

14. Anwendungen der Differentialrechnung

14.1 Extremwertberechnung

Gegeben sei eine C^2 -Funktion $f : \mathbb{R}^n \supset D \rightarrow \mathbb{R}$, wobei D offen sei. Wir interessieren uns für die Bestimmung der Extrema (Maxima und Minima) von f .

Satz (14.1.1) (Notwendige Bedingungen)

Ist $\mathbf{x}^0 \in D$ ein lokales Extremum von f , so gelten:

- a) $\nabla f(\mathbf{x}^0) = 0$.
- b) Ist \mathbf{x}^0 ein lokales Minimum von f , so ist die Hesse-Matrix $\nabla^2 f(\mathbf{x}^0)$ positiv semidefinit; ist \mathbf{x}^0 ein lokales Maximum, so ist $\nabla^2 f(\mathbf{x}^0)$ negativ semidefinit.

Satz (14.1.2) (Hinreichende Bedingungen)

Für $\mathbf{x}^0 \in D$ gelten:

- a) Ist $\nabla f(\mathbf{x}^0) = \mathbf{0}$ und ist $\nabla^2 f(\mathbf{x}^0)$ positiv (negativ) definit, so ist \mathbf{x}^0 ein strenges lokales Minimum (Maximum) von f .
- b) Ist $\nabla f(\mathbf{x}^0) = \mathbf{0}$ und ist $\nabla^2 f(\mathbf{x}^0)$ indefinit, so ist \mathbf{x}^0 ein Sattelpunkt von f .

Bemerkungen (14.1.3)

a) Punkte $\mathbf{x}^0 \in D^0$ mit $\nabla f(\mathbf{x}^0) = \mathbf{0}$ heißen **stationäre Punkte** von f . Dies sind nicht zwingend lokale Extrema; in einem stationären Punkt sind durch die EVen zu positiven bzw. negativen EVen der Hesse-Matrix Richtungen gegeben, in denen f wächst bzw. fällt.

b) Ein stationärer Punkt mit $\det \nabla^2 f(\mathbf{x}^0) = 0$, d.h. $\lambda = 0$ ist ein EW von $\nabla^2 f(\mathbf{x}^0)$, heißt **ausgeartet**. Ist \mathbf{x}^0 nichtausgearteter stationärer Punkt, so trifft genau einer der drei in (14.1.2) genannten Fälle zu.

c) \mathbf{x}^0 lokales Minimum $\Leftrightarrow \mathbf{x}^0$ strenges lok. Minimum
 \Downarrow $\nabla^2 f(\mathbf{x}^0)$ pos. semidef. $\Leftrightarrow \nabla^2 f(\mathbf{x}^0)$ pos. definit \Uparrow

Man beachte, dass hierbei in keinem Fall die Umkehrung der Implikation gilt!

d) Ist f C^3 -Funktion, \mathbf{x}^0 stationärer Punkt mit $\nabla^2 f(\mathbf{x}^0)$ pos. definit, so gibt es eine Umgebung von \mathbf{x}^0 und ein $C > 0$, so dass dort gilt:

$$f(\mathbf{x}) - f(\mathbf{x}^0) \geq C \|\mathbf{x} - \mathbf{x}^0\|^2.$$

Beispiel (14.1.4)

Gegeben sei $z = f(x, y) := y^2 (x - 1) + x^2 (x + 1)$.

Wir berechnen: $\nabla f(x, y) = (y^2 + 3x^2 + 2x, 2y(x - 1))^T$.

Damit hat f die beiden stationären Punkte:

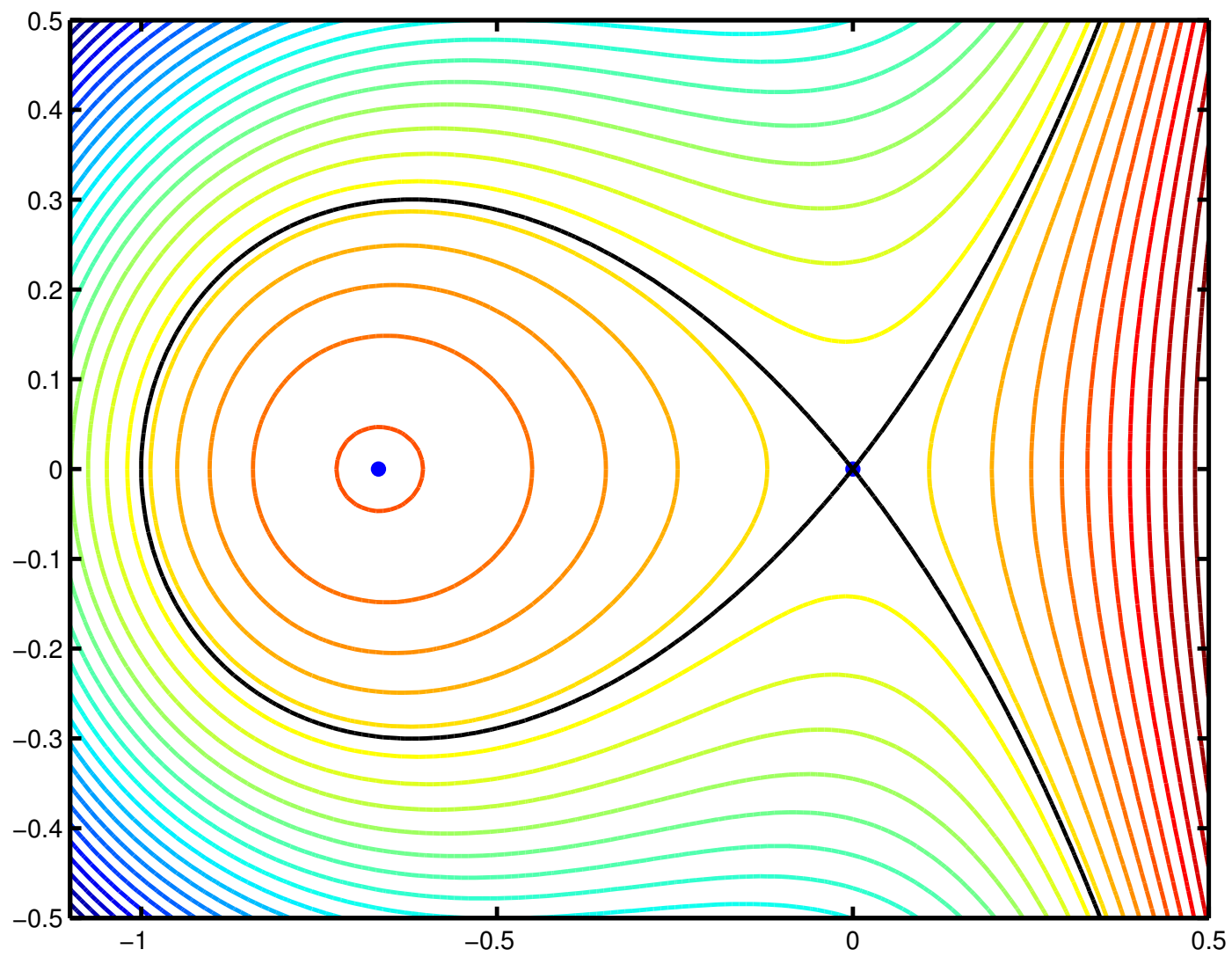
$$\mathbf{x}^0 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{x}^1 = \begin{pmatrix} -2/3 \\ 0 \end{pmatrix}$$

mit den Funktionswerten $f(\mathbf{x}^0) = 0$, $f(\mathbf{x}^1) = 4/27$.

Mit $f_{xx} = 6x + 2$, $f_{xy} = 2y$, $f_{yy} = 2(x - 1)$, ergibt sich:

$\nabla^2 f(\mathbf{x}^0) = \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}$ indefinit, \mathbf{x}^0 also Sattelpunkt,

$\nabla^2 f(\mathbf{x}^1) = \begin{pmatrix} -2 & 0 \\ 0 & -10/3 \end{pmatrix}$ neg. definit, \mathbf{x}^1 also strenges lokales Maximum.



Beispiel (14.1.5)

Gegeben sei $z = f(x, y) := yx^2(4 - x + y)$.

Wir berechnen: $\nabla f(x, y) = (xy(8 - 3x - 2y), x^2(4 - x - 2y))^T$.

Damit hat f die folgenden stationären Punkte:

$$\mathbf{x}^0 = \begin{pmatrix} 0 \\ y \end{pmatrix}, \quad \mathbf{x}^1 = \begin{pmatrix} 4 \\ 0 \end{pmatrix}, \quad \mathbf{x}^2 = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

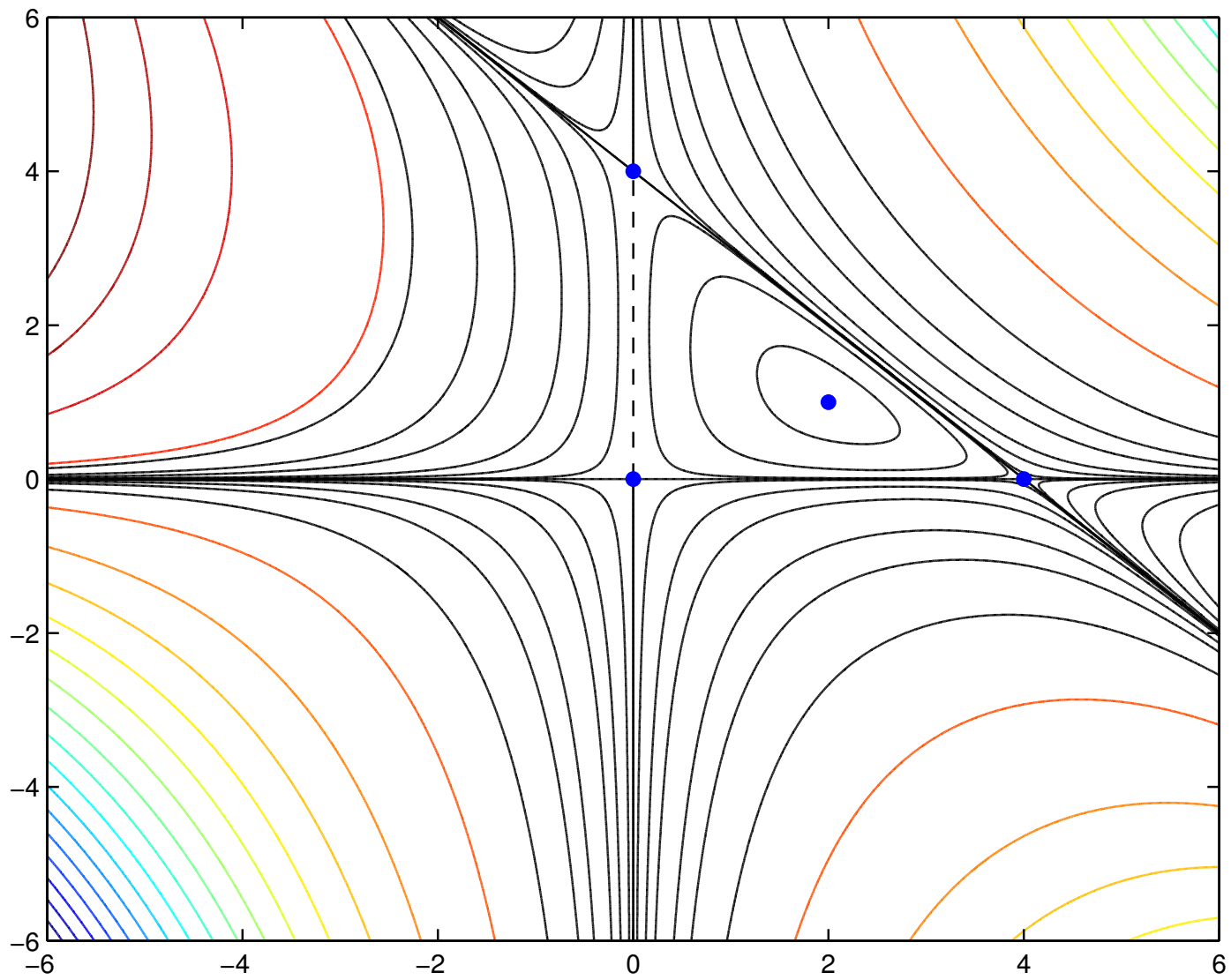
mit den Funktionswerten $f(\mathbf{x}^0) = f(\mathbf{x}^1) = 0$, $f(\mathbf{x}^2) = 4$.

Mit $f_{xx} = y(8 - 6x - 2y)$, $f_{xy} = x(8 - 3x - 4y)$, $f_{yy} = -2x^2$ ergibt sich: $\nabla^2 f(\mathbf{x}^0)$ singular, \mathbf{x}^0 ist also ausgeartet,

$\nabla^2 f(\mathbf{x}^1) = \begin{pmatrix} 0 & -16 \\ -16 & -32 \end{pmatrix}$ indefinit, \mathbf{x}^1 also Sattelpunkt,

$\nabla^2 f(\mathbf{x}^2) = \begin{pmatrix} -6 & -4 \\ -4 & -8 \end{pmatrix}$ neg. definit, \mathbf{x}^2 also strenges lokales

Maximum.



14.2 Implizit definierte Funktionen

Wir fragen, ob man ein gegebenes Gleichungssystem $g(\mathbf{x}, \mathbf{y}) = 0$ bestehend aus m Gleichungen in den $(n + m)$ Unbekannten (\mathbf{x}, \mathbf{y}) nach m dieser Unbekannten (hier o.B.d.A. mit \mathbf{y} bezeichnet) auflösen kann:

$$g(\mathbf{x}, \mathbf{y}) = 0 \quad \Leftrightarrow \quad \mathbf{y} = \mathbf{f}(\mathbf{x}), \quad \mathbf{y} \in \mathbb{R}^m, \quad \mathbf{x} \in \mathbb{R}^n. \quad (14.2.1)$$

Wir sagen in diesem Fall, die Funktion $\mathbf{f} : \mathbb{R}^n \supset D \rightarrow \mathbb{R}^m$ wird durch das Gleichungssystem $g(\mathbf{x}, \mathbf{y}) = 0$ **implizit definiert**. Die Lösungsmenge des Gleichungssystems lässt sich dann (lokal) als Graph der Funktion \mathbf{f} interpretieren.

Beispiele (14.2.2)

a) Die Kreisgleichung $g(x, y) = x^2 + y^2 - r^2 = 0$, $r > 0$, definiert implizit die vier Funktionen

$$y = \pm\sqrt{r^2 - x^2}, \quad x = \pm\sqrt{r^2 - y^2}, \quad -r \leq x, y \leq r.$$

b) Ist g in (14.2.1) eine affin-lineare Funktion, so hat man bei geeigneter Aufteilung der Variablen:

$$g(\mathbf{x}, \mathbf{y}) := \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{b} = \mathbf{0}, \quad \mathbf{A} \in \mathbb{R}^{(m,n)}, \quad \mathbf{B} \in \mathbb{R}^{(m,m)}, \quad \mathbf{b} \in \mathbb{R}^m.$$

Das Gleichungssystem ist nach den Variablen \mathbf{y} auflösbar, falls die Matrix \mathbf{B} regulär ist. Dann gilt

$$g(\mathbf{x}, \mathbf{y}) = \mathbf{0} \Leftrightarrow \mathbf{y} = -\mathbf{B}^{-1}(\mathbf{A}\mathbf{x} + \mathbf{b}) =: \mathbf{f}(\mathbf{x}).$$

Insbesondere ergibt sich für die Jacobi-Matrix der Funktion \mathbf{f} :

$$\mathbf{Jf}(x) = -\mathbf{B}^{-1}\mathbf{A} = -\left(\frac{\partial g}{\partial \mathbf{y}}(\mathbf{x}, \mathbf{y})\right)^{-1} \left(\frac{\partial g}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{y})\right), \quad \mathbf{y} = \mathbf{f}(\mathbf{x}).$$

Zum Nachweis der Auflösbarkeit von (14.2.1) nach \mathbf{y} schreiben wir das Gleichungssystem als Fixpunktgleichung

$$g(\mathbf{x}, \mathbf{y}) = \mathbf{0} \Leftrightarrow \mathbf{h}(\mathbf{x}, \mathbf{y}) := \mathbf{y} - \mathbf{B}^{-1}g(\mathbf{x}, \mathbf{y}) = \mathbf{y}. \quad (14.2.3)$$

Dabei wird \mathbf{x} als Parameter der Fixpunktgleichung betrachtet und \mathbf{B} ist eine reguläre (m, m) -Skalierungsmatrix.

Satz (14.2.4) (Parameterabhängiger Fixpunktsatz)

Seien $U := \overline{K}_{\varepsilon_1}(\mathbf{x}^0)$ und $V := \overline{K}_{\varepsilon_2}(\mathbf{y}^0)$ abgeschlossene Kugeln mit Radien $\varepsilon_1, \varepsilon_2 > 0$ um die Mittelpunkte $\mathbf{x}^0 \in \mathbb{R}^n$ und $\mathbf{y}^0 \in \mathbb{R}^m$.

Die Funktion $\mathbf{h} : U \times V \rightarrow \mathbb{R}^m$ sei stetig, $L \in]0, 1[$ und es mögen die folgenden Voraussetzungen erfüllt sein:

- a)** $\mathbf{h}(\mathbf{x}^0, \mathbf{y}^0) = \mathbf{y}^0$,
- b)** $\|\mathbf{h}(\mathbf{x}, \mathbf{y}_1) - \mathbf{h}(\mathbf{x}, \mathbf{y}_2)\| \leq L \|\mathbf{y}_1 - \mathbf{y}_2\|, \forall \mathbf{x} \in U, \mathbf{y}_1, \mathbf{y}_2 \in V$,
- c)** $\|\mathbf{h}(\mathbf{x}, \mathbf{y}^0) - \mathbf{y}^0\| \leq \varepsilon_2 (1 - L)$, für alle $\mathbf{x} \in U$.

Dann existiert eine *stetige* Funktion $\mathbf{f} : U \rightarrow V$ mit $\mathbf{h}(\mathbf{x}, \mathbf{f}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$ für alle $\mathbf{x} \in U$.

Satz (14.2.5) (Satz über implizite Funktionen)

Sei $g : D \rightarrow \mathbb{R}^m$ eine C^1 -Funktion mit $D \subset \mathbb{R}^{n+m}$ offen.
 $(\mathbf{x}^0, \mathbf{y}^0) \in D$ sei ein Lösungspunkt, also $g(\mathbf{x}^0, \mathbf{y}^0) = \mathbf{0}$. Ferner
sei die Jacobi-Matrix $\frac{\partial \mathbf{g}}{\partial \mathbf{y}}(\mathbf{x}^0, \mathbf{y}^0)$ **regulär**.

Dann gibt es Umgebungen U von $\mathbf{x}^0 \in \mathbb{R}^n$ und V von $\mathbf{y}^0 \in \mathbb{R}^m$, mit $U \times V \subset D$, sowie eine eindeutig bestimmte stetige Funktion $f : U \rightarrow V$ mit

$$f(\mathbf{x}^0) = \mathbf{y}^0 \quad \text{und} \quad g(\mathbf{x}, f(\mathbf{x})) = \mathbf{0} \quad (\forall \mathbf{x} \in U).$$

Die Funktion f ist sogar *stetig differenzierbar* auf U und für die Jacobi-Matrix von f gilt $(\forall \mathbf{x} \in U)$

$$\mathbf{J}f(\mathbf{x}) = - \left(\frac{\partial \mathbf{g}}{\partial \mathbf{y}}(\mathbf{x}, f(\mathbf{x})) \right)^{-1} \left(\frac{\partial \mathbf{g}}{\partial \mathbf{x}}(\mathbf{x}, f(\mathbf{x})) \right).$$

Die obige Formel für $\mathbf{Jf}(\mathbf{x})$ erhält man durch **implizite Differentiation** des Gleichungssystems $g(\mathbf{x}, \mathbf{y}) = 0$. Mit der Kettenregel folgt nämlich:

$$\frac{d}{d\mathbf{x}}g(\mathbf{x}, \mathbf{f}(\mathbf{x})) = \frac{\partial g}{\partial \mathbf{x}}(\mathbf{x}, \mathbf{f}(\mathbf{x})) + \frac{\partial g}{\partial \mathbf{y}}(\mathbf{x}, \mathbf{f}(\mathbf{x})) \cdot \mathbf{Jf}(\mathbf{x}) = 0$$

Beispiele (14.2.6)

a) Die Kreisgleichung $g(x, y) = x^2 + y^2 - r^2 = 0$, $r > 0$, erfüllt im Punkt $(x^0, y^0) := (0, r)$ die Voraussetzung des Satzes über implizite Funktionen, da $g_y(0, r) = 2r \neq 0$. In einer Umgebung von $x^0 = 0$ ist die Kreisgleichung also nach $y = f(x)$ auflösbar und dies auch eindeutig, wenn man $f(0) = y^0 = +r$ verlangt.

(Wir wissen ja schon: $y = +\sqrt{r^2 - x^2}$, $-r \leq x \leq r$).

Die Ableitungen von f lassen sich nun implizit berechnen:

$$2x + 2yy' = 0 \Rightarrow y' = -x/y = -x/\sqrt{r^2 - x^2}.$$

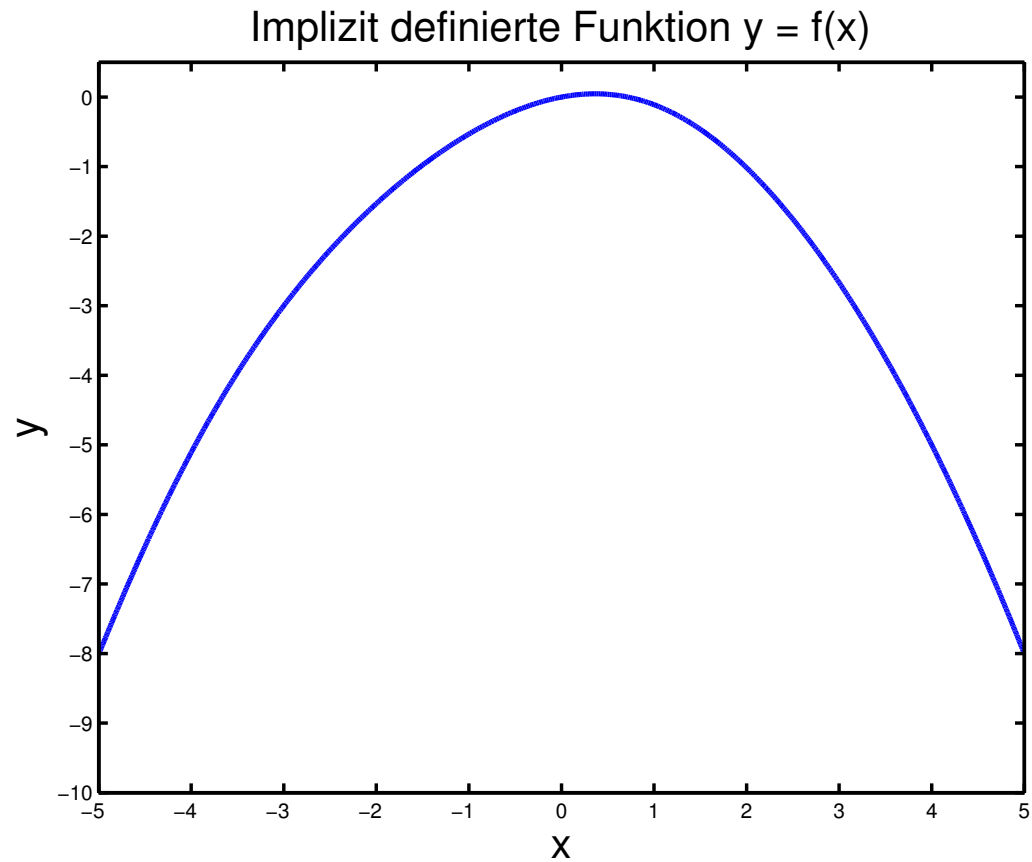
b) Die Gleichung $g(x, y) := e^{y-x} + 3y + x^2 - 1 = 0$ ist für jedes $x \in \mathbb{R}$ eindeutig nach $y := f(x)$ auflösbar. Wegen $g_y = e^{y-x} + 3 > 0$ ist der Satz über implizite Funktionen anwendbar. Implizite Differentiation der Gleichung nach x ergibt:

$$e^{y-x} (y' - 1) + 3y' + 2x = 0 \Rightarrow y' = \frac{e^{y-x} - 2x}{e^{y-x} + 3}.$$

Damit ist f sogar eine C^2 -Funktion. Nochmaliges implizites Differenzieren ergibt:

$$\begin{aligned} e^{y-x} y'' + e^{y-x} (y' - 1)^2 + 3y'' + 2 &= 0 \\ \Rightarrow y'' &= -\frac{2 + e^{y-x} (y' - 1)^2}{e^{y-x} + 3}. \end{aligned}$$

Man beachte, dass hier eine *explizite* Auflösung nach y nicht möglich ist. Die Funktion f kann also ebenso wie die Ableitungen nur numerisch berechnet werden.



Aufgabe: Zeigen Sie, dass f genau ein Extremum besitzt und bestimmen Sie dieses mit einer geeigneten Fixpunkt-Iteration.

Implizite Darstellung ebener Kurven

Die Lösungsmenge einer skalaren Gleichung $g(x, y) = 0$ (g hinreichend glatt) ist nach dem Satz über implizite Funktionen im Fall $\nabla g = (g_x, g_y)^T \neq 0$ lokal als Graph einer Funktion $y = f(x)$ oder $x = f(y)$ gegeben.

Lösungspunkte $\mathbf{x}^0 = (x^0, y^0)$ mit $\nabla g(\mathbf{x}^0) \neq 0$ heißen **regulär**, solche mit $\nabla g(\mathbf{x}^0) = 0$ heißen **singuläre Lösungspunkte**.

Mittels der Taylor-Terme zweiter Ordnung lassen sich die singulären Lösungen von $g(x, y) = 0$ wie folgt klassifizieren:

$$\det \nabla^2 g(\mathbf{x}^0) > 0 : \mathbf{x}^0 \quad \text{isolierter Punkt}$$

$$\det \nabla^2 g(\mathbf{x}^0) < 0 : \mathbf{x}^0 \quad \text{Doppelpunkt}$$

$$\det \nabla^2 g(\mathbf{x}^0) = 0 : \mathbf{x}^0 \quad \text{Rückkehrpunkt (Spitze)}$$

(Letzteres stimmt im Allg. nur dann mit der Anschauung überein, wenn die Terme dritter Ordnung nicht verschwinden.)

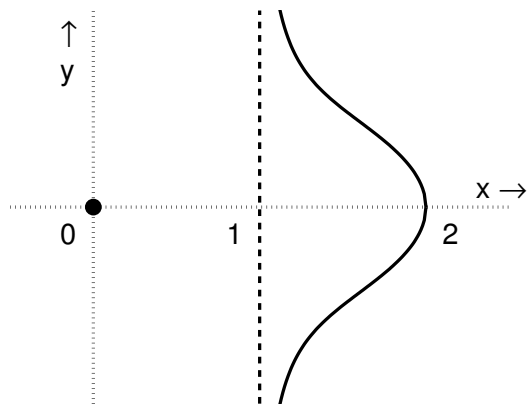
Beispiele (14.2.7)

In folgenden Beispielen ist $x^0 = 0$ ein singulärer Lösungspunkt.

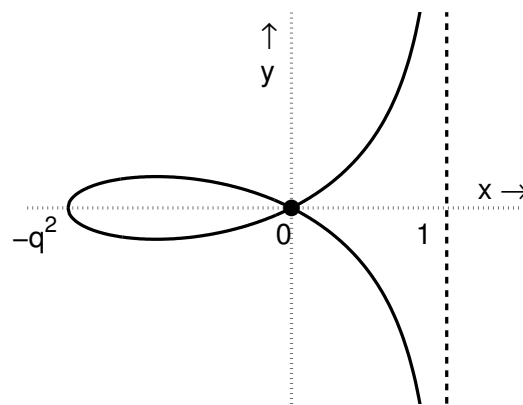
a) $g(x, y) = y^2(x - 1) + x^2(x - 2)$

b) $g(x, y) = y^2(x - 1) + x^2(x + q^2)$

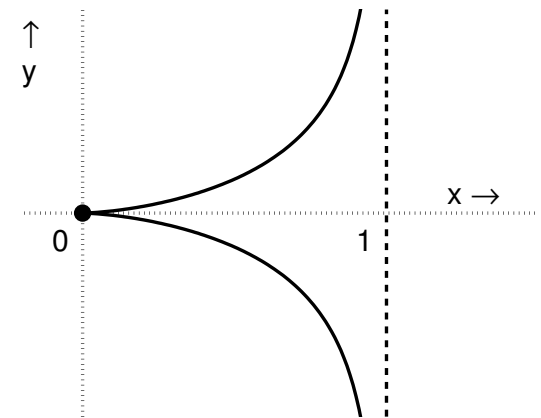
c) $g(x, y) = y^2(x - 1) + x^3$



a)



b)



c)

Implizite Darstellung von Flächen im \mathbb{R}^3

Die Lösungsmenge einer skalaren Gleichung $g(x, y, z) = 0$ (g hinreichend glatt) ist nach dem Satz über implizite Funktionen im Fall $\nabla g = (g_x, g_y, g_z)^\top \neq \mathbf{0}$ lokal als Graph einer Funktion $z = f(x, y)$ oder $y = f(x, z)$ oder $x = f(y, z)$ gegeben.

Die **Tangentialebene** an diese Fläche in einem regulären Lösungspunkt \mathbf{x}^0 lautet dann nach dem Taylorschen Satz (nur Terme erster Ordnung, implizite Darstellung):

$$\nabla g(\mathbf{x}^0)^\top (\mathbf{x} - \mathbf{x}^0) = 0. \quad (14.2.8)$$

Man findet also wiederum, dass der Gradient $\nabla g(\mathbf{x}^0)$ auf der Fläche $g(x, y, z) = 0$ (Niveaumenge oder Äquipotentialfläche) senkrecht steht; vgl. auch (13.1.15).

Das Umkehrproblem

Sei $f : D \rightarrow \mathbb{R}^n$ eine C^1 -Funktion und $D \subset \mathbb{R}^n$ offen. Wir fragen nach der (lokalen) Invertierbarkeit von f . Dazu bringen wir das Gleichungssystem $y = f(x)$ in die Form $g(x, y) := f(x) - y = 0$ und wenden auf g den Satz über implizite Funktionen an.

Satz (14.2.9) (Umkehrsatz)

Ist die Jacobi-Matrix $Jf(x^0)$ für ein $x^0 \in D$ regulär, so gibt es offene Umgebungen U von x^0 und V von $y^0 := f(x^0)$, so dass die Funktion f den Bereich U *bijektiv* auf V abbildet.

Ferner ist die Umkehrabbildung $f^{-1} : V \rightarrow U$ ebenfalls eine C^1 -Funktion, und es gilt für alle $x \in U$:

$$Jf^{-1}(y) = (Jf(x))^{-1}, \quad y = f(x).$$

Beispiel (14.2.10)

Die Transformation des \mathbb{R}^2 mittels Polarkoordinaten

$$\mathbf{f} :]0, \infty[\times \mathbb{R} \rightarrow \mathbb{R}^2, \quad \mathbf{x} = \mathbf{f}(r, \varphi) := (r \cos(\varphi), r \sin(\varphi))^T$$

ist wegen $\det \mathbf{Jf}(r, \varphi) = r$ an jeder Stelle (r_0, φ_0) mit $r_0 > 0$ lokal invertierbar. Die Umkehrfunktion ist eine C^1 -Funktion mit der Jacobi-Matrix

$$\mathbf{Jf}^{-1}(\mathbf{x}) = \begin{pmatrix} \cos \varphi & -\frac{1}{r} \sin \varphi \\ \sin \varphi & \frac{1}{r} \cos \varphi \end{pmatrix}, \quad \mathbf{x} = \mathbf{f}(r, \varphi).$$

Bemerkung (14.2.11)

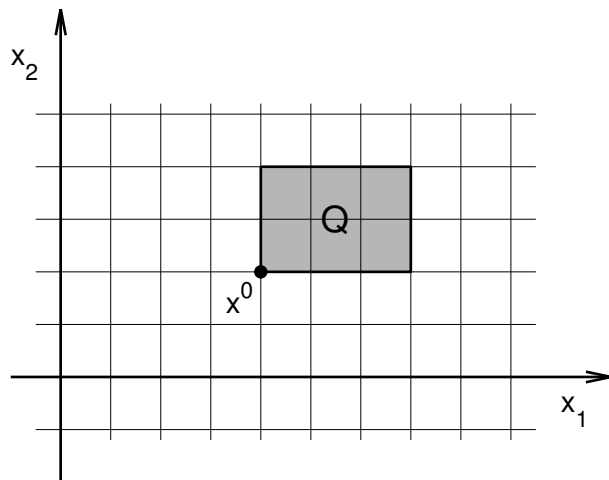
Bijektive C^1 -Funktionen, deren Umkehrungen ebenfalls C^1 -Abbildungen sind, werden auch **C^1 -Diffeomorphismen** genannt.

Transformation von Volumina

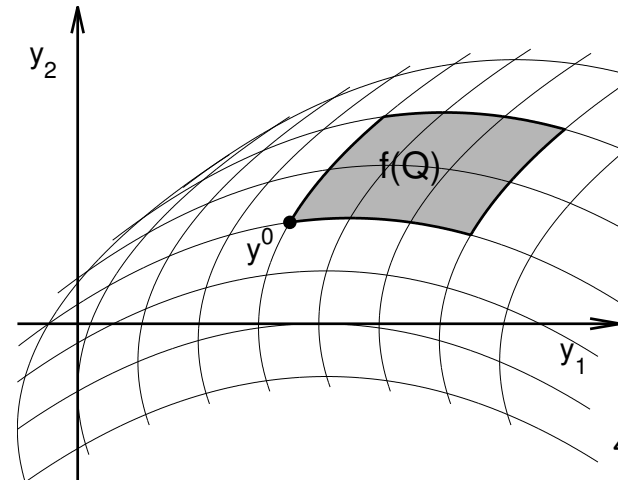
Sei $f : \mathbb{R}^n \supset D \rightarrow \mathbb{R}^n$ ein auf (offenen) Umgebungen U_1 von $\mathbf{x}^0 \in D$ und U_2 von $\mathbf{y}^0 := f(\mathbf{x}^0)$ definierter C^1 -Diffeomorphismus.

Ein „kleines“ Volumen V bei \mathbf{x}^0 wird dann auf ein „kleines“ Volumen \tilde{V} bei \mathbf{y}^0 abgebildet. Für die (Vorzeichen-behafteten) Volumina gilt näherungsweise

$$\text{vol}(\tilde{V}) \approx \det(\mathbf{J}f(\mathbf{x}^0)) \cdot \text{vol}(V). \quad (14.2.12)$$



$$\mathbf{y} = f(\mathbf{x})$$



14.3 Gleichungsrestringierte Optimierung

Gesucht sind die Extremwerte einer Funktion $f : D \rightarrow \mathbb{R}$ ($D \subset \mathbb{R}^n$ offen) auf einer Teilmenge G des Definitionsbereichs, die durch Gleichungen der Form $g_i(\mathbf{x}) = 0$ definiert ist:

$$G := \{\mathbf{x} \in D : \mathbf{g}(\mathbf{x}) = \mathbf{0}\}.$$

Dabei sei $\mathbf{g} := (g_1, \dots, g_m)^\top : D \rightarrow \mathbb{R}^m$ ($m < n$) eine C^1 -Abbildung.

Grundidee von Lagrange: Ankopplung der Nebenbedingungen $g_i(\mathbf{x}) = 0$, mittels (unbekannter) **Lagrange-Multiplikatoren** $\lambda_i \in \mathbb{R}$ an die zu maximierende (minimierende) Funktion f . Man betrachtet also – anstelle von f – die **Lagrange-Funktion**

$$L(\mathbf{x}, \boldsymbol{\lambda}) := f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{g}(\mathbf{x}) \quad (14.3.1)$$

und maximiert (bzw. minimiert) diese für festes $\boldsymbol{\lambda}$ über dem gesamten Bereich $\mathbf{x} \in D$ (unrestringiert!).

Satz (14.3.2) (Lagrange–Lemma)

Minimiert bzw. maximiert $\mathbf{x}^0 \in D$ die Lagrange–Funktion $L(\mathbf{x}, \boldsymbol{\lambda})$ für ein festes $\boldsymbol{\lambda}$ über D und gilt $\mathbf{g}(\mathbf{x}^0) = \mathbf{0}$, so ist \mathbf{x}^0 zugleich ein globales Minimum bzw. Maximum von f über der zulässigen Menge $G := \{\mathbf{x} \in D : \mathbf{g}(\mathbf{x}) = \mathbf{0}\}$.

Bemerkung (14.3.3)

Sind f und \mathbf{g} C^1 -Funktionen, so lautet die notwendige Bedingung für eine Extremalstelle \mathbf{x}^0 von L :

$$\nabla_{\mathbf{x}}L(\mathbf{x}, \boldsymbol{\lambda}) = \nabla f(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{x}) = \mathbf{0}. \quad (14.3.4)$$

(14.3.4) bildet zusammen mit $\mathbf{g}(\mathbf{x}) = \mathbf{0}$ ein Gleichungssystem für die $(n + m)$ Unbekannten \mathbf{x} und $\boldsymbol{\lambda}$. Die Lösungen $(\mathbf{x}^0, \boldsymbol{\lambda}^0)$ sind gerade die Kandidaten für die gesuchten Extremalstellen (**notwendige Bedingung**).

Satz (14.3.5) (Lagrange–Multiplikatoren–Regel)

Sind $f, g \in C^1$ -Funktionen auf der offenen Menge $D \subset \mathbb{R}^n$, ist $\mathbf{x}^0 \in D$ ein lokales Extremum von f unter der Nebenbedingung $g(\mathbf{x}) = 0$ und gilt die **Regularitätsbedingung (constraint qualification)**

$$\text{Rang} \left(\mathbf{J}g(\mathbf{x}^0) \right) = m,$$

so existieren Lagrange-Multiplikatoren $\lambda_1, \dots, \lambda_m$, so dass für die Lagrange–Funktion $L = f + \boldsymbol{\lambda}^\top g$ die **notwendige Bedingung erster Ordnung** (14.3.4) gilt: $\nabla_{\mathbf{x}} L(\mathbf{x}^0, \boldsymbol{\lambda}) = \mathbf{0}$.

Beispiel (14.3.6)

Gesucht ist ein Quader mit maximalem Volumen bei vorgegebener Oberfläche. Bezeichnen $x, y, z > 0$ die Kantenlängen, so ist die Funktion (Volumen) $f(x, y, z) := xyz$ zu maximieren unter der Nebenbedingung $g(x, y, z) := xy + xz + yz - a = 0$, $a > 0$.

Da die Regularitätsbedingung für zulässige Punkte erfüllt ist, lässt sich (14.3.5) anwenden. Mit der Lagrange-Funktion

$$L(x, y, z, \lambda) = x y z + \lambda (x y + x z + y z - a)$$

erhält man:

$$\nabla_{\mathbf{x}} L = \begin{pmatrix} y z + \lambda (y + z) \\ x z + \lambda (x + z) \\ x y + \lambda (x + y) \end{pmatrix} = \mathbf{0}.$$

Für zulässige Punkte $x, y, z > 0$ ist auch $\lambda \neq 0$ und man erhält die (eindeutig bestimmte) Lösung

$$x_0 = y_0 = z_0 = \sqrt{a/3}, \quad \lambda = -0.5 \sqrt{a/3}, \quad f(\mathbf{x}_0) = (a/3)^{3/2}.$$

Beispiel (14.3.7)

Gesucht sind die Punkte auf der Parabel $y = x^2 - a$, die vom Ursprung den kürzesten Abstand haben.

Zu minimieren ist also die Funktion $f(x, y) = x^2 + y^2$ unter der Nebenbedingung $g(x, y) = y - x^2 + a = 0$.

Mit der Lagrange-Funktion $L(x, y, \lambda) = x^2 + y^2 + \lambda(y - x^2 + a)$ ergibt sich die notwendige Bedingung

$$\nabla_{\mathbf{x}}L = \begin{pmatrix} 2x - 2\lambda x \\ 2y + \lambda \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

Hieraus ergeben sich die Lösungen:

- (i) $x = 0, \quad y = -a, \quad \lambda = 2a$
- (ii) $x = \pm\sqrt{a - 0.5}, \quad y = -0.5, \quad \lambda = 1.$

Die Lösungen (ii) sind nur für $a \geq 0.5$ erklärt. Durch (i) wird die Lösung für $a \leq 0.5$ beschrieben.

Die Beziehung (14.3.4) ist die „notwendige Bedingung erster Ordnung“ für gleichungsrestringierte Optimierungsaufgaben. Daneben lassen sich auch notwendige und hinreichende Bedingungen „zweiter Ordnung“ beweisen. Dazu nehmen wir an, dass die Funktionen f, g sogar zweifach stetig differenzierbar sind.

Satz (14.3.8)

a) Notwendige Bedingung zweiter Ordnung:

Ist $\mathbf{x}^0 \in D$ lokales Minimum von f auf G , ist die Regularitätsbedingung aus (14.3.5) erfüllt, so ist die Hesse-Matrix $\nabla_{\mathbf{x}}^2 L(\mathbf{x}^0, \boldsymbol{\lambda})$ der Lagrange-Funktion positiv semidefinit auf dem Tangentialraum der zulässigen Menge G

$$T_G(\mathbf{x}^0) := \left\{ \mathbf{y} \in \mathbb{R}^n : \nabla g_i(\mathbf{x}^0)^\top \mathbf{y} = 0, i = 1, \dots, m \right\}.$$

Dies bedeutet: $\forall \mathbf{y} \in T_G(\mathbf{x}^0) : \mathbf{y}^\top \nabla_{\mathbf{x}}^2 L(\mathbf{x}^0, \boldsymbol{\lambda}) \mathbf{y} \geq 0$.

b) Hinreichende Bedingung zweiter Ordnung:

Ist $\mathbf{x}^0 \in G$ und existieren Lagrange-Parameter $\boldsymbol{\lambda}$, so dass \mathbf{x}^0 stationärer Punkt der Lagrange-Funktion ist und deren Hesse-Matrix $\nabla_{\mathbf{x}}^2 L(\mathbf{x}^0, \boldsymbol{\lambda})$ positiv definit auf $T_G(\mathbf{x}^0)$ ist, so ist \mathbf{x}^0 ein strenges lokales Minimum von f auf G .

Beispiel (14.3.9)

Zu minimieren sei die Funktion $f(x_1, x_2) := x_1 x_2$ unter der Nebenbedingung $g(x_1, x_2) := x_1^2 + x_2^2 - 1 = 0$.

Mit der Lagrange-Funktion $L(\mathbf{x}, \lambda) = x_1 x_2 + \lambda (x_1^2 + x_2^2 - 1)$ und den notwendigen Bedingungen $L_{x_1} = L_{x_2} = g = 0$ ergeben sich die folgenden vier Lösungskandidaten:

$$\begin{aligned} \lambda_{1,2} = 0.5 & \quad : \quad \mathbf{x}^1 = (\sqrt{0.5}, -\sqrt{0.5})^\top, \quad \mathbf{x}^2 = -\mathbf{x}^1 \\ \lambda_{3,4} = -0.5 & \quad : \quad \mathbf{x}^3 = (\sqrt{0.5}, \sqrt{0.5})^\top, \quad \mathbf{x}^4 = -\mathbf{x}^3. \end{aligned}$$

Wegen $f(\mathbf{x}^1) = f(\mathbf{x}^2) = -0.5$ und $f(\mathbf{x}^3) = f(\mathbf{x}^4) = 0.5$ sind $\mathbf{x}^1, \mathbf{x}^2$ globale Minima und $\mathbf{x}^3, \mathbf{x}^4$ globale Maxima von f auf G .

Wir untersuchen die Bedingungen zweiter Ordnung für \mathbf{x}^1 :

Mit $\nabla g = 2\mathbf{x}$ ergibt sich der Tangentialraum zu $T_G(\mathbf{x}^1) = \{\mathbf{y} \in \mathbb{R}^2 : (\mathbf{x}^1)^\top \mathbf{y} = 0\}$. Ferner ist

$$\nabla_{\mathbf{x}}^2 L(\mathbf{x}^1, \lambda_1) = \begin{pmatrix} 2\lambda_1 & 1 \\ 1 & 2\lambda_1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}.$$

Die Hesse-Matrix der Lagrange-Funktion ist also *nicht* positiv definit auf \mathbb{R}^2 , allerdings gilt für $\mathbf{y} \in T_G(\mathbf{x}^1) \setminus \{\mathbf{0}\}$:

$$\mathbf{y}^\top \nabla_{\mathbf{x}}^2 L(\mathbf{x}^1, \lambda_1) \mathbf{y} = (y_1 + y_2)^2 = 4y_1^2 > 0.$$

Damit ist die hinreichende Bedingung zweiter Ordnung erfüllt und \mathbf{x}^1 ist ein strenges lokales (und sogar ein globales) Minimum von f auf G .

14.4 Das Newton–Verfahren

Gegeben sei eine C^1 –Funktion $f : D \rightarrow \mathbb{R}^n$ mit $D \subset \mathbb{R}^n$ offen. Gesucht ist eine Nullstelle von f , also eine Lösung $\mathbf{x}^* \in D$ des Gleichungssystems $f(\mathbf{x}) = \mathbf{0} \in \mathbb{R}^n$.

Die Taylor–Entwicklung von f um einen Startwert $\mathbf{x}^0 \in D$ lautet

$$f(\mathbf{x}) = f(\mathbf{x}^0) + \mathbf{J}f(\mathbf{x}^0)(\mathbf{x} - \mathbf{x}^0) + o(\|\mathbf{x} - \mathbf{x}^0\|).$$

Setzt man hierin für \mathbf{x} die gesuchte Nullstelle \mathbf{x}^* ein und vernachlässigt man den Fehlerterm $o(\|\mathbf{x}^* - \mathbf{x}^0\|)$, so erhält man

$$\mathbf{J}f(\mathbf{x}^0)(\mathbf{x}^* - \mathbf{x}^0) \approx -f(\mathbf{x}^0).$$

Als neue Näherung für \mathbf{x}^* wählt man daher die Lösung \mathbf{x}^1 des *linearen Gleichungssystems*

$$\mathbf{J}f(\mathbf{x}^0)(\mathbf{x}^1 - \mathbf{x}^0) = -f(\mathbf{x}^0). \quad (14.4.1)$$

Sofern die Jacobi-Matrix $\mathbf{Jf}(\mathbf{x}^0)$ regulär ist, besitzt (14.4.1) eine eindeutig bestimmte Lösung.

Algorithmus (14.4.2)

Wähle einen Startvektor $\mathbf{x}^0 \in D$, $\text{tol} > 0$

für $k = 0, 1, 2, \dots$

Berechne $\Delta \mathbf{x}^k$ aus $\mathbf{Jf}(\mathbf{x}^k) \Delta \mathbf{x}^k = -\mathbf{f}(\mathbf{x}^k)$,

$\mathbf{x}^{k+1} := \mathbf{x}^k + \Delta \mathbf{x}^k$,

Abbruch, falls $\|\mathbf{x}^{k+1} - \mathbf{x}^k\| < \text{tol} \cdot \|\mathbf{x}^k\|$

end k

Satz (14.4.3) (Skalierungsinvarianz)

Das Newton-Verfahren ist invariant unter linearen Transformationen (Skalierungen) der Form

$$\mathbf{f}(\mathbf{x}) \rightarrow \mathbf{g}(\mathbf{y}) := \mathbf{A} \cdot \mathbf{f}(\mathbf{B} \mathbf{y}),$$

mit beliebigen (festen) regulären Matrizen $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{(n,n)}$.

Satz (14.4.4) (Konvergenzsatz)

Sei $f : \mathbb{R}^n \supset D \rightarrow \mathbb{R}^n$ C^1 -Funktion, D offen, und $\mathbf{x}^* \in D$ eine Nullstelle von f . Die Jacobi-Matrix $\mathbf{J}f(\mathbf{x})$ sei regulär für alle $\mathbf{x} \in D$, und es gelte eine Lipschitz-Bedingung ($L > 0$)

$$\forall \mathbf{x}, \mathbf{y} \in D : \left\| (\mathbf{J}f(\mathbf{x}))^{-1} (\mathbf{J}f(\mathbf{y}) - \mathbf{J}f(\mathbf{x})) \right\| \leq L \|\mathbf{y} - \mathbf{x}\|.$$

Dann ist das Newton-Verfahren (14.4.2) für alle Startwerte $\mathbf{x}^0 \in D$ mit $\|\mathbf{x}^0 - \mathbf{x}^*\| < 2/L =: r$ und $K_r(\mathbf{x}^*) \subset D$ wohldefiniert. Alle Newton-Iterierten \mathbf{x}^k liegen in der Kugel $K_r(\mathbf{x}^*)$ und die Folge konvergiert **quadratisch** gegen \mathbf{x}^* :

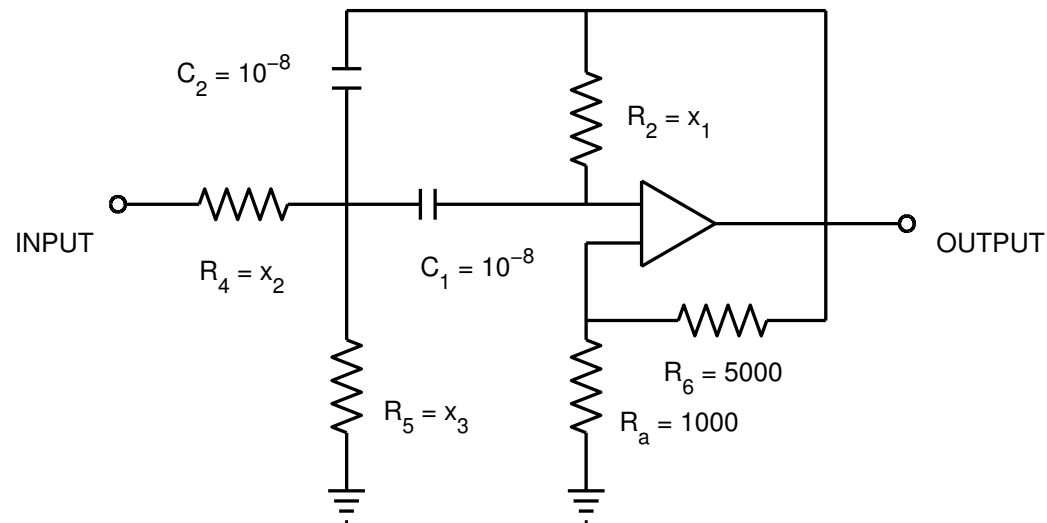
$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq (L/2) \|\mathbf{x}^k - \mathbf{x}^*\|^2.$$

Ferner ist \mathbf{x}^* die einzige Nullstelle von f innerhalb der Kugel $K_r(\mathbf{x}^*)$.

Beispiel (14.4.5)

Die Transformation vom Eingang- zum Ausgangssignal einer elektrischen Schaltung (aktiver Filter) wird beschrieben durch

$$T(s) = \frac{a_1 s}{s^2 + b_1 s + b_0}, \quad a_1 = 1.2 \cdot 10^5 / x_2,$$
$$b_0 = \frac{(x_2 + x_3) \cdot 10^{10}}{x_1 x_2 x_3}, \quad b_1 = \frac{2 \cdot 10^5}{x_1} - \frac{2 \cdot 10^4 (x_2 + x_3)}{x_2 x_3}$$



Die Widerstände x_1, x_2, x_3 sollen so angepasst werden, dass $a_1 = 1$, $b_0 = 40$ und $b_1 = 1$ gilt. Damit ergibt sich das folgende nichtlineare Gleichungssystem

$$\mathbf{f}(\mathbf{x}) := \begin{pmatrix} 1.2 \cdot 10^5 / x_2 - 1 \\ \frac{(x_2 + x_3) \cdot 10^{10}}{x_1 x_2 x_3} - 40 \\ \frac{2 \cdot 10^5}{x_1} - \frac{2 \cdot 10^4 (x_2 + x_3)}{x_2 x_3} - 1 \end{pmatrix} = \mathbf{0}.$$

Die Anwendung des Newton-Verfahrens mit den Startwerten $x_1 = x_2 = x_3 = 10^3$ liefert nach 14 Iterationen eine Lösung mit der relativen Genauigkeit von etwa 10 Dezimalziffern:

$$\mathbf{x}^* = (4.413911092 \cdot 10^4, 1.200000000 \cdot 10^5, 5.944486448 \cdot 10^3)^\top.$$

Anhand des Lösungsverlaufs ($\|\Delta \mathbf{x}^k\|$) lässt sich die quadratische Konvergenz bestätigen.

Das gedämpfte Newton-Verfahren

Zur Verbesserung der (globalen) Konvergenzeigenschaften des Newton-Verfahrens reduziert man „weit weg“ von der Lösung \mathbf{x}^* die Länge der Newton-Korrektur $\Delta \mathbf{x}^k$, verwendet also anstelle von (14.4.2) die folgende modifizierte Newton-Iteration:

$$\begin{aligned} \mathbf{Jf}(\mathbf{x}^k) \Delta \mathbf{x}^k &= -\mathbf{f}(\mathbf{x}^k) \\ \mathbf{x}^{k+1} &= \mathbf{x}^k + \lambda_k \Delta \mathbf{x}^k . \end{aligned} \tag{14.4.6}$$

Der **Dämpfungsparameter** $\lambda_k \in]0, 1]$ wird hierbei so gewählt, dass eine geeignete **Testfunktion** $T(\mathbf{x})$, etwa $T(\mathbf{x}) := \|\mathbf{f}(\mathbf{x})\|$, in jedem Iterationsschritt fällt. Zur algorithmischen Bestimmung von λ_k kann man eine einfache Halbierungsstrategie verwenden. Ferner lässt sich zeigen, dass man hiermit stets ein $\lambda_k > 0$ finden kann, für das die Testfunktion fällt.

Algorithmus (14.4.7) (Schrittweitenbestimmung)

$$\mu := \lambda_{-1} := 1$$

für $k = 0, 1, 2, \dots$

Wähle $\lambda_k \in \{\mu, \mu/2, \mu/4, \mu/8, \dots\}$ maximal, mit

$$T(\mathbf{x}^k + \lambda_k \Delta \mathbf{x}^k) < T(\mathbf{x}^k)$$

Abbruch, falls $\lambda_k < \lambda_{\min}$

$$\mathbf{x}_{k+1} := \mathbf{x}^k + \lambda_k \Delta \mathbf{x}^k$$

$$\mu := \begin{cases} \min(1, 2 \lambda_k), & \text{falls } \lambda_k = \lambda_{k-1} \\ \lambda_k, & \text{sonst} \end{cases}$$

end k

In der Literatur findet man ausgefeiltere Strategien zur Schrittweitenbestimmung, insbesondere ist es günstig, skalierungsinvariante Testfunktionen zu verwenden.

Nichtlineare Ausgleichsprobleme

Ein häufig auftretendes Problem ist die Anpassung einer parameterabhängigen Funktion an vorgegebene Messdaten. Gegeben seien Messdaten $(t_i, y_i) \in \mathbb{R}^2$, $i = 1, \dots, m$, sowie eine **Ansatzfunktion** $y = f(t; \mathbf{x})$, wobei $\mathbf{x} = (x_1, \dots, x_n)^\top$ die Parameter beschreibt. Es wird angenommen, dass $n \ll m$ ist, also mehr Messdaten zur Verfügung stehen als Parametern.

Aufgabe ist es, die Funktion $g(\mathbf{x}) := \|\mathbf{F}(\mathbf{x}) - \mathbf{y}\|_2^2$, $\mathbf{x} \in \mathbb{R}^n$, zu minimieren. Dabei ist

$$\mathbf{F}(\mathbf{x}) := (f(t_1; \mathbf{x}), \dots, f(t_m; \mathbf{x}))^\top, \quad \mathbf{y} := (y_1, \dots, y_m)^\top.$$

Die Grundidee des Newton–Verfahrens, ein nichtlineares Gleichungssystem in der Nähe einer Iterierten \mathbf{x}^k durch Linearisierung näherungsweise zu lösen, lässt sich auf nichtlineare Ausgleichsprobleme übertragen. Dieses Verfahren geht auf Gauß zurück und wird **Gauß–Newton–Verfahren** genannt.

Sei $\mathbf{x}^k \in \mathbb{R}^n$ Näherung für eine Lösung \mathbf{x}^* des Ausgleichsproblems. Wir ersetzen \mathbf{F} durch das lineare Taylor–Polynom zum Entwicklungspunkt \mathbf{x}^k :

$$g(\mathbf{x}) \approx \tilde{g}(\mathbf{x}) := \left\| \mathbf{F}(\mathbf{x}^k) + \mathbf{JF}(\mathbf{x}^k) (\mathbf{x} - \mathbf{x}^k) - \mathbf{y} \right\|_2^2.$$

Die Minimierung von \tilde{g} (anstelle von g) bildet nun eine *lineare* Ausgleichsaufgabe, deren Lösung \mathbf{x}^{k+1} z.B. mittels QR-Zerlegung (Householder/Givens) berechnet werden kann.

Algorithmus (14.4.8) (Gauß–Newton–Verfahren)

Wähle einen Startvektor $\mathbf{x}^0 \in D$, $\text{tol} > 0$

für $k = 0, 1, 2, \dots$

Berechne $\Delta \mathbf{x}^k$ aus der Minimierung von $\left\| \mathbf{F}(\mathbf{x}^k) + \mathbf{JF}(\mathbf{x}^k) \Delta \mathbf{x}^k - \mathbf{y} \right\|_2^2$,

Berechne eine Schrittweite $\lambda_k > 0$,

$$\mathbf{x}^{k+1} := \mathbf{x}^k + \lambda_k \Delta \mathbf{x}^k,$$

Abbruch, falls $\|\Delta \mathbf{x}^k\| < \text{tol} \cdot \|\mathbf{x}^k\|$

end k .

In **MATLAB** steht professionelle Software zur Lösung nichtlinearer Ausgleichsaufgaben zur Verfügung. Standardaufruf ist

$$[\mathbf{x}, \text{resnorm}, \mathbf{residual}] = \text{lsqnonlin}(@\text{fun}, \mathbf{x0})$$

Dabei ist `fun` eine *function* (Unterprogramm) zur Berechnung des Residuums $\mathbf{F}(\mathbf{x}) - \mathbf{y}$, `x0` ist ein vorzugebender Startvektor. Der Lösungsvektor ist durch `x` gegeben, `residual` gibt den optimalen Residuenvektor $\mathbf{F}(\mathbf{x}^*) - \mathbf{y}$ an und `resnorm` dessen Euklidische Norm.

Beispiel (14.4.9) (Blutspiegelmessungen)

Die Konzentration eines Arzneistoffes im Blutplasma (der Blutspiegel) wird nach der Einnahme (Zeit $t = 0$) vereinfachend durch die so genannte Bateman-Funktion* beschrieben:

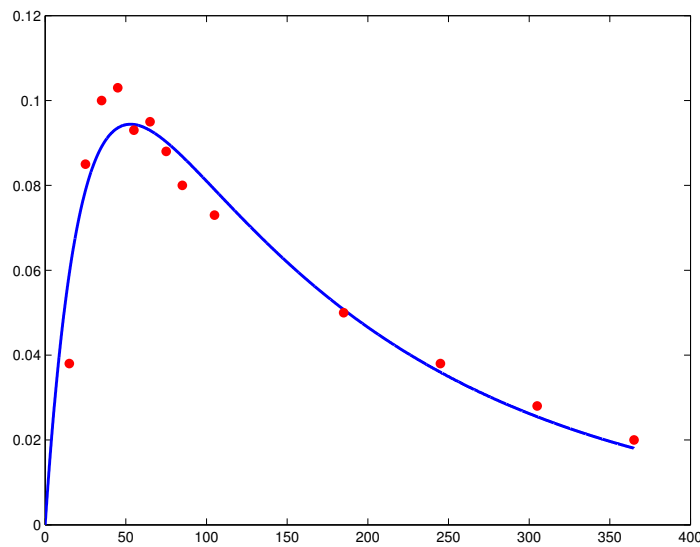
*benannt nach dem brit. Mathematiker Harry Bateman (1882–1946)

$$y(t; \mathbf{x}) := x_3 [\exp(-x_1 t) - \exp(-x_2 t)]$$

Aufgabe ist es, die Parameter x_1, x_2, x_3 so zu bestimmen, dass diese Funktion möglichst gut die folgenden Messdaten wiedergibt:

t_i :	15	25	35	45	55	65	75	85	105	185
y_i :	0.038	0.085	0.1	0.103	0.093	0.095	0.088	0.08	0.073	0.05

t_i :	245	305	365
y_i :	0.038	0.028	0.02



$$x_1^* = 0.0057$$

$$x_2^* = 0.0443$$

$$x_3^* = 0.1470$$