

Anschaulich bedeutet dies, daß der Zustand eines Systems, das durch diese Differentialgleichung beschrieben wird, für $\vec{y}(t) = \vec{y}_0$ zeitlich konstant ist, das System befindet sich also im Gleichgewicht.

Da die Ableitung einer konstanten Funktion verschwindet, sind die Fixpunkte des Differentialgleichungssystems $\vec{y}(t) = f(\vec{y}(t), t)$ gerade die Lösungen des Gleichungssystems

$$f(\vec{y}_0, t) = 0 \quad \text{für alle } t \in [t_0, t_1].$$

Bei einem nichtlinearen Differentialgleichungssystem ist das ein nichtlineares Gleichungssystem, man wird sich daher oft mit Näherungslösungen begnügen müssen. (Die Variable t tritt natürlich nur bei nichtautonomen Systemen auf; bei den in Naturwissenschaft und Technik häufigen autonomen Systemen haben wir ein Gleichungssystem, in dem nur die Komponenten von \vec{y}_0 vorkommen.)

Ein klassisches Beispiel, bei dem sich die Fixpunkte leicht ausrechnen lassen, ist das Raubtier-Beutetier-Modell, das 1925 von LOTKA und VOLTERRA vorgeschlagen wurde: In einem Gebiet gebe es eine Population von Raubtieren, die sich von genau einer Art von Beutetieren ernähren. Für die Beutetiere sei genügend Nahrung vorhanden, so daß diese sich, falls es keine Raubtiere gäbe, beliebig vermehren könnten. Wenn wir die Populationsstärke zum Zeitpunkt t mit $x(t)$ bezeichnen, können wir also annehmen, daß die Beutetiere bei Abwesenheit der Raubtiere eine konstanter Wachstumsrate hätten, d.h.

$$\dot{x}(t) = \alpha x(t) \quad \text{mit } \alpha > 0.$$

Die Raubtiere, deren Bestand zum Zeitpunkt t wir mit $y(t)$ bezeichnen wollen, würden in Abwesenheit der Beutetiere relativ schnell verhungern und somit aussterben, was wir durch eine negative Wachstumsrate modellieren können:

$$\dot{y}(t) = -\gamma y(t) \quad \text{mit } \gamma > 0.$$

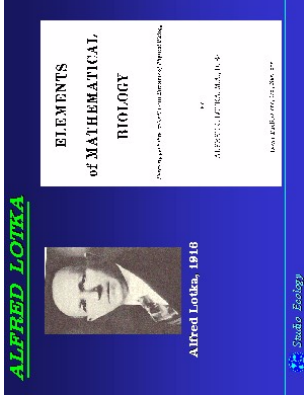
Nun sind aber die Raubtiere und die Beutetiere nicht isoliert voneinander, sondern es kommt zu Begegnungen zwischen den beiden Populationen. Deren Häufigkeit ist etwa proportional zum Produkt der beiden Populationsstärken, und die Auswirkung einer solchen Begegnung ist

positiv für die Wachstumsrate der Raubtiere, aber negativ für die der Beutetiere. Unser Modell läßt sich somit beschreiben durch das System

$$\dot{x}(t) = \alpha x(t) - \beta x(t)y(t)$$

$$\dot{y}(t) = -\gamma y(t) + \delta x(t)y(t)$$

mit positiven reellen Zahlen $\alpha, \beta, \gamma, \delta$.



Der amerikanische Wissenschaftler ALFRED LOTKA (1880–1949) war von der Ausbildung her ein Mathematiker, interessierte sich aber Zeit seines Lebens stark für Physik, insbesondere Thermodynamik, und war einer der ersten, der die Evolution unter physikalischen Gesichtspunkten betrachtete. Er zählt zu den Pionieren der Selbstorganisation, der Bioenergetik und (auch für eine Versickerung) der Demographie.



VITO VOLTERRA (1860–1940) wurde in Ancona im damaligen Kirchenstaat geboren. Er studierte bereits als Elfjähriger mathematische Literatur, promovierte dann aber in Physik über ein Thema aus der Hydrodynamik. Er hatte Lehrstühle für Mechanik und für mathematische Physik in Pisa, Turin und Rom. Seine wichtigsten Arbeiten beschäftigen sich mit partiellen Differentialgleichungen und vor allem Integralgleichungen. Ab 1922 kämpfte er im italienischen Parlament gegen den Faschismus und verlor deshalb 1931 nach Auflösung des Parlaments seinen Lehrstuhl in Rom. Den Rest seines Lebens verbrachte er größtenteils im Exil.

Zur Bestimmung der Gleichgewichtslösungen müssen wir für $x(t)$ und $y(t)$ Konstanten einsetzen; dies führt auf die Gleichungen

$$0 = \alpha x_0 - \beta x_0 y_0 = x_0(\alpha - \beta y_0)$$

$$0 = -\gamma y_0 + \delta x_0 y_0 = y_0(\gamma - \delta x_0).$$

Es gibt somit genau zwei Gleichgewichtslösungen: Einmal die uninteressante Lösung $x(t) = y(t) \equiv 0$, die im wesentlichen besagt, daß ohne Raub- und Beutetiere in diesem System nichts passiert, und dann noch

die Lösung

$$x(t) \equiv \frac{\gamma}{\delta} \quad \text{und} \quad y(t) \equiv \frac{\alpha}{\beta}.$$

Falls die beiden Populationen diese Stärken haben, fressen also die Raubtiere genau so viele Beutetiere weg, wie nachwachsen; umgekehrt ernähren die Beutetiere gerade aus, um die Raubtierpopulation zu ernähren. Was passiert, wenn die Populationen nicht im Gleichgewicht sind? Wir haben offensichtlich kaum Chancen, das Differentialgleichungssystem explizit zu lösen, aber wir können trotzdem versuchen, etwas über die Lösungskurven in Erfahrung zu bringen.

Wenn wir y als Funktion von x betrachten, ist

$$y'(x) = \frac{dy}{dx} = \frac{\dot{y}(t)}{\dot{x}(t)} = \frac{-\gamma y + \delta x y}{\alpha x - \beta x y} = \frac{y}{\alpha - \beta y} \cdot \frac{\delta x - \gamma}{x},$$

wir haben also eine Differentialgleichung mit getrennten Veränderlichen. Trennung der Variablen führt auf

$$\int \left(\frac{\alpha}{y} - \beta \right) dy = \int \left(\delta - \frac{\gamma}{x} \right) dx$$

oder

$$\alpha \ln y - \beta y = \delta x - \gamma \ln x + C.$$

Anwendung der Exponentialfunktion macht daraus

$$\frac{y^\alpha}{e^{\beta y}} = \frac{e^{\delta x}}{x^\gamma} \cdot e^C.$$

Diese Gleichung können wir zwar weder nach y noch nach x auflösen, aber eine einfache Kurvendiskussion der Funktionen

$$f(y) = \frac{y^\alpha}{e^{\beta y}} \quad \text{und} \quad g(x) = \frac{e^{\delta x}}{x^\gamma}$$

zeigt, daß die Ableitung in beiden Fällen außer im Nullpunkt noch in genau einem weiteren Punkt verschwindet, nämlich dort wo x bzw. y gleich der entsprechenden Koordinate des nichttrivialen Gleichgewichtspunkts ist. f hat in diesem Punkt ein Maximum, g ein Minimum, und für $t \rightarrow 0$ oder $t \rightarrow \infty$ geht f gegen null und g gegen unendlich.

Da beide Funktionen im positiven Bereich der reellen Achse nur positive Werte annehmen, gibt es für eine vorgegebene positive Zahl c somit höchstens zwei Werte, an denen sie von f bzw. g angenommen wird; für zu große c , gibt es kein y mehr mit $f(y) = c$, und für zu kleine c kein x mit $g(x) = c$.

Daraus folgt nach kurzer Überlegung, daß die Lösungskurven (abgesehen von den beiden Fixpunkten) auf geschlossenen Kurven um den nichttrivialen Fixpunkt liegen. Da kein Punkt auf einer solchen Kurve Fixpunkt ist, kann keine Lösungskurve für $t \rightarrow \infty$ gegen einen Punkt einer solchen Kurve konvergieren, die Lösungskurven müssen also den nichttrivialen Gleichgewichtspunkt permanent umrunden.

Betrachten wir eine konkrete Lösung $(x(t), y(t))$ des Differentialgleichungssystems, das wir der Einfachheit halber kurz als

$$\begin{pmatrix} \dot{x}(t) \\ \dot{y}(t) \end{pmatrix} = F(x(t), y(t))$$

schreiben wollen, und fixieren wir einen Zeitpunkt t_0 ; für diesen sei $x(t_0) = a$ und $y(t_0) = b$. Dann muß es nach obiger Diskussion ein kleinste Zeitspanne T geben, so daß auch

$$x(t_0 + T) = a \quad \text{und} \quad y(t_0 + T) = b$$

ist. Für die beiden Funktionen

$$u(t) \stackrel{\text{def}}{=} x(t+T) \quad \text{und} \quad v(t) \stackrel{\text{def}}{=} y(t+T)$$

ist dann $u(t_0) = a$ und $v(t_0) = b$; außerdem ist

$$\begin{pmatrix} \dot{u}(t) \\ \dot{v}(t) \end{pmatrix} = \begin{pmatrix} \dot{x}(t+T) \\ \dot{y}(t+T) \end{pmatrix} = F(x(t+T), y(t+T)) = F(u(t), v(t)),$$

$(x(t), y(t))$ und $(u(t), v(t))$ lösen also dasselbe Anfangswertproblem. Falls wir zeigen können, daß F eine LIPSCHITZ-Bedingung erfüllt, müssen die beiden Funktionen also übereinstimmen.

Im betrachteten Beispiel ist

$$F(x, y) = \begin{pmatrix} \alpha x - \beta x y \\ -\gamma y + \delta x y \end{pmatrix},$$

also

$$\|F(x_1, y_1) - F(x_2, y_2)\| = \left\| \begin{pmatrix} \alpha(x_1 - x_2) - \beta(x_1 y_1 - x_2 y_2) \\ -\gamma(y_1 - y_2) + \delta(x_1 y_1 - x_2 y_2) \end{pmatrix} \right\|.$$

Im Quadrat $-R \leq x, y \leq R$ ist

$$|\alpha(x_1 - x_2) - \beta(x_1 y_1 - x_2 y_2)| \leq |\alpha(x_1 - x_2)| + |\beta(x_1 y_1 - x_2 y_2)|$$

und

$$\begin{aligned} |(x_1 y_1 - x_2 y_2)| &= |x_1 y_1 - x_1 y_2 + x_1 y_2 - x_2 y_2| \\ &= |x_1(y_1 - y_2) + y_2(x_1 - x_2)| \\ &\leq |x_1(y_1 - y_2)| + |y_2(x_1 - x_2)| \\ &\leq R(|y_1 - y_2| + |x_1 - x_2|), \end{aligned}$$

also ist

$$\begin{aligned} &|\alpha(x_1 - x_2) - \beta(x_1 y_1 - x_2 y_2)| \\ &\leq \alpha |x_1 - x_2| + \beta R(|y_1 - y_2| + |x_1 - x_2|). \end{aligned}$$

Analog folgt die Ungleichung

$$\begin{aligned} &|\gamma(y_1 - y_2) + \delta(x_1 y_1 - x_2 y_2)| \\ &\leq \gamma |y_1 - y_2| + \delta R(|y_1 - y_2| + |x_1 - x_2|). \end{aligned}$$

Wir arbeiten hier mit der Maximumnorm von Vektoren, $\left\| \begin{pmatrix} x \\ y \end{pmatrix} \right\|$ ist also das Maximum von $|x|$ und $|y|$, und entsprechend ist

$$\left\| \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} - \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} \right\| = \max\{|x_1 - x_2|, |y_1 - y_2|\}.$$

Mit

$$L = \max\{\alpha + 2\beta R, \gamma + 2\delta R\}$$

ist somit

$$\|F(x_1, y_1) - F(x_2, y_2)\| \leq L \left\| \begin{pmatrix} x_1 \\ y_1 \end{pmatrix} - \begin{pmatrix} x_2 \\ y_2 \end{pmatrix} \right\|;$$

F erfüllt also eine LIPSCHITZ-Bedingung, so daß wir aus dem Satz von PICARD-LINDELOF folgern können, daß das Anfangswertproblem in jedem abgeschlossenen Quadrat eindeutig lösbar ist.

Da jede Lösungskurve in einem abgeschlossenen Quadrat liegt (sonst müßte sie irgendwo gegen unendlich gehen), ist also

$$u(t) = x(t) \quad \text{und} \quad v(t) = y(t) \quad \text{für alle } t \geq t_0$$

d.h.

$$x(t+T) = x(t) \quad \text{und} \quad y(t+T) = y(t) \quad \text{für alle } t \geq t_0.$$

Damit wissen wir, daß alle Lösungsfunktionen periodisch sind.

In der unmittelbaren Umgebung der nichttrivialen Gleichgewichtslösung können wir sogar noch etwas mehr sagen: Durch TAYLOR-Entwicklung der oben betrachteten Funktionen f und g überzeugt man sich leicht davon, daß die Lösungskurven dort näherungsweise als Ellipsen aufgefaßt werden können.

Nach all diesen Vorbereitungen sollten wir uns endlich eine konkrete Lösungskurve anschauen, d.h. wir sollten das Problem in einem Spezialfall numerisch lösen.

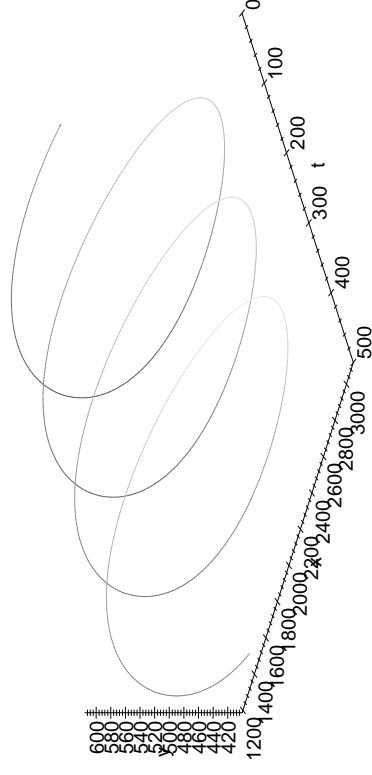


Abb. 40: Numerische Simulation des Raubtier-Beutetier-Systems

Abbildung 40 zeigt das Ergebnis; die spiralförmige Kurve entspricht genau unseren Erwartungen.

Besser können wir diese überprüfen, wenn wir eine Reihe von Lösungskurven in der xy -Ebene betrachten; Abbildung 41 zeigt solche Kurven zu verschiedenen Anfangsbedingungen. Abgesehen von der Gleichgewichtslösung, die einfach ein Punkt ist, sieht man die vorhergesagten geschlossenen Kurven, die das Gleichgewicht umrunden.

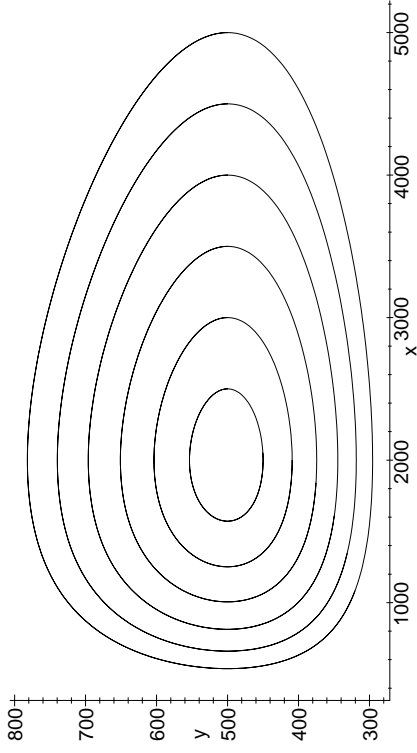


Abb. 41: Lösungskurven in der xy -Ebene

Man kann sich leicht klarmachen, was der Umlauf auf so einer Lösungskurve biologisch bedeutet: Im jeweils untersten Punkt ist die Raubtierpopulation minimal, so daß sich die Beutetiere stark vermehren können; dadurch verbessert sich die Nahrungsgrundlage für die Raubtiere, was nun auch zu deren Vermehrung führt, so daß die Kurve auf ihren am weitesten rechts gelegenen Punkt zusteuert, in dem die Beutetiere ihre maximale Populationsstärke erreichen. Die gestiegene Raubtierpopulation frisst nun, nachdem sie ihren Gleichgewichtswert überschritten hat, mehr Beutetiere als nachwachsen, kann sich aber wegen der großen Anzahl vorhandener Beutetiere weiterhin vermehren auf ein Maximum hin, das am obersten Punkt der Lösungskurve erreicht ist. Danach reicht der bereits gesunkene Bestand an Beutetieren nicht mehr aus als Nahrungsgrundlage für die Raubtiere, ihre Population geht also zurück, reicht aber immer noch aus, um die Beutetiere weiter zu dezimieren. Im Punkt links außen hat deren Bestand schließlich sein Minimum erreicht; die

weiterhin sinkende Raubtierpopulation frißt nun weniger Beutetiere als nachwachsen und leidet trotzdem weiter an Nahrungsmangel. Sobald sie ihr Minimum erreicht hat, schließt sich der Kreis, und der gleiche Zyklus beginnt von vorne.

Auch wenn Raubtiere und Beutetiere in der Technischen Informatik keine große Rolle spielen, sollten wir uns doch zumindest kurz fragen, ob die mathematische Lösung irgendetwas mit der biologischen Realität zu tun hat – der Zusammenhang zwischen idealisierten mathematischen Modellen und realen Systemen ist schließlich auch in der Technischen Informatik von Bedeutung.

Wie in vielen praktischen Anwendungen der Mathematik sind die Annahmen des Modells auch hier viel zu einfach: Es gibt kaum je zwei Arten, die völlig isoliert vom Rest der Welt leben. Trotzdem wurden die vorhergesagten Zyklen schon beobachtet: Die Hudson Bay Company sammelte rund hundert Jahre lang Daten über gekaufte Felle von Luchsen (Raubtieren) und Schneehasen (deren Beute); da die Trapper kaum beeinflussen können, was in ihre Fallen läuft, sollten diese Anzahlen ungefähr proportional sein zu den jeweiligen Populationszahlen. Abbildung 42 zeigt ungefähr die vorhergesagten zyklischen Schwankungen – auch wenn die Schwankungen zwischen 1875 und 1905 in die falsche Richtung gehen: Dort wird der Gleichgewichtspunkt nicht gegen den Uhrzeigersinn umrundet, sondern im Uhrzeigersinn.

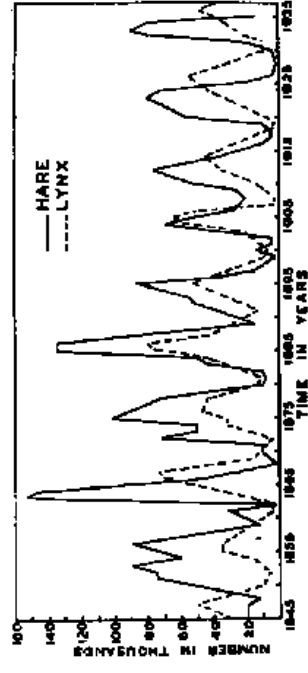


Abb. 42: Luchse und Schneehasen

Über den genauen Grund dafür gibt es immer noch viele Spekulationen; der Grund, warum Abweichungen vom Modell auftreten müssen ist aber einfach zu verstehen: Die Realität deutlich ist komplizierter als das extrem vereinfachte Modell von LOTKA und VOLTERRA.

Das vorliegende Beispiel ist natürlich, wie fast alle Beispiele in einer Anfängervorlesung, viel zu elementar: Im allgemeinen kann man einer Differentialgleichung nicht auf so einfache Weise so viele Eigenschaften der Lösungen ansehen.

Die qualitative Theorie der Differentialgleichungen wendet denn auch viele Methoden an, die weit jenseits des Stoffs dieser Vorlesung liegen, und selbst damit kann sie in komplizierteren Fällen nur deutlich weniger Information aus der Differentialgleichung extrahieren als in diesem Beispiel.

Einen ersten Überblick über das Verhalten der Lösungen einer *autonomen* Differentialgleichung in nur zwei Variablen liefert auch bei beliebig komplizierten Systemen ein graphisches Verfahren: Bei einer autonomen Differentialgleichung

$$\dot{\vec{y}}(t) = F(\vec{y}(t))$$

definiert die Funktion F ein Vektorfeld, wie wir es aus [HMI], Kapitel 2, kennen.

Spätestens an dieser Stelle wird klar, daß wir in diesem Kapitel bei der Unterscheidung von Punkten und Vektoren geschluppt haben: Für die linearen homogenen Differentialgleichungssysteme, die den Hauptinhalt dieses Kapitels bilden, war es völlig in Ordnung, nur von Vektoren zu reden: Dort gibt es einen wohldefinierten Nullpunkt, so daß sich Punkte und Vektoren in kanonischer Weise entsprechen.

Zur geometrischen Interpretation von $\dot{\vec{y}}(t) = F(\vec{y}(t))$ ist es aber sinnvoller, das Argument \vec{y} von F als Punkt \mathbf{y} aufzufassen und den Funktionswert als Tangentenvektor in diesem Punkt zu interpretieren. Eine Lösungskurve der Differentialgleichung ist dann eine Kurve, die in jedem Punkt \mathbf{y} den Vektor $F(\mathbf{y})$ als Tangentenvektor hat.

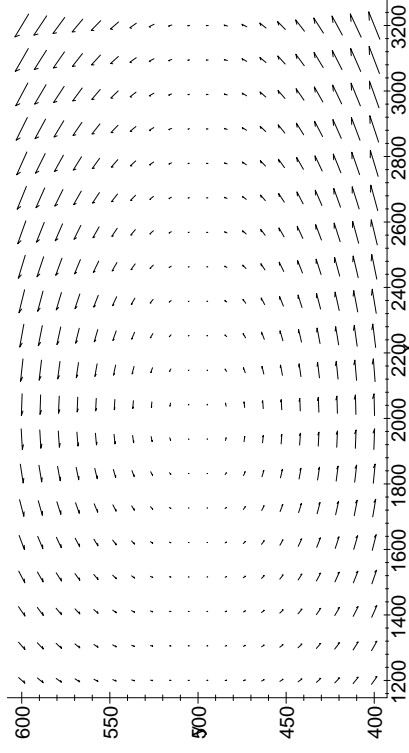


Abb. 43: Das Vektorfeld zur Raubtier-Beutefier-Gleichung

Abbildung 43 zeigt das Vektorfeld im betrachteten Beispiel; es legt zumindest die Vermutung nahe, daß die Lösungen zyklisch um einen Punkt rotieren. *Genau* können das wir freilich aufgrund der graphischen Information nicht sagen: Eine visuell nicht wahrnehmbare Richtungsänderung der Vektoren gehört zu einer Lösung die sich spiralförmig auf den Gleichgewichtspunkt zusammenzieht oder aber spiralförmig ins Unendliche geht.

Schon bei der Visualisierung von Vektorfeldern haben wir gesehen, daß es gelegentlich übersichtlicher ist, auf die Längeninformation zu verzichten und nur die Richtung zu betrachten. Bei Differentialgleichungen, bei denen es bei einer graphischen Lösung praktisch nur auf die *Richtung* des Vektorfelds in jedem Punkt ankommt, gilt dies umso mehr; oft versucht man daher die Lösungskurve durch ein auf Einheitslänge normiertes Vektorfeld zu führen. Abbildung 44 zeigt, wie dies im vorliegenden Beispiel aussieht.

g) Stabilitätsfragen

Die Beschreibung eines realen Systems durch ein mathematisches Modell wie ein Differentialgleichungssystem ist abgesehen von einigen ganz einfachen Fällen immer mit einer Idealisierung verbunden; das reale System verhält sich daher nicht *exakt* so wie das Modell es vorhersagt.

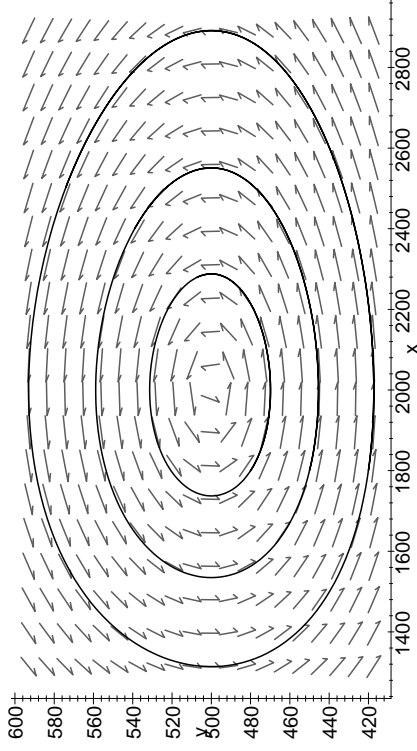


Abb. 44: Anpassung von Lösungskurven an das Vektorfeld

Auch die Anfangsbedingungen des Modells, die dem Zustand des realen Systems zu einem vorgegebenen Anfangszeitpunkt entsprechen, lassen sich nur durch fehlerbehaftete Messungen bestimmen. Hinzu kommt, daß man bei der Auswertung des mathematischen Modells nur selten wirklich mit realen Zahlen rechnet; meistens rechnet man per Computer und somit (falls man keine sehr spezialisierte Mathematiksoftware benutzt) mit rundungsfehlerbehafteten Gleitkommaoperationen.

Aus einem mathematischen Modell abgeleiteten Aussagen können daher nur dann nützlich für die Vorhersage von realen Systemen sein, wenn sie stabil sind gegenüber kleineren Änderungen von Koeffizienten und Anfangsbedingungen.

Betrachten wir dazu ein Beispiel: Eine Größe $y(t)$ sei beschrieben durch das Anfangswertproblem

$$\dot{y}(t) = y(t) + 2e^{-t} \quad \text{mit} \quad y(0) = 1.$$

Wie man sich sofort durch Einsetzen überzeugt, ist $y(t) = e^{-t}$ eine Lösung. Die rechte Seite

$$F(y, t) = y - e^{-t}$$

genügt offensichtlich auf ganz \mathbb{R}^2 einer LIPSCHITZ-Bedingung mit Kon-

stante eins, denn $F_y(y, t) \equiv 1$, und man sieht auch direkt, daß

$$|F(y_1, t) - F(y_2, t)| = |y_1 - y_2| \leq 1 \cdot |y_1 - y_2|$$

ist. Somit ist $y(t) = e^{-t}$ die *einzige* Lösung des Anfangswertproblems. Trotzdem ist diese Lösung für alle praktischen Zwecke völlig wertlos:

$$\dot{y}(t) = y(t) - 2e^{-t}$$

ist eine inhomogene lineare Differentialgleichung, deren zugehörige homogene Gleichung

$$\dot{y}(t) = y(t)$$

die allgemeine Lösung

$$y(t) = \lambda e^t$$

hat. Die allgemeine Lösung der inhomogenen Gleichung ist daher

$$y(t) = e^{-t} + \lambda e^t \quad \text{und} \quad y(0) = 1 + \lambda.$$

Sobald also die Anfangsbedingung auch nur minimal gestört wird, geht die Lösung für große t nicht mehr gegen null, sondern je nach Vorzeichen von λ gegen $\pm\infty$. Bei einem solchen *schlecht gestellten* oder *strukturell instabilen* Problem läßt sich also mathematisch nichts vorhersagen.

Auch eine numerische Lösung des Anfangswertproblems wird wegen allfälliger Rundungsfehler über kurz oder lang die exakte Lösungskurve $y(t) = e^{-t}$ verlassen und auf eine der zumindest anfänglich benachbarten anderen Kurven überwechseln, so daß auch sie für $t \rightarrow \infty$ divergiert. Abbildung 45 zeigt eine mit einem RUNGE-KUTTA-Verfahren der Ordnung vier/fünf berechnete numerische Lösung; wie man sieht, hat sie ab etwa $t = 14$ nichts mehr mit der korrekten Lösung zu tun.

Anders sieht es aus beim Anfangswertproblem

$$\dot{y}(t) = -y(t) + 1 \quad \text{mit} \quad y(0) = 1.$$

Hier überzeugt man sich leicht, daß $y(t) = 1$ die einzige Lösung ist, aber jetzt hat die zugehörige homogene Differentialgleichung die allgemeine Lösung $y(t) = \lambda e^{-t}$; die allgemeine Lösung der Differentialgleichung

$$\dot{y}(t) = -y(t) + 1$$

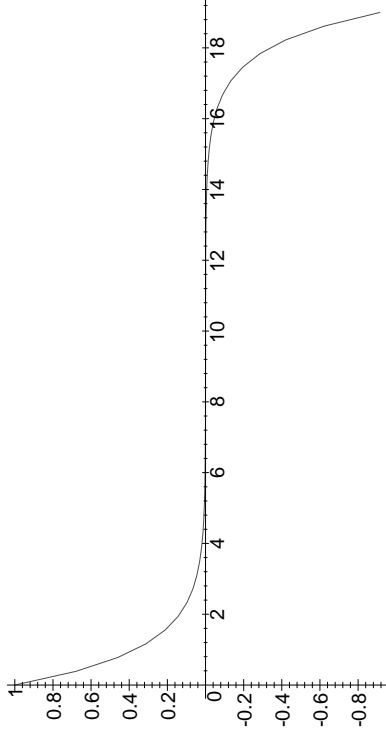


Abb. 45: Divergenz einer numerisch berechneten Lösungskurve

ist also

$$y(t) = 1 + \lambda e^{-t} \quad \text{mit} \quad y(0) = 1 + \lambda.$$

Kleine Störungen λ des Anfangswerts werden jetzt durch den Faktor e^{-t} weggedämpft; unabhängig von solchen Störungen bleibt also die Lösung $y \equiv 1$ stabil.

Um allgemeiner zu sehen, was in der Umgebung eines Gleichgewichts passieren kann, versuchen wir, die Gleichung in der Umgebung des Gleichgewichts anzunähern durch die einzige Klasse von Differentialgleichungen, die wir wirklich beherrschen, die linearen homogenen Differentialgleichungen mit konstanten Koeffizienten.

Dazu erinnern wir uns an die Definition einer differenzierbaren Funktion mehrerer Veränderlicher: $F: \mathbb{R}^n \rightarrow \mathbb{R}_n$ ist im Punkt $\mathbf{x} \in \mathbb{R}^n$ differenzierbar, wenn in einer Umgebung des Punktes gilt

$$F(\mathbf{x} + \vec{h}) = F(\mathbf{x}) + J_F(\mathbf{x})\vec{h} + o(\|\vec{h}\|),$$

wobei $J_F(\mathbf{x})$ die JACOBI-Matrix von F in \mathbf{x} ist.

Wir betrachten ein autonomes Differentialgleichungssystem

$$\dot{\vec{y}}(t) = F(\vec{y}(t))$$

mit einer differenzierbaren Funktion F mit Fixpunkt \vec{y}_0 . In der Umgebung des Fixpunkts ist dann

$$F(\vec{y}_0 + \vec{h}) = F(\vec{y}_0) + J_F(\vec{y}_0)\vec{h} + o(\|\vec{h}\|) = \vec{y}_0 + J_F(\vec{y}_0)\vec{h} + o(\|\vec{h}\|);$$

falls wir den Fehlerterm $o(\|\vec{h}\|)$ vernachlässigen, genügt also die Differenz

$$\vec{u}(t) \stackrel{\text{def}}{=} \vec{y}(t) - \vec{y}_0$$

zwischen der Gleichgewichtslösung und einer nahe benachbarten Lösung näherungsweise einer linearen homogenen Differentialgleichung

$$\dot{\vec{u}}(t) = J_F(\vec{y}_0)\vec{u}(t).$$

Das Langzeitverhalten von deren Lösungen hängt, wie wir aus §3c) wissen, von den Eigenwerten der Matrix $J_F(\vec{y}_0)$ ab: Falls diese allesamt negativen Realteil haben, konvergiert jede Lösung \vec{u} für $t \rightarrow \infty$ gegen den Nullpunkt; falls alle Eigenwerte positiven Realteil haben, divergiert jede Lösung außer der Null ins Unendliche. Im ersten Fall sprechen wir von einem *stabilen* oder *anziehenden* Fixpunkt, im zweiten von einem *instabilen* oder *abstoßenden*. Falls einige Eigenwerte positiven und andere negativen Realteil haben, reden wir von einem *Sattelpunkt*; hier hängt es von der Richtung ab, ob eine Störung weggedämpft wird oder nicht, allerdings wird in der Praxis fast jede Störung zur Divergenz führen, denn nur im linearen Unterraum, der von den Eigenvektoren zu den Eigenwerten mit negativem Realteil aufgespannt wird, werden die Störungen weggedämpft. Eine zufällige Störung wird aber meist sehr schnell aus diesem Unterraum herausführen, so daß dann auch Eigenwerte mit positivem Realteil eine Rolle spielen.

Bei Eigenwerten mit Realteil null reicht die Linearisierung nicht aus, um zu Aussagen über das Stabilitätsverhalten zu kommen, da dann die Terme höherer Ordnung das Geschehen dominieren. Es kann sehr schwer sein, in so einem Fall die Dynamik vorherzusagen.

Abbildung 46 zeigt das Vektorfeld in der Nähe eines stabilen Fixpunkts; alle Lösungskurven laufen auf diesen Punkt zu. Bei einem instabilen Fixpunkt hätten wir dasselbe Bild, nur daß dann alle Pfeile in Gegenrichtung zeigen würden.

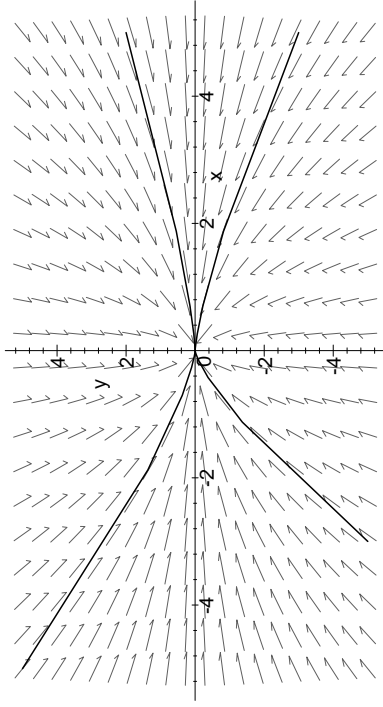


Abb. 46: Die Umgebung eines stabilen Fixpunkts

Auch in Abbildung 47 ist ein stabiler Fixpunkt zu sehen; hier hat aber die JACOBI-Matrix zwei konjugiert komplexe Eigenwerte, so daß sich benachbarte Lösungen spiralförmig auf den Fixpunkt zusammenziehen. Bei einem abstoßenden Fixpunkt hätten wir wieder im wesentlichen dasselbe Bild, aber mit umgedrehten Pfeilen.

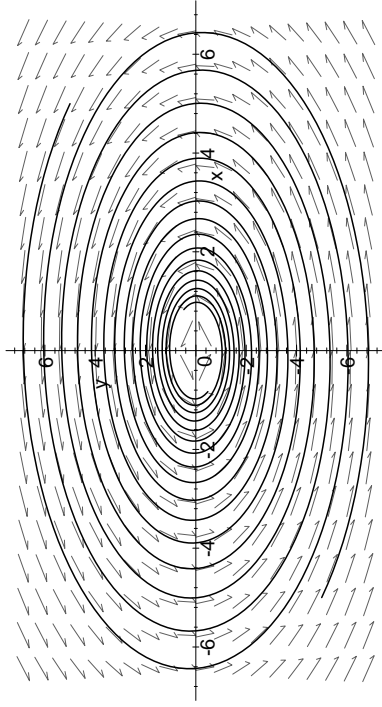


Abb. 47: Zwei konjugiert komplexe Eigenwerte mit negativem Realteil

Abbildung 48 zeigt die Umgebung eines Sattelpunkts; hier haben wir Lösungskurven, die sich zwar asymptotisch der y -Achse annähern, auf dieser aber gegen plus oder minus unendlich gehen.

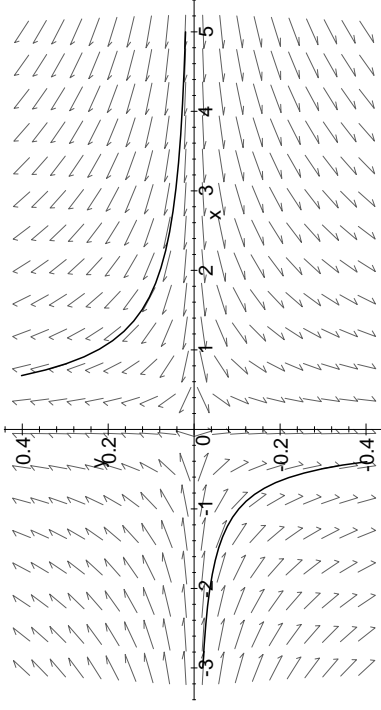


Abb. 48: Umgebung eines Sattelpunkts

Als Beispiel wollen wir die Lösungen der LORENZ-Gleichungen

$$\dot{x}(t) = p(y(t) - x(t))$$

$$\dot{y}(t) = r x(t) - y(t) - x(t)z(t)$$

$$\dot{z}(t) = -bz(t) + x(t)y(t)$$

untersuchen. Dieses Differentialgleichungssystem ist eine extreme Vereinfachung der sogenannten NAVIER-STOKES-Gleichung, einer partiellen Differentialgleichung, die Strömungsphänomene beschreibt. Für technische Informatiker interessanter ist wohl, daß dasselbe System nach HAKEN (Phys. Lett. A53 (1975), 77–78) auch das Verhalten von Lasern beschreiben kann.

Die Funktionen $x(t)$, $y(t)$ und $z(t)$ verlieren im Vereinfachungsprozeß ihre unmittelbare physikalische Bedeutung; die Parameter lassen sich allerdings physikalisch interpretieren: Für die atmosphärische Konvektion ist nach LORENZ

$$p = 10, \quad r = 28 \quad \text{und} \quad b = \frac{8}{3}$$

eine sinnvolle Wahl.



EDWARD NORTON LORENZ wurde 1917 im US-Bundesstaat Connecticut geboren; er studierte Mathematik in Dartmouth College (A.B. 1938) und Harvard (M.A. 1940). Nach seinem Kriegsdienst ging er ans MIT, wo er 1948 über Meteorologie promovierte. Sowohl dem MIT, wo er 1981 als Professor emeritiert wurde, als auch der Meteorologie blieb er fortan treu. Zu seinen vielen Auszeichnungen gehört unter anderem der Kyoto-Preis von 1991, der wohl höchstdotierte Wissenschaftspreis.

Da uns der erste Augenschein nichts über das Verhalten der Lösungen zeigt, empfiehlt es sich, daß wir uns durch numerische Simulation einen ersten Eindruck verschaffen. Abbildung 49 zeigt die Lösung des Anfangswertproblems mit $x(0) = 2$ und $y(0) = z(0) = 10$; die meisten werden ähnliche Bilder wohl schon gesehen haben.

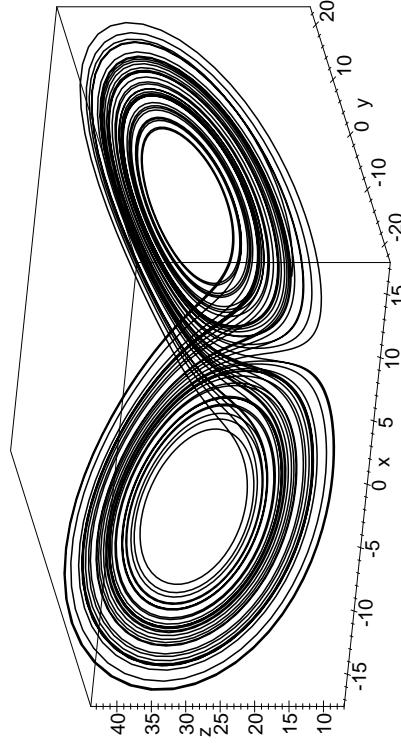


Abb. 49: Eine Bahnkurve des LORENZ-Systems

Leider ist dieses Bild einerseits etwas unübersichtlich, andererseits zeigt es nur eine einzige Lösungskurve. Um besser zu verstehen, was hier passiert, beschränken wir uns auf die Funktion $x(t)$ und betrachten diese für zwei Lösungskurven; Abbildung 50 zeigt die für die Anfangsbedingungen

$$x(0) = 2, \quad y(0) = z(0) = 10 \quad \text{und} \quad x(0) = 2,01, \quad y(0) = z(0) = 10.$$

Wie man sieht, sind die beiden Lösungskurven bis etwa zum Zeitpunkt $t = 6,5$ praktisch ununterscheidbar, danach gehen sie aber recht schnell auseinander und haben ab etwa $t = 10$ nichts mehr miteinander zu tun.

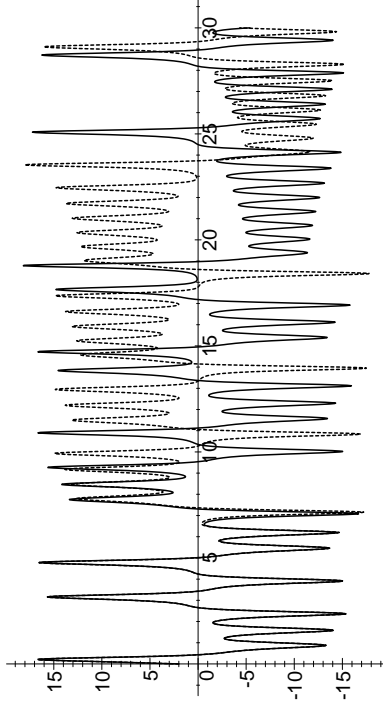


Abb. 50: Die x -Koordinaten zweier Lösungen mit benachbarten Anfangswerten

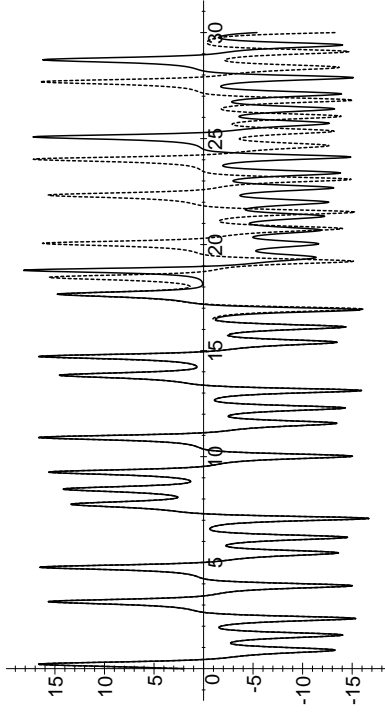
Fast das gleiche Bild ergibt sich, wenn wir die gestrichelte Kurve nicht mit $x(0) = 2,01$ anfangen lassen, sondern – bei sonst unveränderten Werten – bei

$$x(0) = 2,0000001 = 2 + 10^{-6}.$$

Jetzt sind die Kurven zwar bis etwa $t = 16$ praktisch ununterscheidbar, aber spätestens ab etwa $t = 20$ haben sie auch hier nichts mehr miteinander zu tun.

Gerade beim zweiten Fall sollte uns das zu denken geben: Wenn wir Differentialgleichungen zur Vorhersage benutzen, stammen die Anfangsbedingungen im allgemeinen aus einer Messung. Man kann aber nur selten so genau messen, daß sich die beiden Werte 2 und 2,000001 unterscheiden ließen; um eine sinnvolle Voraussage über das Verhalten der Lösung zum Zeitpunkt $t = 20$ zu machen, muß man aber nach Abbildung 51 den Wert $x(0)$ mit dieser Genauigkeit kennen.

Es kommt noch schlimmer: In Abbildung 52 ist die dick ausgezogene Kurve wieder eine numerische Simulation der Lösung zu den Anfangs-

Abb. 51: Effekt einer Störung des Anfangswerts um 10^{-6}

bedingungen $x(0) = 2$ und $y(0) = z(0) = 10$, die gestrichelte Kurve allerdings auch! Die beiden Kurven unterscheiden sich nur dadurch, daß die numerische Simulation bei der dick ausgezogenen Kurve (wie auch bei allen anderen bisherigen Kurven) mit Schrittweite 0,02 arbeitete, wohingegen die Schrittweite für die gestrichelte Kurve mit 0,01 nur halb so groß war. Auch das reicht schon, daß die Kurven ab etwa $t = 10$ nichts mehr miteinander zu tun haben, und damit dürfte wohl auch klar sein, daß keine der bislang betrachteten Kurven für größere Werte von t irgendetwas mit der „wahren“ Lösungsfunktion $x(t)$ zu tun hatte.

LORENZ mußte dieses Verhalten der Lösungen auf die harte Weise lernen: Er fand zu seinem großen Erstaunen, daß sich seine Rechenergebnisse nicht reproduzieren ließen, wenn er Zwischenergebnisse für Kontrollrechnungen nur in gerundeter Form eintippte: Der geringe Rundungsfehler bei der Eingabe des Startwerts reichte bereits aus, um das Langzeitverhalten des Systems grundlegend zu verändern.

Falls das gleiche Phänomen auch in der wirklichen Atmosphäre auftritt, können also minimale Veränderungen etwa des Luftdrucks oder der Temperatur auf längere Sicht zu einer dramatisch anderen Entwicklung des Wetters führen – eine Idee, die vielen Meteorologen damals als zu phantastisch erschien um ernstgenommen zu werden: Am 22. Januar

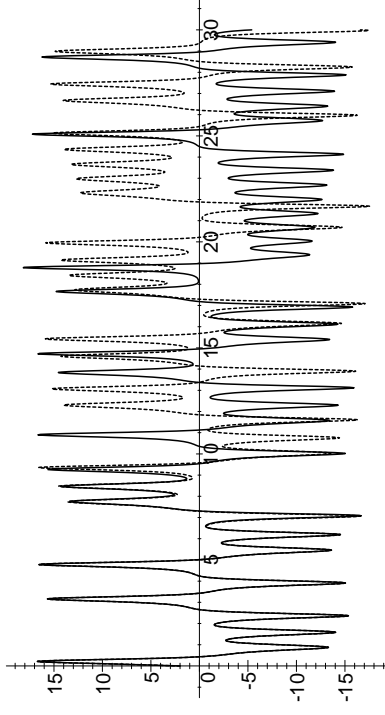


Abb. 52: Effekt einer Schrittweithalbierung bei der numerischen Simulation

1963 berichtete LORENZ vor der New York Academy of Sciences über seine Ergebnisse (*Trans. N.Y. Acad. Sci.* **25** (1963), 409–432) und schloß seinen Vortrag mit den Worten:

Als die Instabilität eines gleichförmigen Flusses gegenüber infinitesimalen Störungen erstmals als Erklärung für das Auftreten von Zykklonen und Antizyklonen in der Atmosphäre vorgeschlagen wurde, war diese Idee nicht allgemein akzeptiert. Ein Meteorologe bemerkte, daß, falls die Theorie korrekt wäre, ein Flügelschlag einer Möwe ausreichen würde, um die Entwicklung des Wetters für immer zu verändern. Die Kontrolle ist noch nicht entschieden, aber die neueste Evidenz scheint für die Möwen zu sprechen.

Inzwischen ist der Sieg der Möwen bekanntlich allgemein anerkannt; man fordert sogar nicht einmal mehr den relativ kräftigen Flügelschlag einer Möwe, um das Wetter permanent zu verändern: Im Dezember 1972 hielt LORENZ vor der American Association for the Advancement of Sciences in Washington, DC, einen Vortrag mit dem Titel *Predictability: Does the Flap of a Butterfly's Wings in Brazil set off a Tornado in Texas*, und seitdem geht das Wort vom *Schmetterlingseffekt* um die Welt.

Auch das Wort *Chaos* wird heute meist auf diese Weise definiert: Kleinste Änderungen bei den Anfangsbedingungen führen zu dramatischen

Veränderungen des Langzeitverhaltens. Allerdings muß man hier aufpassen: Bei der Differentialgleichung

$$\dot{y}(t) = y(t)$$

mit Anfangsbedingungen

$$y(0) = 1 \quad \text{und} \quad y'(0) = 1 + \varepsilon$$

unterscheiden sich die Lösungen $y(t) = e^t$ und $y(t) = (1 + \varepsilon) \cdot e^t$ zur Zeit t um $\varepsilon \cdot e^t$, was auch bei kleinsten ε -Werten sehr schnell eine sehr große Zahl wird: bei $\varepsilon = 10^{-6}$ und $t = 50$ etwa ist die Differenz bereits größer als $5 \cdot 10^{15}$. Trotzdem wird hier niemand von Chaos reden, denn beide Lösungen gehen sehr schnell gegen den „Gleichgewichtspunkt“ unendlich. Von „echtem“ Chaos verlangt man daher auch noch, daß die Lösungen nicht gegen eine (endliche oder unendliche) Gleichgewichtslösung konvergieren und auch nicht gegen eine periodische Lösung. Chaos in diesem Sinne ist sehr schwer nachzuweisen; für die LORENZ-Gleichung mit den klassischen Parameterwerten wurde erst Ende September 1999 ein *preprint* veröffentlicht, in dem dies (mit großem theoretischen wie auch rechnerischem Aufwand) gezeigt wird; siehe <http://www.math.gatech.edu/~mischalk/papers/tor3.ps>.

Chaos heißt nun allerdings nicht, daß wir dann überhaupt nichts über das Verhalten der Lösungen aussagen können. Beispielsweise können wir bereits mit unseren einfachen Mitteln zeigen, daß das Bild in Abbildung 49 zumindest qualitativ das Verhalten der Lösungskurven korrekt wiedergibt – quantitativ ist natürlich ab spätestens etwa $t = 10$ alles falsch.

Dazu berechnen wir zunächst die Gleichgewichtslösungen: Im Gleichungssystem

$$\begin{aligned} 0 &= p(y - x) \\ 0 &= rx - y - xz \\ 0 &= -bz + xy \end{aligned}$$

zeigt die erste Gleichung, falls wir den uninteressanten Fall $p = 0$ ausschließen, daß die x -Koordinate und die y -Koordinate eines jeden Fixpunkts übereinstimmen müssen.

Falls beide Koordinaten verschwinden, zeigt die dritte Gleichung ($b \neq 0$ vorausgesetzt), daß dann auch die z -Koordinate verschwinden muß; Einsetzen in die Gleichungen zeigt, daß der Nullpunkt in der Tat ein Fixpunkt ist.

Im Fall $x \neq 0$ können wir y in der zweiten Gleichung durch x ersetzen und dann durch x dividieren; dies ergibt die z -Koordinate

$$z = r - 1.$$

Damit zeigt die dritte Gleichung, daß es für $r \neq 1$ noch zwei weitere Fixpunkte gibt mit

$$x = y = \pm \sqrt{b(r-1)} \quad \text{und} \quad z = r - 1.$$

Die Untersuchung des relativ uninteressanten Nullpunkts sei dem Leser als Übungsaufgabe überlassen; hier seien nur die beiden anderen Fixpunkten betrachtet. Für

$$x = y = \pm \sqrt{b(r-1)} \quad \text{und} \quad z = r - 1$$

führen wir, wie oben im allgemeinen Fall, neue Variablen u , v und w ein, die den Abstand zum Fixpunkt beschreiben, d.h.

$$x = \pm \sqrt{b(r-1)} + u, \quad y = \pm \sqrt{b(r-1)} + v \quad \text{und} \quad z = r - 1 + w.$$

Zur Linearisierung in der Nähe des Gleichgewichtspunkt vernachlässigen wir alle nichtlinearen Terme in $u(t)$, $v(t)$ und $w(t)$; das entstehende lineare Differentialgleichungssystem hat dann die JACOBI-Matrix im Fixpunkt als Matrix, ist also

$$\begin{pmatrix} \dot{u}(t) \\ \dot{v}(t) \\ \dot{w}(t) \end{pmatrix} = A \begin{pmatrix} u(t) \\ v(t) \\ w(t) \end{pmatrix}$$

mit

$$A = \begin{pmatrix} -p & p & 0 \\ 1 & -1 & \pm \sqrt{b(r-1)} \\ \pm \sqrt{b(r-1)} & \pm \sqrt{b(r-1)} & -b \end{pmatrix}.$$

Das charakteristische Polynom

$$\det(A - \lambda E) = -\lambda^3 - (b+1+p)\lambda^2 - b(r-p)\lambda - 2pb(r-1)$$

ist für beide Fixpunkte dasselbe, verleiht aber nicht dazu, es allgemein lösen zu wollen. Wir setzen daher die von LORENZ vorgeschlagenen speziellen Parameterwerte ein und erhalten

$$-\lambda^3 - \frac{41}{3}\lambda^2 - \frac{304}{3}\lambda - 1440,$$

was immer noch so schlimm ist, daß wir es besser numerisch lösen. Die drei Lösungen ergeben sich näherungsweise als

$$-13,85457791 \quad \text{und} \quad 0,09395562396 \pm 10,19450522i.$$

Es gibt also einen negativen Eigenwert und zwei konjugiert komplexe Eigenwerte mit positivem Realteil. Damit ist klar, wie Lösungskurven des linearisierten Systems in der Nähe der beiden Fixpunkte aussehen: Der negative Eigenwert sorgt dafür, daß die Lösungen asymptotisch in die Ebene gedrückt werden, die von den Eigenvektoren zu den beiden anderen Eigenwerten aufgespannt wird, und die beiden komplexen Eigenwerte sorgen dafür, daß sie dort spiralförmig nach außen gehen.

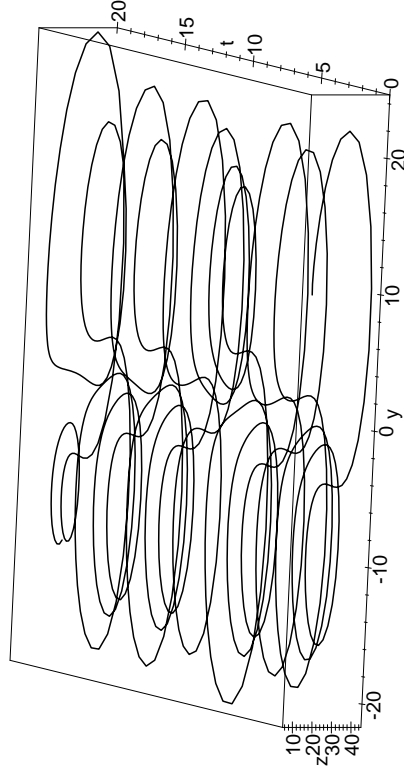


Abb. 53: y - und z -Koordinate als Funktion von t

Damit wird das Verhalten des LORENZ-Systems klar: Wir haben zwei Ebenen, die jeweils einen Sattelpunkt enthalten; kommt eine Lösung in die Nähe eines solchen Sattelpunkts, wird sie von der entsprechenden

Ebenen eingefangen und geht dort spiralförmig nach außen. Wenn sie sich hinreichend weit vom Sattelpunkt entfernt hat, sind die Voraussetzungen für die obige Linearisierung nicht mehr gegeben; die Lösung kann daher der Ebenen entkommen, wird aber über kurz oder lang von der Ebenen des anderen Sattelpunkts eingefangen und so weiter. Abbildung 53 zeigt dieses Verhalten etwas klarer als Abbildung 49: Hier sind die die y - und die z -Koordinate der Lösungskurve über der Zeit aufgetragen.

Für zweimal stetig differenzierbare Funktionen gibt es bekanntlich auch eine hinreichende Bedingung sowie die Möglichkeit, Maxima und Minima voneinander zu unterscheiden: Falls $f'(x_0)$ verschwindet und $f''(x_0)$ negativ ist, hat f im Punkt x_0 ein Maximum; bei positivem $f''(x_0)$ liegt ein Minimum vor. Auch hier folgt alles sofort aus der Definition der zweimaligen Differenzierbarkeit: Wegen

$$\begin{aligned} f(x_0 + h) &= f(x_0) + h f'(x_0) + \frac{h^2}{2} f''(x_0) + o(h^2) \\ &= f(x_0) + \frac{h^2}{2} f''(x_0) + o(h^2) \end{aligned}$$

sieht der Graph von f in diesen Fällen in der unmittelbaren Umgebung von x_0 aus wie eine nach unten bzw. oben geöffnete Parabel.

b) Verallgemeinerung aufs Mehrdimensionale

Nun betrachten wir eine stetig differenzierbare Funktion $f: D \rightarrow \mathbb{R}$ auf einer offenen Teilmenge $D \subset \mathbb{R}^n$. Dann bedeutet Differenzierbarkeit bekanntlich, daß es in jedem Punkt $\mathbf{x}_0 \in D$ einen Vektor

$$\nabla f(\mathbf{x}_0) = \text{grad } f(\mathbf{x}_0) \in \mathbb{R}^n$$

gibt, den Gradienten, so daß für hinreichend kleine Vektoren $\vec{h} \in \mathbb{R}^n$ gilt

$$f(\mathbf{x}_0 + \vec{h}) = f(\mathbf{x}_0) + \text{grad } f(\mathbf{x}_0) \cdot \vec{h} + o(|\vec{h}|).$$

Hier muß also für jeden Extremwert $\text{grad } f(\mathbf{x}_0)$ gleich dem Nullvektor sein, denn setzt man für \vec{h} ein kleines Vielfaches $t \cdot \text{grad } f(\mathbf{x}_0)$ des Gradienten ein, wäre sonst

$$f(\mathbf{x}_0 + \vec{h}) = f(\mathbf{x}_0) + t(\text{grad } f(\mathbf{x}_0) \cdot \text{grad } f(\mathbf{x}_0)) + o(|\vec{h}|)$$

für kleine positive t größer als $f(\mathbf{x}_0)$ und für kleine negative t kleiner.

Die Frage, welche Nullstellen des Gradienten wirklich Extremwerten entsprechen, ist schwieriger; in der Praxis wird es oft am einfachsten sein, sich die Umgebung des betreffenden Punktes mit irgendwelchen *ad hoc*-Methoden genauer anzusehen und dann zu entscheiden.

Klassisches Beispiel eines Punktes, in dem der Gradient verschwindet, ohne daß ein Extremwert vorliegt, ist der in Abbildung 54 gezeigte

Kapitel 5 Optimierung, Fehlerrechnung und Statistik

In der Schule werden Ableitungen hauptsächlich benutzt, um die Extremwerte einer Funktion zu bestimmen; ein Gesichtspunkt, der im letzten Semester bei der Differentialrechnung mehrerer Veränderlicher keine Rolle spielte. In diesem letzten Kapitel der Vorlesung soll dies nachgeholt werden, wobei insbesondere die Anwendungen auf die Fehler- und Ausgleichsrechnung wichtige Beispiele liefern. Zu deren besseren Verständnis sollen auch einige Grundbegriffe der Statistik erörtert werden.

§1: Extrema von Funktionen mehrerer Veränderlicher

a) Der eindimensionale Fall

Erinnern wir uns an die Schule: Wenn die stetig differenzierbare Funktion $f: (a, b) \rightarrow \mathbb{R}$ im Punkt $x_0 \in (a, b)$ ein Extremum annimmt, verschwindet dort die Ableitung $f'(x_0)$. Der Grund ist klar: Nach Definition der Differenzierbarkeit ist

$$f(x_0 + h) = f(x_0) + h f'(x_0) + o(h);$$

falls $f'(x_0)$ nicht verschwindet, ist $f(x_0 + h)$ für kleine h mit demselben Vorzeichen wie $f'(x_0)$ größer und für solche mit entgegengesetztem Vorzeichen kleiner als $f(x_0)$. In x_0 kann f somit weder ein Maximum noch ein Minimum annehmen.

Die Umkehrung gilt nicht: Standardbeispiel ist die Funktion $f(x) = x^3$, für die $f'(0)$ verschwindet, ohne daß im Nullpunkt ein Maximum oder Minimum wäre.

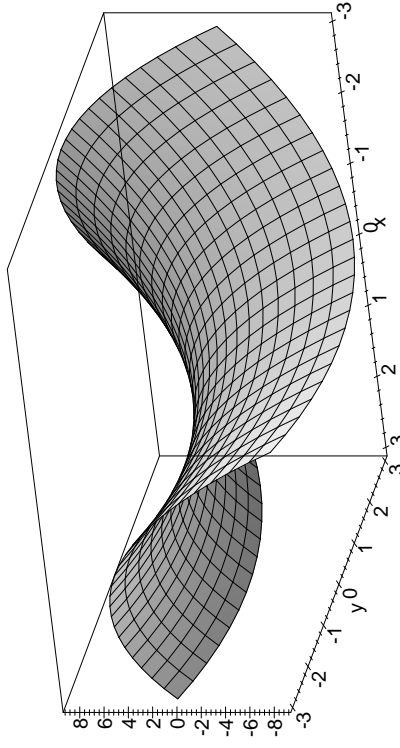


Abb. 54: Graph der Funktion $f(x, y) = x^2 - y^2$

Sattelpunkt, hier dargestellt als Funktionswert über dem Punkt $(0, 0)$ für die Funktion $f(x, y) = x^2 - y^2$.

Für zweifach stetig differenzierbare Funktionen kann man genau wie im eindimensionalen Fall ein hinreichendes Kriterium finden, das nur von der zweiten Ableitung im Punkt \mathbf{x}_0 abhängt:

Die zweite Ableitung von $f \in C^2(D, \mathbb{R})$ im Punkt $\mathbf{x}_0 \in D$ ist bekanntlich gegeben durch die HESSE-Matrix

$$H_f(\mathbf{x}_0) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

und zweimalige Differenzierbarkeit bedeutet, daß

$$f(\mathbf{x}_0 + \vec{h}) = f(\mathbf{x}_0) + \text{grad } f(\mathbf{x}_0) \cdot \vec{h} + \frac{1}{2} \vec{h}^T H_f(\mathbf{x}_0) \vec{h} + o(|\vec{h}|^2)$$

ist für kleine \vec{h} .

Wenn $\text{grad } f(\mathbf{x}_0)$ verschwindet, hängt also das Verhalten von f in der Umgebung von \mathbf{x}_0 ab von der quadratischen Form

$$\vec{h} \mapsto \vec{h}^T H_f(\mathbf{x}_0) \vec{h}.$$

Definition: a) Eine Matrix $A \in \mathbb{R}^{n \times n}$ heißt *positiv definit*, wenn für alle Vektoren $\vec{v} \neq \vec{0}$ aus \mathbb{R}^n gilt:

$${}^t \vec{v} A(\mathbf{x}_0) \vec{v} > 0.$$

b) A heißt *negativ definit*, wenn für alle $\vec{v} \neq \vec{0}$ aus \mathbb{R}^n gilt:

$${}^t \vec{v} A(\mathbf{x}_0) \vec{v} < 0.$$

c) A heißt *indefinit*, wenn es Vektoren $\vec{v}, \vec{w} \in \mathbb{R}^n$ gibt mit

$${}^t \vec{v} A(\mathbf{x}_0) \vec{v} > 0 \quad \text{und} \quad {}^t \vec{w} A(\mathbf{x}_0) \vec{w} < 0.$$

Mit dieser Terminologie ist das folgende Lemma klar:

Lemma: Wenn die differenzierbare Funktion $f \in C^1(D, \mathbb{R})$ im Punkt $\mathbf{x}_0 \in D$ ein lokales Extremum hat, ist dort ihr Gradient gleich dem Nullvektor.

Falls umgekehrt für $f \in C^2(D, \mathbb{R})$ der Gradient im Punkt $\mathbf{x} \in D$ verschwindet, gilt:

- a) Falls die HESSE-Matrix $H_f(\mathbf{x}_0)$ positiv definit ist, hat f im Punkt \mathbf{x}_0 ein Minimum.
- b) Falls $H_f(\mathbf{x}_0)$ negativ definit ist, hat f im Punkt \mathbf{x}_0 ein Maximum.
- c) Falls $H_f(\mathbf{x}_0)$ indefinit ist, hat f im Punkt \mathbf{x}_0 kein Extremum. ■

Damit uns das etwas nützt, brauchen wir jetzt nur noch ein Kriterium, mit dem wir feststellen können, welche Definitheitseigenschaften die HESSE-Matrix hat. Dazu erinnern wir uns daran, daß die HESSE-Matrix symmetrisch ist, und daß nach Kapitel 4, §2d) jede symmetrische Matrix diagonalisierbar ist.

Für eine Diagonalmatrix A mit Einträgen $\lambda_1, \dots, \lambda_n$ und einen Vektor \vec{v} mit Komponenten v_1, \dots, v_n wird obige quadratische Form zu

$$(v_1, v_2, \dots, v_n) \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} = \lambda_1 v_1^2 + \dots + \lambda_n v_n^2;$$

eine Diagonalmatrix ist also genau dann positiv definit, wenn alle Diagonaleinträge positiv sind und genau dann negativ definit, wenn sie alle

negativ sind. Falls es sowohl positive als auch negative Diagonaleinträge gibt, ist die Matrix indefinit.

Nun ist es für den Wertebereich einer Funktion irrelevant, bezüglich welches Koordinatensystems wir die Argumente ausdrücken; wir können eine symmetrische Matrix also bezüglich einer Basis aus Eigenvektoren betrachten, wo sie zur Diagonalmatrix wird mit den Eigenwerten als Einträgen. Daher gilt:

Lemma: Eine symmetrische Matrix ist genau dann positiv definit, wenn alle ihre Eigenwerte positiv sind und genau dann negativ definit, wenn alle ihre Eigenwerte negativ sind. Falls es sowohl positive als auch negative Eigenwerte gibt, ist sie indefinit. ■

Da die Determinante einer Matrix gleich dem Produkt ihrer Eigenwerte ist, folgt, daß eine Matrix nur dann positiv definit sein kann, wenn ihre Determinante positiv ist; für negativ definite $n \times n$ -Matrizen muß die Determinante bei geradem n ebenfalls positiv sein, bei ungeradem negativ.

Für symmetrische 2×2 -Matrizen läßt sich daraus leicht ein notwendiges und hinreichendes Kriterium machen: Das charakteristische Polynom von

$$A = \begin{pmatrix} a & b \\ b & d \end{pmatrix}$$

mit Eigenwerten λ_1 und λ_2 ist

$$\lambda^2 - (a+d)\lambda + (ad - b^2) = (\lambda - \lambda_1)(\lambda - \lambda_2);$$

daher ist

$$\lambda_1 + \lambda_2 = a + d.$$

(In der Tat rechnet man auf genau die gleiche Weise leicht nach, daß für jede $n \times n$ -Matrix die Summe der n Eigenwerte gleich der Summe der n Diagonaleinträge ist, die sogenannte *Spur* der Matrix.)

Wenn $\det A = ad - b^2$ positiv ist, haben nicht nur λ_1 und λ_2 , sondern auch a und d dasselbe Vorzeichen, das somit gleich dem von $a + d = \lambda_1 + \lambda_2$ ist. Also ist A genau dann positiv definit, wenn $\det A > 0$ und $a > 0$

ist, negativ definit, wenn $\det A > 0$ und $a < 0$ ist, und indefinit wenn $\det A < 0$ ist. (Anstelle von a könnte hier natürlich überall auch d stehen.)

Beispielsweise ist die Matrix $\begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix}$ positiv definit, denn sie hat Determinante eins und positive Diagonaleinträge. Im obigen Beispiel des Sattelpunkts mit $f(x, y) = x^2 - y^2$ ist

$$H_f(0, 0) = \begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}$$

offensichtlich indefinit, was man nicht nur an der negativen Determinanten sieht.

§2: Maxima und Minima unter Nebenbedingungen

Bei einem realen physikalischen oder technischen Prozeß können sich die Variablen selten frei im gesamten \mathbb{R}^n bewegen: Physikalisch sinnvoll ist meist nur eine beschränkte Teilmenge. Im Gegensatz zur Dimension eins, wo diese Teilmenge praktisch immer ein Intervall ist, gibt es aber im Mehrdimensionalen keinen Grund, warum diese Teilmenge offen oder zumindest der Abschluß einer offenen Teilmenge sein sollte: Im \mathbb{R}^3 kann man sich beispielsweise auch interessieren für das Maximum oder Minimum der Ladungsdichte auf einer Kugeloberfläche oder die elektrische Feldstärke oder Temperaturverteilung auf der Innenhaut eines Reaktordruckbehälters.

Diese Maxima oder Minima sind im allgemeinen keine lokalen Maxima oder Minima der betrachteten Funktion: Wenn man die jeweilige Fläche verläßt, läßt sich der Funktionswert selbst für einen solchen Extremwert meist noch – je nach Richtung – sowohl vergrößern als auch verkleinern. Dementsprechend können die Methoden, die wir in §1 diskutiert haben, solche Extremwerte üblicherweise nicht finden; wir brauchen weitere Werkzeuge, die in diesem Paragraphen bereitgestellt werden sollen.

Die Situation, um die es hier geht, ist typischerweise die folgende: Gegeben ist eine Funktion $f: \mathbb{R}^n \rightarrow \mathbb{R}$, möglicherweise auch nur auf einer Teilmenge $D \subset \mathbb{R}^n$ definiert, deren Extremwerte nicht auf \mathbb{R}^n oder D

gesucht werden, sondern nur auf einer Teilmenge, die beispielsweise durch das Verschwinden einer weiteren Funktion $g: \mathbb{R}^n \rightarrow \mathbb{R}$ gegeben ist. Falls wir uns für Extremwerte auf einer Kugel vom Radius r um den Nullpunkt interessieren, wäre dies etwa die Funktion

$$g: \begin{cases} \mathbb{R}^3 & \rightarrow \mathbb{R} \\ (x, y, z) & \mapsto x^2 + y^2 + z^2 - r^2. \end{cases}$$

Eine mögliche Strategie zur Lösung solcher Probleme besteht darin, die Gleichung $g = 0$ nach einer der Variablen aufzulösen, diese dann in f einzusetzen und sodann eine gewöhnliche Extremwertaufgabe zu lösen. Diese Auflösung ist *explizit* nur in sehr einfachen Fällen möglich, aber selbst wenn wir nur wissen, daß eine solche Auflösung *existiert*, können wir doch damit argumentieren und Kriterien ableiten.

Unter Maxima und Minima sollen hier *lokale* Extrema verstanden werden, so daß wir die üblichen Kriterien anwenden können:

Definition: Wir sagen, die Funktion $f: D \rightarrow \mathbb{R}$ auf einer Teilmenge $D \subseteq \mathbb{R}^n$ habe im Punkt $\mathbf{a} \in D$ ein lokales $\left\{ \begin{array}{l} \text{Maximum} \\ \text{Minimum} \end{array} \right\}$ unter der Nebenbedingung $g = 0$, wobei $g: D \rightarrow \mathbb{R}$ eine weitere Funktion ist, wenn $g(\mathbf{a}) = 0$ ist und es eine Umgebung U von \mathbf{a} gibt, so daß für alle $\mathbf{x} \in U$ gilt: Ist $g(\mathbf{x}) = 0$, so ist $f(\mathbf{x}) \begin{cases} \leq \\ \geq \end{cases} f(\mathbf{a})$.

Als Einstiegsbeispiel betrachten wir eine beliebige Schulbuchaufgabe zur Minimumbestimmung: Eine Konservendose soll bei einem vorgegebenen Volumen von 100 cm^3 möglichst wenig Blech benötigen, d.h. ihre Oberfläche soll minimal sein.

Die Oberfläche eines Zylinders der Höhe h mit einer Grundfläche vom Radius r ist

$$f(r, h) = 2\pi r^2 + 2\pi r \cdot h;$$

die Nebenbedingung für das Volumen $V = \pi r^2 h$ besagt, daß

$$g(r, h) = \pi r^2 h - 100 = 0$$

sein soll.

Hier läßt sich natürlich die Nebenbedingung sofort nach h auflösen:

$$h = \frac{100}{\pi r^2},$$

und wir müssen nur noch die Funktion

$$F(r) = f\left(r, \frac{100}{\pi r^2}\right) = 2\pi r^2 + \frac{200}{r}$$

minimieren. Für diese ist

$$F'(r) = 4\pi r - \frac{200}{r^2},$$

und dies verschwindet genau dann, wenn

$$4\pi r^3 = 200 \quad \text{oder} \quad r = \sqrt[3]{\frac{50}{\pi}}$$

ist.

In diesem einfachen Fall kann man solche Aufgaben also zurückführen auf gewöhnliche Extremwertaufgaben, indem man die Nebenbedingung nach einer der Variablen auflöst und diese dann in f einsetzt; in anderen Fällen kann man gelegentlich die Nebenbedingung durch geeignete Parameterwahl oder Wahl eines angepaßten Koordinatensystems berücksichtigen. Im allgemeinen wird aber beides nicht möglich sein, so daß wir andere Methoden brauchen.

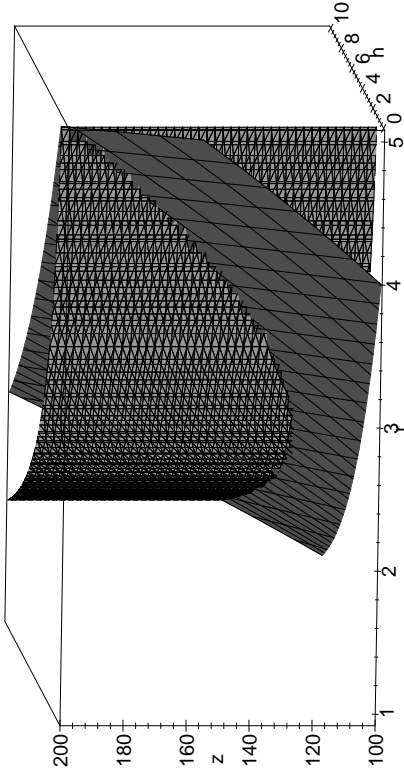


Abb. 55: Oberfläche einer Konservendose mit festem Volumen

Unser bisherige Theorie für lokale Extrema ist in dieser Situation nicht anwendbar, denn die lokalen Extrema von f werden nur in den seltensten Fällen die Nebenbedingung $g = 0$ erfüllen; im obigen Beispiel zeigt Abbildung 55 die Nebenbedingung als eng schraffierte Fläche dargestellt und der Graph von f als weiter schraffierte; wie man sieht, läßt sich der Wert von f problemlos verkleinern, wenn man nur die Fläche $g = 0$ verläßt, und in der Tat ist auch ohne jede Mathematik sofort klar, daß man mit weniger Blech auskommt, wenn man die Konservendose einfach schmaler oder kürzer macht.

Die Grundidee für ein alternatives Verfahren wird klar bei der Betrachtung der Niveaulinien in Abbildung 56: Die Niveaulinie für $g = 0$ ist gestrichelt eingezeichnet, verschiedene Niveaulinien von f als durchgezogene Kurven.

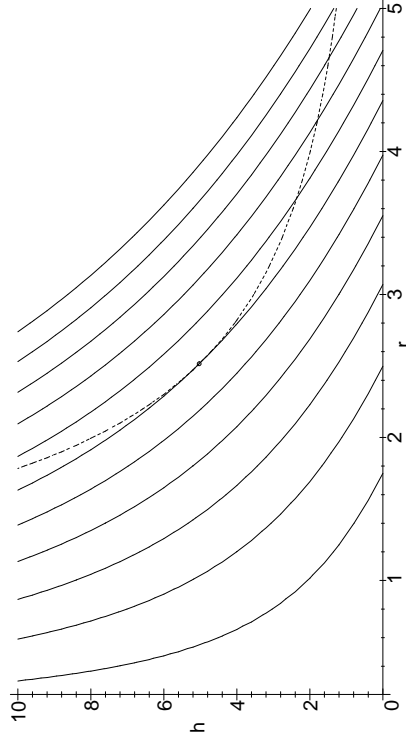


Abb. 56: Niveaulinien für Oberfläche und Volumen

Wie man sieht, schneiden einige dieser Niveaulinien die gestrichelte Kurve überhaupt nicht: Wenn man zu wenig Blech hat, kann man keine Dose mit 100 cm^3 Inhalt zusammenlöten. Wenn es dagegen genug Blech gibt, gibt es gleich zwei Schnittpunkte: Die Dose kann entweder eher höher oder eher breiter gemacht werden. In einem solchen Fall kann man die Niveaulinie durch eine zu einem etwas niedrigeren Niveau ersetzen,

die im allgemeinen auch wieder Schnittpunkte haben wird, so daß das Niveau noch nicht minimal sein kann. Erst wenn man im Minimum ist, fallen die beiden Schnittpunkte zusammen; wenn man nun das Niveau noch weiter erniedrigt, gibt es keine Schnittpunkte mehr.

Da somit im Minimum zwei Schnittpunkte zusammenfallen, berühren sich dort die Niveaulinien von f und von g , d.h. sie haben eine gemeinsame Tangente. Da der Gradient, wie wir wissen, senkrecht auf der Tangenten der Niveaulinien steht (die Richtungsableitung entlang einer Niveaulinie ist schließlich null), sind somit die Gradienten von f und g im Minimum zueinander parallel, d.h. der eine ist ein Vielfaches des anderen.

Dies gilt nicht nur im vorliegenden Beispiel, sondern allgemein:

Satz: $D \subseteq \mathbb{R}^n$ sei eine offene Menge und $f, g \in C^1(D, \mathbb{R})$ seien stetig differenzierbare Funktionen auf D . Falls f im Punkt $\mathbf{a} \in D$ ein Extremum hat unter der Nebenbedingung $g(\mathbf{x}) = 0$, so sind $\text{grad } f(\mathbf{a})$ und $\text{grad } g(\mathbf{a})$ linear abhängig.

Beweis: Die Grundidee ist einfach: Auch wenn wir die Nebenbedingung nicht *explizit* nach einer der Variablen auflösen können, sagt uns der Satz über implizite Funktionen in vielen Fällen dennoch, daß zumindest lokal eine Auflösung existiert. Diese Auflösung kennen wir zwar nicht, aber wir können mit ihr argumentieren und, zumindest formal, auch rechnen.

Falls $\text{grad } g(\mathbf{a})$ der Nullvektor ist, gibt es nichts mehr zu beweisen, denn jede Menge, die den Nullvektor enthält, ist linear abhängig.

Wir können daher annehmen, daß $\text{grad } g(\mathbf{a})$ mindestens eine von Null verschiedene Komponente hat, und durch Ummummern der Koordinaten können wir o.B.d.A. annehmen, daß dies die n -te Komponente ist, d.h. $g_{x_n}(\mathbf{a}) \neq 0$.

Dann gibt es nach dem Satz über implizite Funktionen ([HJM I], Kap. 2, §3d) eine Umgebung U von (a_1, \dots, a_{n-1}) und eine Funktion $h: U \rightarrow \mathbb{R}$ mit $h(a_1, \dots, a_{n-1}) = a_n$, so daß

$$g(x_1, \dots, x_{n-1}, h(x_1, \dots, x_{n-1})) = 0 \quad \text{für alle } (x_1, \dots, x_{n-1}) \in U.$$

Nachdem f in \mathbf{a} ein lokales Extremum unter der Nebenbedingung $g = 0$ hat, nimmt die Funktion

$$F(x_1, \dots, x_{n-1}) \stackrel{\text{def}}{=} f(x_1, \dots, x_{n-1}, h(x_1, \dots, x_{n-1}))$$

in (a_1, \dots, a_{n-1}) ein lokales Extremum im üblichen Sinne an, d.h. der Gradient von F verschwindet dort.

Nach der Kettenregel ist für $i = 1, \dots, n-1$

$$F_{x_i}(a_1, \dots, a_{n-1}) = f_{x_i}(\mathbf{a}) + f_{x_n}(\mathbf{a}) \cdot h_{x_i}(a_1, \dots, a_{n-1}),$$

und nach dem Satz über implizite Funktionen ist $h_{x_i} = -g_{x_i}/g_{x_n}$, d.h.

$$F_{x_i}(a_1, \dots, a_{n-1}) = f_{x_i}(\mathbf{a}) - f_{x_n}(\mathbf{a}) \frac{g_{x_i}(\mathbf{a})}{g_{x_n}(\mathbf{a})}.$$

Da die linke Seite verschwindet, gilt dasselbe auch für die rechte. Die rechte Seite ist im Gegensatz zur linken auch für $i = n$ definiert und verschwindet aus trivialen Gründen; also ist für alle i

$$f_{x_i}(\mathbf{a}) - \frac{f_{x_n}(\mathbf{a})}{g_{x_n}(\mathbf{a})} g_{x_i}(\mathbf{a}) = 0$$

oder, anders ausgedrückt,

$$\text{grad } f(\mathbf{a}) - \frac{f_{x_n}(\mathbf{a})}{g_{x_n}(\mathbf{a})} \text{grad } g(\mathbf{a}) = \vec{0}.$$

Damit sind die beiden Gradienten in der Tat linear abhängig. ■

Falls der Gradient von g im Punkt \mathbf{a} nicht verschwindet, gibt es somit eine Zahl $\lambda \in \mathbb{R}$, so daß

$$\text{grad } f(\mathbf{a}) - \lambda \text{grad } g(\mathbf{a}) = \vec{0}$$

ist, nämlich

$$\lambda = \frac{f_{x_n}(\mathbf{a})}{g_{x_n}(\mathbf{a})}.$$

Diese Zahl bezeichnet man als LAGRANGESCHEN Multiplikator; mit seiner inhaltlichen Interpretation werden wir uns in Kürze beschäftigen.



JOSEPH-LOUIS LAGRANGE (1736–1813) wurde als GIUSEPPE LODOVICO LAGRANGIA in Turin geboren und studierte dort zunächst Latein. Erst eine alte Arbeit von HALLEY über algebraische Methoden in der Optik weckte sein Interesse an der Mathematik, woraus ein ausgedehnter Briefwechsel mit EULER entstand. In einem Brief vom 12. August 1755 berichtete er diesem unter anderem über seine Methode zur Berechnung von Maxima und Minima; 1756 wurde er auf EULERS Vorschlag, Mitglied der Berliner Akademie; zehn Jahre später zog er nach Berlin und wurde dort EULERS Nachfolger als mathematischer Direktor der Akademie. 1787 wechselte er an die Pariser Académie des Sciences, wo er bis zu seinem Tod blieb und unter anderem an der Einführung des metrischen Systems beteiligt war. Seine Arbeiten umspannen weite Teile der Analysis, Algebra und Geometrie.

Zur praktischen Bestimmung von Extremwerten unter Nebenbedingungen geht man wie folgt vor: Über die Punkte, in denen der Gradient von g verschwindet, macht obiger Satz keine verwertbare Aussage; diese Punkte müssen also vorab berechnet und untersucht werden.

Danach müssen die Punkte gefunden werden, in denen es ein $\lambda \in \mathbb{R}$ gibt, so daß

$$\begin{aligned} f_{x_1}(\mathbf{x}) - \lambda g_{x_1}(\mathbf{x}) &= 0 \\ &\vdots \\ f_{x_n}(\mathbf{x}) - \lambda g_{x_n}(\mathbf{x}) &= 0 \\ g(\mathbf{x}) &= 0 \end{aligned}$$

ist. Dies ist ein System von $n+1$ Gleichungen für die $n+1$ Unbekannten, allerdings ist dieses Gleichungssystem nur selten linear und damit oft nicht mit bekannten Methoden lösbar. Manchmal kann man das Gleichungssystem durch geeignete Umformungen und Fallunterscheidungen vollständig lösen, in anderen Fällen helfen nur die aus der Numerik bekannten Näherungsverfahren wie etwa die Methode von NEWTON-RAPHSON.

Falls alle Gleichungen Polynomgleichungen sind (oder durch Einführung geeigneter zusätzlicher Variablen auf Polynomgleichungen zurückgeführt werden können), kann man im Falle einer endlichen Lösungsmenge diese auch exakt bestimmen: Genau wie der GAUSS-Algorithmus zur Lösung eines linearen Gleichungssystems dieses auf eine

Treppengestalt bringt, aus der man die Lösungen einfach ermitteln kann, gibt es in der Computeralgebra einen Algorithmus, der dasselbe für beliebige Systeme von Polynomgleichungen versucht; die Gleichungen, die dieser Algorithmus liefert, bezeichnet man als GRÖBNER-Basis oder Standardbasis. Zum Verständnis dieses Algorithmus, den man als eine Art Synthese aus EUKLIDISCHEN Algorithmus und GAUSS-Algorithmus ansehen kann, sind Kenntnisse der kommutativen Algebra erforderlich, für die die Zeit in dieser Vorlesung nicht ausreicht; bei einigen Implementierungen werden zusätzlich auch noch Algorithmen aus der Informatik eingesetzt, die typischerweise nicht in Grundvorlesungen behandelt werden. Deshalb sei hier nur darauf hingewiesen, daß die gängigen universellen Computeralgebrasysteme wie Maple, Mathematica, MuPad allesamt entsprechende Routinen enthalten, mit denen man auch dann experimentieren kann, wenn man die dahinterstehende Theorie nicht versteht.

Als Beispiel, wie gelegentlich auch ein nichtlineares Gleichungssystem elementar gelöst werden kann, betrachten wir eine Anwendung aus den Wirtschaftswissenschaften: Die Gesamtproduktion eines Unternehmens oder eines Staats in Abhängigkeit von n eingesetzten Ressourcen x_1, \dots, x_n wird oft modelliert durch eine sogenannte COBB-DOUGLAS-Funktion der Form

$$P(x_1, \dots, x_n) = \alpha x_1^{\epsilon_1} \dots x_n^{\epsilon_n},$$

benannt nach den beiden Wissenschaftlern, die dieses Modell 1928 für die amerikanische Gesamtproduktion in Abhängigkeit von Kapital und Arbeit in den Jahren 1899 bis 1922 entwickelten. (Sie fanden $P \approx 1,01A^{3/4}K^{1/4}$ mit $A = \text{Anzahl der Beschäftigten}$ und $K = \text{Kapitaleinsatz}$.)

Betrachten wir stattdessen die Produktion eines Wirtschaftsguts aus zwei Ressourcen x, y gemäß der Funktion

$$f(x, y) = P(x, y) = x^{1/2} y^{1/4}.$$

Falls wir der Einfachheit halber annehmen, daß die Kosten pro Einheit für x und y gleich sind und die Gesamtkosten höchstens gleich zwölf sein dürfen, müssen wir f maximieren unter der Nebenbedingung

$$x + y \leq 12.$$

Nun ist aber f eine monoton wachsende Funktion sowohl von x als auch von y , d.h. die maximale Produktion wird sicherlich erreicht in einem Punkt, für den $x + y = 12$ ist, denn für jeden anderen Punkt

(x, y) mit $x + y < 12$ ist $f(x, y) < f(x, 12 - x)$. Daher können wir die Nebenbedingung in der gewohnten Form

$$g(x, y) = x + y - 12 = 0$$

schreiben. Diese Nebenbedingung sowie die zu maximierende Funktion sind in Abbildung 57 dargestellt.

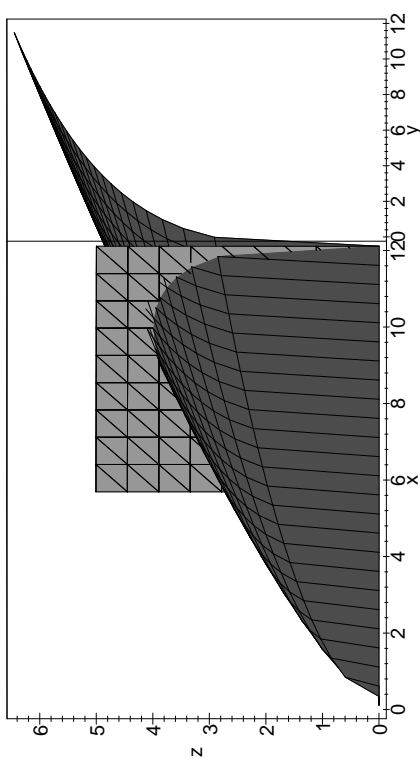


Abb. 57: Maximierung einer Produktionsfunktion bei festem Kapitaleinsatz

Ableitung beider Funktionen zeigt, daß

$$\text{grad } g = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \text{und} \quad \text{grad } f = \begin{pmatrix} y^{1/4}/2x^{1/2} \\ x^{1/2}/4y^{3/4} \end{pmatrix}$$

ist; das zu lösende Gleichungssystem wird also zu

$$\begin{aligned} \frac{y^{1/4}}{2x^{1/2}} - \lambda &= 0 \\ \frac{x^{1/2}}{4y^{3/4}} - \lambda &= 0. \end{aligned}$$

$$x + y - 12 = 0$$

(Die Nenner brauchen uns nicht zu stören, denn da $f(0, y) = f(x, 0) = 0$ ist, kommen Lösungen mit $x = 0$ oder $y = 0$ für das Maximum ohnehin nicht in Frage; wir können sie also getrost ausschließen.)

Als Ansatz zu einer möglichen Lösung können wir ausnutzen, daß λ in den beiden ersten Gleichungen isoliert steht; wenn wir danach auflösen und gleichsetzen, erhalten wir die Gleichung

$$\frac{y^{1/4}}{2x^{1/2}} = \frac{x^{1/2}}{4y^{3/4}}.$$

Multiplikation mit dem Hauptnenner macht daraus

$$4y^{1/4} y^{3/4} = 2x^{1/2} x^{1/2} \quad \text{oder} \quad 2y = x.$$

Einsetzen in die dritte Gleichung ergibt $3y = 12$, also ist

$$y = 4 \quad \text{und} \quad x = 8;$$

der Maximalwert von f ist

$$f(8, 4) = 8^{1/2} \cdot 4^{1/4} = 2\sqrt{2} \cdot \sqrt{2} = 4.$$

Auch den LAGRANGESCHEN Multiplikator λ können wir noch ausrechnen:

$$\lambda = \frac{y^{1/4}}{2x^{1/2}} = \frac{4^{1/4}}{2 \cdot 8^{1/2}} = \frac{\sqrt{2}}{2 \cdot 2\sqrt{2}} = \frac{1}{4}.$$

Die Berechnung von λ war für die Bestimmung des Optimums eigentlich überflüssig; λ ist nur eine Hilfsgröße zur Berechnung des Extremums. Wir wollen uns als nächstes überlegen, daß wir λ auch inhaltlich interpretieren können: Dazu betrachten wir eine Nebenbedingung

$$g(x_1, \dots, x_n) = c$$

mit *variabler* rechter Seite c und ein Extremum der Funktion

$$f(x_1, \dots, x_n).$$

Dieses Extremum wird natürlich von c abhängen; wir schreiben es in der Form

$$(x_1(c), \dots, x_n(c))$$

und nehmen an, daß die Funktionen $x_i(c)$ stetig differenzierbar seien. (Ein interessierter Leser kann sich anhand des Satzes über implizite Funktionen überlegen, welche Bedingungen f und g erfüllen müssen,

damit dies garantiert ist.) Der Optimalwert von f in Abhängigkeit von c ist dann

$$F(c) \stackrel{\text{def}}{=} f(x_1(c), \dots, x_n(c)).$$

Nach der Kettenregel aus [HM I], Kapitel 2, §3c) ist

$$F'(c) = \sum_{i=1}^n \frac{\partial f}{\partial x_i} \frac{dx_i(c)}{dc}.$$

Genauso können wir

$$G(c) \stackrel{\text{def}}{=} g(x_1(c), \dots, x_n(c))$$

betrachten und erhalten

$$G'(c) = \sum_{i=1}^n \frac{\partial g}{\partial x_i} \frac{dx_i(c)}{dc}.$$

Da $(x_1(c), \dots, x_n(c))$ ein Optimum ist, sind dort die Gradienten von f und $g - c$ proportional mit Proportionalitätsfaktor λ . Da wir bei der Gradientenbildung nur nach den x_i ableiten, von denen die rechte Seite c nicht abhängt, ist der Gradient von $g - c$ gleich dem von g selbst, d.h.

$$\frac{\partial f}{\partial x_i} = \lambda \frac{\partial g}{\partial x_i} \quad \text{für alle } i.$$

Somit ist $F'(c) = \lambda G'(c)$. Da der Punkt $(x_1(c), \dots, x_n(c))$ die Nebenbedingung mit rechter Seite c erfüllt, ist aber $G(c) = c$ und damit $G'(c) \equiv 1$. Also ist $\lambda = F'(c)$ die Wachstumsrate für das Optimum bei Änderung der rechten Seite der Nebenbedingung.

Im obigen Beispiel steigt also die Maximalmenge $f(x, y)$, die mit Kapitaleinsatz 12 produziert werden kann, für kleines h ungefähr um $h/4$, wenn wir den Kapitaleinsatz auf $12 + h$ erhöhen. Die Erhöhung des Kapitaleinsatzes lohnt sich, wenn für das fertige Produkt ein Preis pro Einheit erzielt werden kann, der größer ist als vier.

Als letztes wollen wir uns noch überlegen, was passiert, wenn wir nicht nur eine, sondern mehrere Nebenbedingungen erfüllen müssen. Es geht

also wieder darum, eine Funktion $f(x_1, \dots, x_n)$ zu optimieren, jetzt aber unter den Nebenbedingungen

$$g_1(x_1, \dots, x_n) \geq 0, \quad \dots \quad g_r(x_1, \dots, x_n) \geq 0.$$

(Es genügt, Bedingungen mit \geq zu betrachten, denn durch Multiplikation mit minus Eins kann man jede Ungleichung mit \leq in eine mit \geq überführen. Auch Gleichungen $g_i = 0$ kann man zumindest formal durch die beiden Ungleichungen $g_i \geq 0$ und $-g_i \geq 0$ ausdrücken.)

Die wichtigsten Beispiele solcher Optimierungsaufgaben sind die Fälle mit linearen Funktionen f und g_i ; hier redet man von *linearen Programmen*. (Das Wort *Programme* in diesem Zusammenhang hat natürlich nichts mit Computerprogrammen zu tun.) Das wichtigste Verfahren zur Lösung solcher Aufgaben, der Simplex-Algorithmus, wird in der Vorlesung *Numerik I* behandelt, so daß wir uns hier auf die *nichtlineare Programmierung* beschränken können.

Man überlegt sich leicht, daß im linearen Fall die Nebenbedingungen ein (endliches oder unendliches) Polyeder im \mathbb{R}^n definieren und eine lineare Funktion, so sie ein endliches Maximum oder Minimum hat, dieses auf dem Rand dieses Polyeders annimmt, und dort sogar in einer Ecke, Man muß daher „nur“ die Ecken dieses Polyeders untersuchen – deren Anzahl allerdings exponentiell mit der Anzahl der Variablen wächst. Trotzdem führt der Simplex-Algorithmus selbst im Fall von Zehntausenden von Variablen im allgemeinen recht schnell ans Ziel, so daß das theoretische Problem der exponentiellen Komplexität im schlimmsten Fall für praktische Anwendungen keine Bedeutung hat.

Bei nichtlinearen Funktionen ist die Situation komplizierter, denn nun kann es auch im Innern Extrema geben: Die Funktion

$$f(x, y) = e^{-x^2 - y^2} \quad \text{mit der Nebenbedingung} \quad x^2 + y^2 \leq 1$$

etwa nimmt ihr Maximum im Punkt $(0, 0)$ an; auf dem Rand des Einheitskreises liegen nur die Minima. Im allgemeinen Fall eines nichtlinearen Programms kann ein Optimum also entweder ganz im Innern liegen oder aber eine beliebige Teilmenge der Nebenbedingungen exakt erfüllen.

Falls wir es mit inneren Punkten zu tun haben, sind diese lokale Maxima oder Minima ohne Nebenbedingungen, und wir haben uns bereits in §1

überlegt, wie man diese bestimmt: In jedem solchen Punkt verschwindet der Gradient der zu optimierenden Funktion.

Im Falle einer einzigen *Gleichung* als Nebenbedingung ist der Gradient von f linear abhängig vom Gradienten der Nebenbedingung; da der Nullvektor von jedem anderen Vektor linear abhängig ist, schließt dies auch den Fall der Optima bei inneren Punkten mit ein. Die naheliegende Verallgemeinerung auf den Fall mehrerer Nebenbedingungen ist daher der

Satz: Die Funktion $f: D \rightarrow \mathbb{R}$ auf $D \subseteq \mathbb{R}^n$ habe im Punkt $\mathbf{a} \in D$ ein Extremum unter den Nebenbedingungen

$$g_1(\mathbf{a}) \geq 0, \quad g_2(\mathbf{a}) \geq 0, \quad \dots, \quad g_r(\mathbf{a}) \geq 0.$$

Dann sind die $r + 1$ Vektoren

$$\text{grad } f(\mathbf{a}), \quad \text{grad } g_1(\mathbf{a}), \quad \text{grad } g_2(\mathbf{a}), \quad \dots, \quad \text{grad } g_r(\mathbf{a})$$

linear abhängig.

Der *Beweis* erfordert keine wesentlich neuen Ideen gegenüber dem Fall einer einzigen Nebenbedingung und sei daher nur kurz skizziert: Falls die Gradienten der g_i im Punkt \mathbf{a} bereits untereinander linear abhängig sind, gibt es nichts mehr zu beweisen; nehmen wir also an, sie seien linear unabhängig. Dann gibt es (mindestens) r verschiedene Variablen x_{j_1} bis x_{j_r} , so daß

$$\frac{\partial g_i}{\partial x_{j_i}}(\mathbf{a}) \neq 0$$

ist. Also kann nach dem Satz über implizite Funktionen jede Nebenbedingung zur Elimination einer anderen Variablen benutzt werden, und im wesentlichen dieselbe Rechnung wie im Fall einer Nebenbedingung zeigt die Behauptung. ■

Die lineare Abhängigkeit der Vektoren

$$\text{grad } f(\mathbf{a}), \quad \text{grad } g_1(\mathbf{a}), \quad \text{grad } g_2(\mathbf{a}), \quad \dots, \quad \text{grad } g_r(\mathbf{a})$$

bezeichnet man als KUHN-TUCKER-Bedingung; sie ist eine offensichtliche Verallgemeinerung der Bedingung von LAGRANGE, ist allerdings

deutlich jünger: Sie erschien 1951 in einer gemeinsamen Arbeit von H.W. KUHN und A.W. TUCKER, vier Jahre, nachdem G. DANTZIG den Simplex-Algorithmus entwickelt hatte, und fast zweihundert Jahre, nachdem LAGRANGE seine Multiplikatoren zur Bestimmung von Extrema unter einer Nebenbedingung eingeführt hatte.

§3: Numerische Verfahren

Wie wir gesehen haben, führt die Methode der LAGRANGESchen Multiplikatoren im allgemeinen auf nichtlineare Gleichungssysteme, die nur in einfachen Fällen explizit lösbar sind. In allen anderen Fällen muß man mit numerischen Methoden arbeiten, und da bietet sich an, das Problem von vornherein ohne den Umweg über LAGRANGESche Multiplikatoren Extrema numerisch zu bearbeiten.

a) Die Gradientenmethode

Für eine differenzierbare Funktion f auf $D \subseteq \mathbb{R}^n$ ist

$$f(\mathbf{x} + \vec{h}) = f(\mathbf{x}) + \text{grad } f(\mathbf{x}) \cdot \vec{h} + o(\|\vec{h}\|);$$

wenn wir ein Maximum (oder Minimum) von f ansteuern wollen, liegt es daher nahe, \vec{h} so zu wählen, daß sich der Funktionswert möglichst stark vergrößert (oder verkleinert).

Nach der CAUCHY-SCHWARZschen Ungleichung ist

$$|\text{grad } f(\mathbf{x}) \cdot \vec{h}| \leq \|\text{grad } f(\mathbf{x})\| \cdot \|\vec{h}\|;$$

wir erhalten also die maximalmögliche Veränderung bei vorgegebener Länge von \vec{h} genau dann, wenn \vec{h} parallel zum Gradienten ist.

Damit bietet sich folgende Strategie an: Wir wählen irgendeinen Ausgangspunkt \mathbf{x}_0 und berechnen dort den Gradienten $\nabla f(\mathbf{x}_0)$. Weiter gehen uns eine Länge ℓ_0 für den Vektor \vec{h} vor, die von der Länge des Gradienten abhängen kann oder auch nicht. Dann setzen wir bei der Suche nach einem Maximum

$$\vec{h}_0 = \frac{\ell_0}{\|\nabla f(\mathbf{x}_0)\|} \nabla f(\mathbf{x}_0);$$

bei der Suche nach Minima nehmen wir das Negative davon.

Als nächstes betrachten wir den Punkt

$$\mathbf{x}_1 \stackrel{\text{def}}{=} \mathbf{x}_0 + \vec{h}_0,$$

berechnen dort den Gradienten $\nabla f(\mathbf{x}_0)$, setzen mit geeignetem ℓ_1

$$\vec{h}_1 = \pm \frac{\ell_1}{\|\nabla f(\mathbf{x}_1)\|} \nabla f(\mathbf{x}_1)$$

(+ für Maxima, – für Minima) zur Definition des nächsten Punkts

$$\mathbf{x}_2 \stackrel{\text{def}}{=} \mathbf{x}_1 + \vec{h}_1$$

und so weiter. In jedem Schritt erhöhen (oder erniedrigen) wir den Funktionswert soweit wie es mit der vorgegebenen Länge ℓ_i nur möglich ist in der Hoffnung, so irgendwann auf ein Maximum (oder Minimum) zu stoßen. Dieses können wir erreichen, wenn wir am Rand des Definitionsbereichs von f angelangt sind, oder aber wenn wir in einem Punkt sind, in dem der Gradient verschwindet: Von dort aus geht es mit diesem Verfahren nicht mehr weiter.

Da wir mit einem numerischen Verfahren nur ein verschwindend geringe Chance haben, exakt in einem Extremum zu enden, zeigt sich hier auch die Notwendigkeit einer intelligenten Wahl der Schrittweiten ℓ_i : Wenn diese zu groß sind, kann es passieren, daß wir endlos um ein Extremum herumoszillieren.

Theoretisch ist auch möglich, daß wir in einem Sattelpunkt landen, aber wenn man sich überlegt, wie die Gradienten in der Umgebung eines Sattelpunktes aussehen, wird schnell klar, daß dies nur sehr selten passiert.

Abbildung 58 zeigt ein einfaches Beispiel für einen mit der Gradientenmethode zurückgelegten Weg; hier wurde in jedem Schritt

$$\begin{pmatrix} h_i \\ k_i \end{pmatrix} = 0,1 \cdot \nabla f(x_i, y_i)$$

gesetzt. Der Weg geht offensichtlich recht zielstrebig auf das Maximum zu.

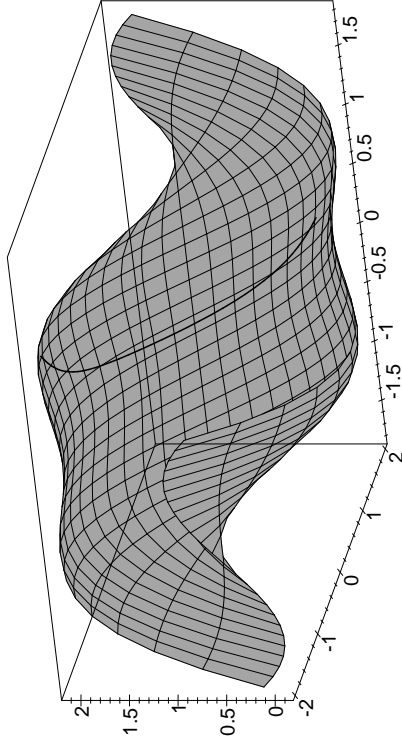


Abb. 58: Eine Anwendung der Gradientenmethode

Abbildung 59 zeigt dasselbe Bild in einen etwas größeren Zusammenhang; hier sehen wir, daß unser Streben nach kurzfristigen Gewinnen langfristig wohl doch nicht so erfolgreich war: Wenn wir vom Startpunkt aus nach rechts in die kleine Mulde abgestiegen wären, hätten wir auf dem gegenüberliegenden Hang deutlich größere Funktionswerte erreicht als im lokalen Maximum, in dem wir schließlich gelandet sind.

Dies ist ein grundsätzliches Problem von Gradientenverfahren: Falls man sie in der Nähe des (absoluten) Optimums starten läßt, führen sie schnell und zuverlässig ans Ziel, ansonsten aber ist die Gefahr sehr groß, daß man in einem nur lokalen Optimum steckenbleibt.

Um von dort wieder weiterzukommen, gibt es verschiedene Strategien. Eine anschaulich recht klare ist die sogenannte „Tunnelung“. Der Name entstand aus der Betrachtung von Minimierungsproblemen; nehmen wir also an, wir wollen das Minimum der Funktion $f(x, y)$ in einem gewissen Bereich finden und ein Gradientenverfahren hat uns in einen Punkt \mathbf{x}_M geführt, von dem aus es nicht mehr weiterkommt. Um zu sehen, ob $z_M = f(\mathbf{x}_M)$ wirklich der kleinste Wert ist, den f im betrachteten Bereich annehmen kann, versuchen wir, eine weitere Lösung der Gleichung

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

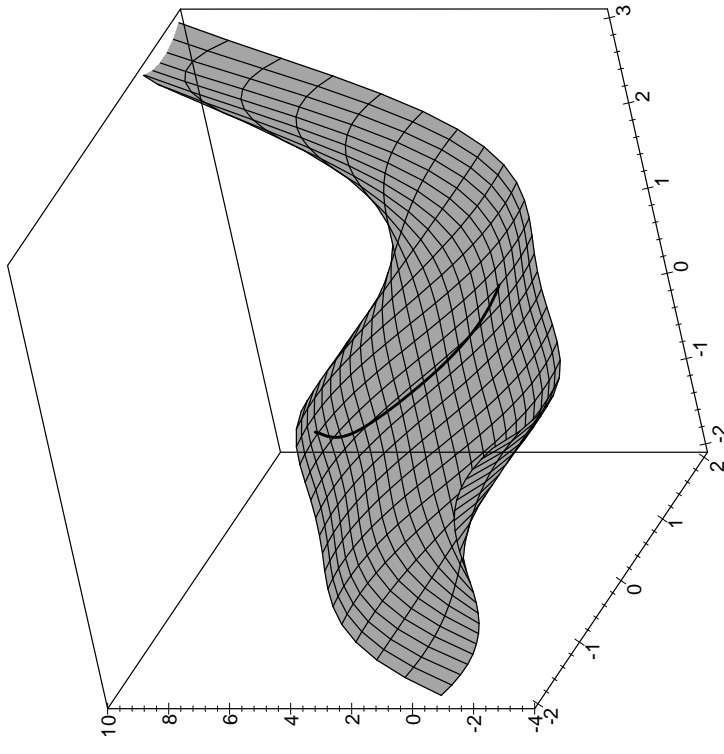


Abb. 59: Der Weg aus Abb. 58 aus einem weiteren Blickwinkel

zu finden. Dafür gibt es eine ganze Reihe numerischer Verfahren, z.B. das Verfahren von NEWTON-RAPHSON, mit denen sich zumindest ein solcher Punkt leicht finden läßt. Leider könnte dieser Punkt unser Ausgangspunkt \mathbf{x}_M sein; deshalb sucht man tatsächlich nicht nach Lösungen der Gleichung $f(\mathbf{x}) = z_M$, sondern nach Lösungen einer leicht abgewandelten Gleichung der Form

$$\tilde{f}(\mathbf{x}) = z_M,$$

wobei \tilde{f} dadurch aus f entsteht, daß man die Funktionswerte in der

unmittelbaren Umgebung von (x_M, y_M) starkt anhebt, so daß das dortige Minimum verschwindet. Dazu kann man beispielsweise eine Funktion der Form

$$\tilde{G}(x, y) = ae^{\frac{(x-x_M)^2+(y-y_M)^2}{b}}$$

mit geeigneten Parametern a, b wählen, wie sie in Abbildung 60 zu sehen ist, und

$$\tilde{f}(x, y) = f(x, y) + \tilde{G}(x, y)$$

setzen.

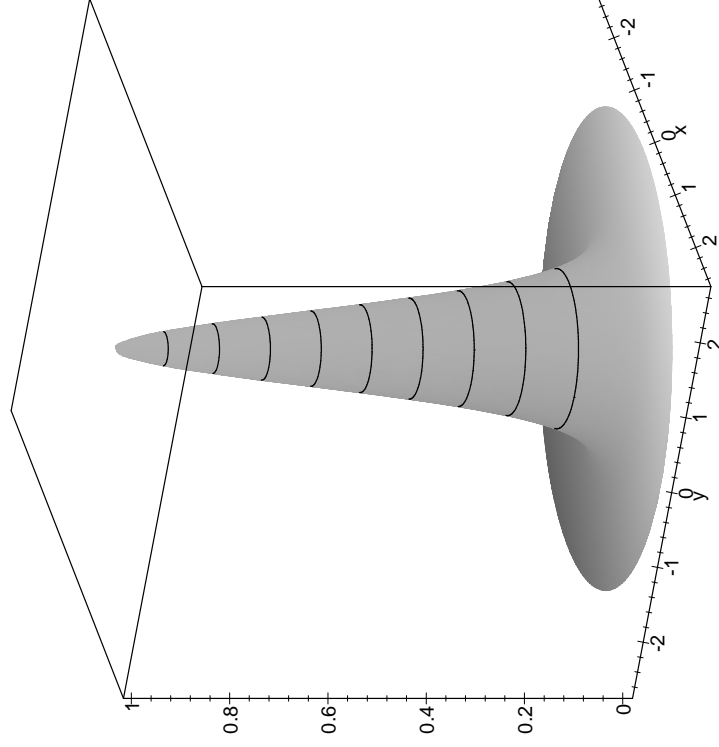


Abb. 60: $G(x, y) = e^{-3(x^2+y^2)}$

Dies bringt das Minimum im Punkt x_M zum Verschwinden und verändert die Funktion praktisch nicht, wenn man nur hinreichend weit entfernt ist von x_M . (Je kleiner b ist, umso lokalisierter ist die Veränderung.) Eine Lösung der Gleichung

$$\tilde{f}(x) = z_M,$$

so es eine gibt, liegt also nicht in der unmittelbaren Umgebung von x_M und ist daher ein guter Ausgangspunkt, um dort die Gradientenmethode noch einmal zu starten bis zum nächsten lokalen Minimum und so weiter. Sobald die Gleichung nicht mehr lösbar ist, können wir ziemlich sicher sein, daß z_M das globale Minimum ist – es sei denn, wir hätten die Parameter a und b sehr dumm gewählt.

Tunnelung ist auch ein wichtiges Konzept in der Physik: Dort versucht ein System bekanntlich stets, sein Energieminimum zu erreichen. Dies kann jedoch daran scheitern, daß es sich in einem lokalen Minimum befindet und nicht genügend Energie aufbringen kann, um den Energie-wall zu überwinden, der es vom absoluten Minimum trennt. Zumindest im Bereich der Quantentheorie gibt es dann auch den sogenannten *Tunneleffekt*, der es einzelnen Teilchen erlaubt, diesen Wall zu tunneln und auf diese Weise einen Zustand niedrigerer Energie zu erreichen.

Im obigen Beispiel geht es nicht um ein Minimum, sondern um ein Maximum, da die Suche danach graphisch besser darstellbar ist. Also graben wir auch keinen Tunnel, sondern spannen ein Hochseil, das irgendwo auf der eingezeichneten Ebenen liegt und uns vom erreichten Zwischenhoch zur Startposition für einen weiteren Anstieg bringt. (Tatsächlich ist die Ebene etwas zu tief eingezeichnet, damit man das alte Maximum noch erkennen kann; das Seil muß also etwas höher hängen.)

b) Der Metropolis-Algorithmus

Eine weitere Idee zur Vermeidung von Zwischenhochs hat ebenfalls viel mit Physik zu tun: Ein Gas erreicht seinen Zustand minimaler Energie dann, wenn die Bewegungsenergie $\frac{1}{2}mv^2$ eines jeden Teilchens gleich null ist, wenn sich also nichts mehr bewegt. Dies geschieht aber höchstens am absoluten Nullpunkt; bei positiven Temperaturen werden die meisten Teilchen positive kinetische Energie haben. Nach LUDWIG

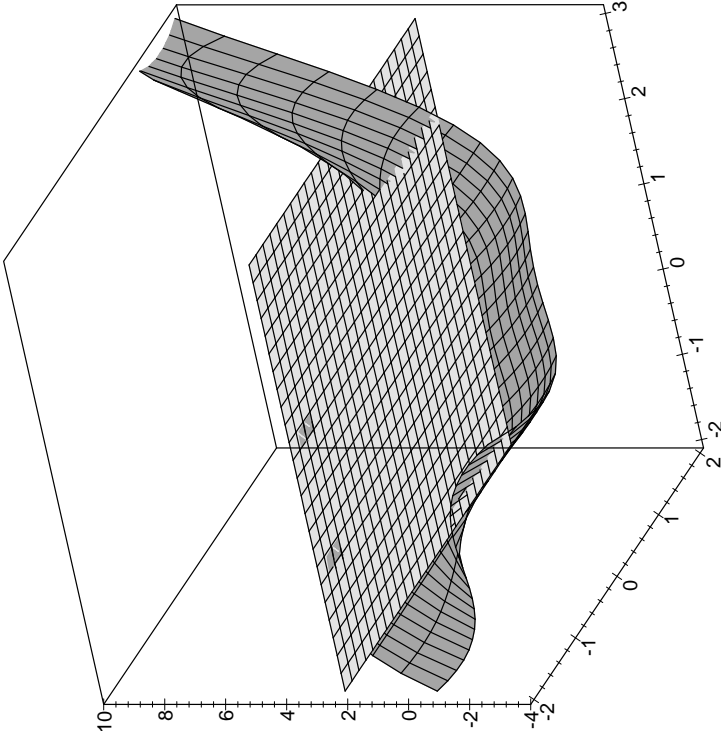


Abb. 61: „Tunnelung“ für Maxima

BOLTZMANN ist dabei die Wahrscheinlichkeit dafür, daß ein Teilchen die Energie $E = \frac{1}{2}mv^2$ hat, bei Temperatur T proportional zu

$$\frac{e^{-E/kT}}$$

mit einer Konstanten $k \approx 1,38066 \cdot 10^{-23} \text{ J/K}$, die heute als BOLTZMANN-Konstante bezeichnet wird.

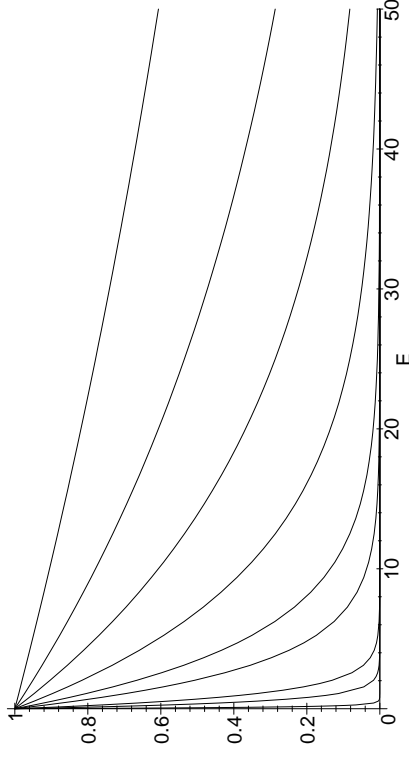


Abb. 62: $e^{-E/kT}$ für $kT = 0, 1, 0,5, 1, 3, 5, 10, 20, 40, 100$



LUDWIG BOLTZMANN (1844–1906) wuchs auf und studierte in Wien; danach lehrte er in Graz, Heidelberg, Berlin, Graz, Wien, Graz, Wien, Leipzig und Wien. Er war Professor für Theoretische Physik, für Mathematik und für Experimentalphysik. Auf seiner letzten Stelle in Wien hielt er eine so erfolgreiche Philosophievorlesung, daß ihn Kaiser Franz Josef in den Palast einlud. Am bekanntesten ist er für die Begründung der statistischen Mechanik, einer damals sehr umstrittene Theorie. Ob die damit verbundenen Anfeindungen zu seinem Selbstmord führten, ist unbekannt.

Bei der *simulierten Abkühlung* oder *BOLTZMANN-Maschine* ahmt man dies nach, indem man mit einer hohen Temperatur startet und der Richtung, in der man weitergeht, einer dieser Temperatur entsprechende Freiheit läßt. Man geht also nicht mehr unbedingt in Richtung des Gradienten, sondern geht zufällig in eine von endlich vielen vorgegebenen Richtungen. Die Wahrscheinlichkeit für den Richtungsvektor \vec{h}_j soll dabei analog zur BOLTZMANN-Verteilung festgelegt werden, d.h. wir ordnen ihm eine „Energie“

$$E_j = \pm(f(\mathbf{x} + \vec{h}_j) - f(\mathbf{x}, y))$$

zu (positiv bei der Suche nach einem Minimum, negativ bei der Suche nach einem Maximum) und die Wahrscheinlichkeit dafür, daß wir in

Richtung \vec{h}_j gehen, soll proportional sein zu $e^{-E_j/kT}$. Sie ist also, falls N Richtungen zur Verfügung stehen, gleich

$$p_j \stackrel{\text{def}}{=} e^{\frac{-E_j/kT}{\sum_{\ell=1}^N e^{-E_\ell/kT}}}$$

Zur Wahl einer Richtung erzeugen wir uns somit eine Zufallszahl Z zwischen null und eins und gehen in Richtung \vec{h}_j , wenn

$$\sum_{\ell=1}^{j-1} p_\ell < Z \leq \sum_{\ell=1}^j p_\ell$$

ist. (Die Frage, wie lang die Richtungsvektoren im wievielten Schritt sein sollen, wollen wir hier ausklammern.)

Bie hohen Temperaturen ist damit die Richtung fast vollständig zufallsbedingt gewählt, während in der Nähe des absoluten Nullpunkts praktisch nur noch die optimale Richtung eine Chance hat. Falls wir bei hoher Temperatur in einem Zwischenextremum landen, sorgt dies also mit recht hoher Wahrscheinlichkeit dafür, daß wir dort nicht steckenbleiben.

Am Ende wollen wir allerdings im absoluten Optimum steckenbleiben, d.h. wir müssen die Temperatur im Verlauf der Rechnung immer weiter senken – daher der Name *simulated annealing* = simulierte Abkühlung. Bei der Anwendung auf Optimierungsprobleme bezeichnet man diese Vorgehensweise als den METROPOLIS-Algorithmus. In welcher Weise man die Temperatur am besten senkt, ist immer noch ein Gebiet aktiver Forschung. Man kann zeigen, daß man statistisch betrachtet praktisch immer im Optimum landet, wenn man mit einer hinreichend hohen Ausgangstemperatur T_1 startet und im r -ten Schritt mit Temperatur $T_1 / \log(r + 1)$ arbeitet, aber bei einer derart langsamen Abkühlung braucht der Algorithmus viel zu lange, um ans Ziel zu kommen.



Nick Metropolis

NICHOLAS METROPOLIS (1915–1999) wuchs auf in Chicago, wo er Physik studierte und 1941 promovierte. Seit 1943 arbeitete er, unterbrochen durch Professuren an der Universität Chicago von 1945–1948 und 1957–1965, in den Los Alamos Laboratories, die ihm im Nachruf als *giant of mathematics and one of the founders of the Information Age* bezeichneten. Sein Ruhm als Mathematiker beruht vor allem auf den von ihm entwickelten Anwendungen statistischer Verfahren auf eine Vielzahl von mathematischen Problemen; zum Pionier des Informationszeitalter macht ihn u. a., daß er einer der ersten Anwender des ersten elektronischen Computers ENIAC war, dessen Nachfolger MANIAC baute und an der Universität Chicago das Institute for Computer Research gründete und bis 1965 leitete.

In Abbildung 63 sieht man, wie sich der Algorithmus bei einer Abkühlungsregel verhält, die im r -ten Schritt mit Temperatur T_1/r arbeitet: Zumindest im gezeigten Fall funktioniert das recht gut.

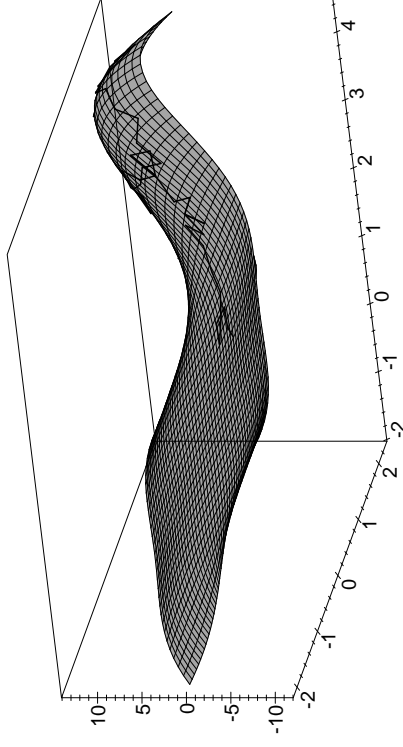


Abb. 63: Der METROPOLIS-Algorithmus für obiges Problem

In anderen Fällen (d.h., wenn andere Zufallszahlen gezogen werden) bleibt man damit aber auch gelegentlich ziemlich lange im Tal hängen; ein Beispiel dafür zeigt Abbildung 64.

Auch hier kommt man aber immerhin in eine gute Startposition, und oft