lange" man in die errechnete Richtung "wandert" hängt unter anderem auch von der gewählten Methode ab.

2.2 Der Ablauf des Gradientenverfahrens

1. Es wird ein beliebiger Startpunkt $P = (x, y)$ ausgewählt.
2. Es wird die Ableitung der Funktion $f(x, y)$, also $\nabla f(x)$ gebildet. Sollte dieser Gradient bereits den Wert null haben sind wir bereits fertig und haben ein lokales Minimum.

²HIER VERWEIS EINFÜGEN

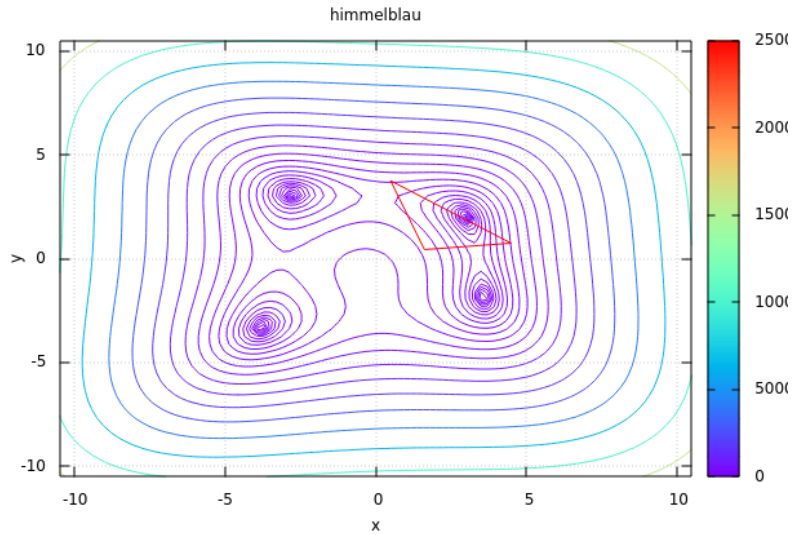


Abbildung 11: Nach dem vierten Transformationsschritt

3. Der Punkt P wird in den Gradienten eingesetzt und der negative Wert berechnet, also $-\nabla f(x_P, y_P)$.
4. Um einen neuen Punkt $P^{[n+1]}$ zu finden stellen wir nun folgende Gleichung auf: $P^{[n+1]} = P^{[n]} + \alpha^{[n]} * s^{[n]}$. Nun haben wir die Richtung in die optimiert werden muss gefunden ($s^{[n]} (= -\nabla f(x_P, y_P))$), jetzt muss noch entschieden werden wie weit wir in diese Richtung "wandern" ($\alpha^{[n]}$, liegt zwischen 0 und 1) (Hinweis: n ist die aktuelle Iteration).

Das kann man einerseits mit Hilfe eines fixen Werts machen, andererseits auch mit einer Folge (wie z.B. $\frac{1}{\sqrt{n}}$) oder einer jedes mal neu errechneten Schrittweite. Zu beachten ist, dass die ausgewählte Folge auf jeden Fall divergent sein muss, da man sonst bis zu einer maximalen Entfernung "wandern" kann und die Wahrscheinlichkeit groß ist das Optimum nie zu erreichen. Sollte man jedes mal einen fixen Wert verliert man zwar den Aufwand für die Berechnung der Schrittweite, jedoch ist die Wahrscheinlichkeit hoch, dass man viele Iterationen benötigt oder in einer endlosen Schleife festhängt.

Es gibt verschiedene Arten wie man die Schrittweite berechnen kann, im Folgenden werde ich die Schrittweitenbestimmung nach Armijo erklären. Dieses Verfahren ist eines der sogenannten "line-search" Verfahren, das von mir gezeigt ist jedoch kein exaktes line-search Verfahren sondern lediglich eine einfache Heuristik.

Die Armijo-Bedingung: $\varphi(\alpha) \leq \varphi(0) + c * \alpha^{[n]} * \varphi(0)'$.

Wobei c eine Konstante zwischen null und eins ist und dazu dient das die Bedingung nicht zu restriktiv ist, $\varphi(\alpha) = f(P^{[n]} + \alpha^{[n]} s^{[n]})$ und $\varphi(0)' = \nabla f(P^n) * s^{[n]}$

5. Für den Anfang wird $\alpha = 1$ gewählt und errechnet ob die Armijo-Bedingung erfüllt ist - sollte sie erfüllt sein haben wir ein passende α gefunden, wenn nicht wird α mit einem β , welches ebenfalls zwischen null und eins liegt, multipliziert und wieder geprüft ob die Ungleichung erfüllt

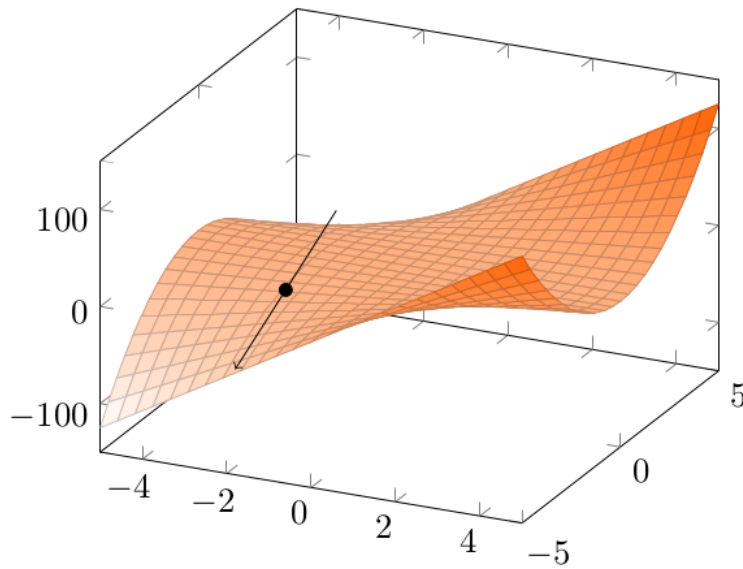


Abbildung 12: Hier zu sehen ist $-\vec{\nabla}f(-2, -2)$ der Funktion $f(x, y) = x * y^2$

ist. (Hinweis: In der Praxis liefern die Werte $c = 0.01$ und $\beta = 0.9$ häufig gute Ergebnisse. ³

6. Wenn man einen neuen Punkt gefunden hat wird bei Punkt 3 fortgesetzt. Da es oft einen großen Rechenaufwand erfordert das genaue Minimum zu finden wird das Verfahren meist beendet wenn der Gradient einen akzeptablen Wert erreicht hat (So kann man als Abbruchbedingung z.B. die Stärke des Anstiegs in einem Punkt (also den Betrag des Gradienten) nehmen).

2.3 Beispiele

2.3.1 Beispiel 1

Zu minimierende Funktion: $f(x, y) = y^4 + 2x^2 - 3xy + 1$, als Abbruchbedingung wählen wir $|\nabla f(x, y)| \leq 0.3$.

Wie man in Abbildung 13 sieht, nähert man sich schrittweise dem Minimum an.

1. Startpunkt $P = (1, 1)$
2. Ableitung: $\nabla f(x, y) = \begin{pmatrix} 4x - 3y \\ 4y^3 - 3x \end{pmatrix}$ (\rightarrow Der Gradient)
3. $-\nabla f(1, 1) = \begin{pmatrix} -4 + 3 \\ -4 + 3 \end{pmatrix} = \begin{pmatrix} -1 \\ -1 \end{pmatrix}$ (\rightarrow Richtung des stärksten Abstiegs)
4. Die Armijo-Bedingung (mit $\alpha = 1$, $c = 0.1$ und $\beta = 0.9$):

$$f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + 1 * \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) \leq f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 0.1 * 1 * \begin{pmatrix} 1 \\ 1 \end{pmatrix} * \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

$$= 1 \leq 0.8$$

$$\rightarrow \alpha^{[2]} = \alpha^{[1]} * \beta = 0.9$$

³<https://www.unibw.de/lrt1/gerdts/lehre/lpnlp/lp-nlp.pdf>

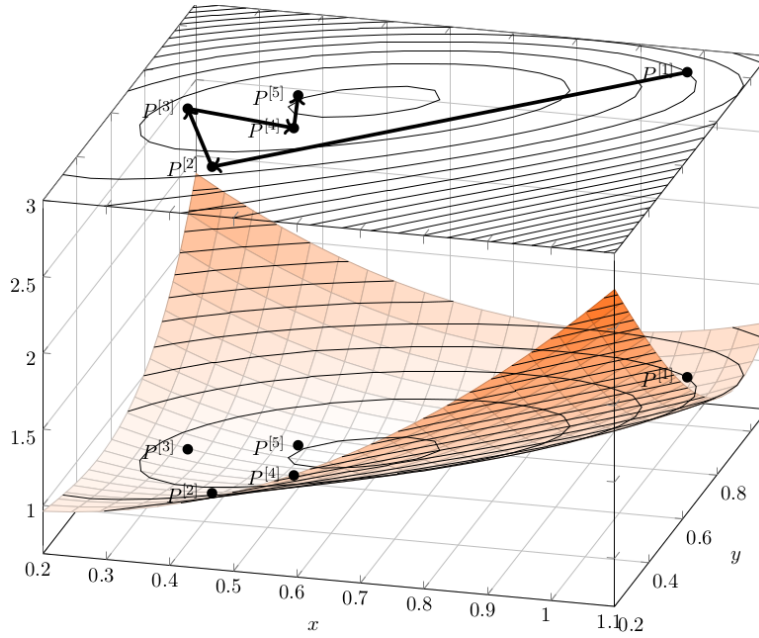


Abbildung 13: Beispiel 1. Zusätzlich sind hier die Höhenschichtlinien eingezeichnet.

$$\begin{aligned}
 f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0.9 * \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) &\leq f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 0.1 * 0.9 * \begin{pmatrix} 1 \\ 1 \end{pmatrix} * \begin{pmatrix} -1 \\ -1 \end{pmatrix} \\
 &= 0.99 \leq 0.82 \\
 \rightarrow \alpha^{[3]} &= \alpha^{[2]} * \beta = 0.81
 \end{aligned}$$

$$\begin{aligned}
 f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0.81 * \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) &\leq f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 0.1 * 0.81 * \begin{pmatrix} 1 \\ 1 \end{pmatrix} * \begin{pmatrix} -1 \\ -1 \end{pmatrix} \\
 &= 0.965 \leq 0.838 \\
 \rightarrow \alpha^{[4]} &= \alpha^{[3]} * \beta = 0.729
 \end{aligned}$$

$$\begin{aligned}
 f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0.729 * \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) &\leq f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 0.1 * 0.729 * \begin{pmatrix} 1 \\ 1 \end{pmatrix} * \begin{pmatrix} -1 \\ -1 \end{pmatrix} \\
 &= 0.932 \leq 0.854 \\
 \rightarrow \alpha^{[5]} &= \alpha^{[4]} * \beta = 0.656
 \end{aligned}$$

$$\begin{aligned}
 f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0.656 * \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) &\leq f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 0.1 * 0.656 * \begin{pmatrix} 1 \\ 1 \end{pmatrix} * \begin{pmatrix} -1 \\ -1 \end{pmatrix} \\
 &= 0.896 \leq 0.869 \\
 \rightarrow \alpha^{[6]} &= \alpha^{[5]} * \beta = 0,59
 \end{aligned}$$

$$\begin{aligned}
 f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0,59 * \begin{pmatrix} -1 \\ -1 \end{pmatrix}\right) &\leq f\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}\right) + 0.1 * 0,59 * \begin{pmatrix} 1 \\ 1 \end{pmatrix} * \begin{pmatrix} -1 \\ -1 \end{pmatrix} \\
 &= 0.86 \leq 0.882
 \end{aligned}$$

Nun haben wir ein passendes α gefunden und können zum nächsten Punkt wandern.

$$5. P^{[2]} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} + 0,59 * \begin{pmatrix} -1 \\ -1 \end{pmatrix} = \begin{pmatrix} 0.41 \\ 0.41 \end{pmatrix}$$

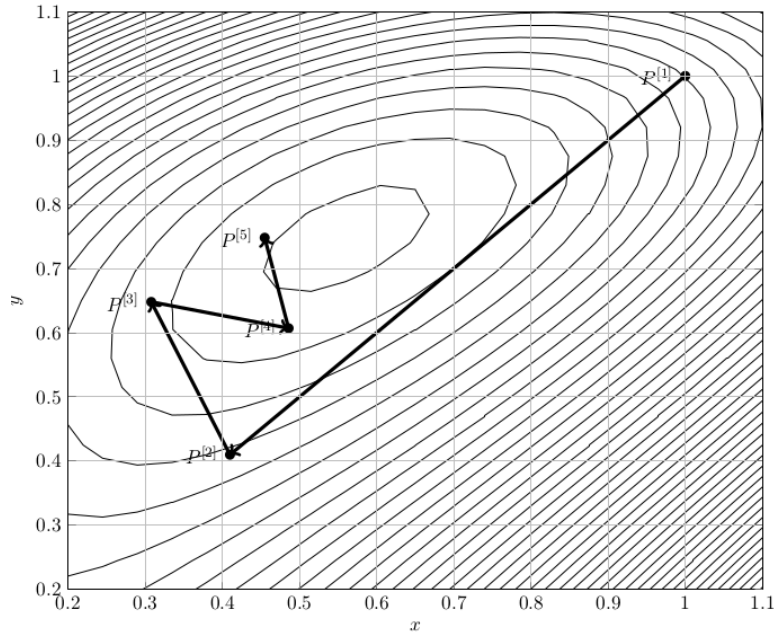


Abbildung 14: Die Höhenschichtlinien genauer gezeigt.

$$6. -\nabla f(0.41, 0.41) = \begin{pmatrix} -0.41 \\ 0.95 \end{pmatrix},$$

$$|\nabla f(0.41, 0.41)| = \sqrt{0.41^2 + 0.95^2} = 1.071$$

7. Die Armijo-Bedingung (mit $\alpha = 1$, $c = 0.1$ und $\beta = 0.5$):

$$f\left(\begin{pmatrix} 0.41 \\ 0.41 \end{pmatrix}\right) + 1 * \begin{pmatrix} -0.41 \\ 0.95 \end{pmatrix} \leq f\left(\begin{pmatrix} 0.41 \\ 0.41 \end{pmatrix}\right) + 0.1 * 1 * \begin{pmatrix} -0.41 \\ 0.95 \end{pmatrix} * \begin{pmatrix} 0.41 \\ -0.95 \end{pmatrix}$$

$$= 4.421 \leq 0.753 \rightarrow \alpha^{[2]} = \alpha^{[1]} * \beta = 0.5$$

$$\rightarrow 1.153 \leq 0.807 \rightarrow \alpha^{[3]} = \alpha^{[2]} * \beta = 0.25$$

$$\rightarrow 0.768 \leq 0.833$$

$$8. P^{[3]} = \begin{pmatrix} 0.41 \\ 0.41 \end{pmatrix} + 0.25 * \begin{pmatrix} -0.41 \\ 0.95 \end{pmatrix} = \begin{pmatrix} 0.308 \\ 0.648 \end{pmatrix}$$

$$9. -\nabla f(0.308, 0.648) = \begin{pmatrix} 0.712 \\ -0.164 \end{pmatrix},$$

$$|\nabla f(0.308, 0.648)| = \sqrt{0.712^2 + 0.164^2} = 0.534$$

10. Die Armijo-Bedingung (mit $\alpha = 1$, $c = 0.1$ und $\beta = 0.5$):

$$f\left(\begin{pmatrix} 0.308 \\ 0.648 \end{pmatrix}\right) + 1 * \begin{pmatrix} 0.712 \\ -0.164 \end{pmatrix} \leq f\left(\begin{pmatrix} 0.308 \\ 0.648 \end{pmatrix}\right) + 0.1 * 1 * \begin{pmatrix} 0.712 \\ -0.164 \end{pmatrix} * \begin{pmatrix} -0.712 \\ 0.164 \end{pmatrix}$$

$$= 1.655 \leq 0.714 \rightarrow \alpha^{[2]} = \alpha^{[1]} * \beta = 0.5$$

$$\rightarrow 0.857 \leq 0.741 \rightarrow \alpha^{[3]} = \alpha^{[2]} * \beta = 0.25$$

$$\rightarrow 0.723 \leq 0.754$$

$$11. P^{[4]} = \begin{pmatrix} 0.308 \\ 0.648 \end{pmatrix} + 0.25 * \begin{pmatrix} 0.712 \\ -0.164 \end{pmatrix} = \begin{pmatrix} 0.486 \\ 0.607 \end{pmatrix}$$

$$12. -\nabla f(0.486, 0.607) = \begin{pmatrix} -0.123 \\ 0.563 \end{pmatrix},$$

$$|\nabla f(0.486, 0.607)| = \sqrt{0.123^2 + 0.563^2} = 0.332$$

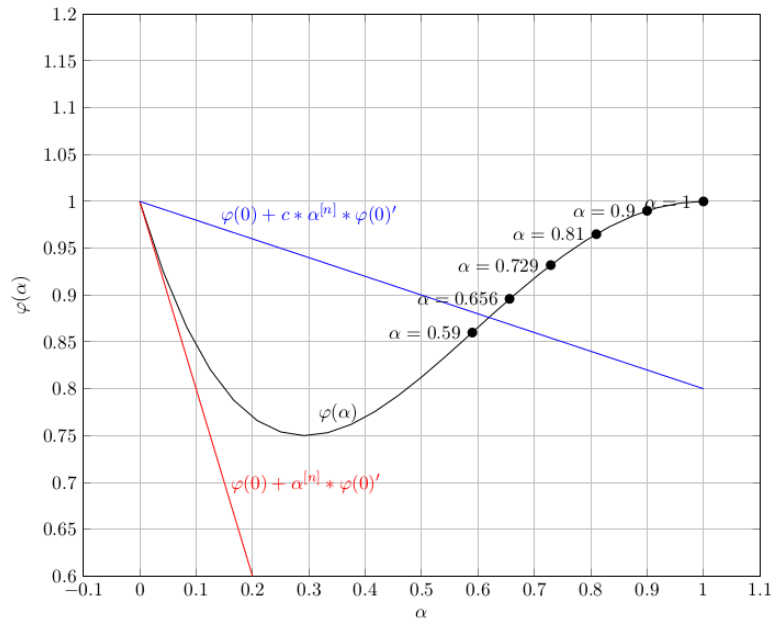


Abbildung 15: Hier sieht man die Funktion $\varphi(\alpha) = (1 - \alpha)^4 + 2 * (1 - \alpha)^2 - 3 * (1 - \alpha) * (1 - \alpha) + 1$, sie ist sozusagen ein Durchschnitt durch die bei Beispiel 1 gegebene Funktion. Diese Funktion wird erzeugt, wenn man vom Punkt $P^{[1]} = (1, 1)$ in Richtung des Gradienten wandert. Die blaue Gerade zeigt die Armijo-Bedingung bei $c = 0.1$, die rote bei $c = 1$ bzw. ohne c .

13. Die Armijo-Bedingung (mit $\alpha = 1$, $c = 0.1$ und $\beta = 0.5$):

$$\begin{aligned} f\left(\begin{pmatrix} 0.486 \\ 0.607 \end{pmatrix}\right) + 1 * \begin{pmatrix} -0.123 \\ 0.563 \end{pmatrix} &\leq f\left(\begin{pmatrix} 0.486 \\ 0.607 \end{pmatrix}\right) + 0.1 * 1 * \begin{pmatrix} -0.123 \\ 0.563 \end{pmatrix} * \begin{pmatrix} 0.123 \\ -0.563 \end{pmatrix} \\ &= 1.863 \leq 0.689 \rightarrow \alpha^{[2]} = \alpha^{[1]} * \beta = 0.5 \\ &\rightarrow 0.852 \leq 0.707 \rightarrow \alpha^{[3]} = \alpha^{[2]} * \beta = 0.25 \\ &\rightarrow 0.706 \leq 0.715 \end{aligned}$$

$$14. P^{[5]} = \begin{pmatrix} 0.486 \\ 0.607 \end{pmatrix} + 0.25 * \begin{pmatrix} -0.123 \\ 0.563 \end{pmatrix} = \begin{pmatrix} 0.455 \\ 0.748 \end{pmatrix}$$

$$15. -\nabla f(0.455, 0.748) = \begin{pmatrix} 0.424 \\ -0.309 \end{pmatrix},$$

$$|\nabla f(0.455, 0.748)| = \sqrt{0.424^2 + 0.309^2} = 0.275$$

Die Abbruchbedingung ist erfüllt und somit ist ein Punkt, welcher nahe genug am Optimum liegt gefunden.

3 Die Koordinatenabstiegsmethode

3.1 Einführung

Bei der Koordinatenabstiegsmethode wird bei jedem Schritt eine Koordinate i gewählt, nach der man optimiert. Es wird bei einem beliebigen Startpunkt begonnen und in jedem Iterationsschritt setzt man alle Koordinaten des aktuellen Punktes in die Funktion ein, nur die Koordinate, nach der man minimieren will, behält man als Variable bei.

Es wird immer in Richtung einer Achse minimiert und die Stelle, an der das Minimum ist, als neue Koordinate beibehalten und so gelangt man mit einer Anzahl an Iterationen zum Minimum der Funktion.

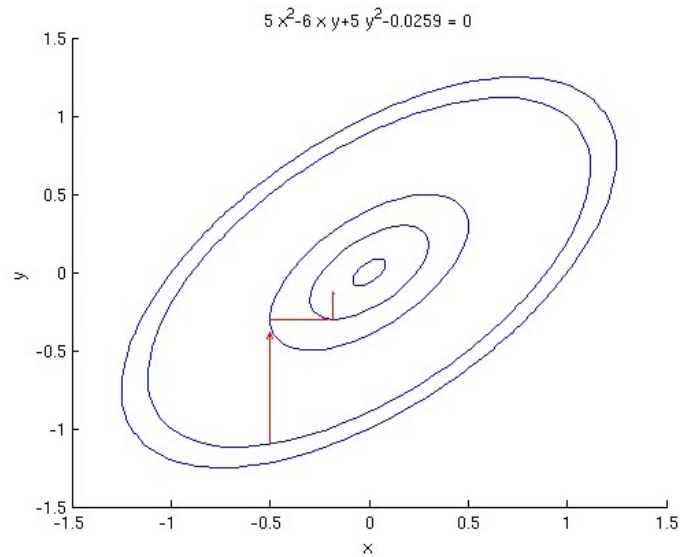


Abbildung 16: Die Koordinatenabstiegsmethode

Wenn zum Beispiel die Funktion $f(x, y) = 2x^2 + 2xy + 1.5y^2$ und der Startpunkt $P = (1, 1)$ gegeben ist, so muss man, wenn entlang der x -Achse minimiert werden soll folgendes eingesetzt werden:

$$\begin{aligned} x &= 1 + 1d \\ y &= 1 + 0d \end{aligned}$$

Soll entlang der y -Achse optimiert werden, so muss folgendes eingesetzt werden:

$$\begin{aligned} x &= 1 + 0d \\ y &= 1 + 1d \end{aligned}$$

3.2 Beispiel

Wir möchten die Funktion $f(x, y) = 2x^2 + 2xy + 1.5y^2$ optimieren und wählen den Startpunkt $P_0 = (1, 1)$. Wir minimieren zuerst entlang der x -Achse und setzen ein

$$\begin{aligned} f(d) &= 2(1 + d)^2 + 2(1 + d) + 1.5 \\ &= 2d^2 + 6d + 5.5 \end{aligned}$$

Diese soeben erhaltene Funktion müssen wir nun noch minimieren.

$$\begin{aligned} f'(d) &= 4d + 6 \\ 4d + 6 &= 0 \\ d &= -\frac{3}{2} \end{aligned}$$

Setzen wir dies in die Gleichung für die x -Koordinate ein und verwenden die alte y -Koordinate erhalten wir den neuen Punkt $P_1 = (-\frac{1}{2}, 1)$. Nun minimieren wir nach y , wir setzen also ein:

$$\begin{aligned} x &= -\frac{1}{2} + 0d \\ y &= 1 + 1d \end{aligned}$$

und erhalten

$$f(d) = 1.5d^2 - d + 6$$

Dies minimieren wir wiederum und erhalten die neue y -Koordinate $\frac{1}{3}$. Damit haben wir den neuen Punkt $(-\frac{1}{2}, \frac{1}{3})$. Dieser Vorgang wird wiederholt, bis sich die errechneten Werte nurnoch marginal unterscheiden. Danach wird abgebrochen.

3.3 Problemfälle

Zeigen Sie, dass die Koordinatenabstiegsmethode für die stetige Funktion

$$f(x, y) = \begin{cases} (x + y - 5)^2 + (x - y - 2)^2 & \text{falls } x \leq y \\ (x + y - 5)^2 + (x - y + 2)^2 & \text{sonst} \end{cases}$$

beginnend beim Startpunkt $(0, 0)$ nicht funktioniert.

Wir setzen ein:

$$\begin{aligned} x &= d \\ y &= 0 \end{aligned}$$

und erhalten die Funktion

$$\begin{aligned} f(d) &= (d - 5)^2 + (d - 2)^2 \\ &= 2d^2 - 14d + 29 \end{aligned}$$

welche ein Minimum an der Stelle $d = \frac{7}{2}$ hat. Dies fortgesetzt für y erhalten wir den neuen Punkt $P = (\frac{7}{2}, \frac{7}{2})$. Bei weiteren Iterationen merken wir, dass d immer 0 beträgt – was bedeutet, dass sich der Punkt nicht weiterbewegt. Wir sitzen auf einem lokalen Minimum fest. Führen wir das mit dem anderen Zweig der Funktion ebenfalls durch, so sehen wir das gleiche Problem.

Die Koordinatenabstiegsmethode hat Probleme bei nicht-glatten* Funktionen. Man könnte bei einem Punkt hängen bleiben, der nicht das Minimum ist, wenn die Contourlinien nicht glatt⁴ sind. [5]

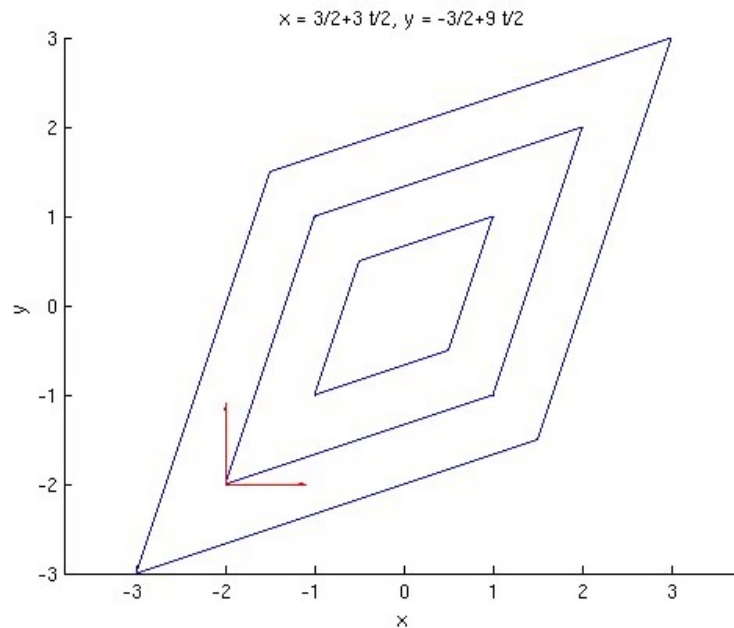


Abbildung 17: Ein Problemfall

Beispiel:

Man stelle sich vor, wir wären am Punkt $(-2, -2)$ angelangt. Die roten Linien zeigen, welche Richtungen man einschlagen kann. Würde man nun eine weitere Iteration des Verfahrens anwenden, würde man auf dem Punkt $(-2, -2)$ festsitzen, da das Wandern entlang der roten Linien den Funktionswert erhöhen würde und somit den derzeitigen Punkt als ideal angesehen wird.

Literatur

- [1] Wikipedia, *Nelder-Mead method* — *Wikipedia, the free encyclopedia*, [Online; accessed 10-December-2015], 2015. Adresse: https://en.wikipedia.org/w/index.php?title=Nelder%E2%80%93Mead_method&oldid=693590045.
- [2] S. Singer und J. Nelder, “Nelder-Mead algorithm”, *Scholarpedia*, Bd. 4, Nr. 7, S. 2928, 2009, revision 91557.
- [3] Unknown, *The Nelder-Mead method (proof)*, [Online excerpt of book; accessed 10-December-2015], 2015. Adresse: <http://mathfaculty.fullerton.edu/mathews/n2003/neldermead/NelderMeadProof.pdf>.
- [4] K. I. M. McKinnon, “Convergence of the nelder-mead simplex method to a non-stationary point”, *SIAM J. Optim, Techn. Ber.*, 1996.
- [5] Wikipedia, *Coordinate descent* — *wikipedia, the free encyclopedia*, [Online; accessed 8-January-2016], 2015. Adresse: https://en.wikipedia.org/w/index.php?title=Coordinate_descent&oldid=696403191.

⁴unendlich oft differenzierbar und stetig (Graph hat keine Ecken)