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.

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

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

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.

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

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

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 $x_0 \in D$ einen Vektor

$$\nabla f(x_0) = \text{grad } f(x_0) \in \mathbb{R}^n$$

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

$$f(x_0 + \vec{h}) = f(x_0) + \text{grad } f(x_0) \cdot \vec{h} + o(|\vec{h}|).$$

Hier muß also für jeden Extremwert $\text{grad } f(x_0)$ gleich dem Nullvektor sein, denn setzt man für \vec{h} ein kleines Vielfaches $t \cdot \text{grad } f(x_0)$ des Gradienten ein, wäre sonst

$$f(x_0 + \vec{h}) = f(x_0) + t(\text{grad } f(x_0) \cdot \text{grad } f(x_0)) + o(|\vec{h}|)$$

für kleine positive t größer als $f(x_0)$ 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

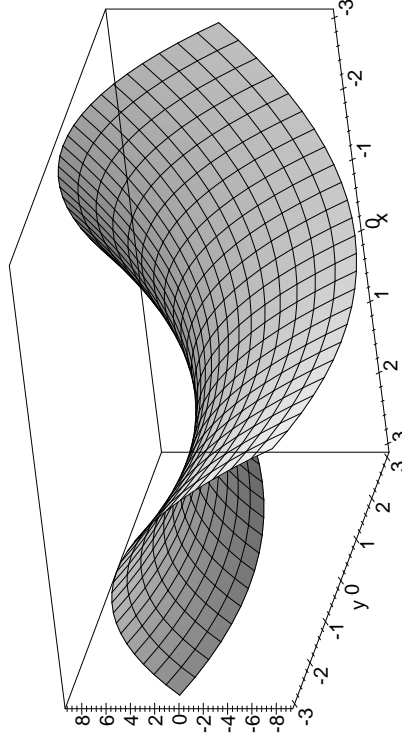


Abb. 54: Graph der Funktion $f(x, y) = x^2 - y^2$

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 x_0 abhängt:

Die zweite Ableitung von $f \in C^2(D, \mathbb{R})$ im Punkt $x_0 \in D$ ist bekanntlich gegeben durch die HESSE-Matrix

$$H_f(x_0) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_n} & \frac{\partial^2 f}{\partial x_2 \partial x_n} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

und zweimalige Differenzierbarkeit bedeutet, daß

$$f(x_0 + \vec{h}) = f(x_0) + \text{grad } f(x_0) \cdot \vec{h} + \frac{1}{2} \vec{h}^t H_f(x_0) \vec{h} + o(|\vec{h}^2|)$$

ist für kleine \vec{h} .

Wenn $\text{grad } f(x_0)$ verschwindet, hängt also das Verhalten von f in der Umgebung von x_0 ab von der quadratischen Form

$$\vec{h} \mapsto \vec{h}^t H_f(x_0) \vec{h}.$$

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:

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

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

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

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

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

Mit dieser Terminologie ist das folgende Lemma klar:

Lemma: Wenn die differenzierbare Funktion $f \in C^1(D, \mathbb{R})$ im Punkt $\vec{x}_0 \in D$ 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 $\vec{x} \in D$ verschwindet, gilt:

- a) Falls die HESSE-Matrix $H_f(\vec{x}_0)$ positiv definit ist, hat f im Punkt \vec{x}_0 ein Minimum.
- b) Falls $H_f(\vec{x}_0)$ negativ definit ist, hat f im Punkt \vec{x}_0 ein Maximum.
- c) Falls $H_f(\vec{x}_0)$ indefinit ist, hat f im Punkt \vec{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 $\begin{cases} \text{Maximum} \\ \text{Minimum} \end{cases}$ 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.

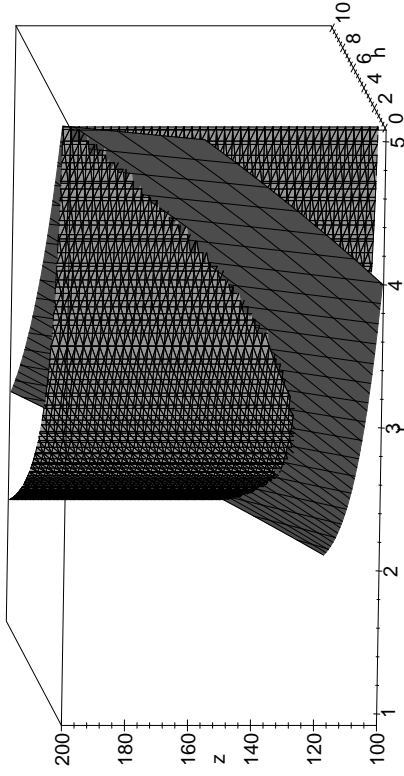


Abb. 55: Oberfläche einer Konservendose mit festem Volumen

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 schmälert 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.

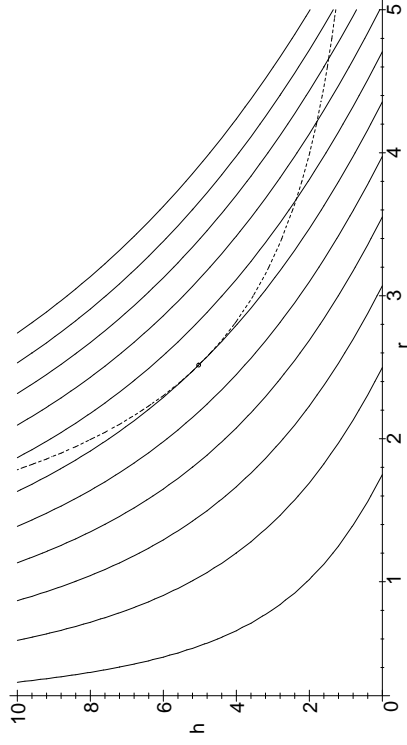


Abb. 56: Niveaulinien für Oberfläche und Volumen

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 (**HIM 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 GRU-SEPPE 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.

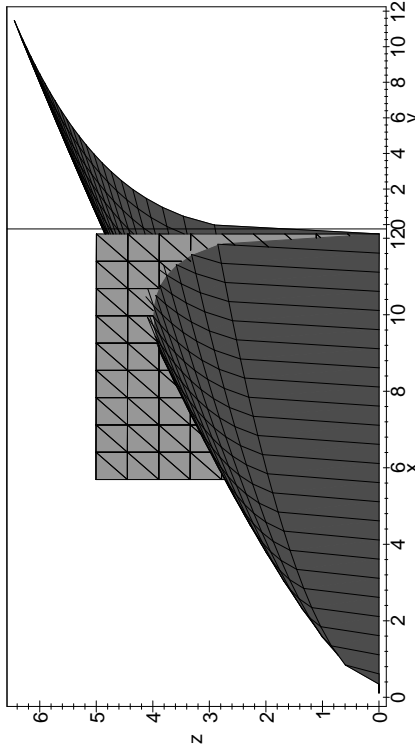


Abb. 57: Maximierung einer Produktionsfunktion bei festem Kapitaleinsatz

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

$$\frac{y^{1/4}}{2x^{1/2}} - \lambda = 0$$

$$\frac{x^{1/2}}{4y^{3/4}} - \lambda = 0$$

$$x + y - 12 = 0$$

(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

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

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

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

der Maximalwert von f ist

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

Auch den LAGRANGESchen Multiplikator λ können wir noch ausrechnen:

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

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 KUHNTUCKER-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

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

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 Punktes

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

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

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

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

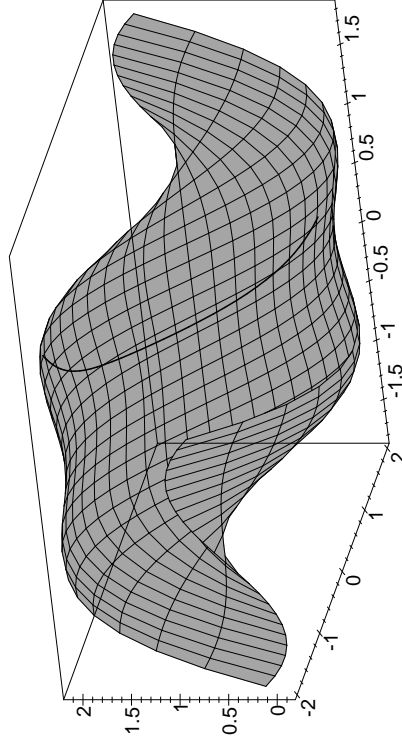


Abb. 58: Eine Anwendung der Gradientenmethode

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

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

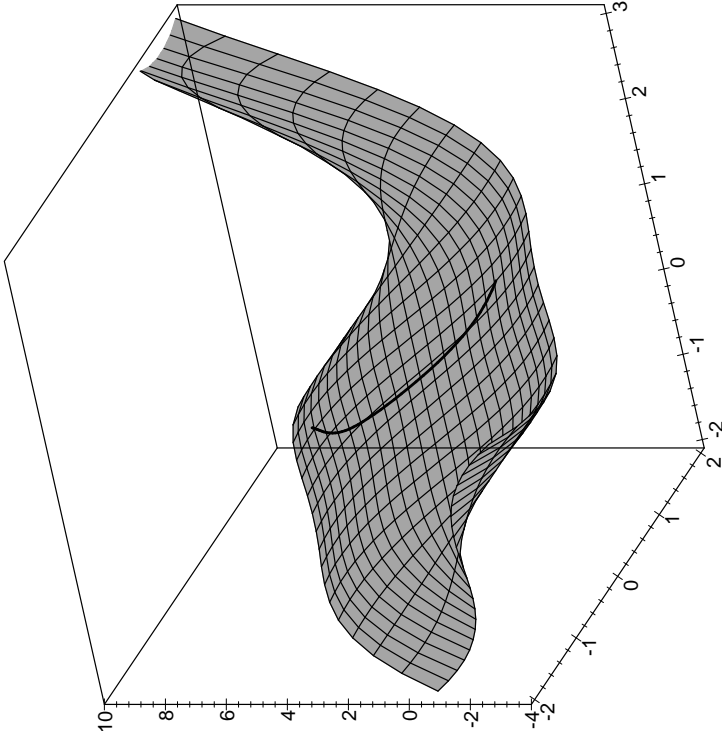


Abb. 59: Der Weg aus Abb. 58 aus einem weiteren Blickwinkel

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 \boldsymbol{x}_M sein; deshalb sucht man tatsächlich nicht nach Lösungen der Gleichung $f(\boldsymbol{x}) = z_M$, sondern nach Lösungen einer leicht abgewandelten Gleichung der Form

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

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

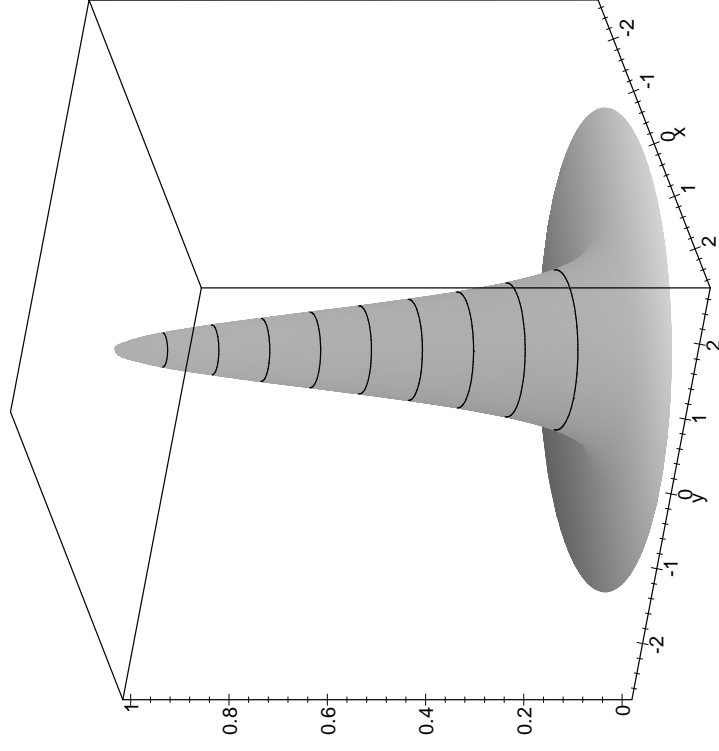
unmittelbaren Umgebung von (x_M, y_M) stark anhebt, so daß das dortige Minimum verschwindet. Dazu kann man beispielsweise eine Funktion der Form

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

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

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

setzen.

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

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

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

so es eine gibt, liegt also nicht in der unmittelbaren Umgebung von \boldsymbol{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

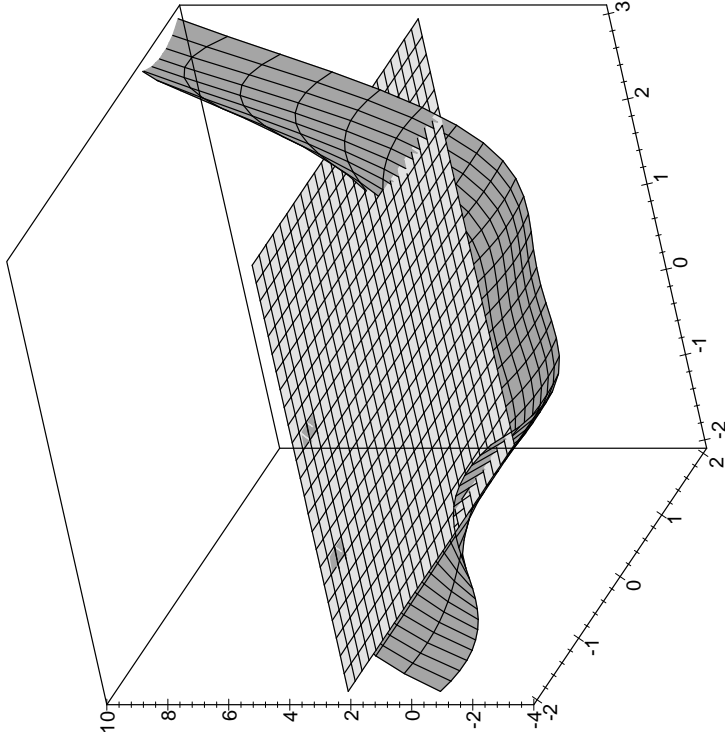


Abb. 61: „Tunnelung“ für Maxima

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

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

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

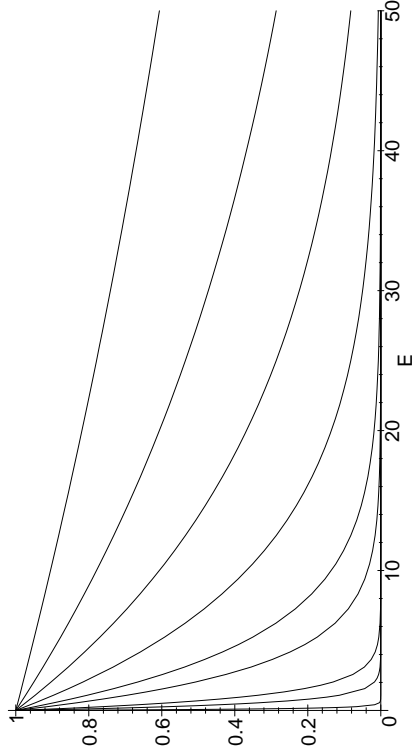


Abb. 62: $e^{-E/kT}$ für $kT = 0, 1, 0, 5, 1, 3, 5, 10, 20, 40, 100$



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.)

Die 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 Laboratorien, die ihn 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 Informationszeitalters 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 Ableitungsregel verhält, die im r -ten Schritt mit Temperatur T_1/r arbeitet: Zumindest im gezeigten Fall funktioniert das recht gut.

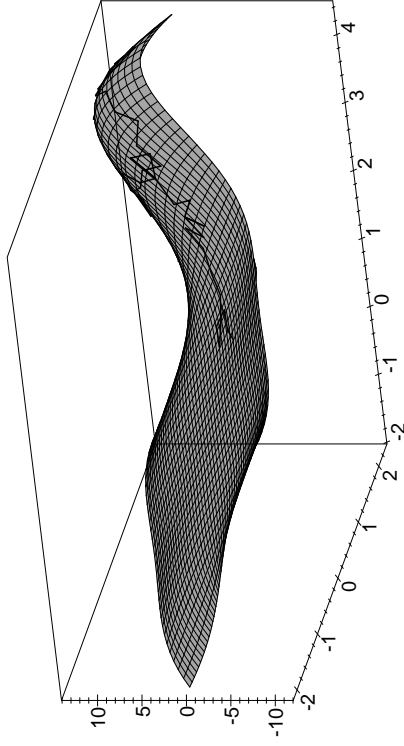


Abb. 63: Der METROPOLIS-Algorithmus für obiges Problem

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

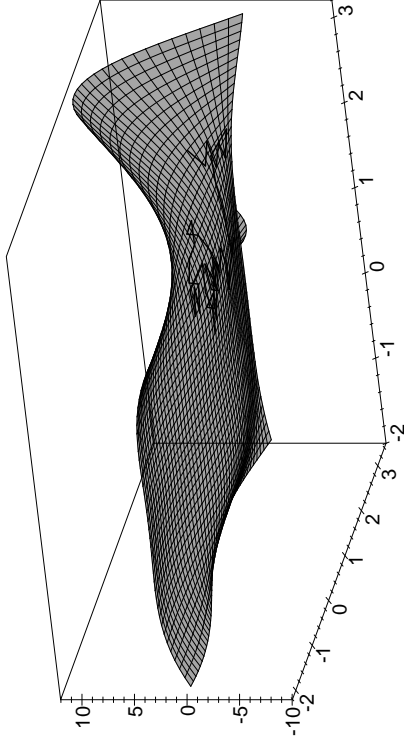


Abb. 64: Ditto mit anderen Zufallszahlen

wird es am besten sein, nach hinreichend vielen METROPOLIS-Schritten einfach ein gewöhnliches Gradientenverfahren zu starten.

Zusammenfassend läßt sich sagen, daß der METROPOLIS-Algorithmus und verwandte Verfahren (die sogenannten Monte-Carlo-Methoden) sehr nützliche Hilfsmittel zur Optimierung sind, falls man so gut wie nichts über die zu optimierende Funktion weiß. Sie funktionieren nicht nur bei kontinuierlichen Problemen, wie den hier betrachteten, sondern auch für diskrete und kombinatorische Optimierungsprobleme.

Sie haben allerdings den Nachteil, daß man nie garantieren kann, daß man ein Optimum erreichen wird, und selbst wenn man eines erreicht, kann die Methode dies nicht erkennen. (Es gibt alternative numerische Methoden, die das können.)

c) Zusammenfassung

Die nichtlineare Optimierung ist ein sehr weites Feld, von dem eine Grundvorlesung wie die *Höhere Mathematik* nur einen kleinen Ausschnitt behandeln kann. Dieser Ausschnitt besteht nicht aus den für die Praxis wichtigsten Verfahren, sondern aus denen, die sich am besten in den Stoff der Vorlesung einordnen. Sie sind zwar (in Kombination mit

dem aus der *Numerik* bekannten Simplex-Verfahren) die Grundbausteine, aus denen die meisten praktisch relevanten Verfahren zusammengesetzt sind, aber für die vielen kleinen Abwandlungen, die dazu führen, daß man ein Problem wirklich effizient lösen kann, müßte man deutlich mehr Zeit aufwenden, als hier zur Verfügung steht. Interessenten seien auf entsprechende Spezialvorlesung aus dem Bereich der Mathematik oder Operations Research verwiesen.

§4: Grundzüge der Fehler- und Ausgleichsrechnung

Physikalische Gesetze machen meist nur dann eine Aussage über ein reales System, wenn alle Umgebungsbedingungen exakt kontrolliert werden können. Das ist in der Praxis natürlich nie möglich. Insbesondere hat man bei der Anwendung physikalischer Prinzipien zur Messung von Daten keine Chance, den exakten Wert der zu messenden Größe zu bestimmen; der gemessene Wert wird immer von zahlreichen kleineren Störungen beeinflusst sein, die man bei einem gut durchgeführten Experiment für alle praktischen Zwecke als zufällig betrachten kann.

Zusätzlich kann die Messung noch durch mehr oder weniger große *systematische* Fehler verfälscht sein; diese können hervorgerufen werden durch ein falsch kalibriertes Meßgerät, Ablesen auf der falschen Skala eines Meßinstruments, durch falsche Anwendung von Meßvorschriften *usw.* Mit diesen systematischen Fehlern wollen wir uns hier nicht beschäftigen; in diesem Paragraphen soll es nur um *Zufallsfehler* gehen.

a) Das Laplacesche Fehlermodell

Der französische Mathematiker PIERRE SIMON, MARQUIS DE LAPLACE (1749–1827), dem wir in dieser Vorlesung bereits mehrfach begegnet sind, entwickelte ein extrem vereinfachtes Modell für das Zustandkommen zufälliger Meßfehler. Trotz seiner unrealistischen Annahmen ist es auch für die Praxis immer noch sehr interessant, da man inzwischen weiß, daß auch sehr viel kompliziertere realistische Fehlerquellen dasselbe Verhalten zeigen, das LAPLACE aus seinem Modell ableitete.

Die Grundannahme des LAPLACESchen Fehlermodells können wir uns so vorstellen, daß eine große Anzahl von „Dämonen“ (oder Fehlerquellen)

unsere Meßergebnisse verfälschen; jeder einzelne dieser „Dämonen“ verursacht einen Fehler derselben Größe ε in positiver oder negativer Richtung, wobei die Wahrscheinlichkeit für $+\varepsilon$ bzw. $-\varepsilon$ für jeden der „Dämonen“ jeweils 50% sein soll und die einzelnen „Dämonen“ unabhängig voneinander handeln sollen.

Im Falle eines einzigen „Dämonen“ wäre der Fehler also mit gleicher Wahrscheinlichkeit $+\varepsilon$ oder $-\varepsilon$, bei zwei „Dämonen“ wäre er in jeweils 25% aller Fälle $+2\varepsilon$ oder -2ε , während sich in 50% der Fälle die beiden Fehler aufheben würden.

Allgemein gibt es bei n „Dämonen“ 2^n gleichwahrscheinliche Möglichkeiten für deren Verhalten; die folgende Tabelle zeigt für $n \leq 5$ jeweils die Anzahl der Fälle, die zu dem in der Kopfzeile angegebenen Gesamtergebnis führen:

$n = 0$											
$n = 1$											
$n = 2$											
$n = 3$											
$n = 4$											
$n = 5$											

Diese dreiecksförmige Anordnung von Zahlen bezeichnet man als *PASCALsches Dreieck*. Offenbar kann man es dadurch rekursiv zeilenweise berechnen, daß man an jede Stelle die Summe der beiden links und rechts darüberstehenden Zahlen schreibt: Die n -te Störung bringt den Fehler genau dann auf $i \cdot \varepsilon$, wenn sie entweder gleich $+\varepsilon$ ist und die ersten $n - 1$ Störungen einen Fehler $(i - 1) \cdot \varepsilon$ produziert haben, oder aber wenn sie gleich $-\varepsilon$ ist und die ersten $n - 1$ Störungen einen Fehler $(i + 1) \cdot \varepsilon$ produziert haben. Entsprechend ist auch klar, daß die Summe aller Zahlen in der n -ten Zeile gleich 2^n ist, denn in der nullten Zeile haben wir Summe eins, und da die jeweils neu hinzukommende Störung genau zwei Möglichkeiten hat, verdoppelt sich die Summe von Zeile zu Zeile. Die Wahrscheinlichkeit dafür, daß sich n Störungen zu $i \cdot \varepsilon$ aufsummieren, ist also gerade gleich der Zahl, die in der n -ten Spalte

unter $i \cdot \varepsilon$ steht (beziehungsweise Null, wenn dort keine Zahl steht), dividiert durch 2^n .

Bekanntlich kann man die Zahlen in diesem Dreieck auch explizit berechnen: An der n -ten Zeile stehen die $n + 1$ Zahlen

$$\binom{n}{i} = \frac{n!}{i!(n-i)!} \quad \text{für } i = 0, \dots, n.$$

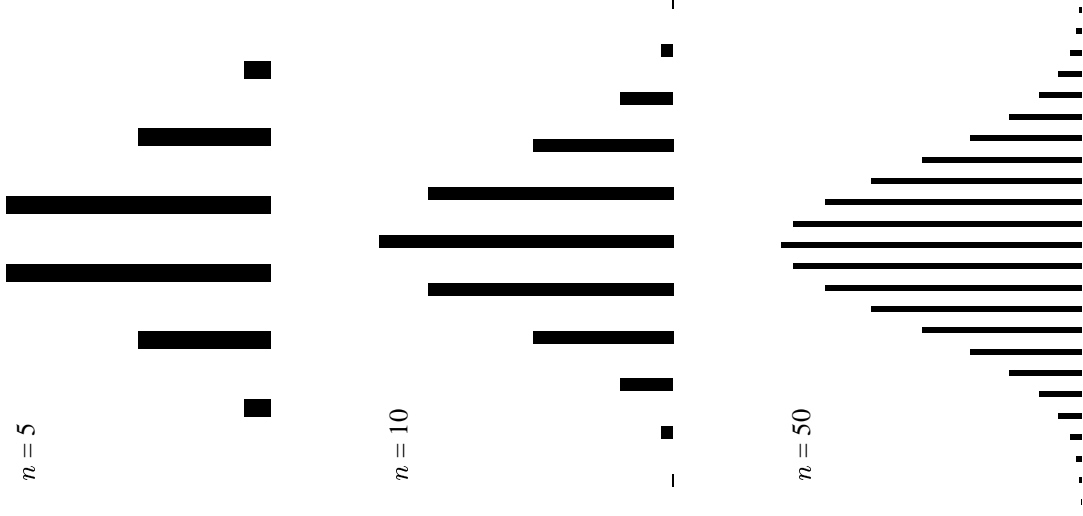
Wer diese Formel nicht kennt, kann sie leicht durch vollständige Induktion beweisen: Für $n = 1$ sowie allgemein für $i = 0$ oder $i = n$ ist alles klar; für $n > 1$ und $0 < i < n$ stehen über $\binom{n}{i}$ die beiden Zahlen $\binom{n-1}{i-1}$ und $\binom{n-1}{i}$, für die in der Tat gilt

$$\begin{aligned} \binom{n-1}{i-1} + \binom{n-1}{i} &= \frac{(n-1)!}{(i-1)!(n-i)!} + \frac{(n-1)!}{i!(n-i-1)!} = \frac{(n-1)!}{i!(n-i)!} \cdot (i + (n-i)) \\ &= \frac{n!}{i!(n-i)!} = \binom{n}{i}. \end{aligned}$$

Mit dieser Formel lassen sich die Fallzahlen für einen bestimmten Fehler im Prinzip berechnen; allerdings ist die Berechnung für große n schnell sehr aufwendig und die Binomialkoeffizienten werden auch schnell sehr groß. Um trotzdem einen Eindruck davon zu bekommen, was für größere n passiert, sind auf den beiden folgenden Seiten die Binomialkoeffizienten für $n = 5, 10, 50, 100, 500$ und 1000 graphisch dargestellt. (Die Tatsache, daß ab $n = 50$ deutlich weniger als $n + 1$ Balken zu sehen sind, erklärt sich daraus, daß die restlichen Binomialkoeffizienten zu klein sind, um noch darstellbar zu sein: Die Diagramme sind so skaliert, daß der größte (mittlere) Balken jeweils eine feste Höhe hat.)

Betrachten wir als nächstes die Größe des Gesamtfehlers. Falls n gerade ist, treten nur Vielfache von 2ε auf und alle diese Vielfachen zwischen $-n\varepsilon$ und $n\varepsilon$ kommen tatsächlich vor; entsprechend sind für ungerades n nur ungeradzahlige Vielfache von ε möglich, und auch hier werden wieder alle solchen Werte zwischen $-n\varepsilon$ und $n\varepsilon$ angenommen. Wir können dies dadurch zusammenfassen, daß in beiden Fällen genau die Werte $(n - 2k)\varepsilon$ mit $k = 0, \dots, n$ angenommen werden, und das PASCALSche Dreieck zeigt, daß der Fehler $(n - 2k)\varepsilon$ in

$$\binom{n}{n-k} = \binom{n}{k}$$

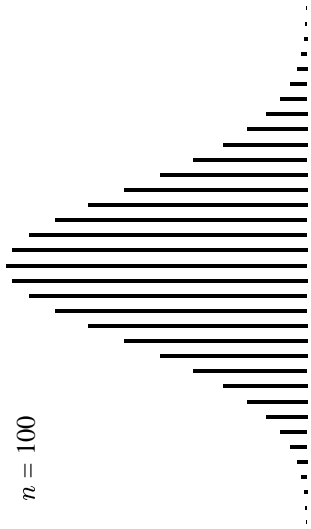


Fällen auftritt. Da n „Dämonen“ insgesamt 2^n Möglichkeiten zur Fehlerzeugung haben, ist die Wahrscheinlichkeit für den Gesamtfehler $(n - 2k)\varepsilon$ also

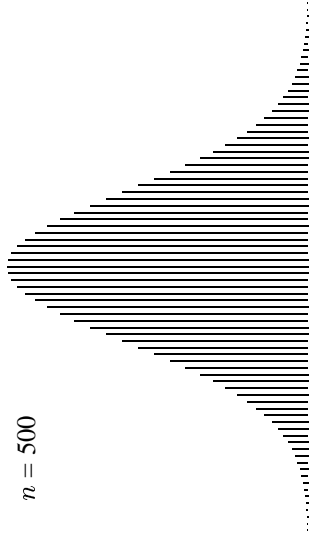
$$\binom{n}{k} \cdot 2^{-n}.$$

Diese Wahrscheinlichkeit sollte für einen festen Fehlerbetrag im wesentlichen unabhängig von n sein: Da wir nicht wirklich an Dämonen glauben, können wir deren Anzahl schließlich nicht in ein realistisches Fehlermodell einfließen lassen.

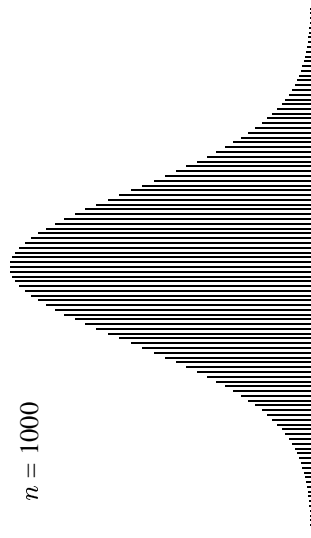
Auch die Balkendiagramme zeigen, daß sich die Verteilung der Fehlerwahrscheinlichkeiten für große n einer festen Kurve annähern sollte, der in Abbildung 65 und auf jedem Zehnmarkstein zu findenden *Glockenkurve* oder GAUSS-Kurve.



$n = 100$



$n = 500$



$n = 1000$

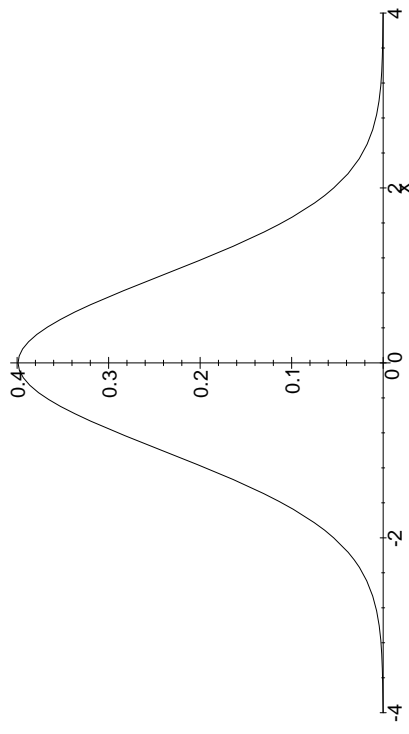


Abb. 65: Die „Glockenkurve“ $y = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

Wenn sich Fehler oder auch beliebige Daten so verteilen, wie es dieser Kurve entspricht, redet man von *normalverteilten* Daten. Damit haben wir also zumindest graphisch gesehen, daß Meßfehler nach dem LAPLACESchen Fehlermodell normalverteilt sind. Das sagt noch nicht unbedingt etwas über die Verteilung realer Meßfehler, da das LAPLACESche

Fehlermodell von unrealistisch einfachen Annahmen ausgeht; nach einem der fundamentalen Gesetze der Statistik, dem *zentralen Grenzwertsatz*, führen aber auch realistischere Annahmen zu genau derselben Verteilung: Sind u_1, \dots, u_n beliebige Quellen von Zufallsfehlern, über deren Verteilung wir (fast) nichts voraussetzen müssen, so ist ihre Summe für hinreichend großes n annähernd normalverteilt. Das eingeklammerte Wort „fast“ ist dabei für praktische Zwecke bedeutungslos, und als „groß“ kann man sich ein n ab etwa dreißig oder vierzig vorstellen.

Im nächsten Paragraphen werden wir uns überlegen, wie man zu einer Gleichung für die Glockenkurve kommt.

b) Statistische Kenngrößen

Die übliche Strategie zum Umgang mit Zufallsfehlern ist wohlbekannt: Man begnügt sich nicht mit einer einzigen Messung, sondern mißt dieselbe Größe mehrmals, so daß man eine ganze Meßreihe

$$x_1, x_2, \dots, x_N$$

erhält. Dann bildet man das *arithmetische Mittel*

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

der Meßreihe in der Hoffnung, daß sich hierbei die Fehler „ausmitteln“, so daß \bar{x} dem theoretisch korrekten Wert \hat{x} nahekommt.

Die Wahl des arithmetischen Mittels läßt sich auch geometrisch begründen: Eine Meßreihe x_1, \dots, x_N für eine Meßgröße mit exaktem Wert \hat{x} definiert einen Vektor im \mathbb{R}^n . Falls es keine Meßfehler gäbe, hätte dieser lauter identische Komponenten \hat{x} . Tatsächlich ist dies natürlich nicht der Fall; wir können aber nach einem Vektor mit identischen Komponenten suchen, der möglichst nahe am Vektor der Meßwerte liegt. Für einen Vektor, dessen sämtliche Komponenten gleich x sind, ist der Euklidische Abstands zum Vektor der Meßwerte gleich

$$d(x) = \sqrt{\sum_{i=1}^N (x - x_i)^2} = \sqrt{N x^2 - 2x \sum_{i=1}^N x_i + \sum_{i=1}^N x_i^2}.$$

Die quadratische Funktion $d(x)^2$ hat ein eindeutig bestimmtes Minimum bei der Nullstelle ihrer Ableitung

$$2N x - 2 \sum_{i=1}^N x_i,$$

also beim arithmetischen Mittel \bar{x} , und dieses ist auch das einzige Minimum von $d(x)$. Wir nehmen daher das arithmetische Mittel \bar{x} als besten verfügbaren Schätzwert für den unbekannt korrekten Wert \hat{x} .

Als Maß für die Schwankungen innerhalb der Meßreihe und damit für die Meßfehler könnte man versucht sein, den Abstand $d(\hat{x})$ zu nehmen; er hat aber den Nachteil, daß er mit steigendem N immer größer wird, d.h. die Schwankungen würden umso größer, je mehr man mißt. Das ist natürlich absurd; daher dividieren wir das Abstandsquadrat noch durch N und definieren die *mittlere quadratische Abweichung* oder *Varianz* der Meßreihe als

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (\hat{x} - x_i)^2.$$

Die (nichtnegative) Quadratwurzel σ hieraus heißt *Standardabweichung* der Meßreihe.

Das Ergebnis einer Messung wird meist angegeben in der Form

$$x = \bar{x} \pm \sigma,$$

man betrachtet also die Standardabweichung der Meßreihe als Maß für den Meßfehler. Da deren Definition allerdings vom (im allgemeinen unbekannt) korrekten Wert \hat{x} abhängt, können wir sie nicht berechnen, sondern müssen im folgenden sehen, wie wir sie zumindest schätzen können.

Als einfachste Möglichkeit bietet sich an, σ^2 durch

$$\frac{1}{N} \sum_{i=1}^N (\bar{x} - x_i)^2$$

zu schätzen, aber das führt sicherlich zu einem zu kleinen Ergebnis: Schließlich ist $d(\bar{x})$ das eindeutig bestimmte Minimum der Abstandsfunktion d , so daß der korrekte Wert $d(\hat{x})$ für $\hat{x} \neq \bar{x}$ notwendigerweise größer sein muß.

In Abschnitt d) werden wir aus dem Fehlerfortpflanzungsgesetz einen besseren Schätzwert für σ herleiten.

Warum betrachten wir eigentlich quadratische Abweichungen und nicht die einfacheren linearen Abweichungen? Nun, der Mittelwert aller Abweichungen ist

$$\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x}) = \frac{1}{N} \sum_{i=1}^N x_i - \frac{1}{N} N \bar{x} = \bar{x} - \bar{x} = 0,$$

also ist dies keine geeignete Maßzahl. Möglich wäre die mittlere *betragsmäßige* Abweichung

$$\frac{1}{N} \sum_{i=1}^N |x_i - \bar{x}|,$$

allerdings wird die im allgemeinen nicht für das arithmetische Mittel \bar{x} minimal, sondern, wie man sich leicht überlegen kann, für jede Zahl \bar{x} mit der Eigenschaft, daß gleichviele Meßwerte größer und kleiner als \bar{x} sind; einen solchen Wert \bar{x} bezeichnet man als *Median* der Meßreihe. Für die Beschreibung wirtschafts- und sozialwissenschaftlicher Daten ist dieser Median meist eine aussagekräftigere Kennzahl als das arithmetische Mittel; in den Naturwissenschaften und der Technik spielt er allerdings keine große Rolle. Im nächsten Paragraphen werden wir sehen, daß auch das LAPLACESche Fehlermodell in natürlicher Weise auf quadratische Abweichungen führt.

c) Das Fehlerfortpflanzungsgesetz

Gegeben seien zwei Größen

$$x = \hat{x} \pm \sigma_x \quad \text{und} \quad y = \hat{y} \pm \sigma_y$$

(die Verallgemeinerung auf mehr als zwei Größen erfordert, wie man sich bei der folgenden Rechnung leicht klarmacht, nur etwas mehr Schreibaufwand; sie ist nicht prinzipiell schwieriger), und eine Größe

$$w = f(x, y),$$

die von diesen beiden abhängt. Um vernünftige Aussagen machen zu können, setzen wir dabei f als stetig differenzierbar voraus.

Für x seien N Meßwerte x_1, \dots, x_N gegeben, und für y entsprechend M Werte y_1, \dots, y_M . Wenn wir echte Zufallsfehler haben, können wir

davon ausgehen, daß die Fehler der x -Werte und die der y -Werte voneinander unabhängig sind, und das wollen wir im folgenden auch annehmen.

Für w haben wir dann NM Werte $w_{ij} = f(x_i, y_j)$, deren Mittelwert die beste Schätzung für den „wahren“ Wert $\hat{w} = f(\hat{x}, \hat{y})$ ist. Dieser Mittelwert ist für komplizierte Funktionen f und/oder große Werte von n und m umständlich auszurechnen; günstiger wäre es, einfach den Mittelwert \bar{x} der x_i und den Mittelwert \bar{y} der y_j zu berechnen, um dann $f(\bar{x}, \bar{y})$ als Schätzung für \hat{w} zu benutzen. Zur Abschätzung des dadurch bedingten Fehler setzen wir

$$x_i = \bar{x} + h_i \quad \text{und} \quad y_j = \bar{y} + k_j;$$

dann ist wegen der Differenzierbarkeit von f

$$\begin{aligned} w_{ij} = f(x_i, y_j) &= f(\bar{x} + h_i, \bar{y} + k_j) \\ &= f(\bar{x}, \bar{y}) + f_x(\bar{x}, \bar{y})h_i + f_y(\bar{x}, \bar{y})k_j + o\left(\sqrt{h_i^2 + k_j^2}\right), \end{aligned}$$

wobei

$$f_x = \frac{\partial f}{\partial x} \quad \text{und} \quad f_y = \frac{\partial f}{\partial y}$$

die partiellen Ableitungen von f bezeichnen. Da die h_i und die k_j als Abweichungen vom Mittelwert die Summe null haben, ist also der Mittelwert der w_{ij} bis auf einen Fehler der Größenordnung $o\left(\sqrt{h^2 + k^2}\right)$ gleich $f(\bar{x}, \bar{y})$, wobei h, k die Betragsmaxima der h_i, k_j sind.

Als nächstes müssen wir den Fehler von \hat{w} berechnen, also den Erwartungswert der $(w_{ij} - \hat{w})^2$. Dazu schreiben wir zunächst

$$x_i = \hat{x} + u_i \quad \text{und} \quad y_j = \hat{y} + v_j,$$

betrachten also anstelle der Abweichungen vom Mittelwert die echten Meßfehler, und erhalten genau wie eben

$$w_{ij} - \hat{w} = f(x_i, y_j) - f(\hat{x}, \hat{y}) \approx u_i \cdot f_x(\hat{x}, \hat{y}) + v_j \cdot f_y(\hat{x}, \hat{y})$$

mit Quadrat

$$\begin{aligned} (w_{ij} - \hat{w})^2 &\approx (u_i \cdot f_x(\hat{x}, \hat{y}) + v_j \cdot f_y(\hat{x}, \hat{y}))^2 \\ &= u_i^2 \cdot f_x(\hat{x}, \hat{y})^2 + v_j^2 \cdot f_y(\hat{x}, \hat{y})^2 + 2u_i \cdot v_j \cdot f_x(\hat{x}, \hat{y}) \cdot f_y(\hat{x}, \hat{y}). \end{aligned}$$

Hier sind die Werte u_i^2, v_j^2 und $u_i v_j$ jeweils Zufallsgrößen, über deren Werte wir nichts sagen können. Wir haben aber gewisse Erwartungen darüber, wie sie sich *im Mittel* verhalten: u_i^2 sollte, da σ_x^2 die mittlere quadratische Abweichung von \hat{x} ist, im Mittel gleich σ_x^2 sein und v_j^2 entsprechend σ_y^2 . Genauso sollten u_i und v_j im Mittel gleich null sein, und wenn wir annehmen, daß die Fehler u_i und v_j voneinander unabhängig sind, sollte auch ihr Produkt im Mittel verschwinden. Diese sogenannten *Erwartungswerte* sind offensichtlich die bestmöglichen Schätzwerte für die jeweiligen Größen; als beste Schätzung für σ_w^2 erhalten wir damit das *GAUSSSCHE Fehlerfortpflanzungsgesetz*

$$\sigma_w^2 = f_x(\hat{x}, \hat{y})^2 \cdot \sigma_x^2 + f_y(\hat{x}, \hat{y})^2 \cdot \sigma_y^2$$

oder

$$\sigma_w = \sqrt{f_x(\hat{x}, \hat{y})^2 \cdot \sigma_x^2 + f_y(\hat{x}, \hat{y})^2 \cdot \sigma_y^2}.$$

Genauso gilt dieses Gesetz auch für Funktionen von mehr als zwei Größen; für $w = f(x_1, \dots, x_n)$ ist

$$\sigma_w = \sqrt{f_{x_1}(\hat{x}_1, \dots, \hat{x}_n)^2 \cdot \sigma_{x_1}^2 + \dots + f_{x_n}(\hat{x}_1, \dots, \hat{x}_n)^2 \cdot \sigma_{x_n}^2}.$$

d) Die Standardabweichung des Mittelwerts und die Schätzung der Varianz

Als einfache Anwendung des Fehlerfortpflanzungsgesetzes betrachten wir die Funktion

$$\bar{x} = f(x_1, \dots, x_N) = \frac{x_1 + \dots + x_N}{N},$$

also den Mittelwert der x_i . Jede Messung x_i sei mit demselben erwarteten Fehler σ behaftet; da alle partiellen Ableitungen von f gleich $1/N$ sind, folgt für den Fehler des Mittelwerts

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{N}}.$$

Dies bestätigt die implizit stets angewandte Regel, daß man durch mehrfaches Messen ein zuverlässigeres Ergebnis erhält; durch 25 Messungen beispielsweise läßt sich der Fehler auf ein Fünftel reduzieren, und für $N \rightarrow \infty$ geht er gegen Null (*Gesetz der großen Zahl*).

Damit wissen wir, wie man aus den Meßwerten auf den Fehler des Mittelwerts schließen kann – sofern man die Fehler der Meßwerte kennt. Wie lassen sich diese schätzen?

Zunächst ist

$$\begin{aligned} \sigma^2 &= \frac{1}{N} \sum_{i=1}^N (\hat{x} - x_i)^2 = \frac{1}{N} \sum_{i=1}^N ((\hat{x} - \bar{x}) + (\bar{x} - x_i))^2 \\ &= \frac{1}{N} \sum_{i=1}^N (\hat{x} - \bar{x})^2 + \frac{1}{N} \sum_{i=1}^N (\bar{x} - x_i)^2 + \frac{2}{N} \sum_{i=1}^N (\hat{x} - \bar{x}) \cdot (\bar{x} - x_i). \end{aligned}$$

Die letzte dieser drei Summen ist

$$2 \cdot \frac{(\hat{x} - \bar{x})}{N} \sum_{i=1}^N (\bar{x} - x_i) = 0,$$

da \bar{x} der Mittelwert der x_i ist. Die zweite Summe ist der Mittelwert der $(\bar{x} - x_i)^2$, also die Varianz der Meßreihe, und von der ersten schließlich wissen wir, daß $(\hat{x} - \bar{x})^2$, das Quadrat des Fehlers des Mittelwerts, den Erwartungswert σ^2/N hat. Die gesamte erste Summe ist somit

$$\frac{1}{N} \cdot N \cdot \frac{\sigma^2}{N} = \frac{\sigma^2}{N},$$

und obige Formel wird zu

$$\sigma^2 = \frac{\sigma^2}{N} + \frac{1}{N} \sum_{i=1}^N (\bar{x} - x_i)^2.$$

Bringt man hier noch den Term σ^2/N auf die linke Seite, so folgt

$$\frac{N-1}{N} \cdot \sigma^2 = \frac{1}{N} \sum_{i=1}^N (\bar{x} - x_i)^2$$

oder

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (\bar{x} - x_i)^2,$$

also

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (\bar{x} - x_i)^2}{N - 1}}.$$

Somit läßt sich auch σ aus den Meßdaten berechnen, der Meßfehler kann also ohne Kenntnis des „wahren“ Werts anhand der gemessenen Werte geschätzt werden.

e) Die Methode der kleinsten Quadrate

Oftmals ist zu gegebenen Beobachtungsdaten grundsätzlich bekannt, welcher *Art* von Gesetz sie genügen sollten, und das Problem besteht „nur“ noch darin, die in diesem Gesetz vorkommenden *Parameter* zu bestimmen. Im einfachsten Fall könnte man etwa an einen Widerstand denken, der dadurch gemessen wird, daß man verschiedene Spannungen U_i anlegt und die zugehörigen Stromstärken I_i mißt. Nach dem Ohmschen Gesetz ist dann $U_i = R \cdot I_i$, aber aufgrund der unvermeidlichen Meßfehler werden die verschiedenen Quotienten U_i/I_i natürlich nicht alle gleich sein. Die Lösung dieses Problems ist klar: Man nimmt den Mittelwert der Quotienten. Schwieriger wird es, wenn mehrere Parameter ins Spiel kommen.

Betrachten wir als Beispiel etwa die Entwicklung der Weltbevölkerung. Da (vor allem auf dem Stand von 1960) vieles für ein stärker als exponentielles Bevölkerungswachstum sprach, entwickelten damals HEINZ VON FOERSTER, PATRICIA M. MORA und LAWRENCE W. AMIOT, allesamt Elektrotechniker an der University of Illinois in Urbana, ein überexponentielles Wachstumsmodell. Die Wachstumsdifferentialgleichung

$$\dot{y}(t) = \beta y(t)^r$$

für die Weltbevölkerung $y(t)$ zum Zeitpunkt t kann durch Trennung der Veränderlichen leicht gelöst werden; die Lösung ist

$$y(t) = \left(\frac{\beta}{t_0 - t} \right)^s \quad \text{mit} \quad s = \frac{1}{r-1}.$$

β , t_0 und r sind Parameter, die aufgrund der bekannten Zahlen zur Bevölkerungsentwicklung geschätzt werden müssen. Besonders interessant ist dabei natürlich der Parameter t_0 , denn zum Zeitpunkt $t = t_0$

wird die Weltbevölkerung unendlich, was man wohl mit dem Weltuntergang gleichsetzen kann.

Für zahlreiche Zeitpunkte t_i sind die Bevölkerungszahlen y_i oder zumindest Schätzungen dafür bekannt, aber natürlich gibt es auch hier keine Parameterwerte, für die *alle* y_i gleich $f(t_i)$ sind; Ziel kann es nur sein, die Parameter so zu wählen, daß die Unterschiede zwischen y_i und $y(t_i)$ nicht zu groß werden.

FOERSTER, MORA und AMIOT taten genau das und kamen auf das Ergebnis

$$\beta = (179 \pm 14) \cdot 10^9$$

$$s = 0,990 \pm 0,009$$

$$t_0 = 2026,87 \pm 5,5 \text{ Jahre}$$

$$r = 2,01 \pm 0,009;$$

der Weltuntergang ist also etwa 2026 fällig. Der beste Schätzwert für das Datum ist, wenn wir 0,87 Jahre in Tage, Stunden und Minuten umrechnen, Freitag, der 13. November gegen 13¹³ uhr – sofern man das bei einer Standardabweichung von fünfeinhalb Jahren überhaupt so genau wissen will. (Für Einzelheiten sei auf die Arbeit *Doomsday: Friday, 13 November, A.D. 2026*, Science **138**, November 1960, Seite 1291–1295, verwiesen.)

Da Weltuntergangsszenarien nicht zum Stoff dieser Vorlesung gehören, wollen wir auf die genaue Behandlung dieses Beispiels verzichten; wir wollen aber zumindest grundsätzlich verstehen lernen, wie man solche Parameterschätzungen durchführen kann, und wir wollen diese Techniken auch auf etwas alltäglichere Probleme anwenden.

Der allgemeine Ansatz ist folgender: Sei $u = f(a, b, \dots; t)$ der von den Parametern a, b, \dots abhängige Zusammenhang zwischen den Größen t und u , und seien $(t_1, u_1), \dots, (t_N, u_N)$ gegebene (Meß-)Werte, an die die Parameter angepaßt werden sollen.

Nehmen wir an, es gäbe Parameter a, b, \dots , für die $u = f(a, b, \dots; t)$ den exakten Zusammenhang zwischen den Größen u und t wieder gibt, daß aber die Werte u_i aufgrund von Beobachtungsfehlern von den

$f(a, b, \dots; t_i)$ abweichen – das ist genau die Situation, mit denen man es in den Naturwissenschaften und in der Technik meist zu tun hat.

Bezüglich der weiteren Vorgehensweise können wir wieder wie beim arithmetischen Mittel und der Standardabweichung argumentieren: Wir versuchen, den EUKLIDISCHEN Abstand zwischen dem Vektor der tatsächlichen Beobachtungswerte und dem der vorhergesagten Werte zu minimieren, d.h. wir wählen die Parameter a, b, \dots so, daß

$$Q(a, b, \dots) \stackrel{\text{def}}{=} \sum_{i=1}^N (u_i - f(a, b, \dots; t_i))^2$$

minimal wird – das ist die von CARL FRIEDRICH GAUSS eingeführte *Methode der kleinsten Quadrate*.

Obiger Ausdruck ist eine Funktion der Parameter a, b, \dots ; falls sie differenzierbar ist, kann sie nur dort ein Minimum annehmen, wo die partiellen Ableitungen

$$\frac{\partial Q}{\partial a}, \quad \frac{\partial Q}{\partial b}, \quad \dots$$

allesamt verschwinden. Dies wird im allgemeinen ein nichtlineares Gleichungssystem definieren, dessen Lösung recht kompliziert sein kann; wir wollen uns hier auf den Fall beschränken, daß die Funktion *f linear* von den Parametern a, b, \dots abhängt. In diesem Fall bezeichnet man die Suche nach den optimalen Parametern als *lineare Regression*; sie hat sowohl in Naturwissenschaften und Technik als auch in den Wirtschafts- und Sozialwissenschaften zahlreiche interessante Anwendungen.

Falls alle Parameter linear in f eingehen, kommen sie in Q höchstens quadratisch vor; die partiellen Ableitungen sind also wieder linear, so daß insgesamt nur ein lineares Gleichungssystem gelöst werden muß.

Wir betrachten hier nur beispielhaft den einfachsten Fall einer Funktion, die von zwei Parametern abhängt: die lineare Funktion

$$u = f(a, b, t) = at + b;$$

die Rechnungen, die wir dafür im folgenden durchführen, funktionieren aber genauso auch für allgemeinere Funktionen f , solange f linear in den Parametern a, b, \dots ist; Nichtlinearität in t ist problemlos.

Gesucht sind also reelle Zahlen a und b , bei denen die Funktion

$$Q(a, b) = \sum_{i=1}^N (u_i - at_i - b)^2$$

minimal wird. Dort müssen dann insbesondere die beiden partiellen Ableitungen $Q_a(a, b)$ und $Q_b(a, b)$ verschwinden. Bevor wir diese ausrechnen, empfiehlt es sich, die Funktion zunächst durch Ausmultiplizieren nach Termen in a und b zu sortieren:

$$\begin{aligned} Q(a, b) &= \sum_{i=1}^N u_i^2 + a^2 \sum_{i=1}^N t_i^2 + Nb^2 - 2a \sum_{i=1}^N u_i t_i \\ &\quad - 2b \sum_{i=1}^N u_i + 2ab \sum_{i=1}^N t_i. \end{aligned}$$

Damit ist

$$Q_a(a, b) = 2a \sum_{i=1}^N t_i^2 - 2 \sum_{i=1}^N u_i t_i + 2b \sum_{i=1}^N t_i$$

und

$$Q_b(a, b) = 2Nb - 2 \sum_{i=1}^N u_i + 2a \sum_{i=1}^N t_i.$$

Division durch zwei ergibt die beiden Gleichungen für das Verschwinden dieser partiellen Ableitungen:

$$a \sum_{i=1}^N t_i^2 + b \sum_{i=1}^N t_i = \sum_{i=1}^N t_i u_i$$

$$\text{und} \quad a \sum_{i=1}^N t_i + bN = \sum_{i=1}^N u_i.$$

Um b zu eliminieren, multiplizieren wir die zweite Gleichung mit $\frac{1}{N} \sum_{i=1}^N t_i$

$$a \cdot \frac{1}{N} \left(\sum_{i=1}^N t_i \right)^2 + b \sum_{i=1}^N t_i = \frac{1}{N} \left(\sum_{i=1}^N t_i \right) \cdot \left(\sum_{i=1}^N u_i \right)$$

und subtrahieren das von der ersten Gleichung:

$$a \left(\sum_{i=1}^N t_i^2 - \frac{1}{N} \left(\sum_{i=1}^N t_i \right)^2 \right) = \sum_{i=1}^N t_i u_i - \frac{1}{N} \left(\sum_{i=1}^N t_i \right) \cdot \left(\sum_{i=1}^N u_i \right),$$

d.h.

$$a = \frac{\sum_{i=1}^N t_i u_i - \frac{1}{N} \left(\sum_{i=1}^N t_i \right) \cdot \left(\sum_{i=1}^N u_i \right)}{\sum_{i=1}^N t_i^2 - \frac{1}{N} \left(\sum_{i=1}^N t_i \right)^2}.$$

Entsprechend läßt sich a eliminieren, wenn wir die erste Gleichung mit $\sum t_i$ und die zweite mit $\sum t_i^2$ multiplizieren; wir erhalten

$$b = \frac{\left(\sum_{i=1}^N t_i \right) \cdot \left(\sum_{i=1}^N t_i u_i \right) - \left(\sum_{i=1}^N t_i^2 \right) \cdot \left(\sum_{i=1}^N u_i \right)}{\left(\sum_{i=1}^N t_i \right)^2 - N \sum_{i=1}^N t_i^2}.$$

Als Beispiel betrachten wir den Zusammenhang zwischen dem Grad der Korruption in einem Land und dessen Bruttozialprodukt pro Einwohner. Jedes Jahr veröffentlicht die Organisation *Transparency International* ihren *corruption perceptions index (CPI)*, in dem jedem Land eine Zahl zwischen null und zehn zugeordnet wird, je nachdem, wie stark Geschäftsleute, Risikospezialisten und die Bevölkerung die Korruption im betreffenden Land einschätzen: Ein Index von zehn bedeutet, daß es praktische keine Korruption gibt, während bei null nichts läuft ohne Bimbos. Die neuesten Daten stammen von Juni 2001 und sind unter <http://www.transparency.org/documents/cpi/2001/cpi2001.de.html> zu finden. Die Zahlen werden als Mittelwerte über die letzten drei Jahren berechnet, so daß singuläre Ereignisse eines Jahres nicht zu sehr ins Gewicht fallen.

Wir vergleichen diese Zahlen mit dem Bruttozialprodukt pro Einwohner, wie es die Weltbank für 1998 festgestellt hat. Dies sind die

neuesten verfügbaren Zahlen; sie sind beispielsweise auf dem Server des Statistischen Bundesamtes unter

<http://www.statistik-bund.de/basis/d/ausl/auslae5.htm>

zu finden. In der folgenden Tabelle sind Staaten aufgelistet, für die sowohl das Bruttozialprodukt pro Einwohner als auch der CPI für 2001 vorliegt; das Bruttozialprodukt pro Einwohner in US-\$ ist kursiv gedruckt, der Korruptionsindex fett.

Ägypten	1290	3,6
Argentinien	8030	3,5
Aserbaidtschan	480	2,0
Australien	20640	8,5
Bangladesch	350	0,4
Belgien	25380	6,6
Bolivien	1010	2,0
Botswana	3070	6,0
Brasilien	4630	4,0
Bulgarien	1220	3,9
Chile	4990	7,5
China (ohne Hongkong & Taiwan)	750	3,5
Costa Rica	2720	4,5
Dänemark	33040	9,5
Deutschland	28570	7,4
Dominikanische Republik	1770	3,1
Ecuador	1520	2,3
Elfenbeinküste	700	2,4
El Salvador	1850	3,6
Estland	3360	5,6
Finnland	24280	9,9
Frankreich	24210	6,7
Ghana	390	3,4
Griechenland	11740	4,2
Großbritannien und Nordirland	21410	8,3
Guatemala	1640	2,9
Honduras	740	2,7
Hongkong	23660	7,9

Indien	440	2,7	Rumänien	1360	2,8
Indonesien	640	1,9	Rußland	2260	2,3
Irland	18710	7,5	Sambia	330	2,6
Island	27830	9,2	Schweden	25580	9,0
Israel	16180	7,6	Schweiz	39980	8,4
Italien	20090	5,5	Senegal	520	2,9
Japan	32350	7,1	Simbabwe	620	2,9
Jordanien	1150	4,9	Singapur	30170	9,2
Kamerun	610	2,0	Slowakei	3700	3,7
Kanada	19170	8,9	Slowenien	9780	5,2
Kasachstan	1340	2,7	Spanien	14100	7,0
Kenia	350	2,0	Südafrika	3310	4,8
Kolumbien	2470	3,8	Taiwan	1233	5,9
Korea, Republik	8600	4,2	Tansania, Vereinigte Republik	220	2,2
Kroatien	4620	3,9	Thailand	2740	3,2
Lettland	2420	3,4	Trinidad und Tobago	4520	5,3
Litauen	2540	4,8	Tschechien	5150	3,9
Luxemburg (BSP von 1996)	45100	8,7	Tunesien	2060	5,3
Malawi	210	3,2	Türkei	3160	3,6
Malaysia	3670	5,0	Uganda	310	1,9
Mauritius	3730	4,5	Ungarn	4510	5,3
Mexiko	3840	3,7	Ukraine	980	2,1
Moldau	380	3,1	Uruguay	6070	5,1
Namibia	1940	5,4	Usbekistan	950	2,7
Neuseeland	14600	9,4	Venezuela	3530	2,8
Nicaragua	370	2,4	Vereinigte Staaten	29240	7,6
Niederlande	24780	8,8	Vietnam	350	2,6
Nigeria	300	1,0			
Norwegen	34310	8,6			
Österreich	26830	7,8			
Pakistan	470	2,3			
Panama	2990	3,7			
Peru	2440	4,1			
Philippinen	1050	2,9			
Polen	3910	4,1			
Portugal	10670	6,3			

Der erste Augenschein zeigt, daß korruptionsärmere Länder oftmals reicher sind: Das weitgehend korruptionsfreie Dänemark hat ein Brutto- sozialprodukt von 33 040 \$ pro Einwohner, das deutlich korruptere Deutschland nur 28 570 \$ und ein stark korruptes Land wie Bangla- desch nur 350 \$. Allerdings gibt es auch Ausnahmen, denn Chile hat, obwohl geringfügig weniger korrupt als Deutschland, nur ein Brutto- sozialprodukt von 4 990 \$ pro Einwohner. Es gibt also sicherlich keinen deterministischen Zusammenhang zwischen Korruption und Wohlstand,

aber es lohnt sich doch, eine Ausgleichsgerade zu berechnen. Obige Formeln führen auf

$$\text{CPI} = 3,18723 + 0,18105 \cdot 10^{-3} \cdot \text{BSP};$$

Abbildung 66 zeigt letztere Gerade zusammen mit den Datenpunkten.

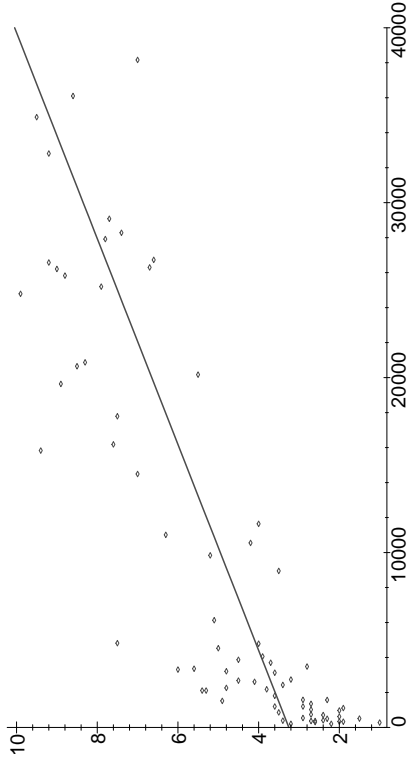


Abb. 66: Korruption 2001 über Bruttonozialprodukt 1998 für 91 Staaten

Natürlich liegen nicht alle Datenpunkte auf der Ausgleichsgeraden: Beispielsweise sind Japan, Belgien und Italien deutlich korrupter, als ihre Wirtschaftskraft vermuten läßt, wohingegen Dänemark, Finnland und Neuseeland deutlich ärmer sind, als man aufgrund der geringen Korruption erwarten sollte. Trotzdem beschreibt die Ausgleichsgerade zumindest einen gewissen Trend.

Für Deutschland ergibt die Formel einen erwarteten Index von 8,05, was zwar praktisch perfekt mit den 1999er CPI von 8,0 übereinstimmt, aber deutlich über dem 2001er CPI von 7,4 liegt. Da 1999 noch kaum jemand das Wort „Bimbes“ kannte, während 2000 „Schwarzgeldaffäre“ zum Wort des Jahres wurde, ist dieser Rückgang nicht weiter verwunderlich und wird wohl auch nächstes Jahr weitergehen, denn der CPI beruht ja auf Studien der vergangenen drei Jahre.

Ähnlich sieht es aus, wenn man das Klausurergebnis der Scheinklausur zur letztjährigen Vorlesung mit den bearbeiteten Übungsaufgaben in

Verbindung bringt: Hier ergibt sich die Gerade

$$\text{Klausurpunkte} = 11,78054131 + 0,05813492738 \cdot \text{Übungsprozente},$$

die in Abbildung 67 zusammen mit den Datenpunkten eingezeichnet ist.

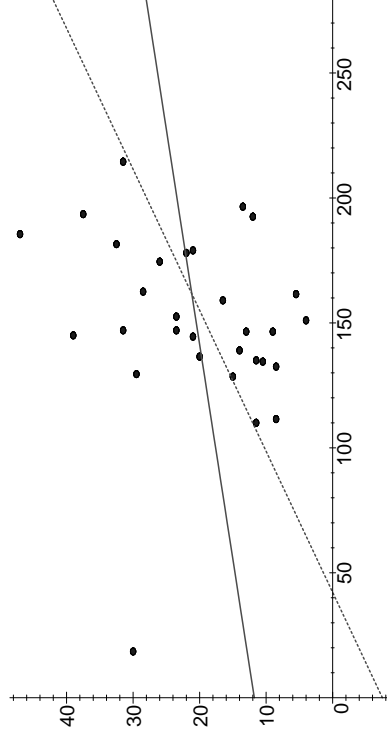


Abb. 67: Zusammenhang zwischen Klausurergebnis und Übungen

Wieder sieht man, daß im großen und ganzen höhere Punktezahlen bei den Übungsaufgaben einem besseren Klausurergebnis entsprechen, aber es gibt auch recht deutliche Ausreißer, deren Interpretation dem Leser überlassen sei.

Wie sehr Ausgleichsgeraden von einzelnen Datenpunkten abhängen können sieht man an der gestrichelt eingezeichneten Geraden

$$\text{Klausurpunkte} = -7,473678952 + 0,1772199512 \cdot \text{Übungsprozente},$$

die fast anhand derselben Daten berechnet wurde: Lediglich der eine Datenpunkt für einen Studenten, der nur ein einziges Übungsblatt bearbeitete und trotzdem die siebtbeste Klausur schrieb, blieb unbeachtet.

Unser nächstes Ziel ist ein Maß, mit dem wir die Qualität einer Ausgleichsgerade $u = at + b$ zu N Datenpaaren (t_i, u_i) beurteilen können.

Zu seiner Definition gehen wir ähnlich vor wie bei der Herleitung des GAUSSschen Fehlerfortpflanzungsgesetzes: \bar{t} sei der Mittelwert der t_i

und \bar{u} der der u_i . Falls die u_i und die t_i voneinander unabhängig sind, erwarten wir, daß

$$\sum_{i=1}^N (t_i - \bar{t})(u_i - \bar{u})$$

verschwinden sollte, denn die Abweichungen $t_i - \bar{t}$ der t_i von ihrem Mittelwert sollten nichts mit den Abweichungen $u_i - \bar{u}$ der u_i von deren Mittelwert zu tun haben, und im Mittel sind beide Abweichungen null.

Tatsächlich ist obiges Argument kein Beweis, sondern die Definition des Wortes „unabhängig“: Wir bezeichnen

$$\text{cov}(\mathbf{t}, \mathbf{u}) \stackrel{\text{def}}{=} \frac{1}{N} \sum_{i=1}^N (t_i - \bar{t})(u_i - \bar{u})$$

als die *Kovarianz* der Meßreihen

$$\mathbf{t} = (t_1, \dots, t_N) \quad \text{und} \quad \mathbf{u} = (u_1, \dots, u_N)$$

und bezeichnen die beiden Meßreihen als statistisch unabhängig, falls diese Kovarianz verschwindet.

Falls im anderen Extrem $u_i = at_i + b$ ist für alle i , so ist $\bar{u} = a\bar{t} + b$ und

$$\begin{aligned} \sum_{i=1}^N (t_i - \bar{t})(u_i - \bar{u}) &= \sum_{i=1}^N (t_i - \bar{t})(at_i - a\bar{t}) \\ &= a \sum_{i=1}^N (t_i - \bar{t})^2 \end{aligned}$$

eine (fast immer) positive Zahl, die eng mit der Standardabweichung der Verteilung der t_i zusammenhängt. Um sie weiter auszurechnen, beachten wir, daß

$$\sum_{i=1}^N (u_i - \bar{u})^2 = a^2 \sum_{i=1}^N (t_i - \bar{t})^2$$

ist und damit

$$\sum_{i=1}^N (t_i - \bar{t})^2 \cdot \sum_{i=1}^N (u_i - \bar{u})^2 = a^2 \left(\sum_{i=1}^N (t_i - \bar{t})^2 \right)^2$$

oder

$$\sqrt{\sum_{i=1}^N (t_i - \bar{t})^2 \cdot \sum_{i=1}^N (u_i - \bar{u})^2} = a \sqrt{\sum_{i=1}^N (t_i - \bar{t})^2}.$$

Damit ist also der Quotient

$$\rho = \frac{\sum_{i=1}^N (t_i - \bar{t})(u_i - \bar{u})}{\sqrt{\sum_{i=1}^N (t_i - \bar{t})^2 \cdot \sum_{i=1}^N (u_i - \bar{u})^2}}$$

je nach Vorzeichen von a gleich ± 1 , wenn alle $u_i = at_i + b$ sind, und $r = 0$, wenn die u_i unabhängig von den t_i sind. Für beliebige Wertepaare (t_i, u_i) sagt uns die CAUCHY-SCHWARZSche Ungleichung aus [HMM], Kap. 2, § 1a), daß stets

$$-1 \leq \rho \leq 1$$

ist.

Diesen Wert ρ nennen wir den *Korrelationskoeffizienten* der beiden Meßreihen \mathbf{t} und \mathbf{u} ; sein Betrag gibt an, wie gut die Beziehung zwischen den u_i und den t_i durch eine Gerade beschrieben werden kann, während sein Vorzeichen angibt, ob diese Gerade ansteigt (die Daten sind *positiv korreliert*) oder sinkt (*negative Korrelation*).

Kürzer läßt er sich auch schreiben als

$$\rho = \frac{\text{cov}(\mathbf{t}, \mathbf{u})}{\sigma_t \sigma_u},$$

wobei σ_t und σ_u die Standardabweichungen der beiden Meßreihen bezeichnen.

Bei den obigen Zahlenbeispielen ergeben sich die Korrelationskoeffizienten zu

$$\rho \approx 0,8483022667$$

für den Zusammenhang zwischen Korruption und Bruttosozialprodukt; die Korrelation ist also recht gut. Bei den Klausurergebnissen dagegen

erhalten wir nur

$$\rho \approx 0,1910296117 \quad \text{bzw.} \quad \rho \approx 0,4218151723,$$

falls der eine Ausreißer nicht mitgerechnet wird – beides nicht gerade sehr berauschend. Letztes Jahr lag der entsprechende Korrelationskoeffizient übrigens bei

$$\rho \approx 0,7465571442;$$

wer will, kann über mögliche Gründe für diesen Unterschied spekulieren.

Zum Abschluß sei noch darauf hingewiesen, daß auch ein Korrelationskoeffizient nahe bei ± 1 nicht unbedingt bedeutet, daß ein *inhaltsreicher* Zusammenhang zwischen den verglichenen Größen besteht: Die Korrelation könnte auch rein formal sein. Ein beliebtes Beispiel, zu dem leider nirgends Zahlenwerte zu finden sind, ist die gute Korrelation zwischen dem Rückgang der Storchpopulationen und der Geburtenrückgang zur Zeit der Industrialisierung: Die meisten Bevölkerungsexperten gehen hier nicht von einem kausalen Zusammenhang aus, sondern eher davon, daß die Industrialisierung den Lebensraum der Störche zerstörte und wegen der gleichzeitigen Einführung der Rentenversicherung die Familien nicht mehr so viele Kinder zur Alterssicherung brauchten.

§5: Die Gaußsche Normalverteilung und die Maximum Likelihood Methode

a) Der Grenzfall des Laplaceschen Fehlermodells

Leider ist in dieser Vorlesung nicht genug Zeit, um auch nur annähernd die notwendigen Grundlagen für einen Beweis des zentralen Grenzwertsatzes bereitzustellen. Wir haben aber immerhin schon am Beispiel des LAPLACESchen Fehlermodell graphisch gesehen, daß für die Verteilungen mit $u_i = \pm \varepsilon$, jeweils mit Wahrscheinlichkeit $\frac{1}{2}$, die Verteilung der Summen gegen eine Normalverteilung konvergiert.

Zumindest in diesem einfachen Fall können wir dies auch rechnerisch einsehen und auf diesem Weg insbesondere auch die Gleichung der

Glockenkurve herleiten. Dazu müssen wir uns zunächst überlegen, was ε tun soll, wenn n gegen unendlich geht.

Wir kennen zwei statistische Kennzahlen zur Beschreibung der Fehlerverteilung: Das arithmetische Mittel und die Varianz. Das arithmetische Mittel ist (z.B. aus Symmetriegründen) null, bleibt also noch die Varianz.

Zu deren Berechnung müssen wir die Fehlerquadrate über die 2^n möglichen Verhaltensweisen der „Dämonen“ summieren, d.h. wir müssen die Streuung

$$\sigma^2 = \frac{1}{2^n} \sum (\varepsilon_1 + \dots + \varepsilon_n)^2$$

berechnen, wobei sich die Summation über alle n -tupel

$$(\varepsilon_1, \dots, \varepsilon_n) \quad \text{mit} \quad \varepsilon_i = \pm \varepsilon$$

erstreckt. Beim Ausmultiplizieren heben sich alle gemischten Terme der Form $\varepsilon_i \varepsilon_j$ gegenseitig weg, denn das Tupel, bei dem nur an der i -ten Stelle das Vorzeichen geändert wurde, liefert einen Summanden $-\varepsilon_i \varepsilon_j$. Also bleiben nur die Quadrate; diese sind alle gleich ε^2 , und es sind pro Summand n Stück. Da die Anzahl der Summanden gleich dem Nenner des Vorfaktors ist, berechnet sich die Varianz daher zu

$$\sigma^2 = n\varepsilon^2$$

Somit müssen wir für $n \rightarrow \infty$ den Einfluß ε jedes einzelnen „Dämonen“ so gegen null gehen lassen, daß $n\varepsilon^2$ konstant bleibt, d.h. wir setzen

$$\varepsilon = \frac{\sigma}{\sqrt{n}}$$

für eine geeignete zu wählende Konstante $\sigma > 0$, die Standardabweichung.

Für festes n kann der Fehler einen der $n+1$ Werte

$$-n\varepsilon, \quad -(n-2)\varepsilon, \quad \dots, \quad (n-2)\varepsilon, \quad n\varepsilon$$

annehmen, was wir in der Form

$$u = (n-2k)\varepsilon = \frac{(n-2k)\sigma}{\sqrt{n}} \quad \text{mit} \quad k = -n, \dots, n$$

schreiben wollen. Dieser Fehler tritt genau dann auf, wenn k der Dämonen den Fehler ε erzeugen und die restlichen $n - k$ den Fehler $-\varepsilon$. Dies geschieht in $\binom{n}{k}$ der 2^n möglichen Fälle; die Wahrscheinlichkeit dafür ist also $\binom{n}{k} 2^{-n}$.

Für $n \rightarrow \infty$ geht dieser Ausdruck gegen null, denn mit n geht schließlich auch die Anzahl der zu betrachtenden Fälle gegen unendlich. Falls es ein Intervall gäbe, in dem die Wahrscheinlichkeit für jeden darin liegenden Fehler größer als irgendein $\alpha > 0$ wäre, ginge allein schon die Summe der Wahrscheinlichkeiten für Fehler aus diesem Teilintervall mit n gegen unendlich, da die Anzahl der dort liegenden möglichen u -Werte wegen der \sqrt{n} im Nenner von u gegen unendlich geht. Da die Summe aller Wahrscheinlichkeiten aber nicht größer als eins werden kann, muß die Wahrscheinlichkeit also in jedem einzelnen Punkt für $n \rightarrow \infty$ gegen null gehen.

Wenn wir n variieren lassen, ist es allerdings ohnehin sinnlos, einen genauen Wert des Fehlers zu betrachten: Für jedes n gibt es nur $n + 1$ mögliche Werte, und bei den meisten größeren Werten von n werden diese Zahlen – abgesehen von der Null – nicht auftreten. Wenn wir etwas von n unabhängiges definieren möchten (und sofern wir nicht an Dämonen glauben, bleibt uns kaum etwas anderes übrig) dürfen wir also nicht den genauen Wert des Fehlers festlegen, sondern müssen ein *Fehlerintervall* betrachten.

Nun wollen wir aber natürlich als Ergebnis keine Funktion, die von einem Intervall abhängt, sondern eine gewöhnliche Funktion von u . Um eine solche zu bekommen, betrachten wir nicht die Wahrscheinlichkeit, sondern die *Wahrscheinlichkeitsdichte*: Für eine kontinuierlich variierende zufällige Größe definieren wir die Wahrscheinlichkeitsdichte $\varphi(u_0)$ im Punkt u_0 als

$$\varphi(u_0) = \lim_{\varepsilon \rightarrow 0} \frac{\text{Wahrscheinlichkeit für } u_0 - \varepsilon \leq u \leq u_0 + \varepsilon}{2\varepsilon},$$

d.h. also als Wahrscheinlichkeit dividiert durch die Intervallbreite. Falls diese Dichte existiert, folgt mehr oder weniger sofort aus der Definition

des RIEMANN-Integrals, daß

$$(\text{Wahrscheinlichkeit für } a \leq u \leq b) = \int_a^b \varphi(u) du$$

ist. Unsere Ziel ist also, diese Wahrscheinlichkeitsdichte φ für $n \rightarrow \infty$ zu berechnen.

Für festes n haben die möglichen Fehlerwerte den Abstand 2ε , wir betrachten daher Intervalle der Länge 2ε mit den Werten $(n - 2k)\varepsilon$ als Mittelpunkten; diese überdecken den möglichen Fehlerbereich lückenlos, und die Wahrscheinlichkeit für einen Fehler in diesem Intervall ist $\binom{n}{k} 2^{-n}$.

Wir interessieren uns daher für den Grenzwert von

$$\frac{\binom{n}{k} 2^{-n}}{2\varepsilon} = \frac{\binom{n}{k} 2^{-n}}{2\sigma/\sqrt{n}} = \binom{n}{k} 2^{-n-1} \frac{\sqrt{n}}{\sigma}$$

für $n \rightarrow \infty$.

b) Die Eulersche Summenformel

Das Problem bei der Berechnung dieses Grenzwerts ist der Binomialkoeffizient

$$\binom{n}{k} = \frac{n!}{k!(n-k)!};$$

um diesen abzuschätzen brauchen wir einen handhabbaren Ausdruck für $n!$.

Dazu schreiben wir

$$\ln n! = \sum_{k=1}^n \ln k$$

und berechnen dies nach einer Methode von EULER, die nicht nur für Summen von Logarithmen anwendbar ist.

Wir betrachten allgemeiner irgendeine reellwertige differenzierbare Funktion f , deren Definitionsbereich das Intervall $[1, n]$ enthält.

Für eine reelle Zahl x bezeichnen wir wie üblich mit $[x]$ die größte ganze Zahl kleiner oder gleich x und mit

$$\{x\} \stackrel{\text{def}}{=} x - [x]$$

den gebrochenen Anteil von x ; ist k eine ganze Zahl, ist somit für $x \in [k, k+1)$

$$\{x\} = x - k.$$

Partielle Integration führt auf die Gleichung

$$\begin{aligned} \int_k^{k+1} (\{x\} - \tfrac{1}{2}) f'(x) dx &= (x - k - \tfrac{1}{2}) f(x) \Big|_k^{k+1} - \int_k^{k+1} f(x) dx \\ &= \frac{f(k+1) + f(k)}{2} - \int_k^{k+1} f(x) dx. \end{aligned}$$

Addition aller solcher Gleichungen von $k = 1$ bis $k = n - 1$ liefert

$$\int_1^n (\{x\} - \tfrac{1}{2}) f'(x) dx = \frac{f(1)}{2} + \sum_{k=2}^{n-1} f(k) + \frac{f(n)}{2} - \int_1^n f(x) dx,$$

womit man die Summe der $f(k)$ berechnen kann:

Satz (EULERSche Summenformel): Für eine differenzierbare Funktion $f: D \rightarrow \mathbb{R}$, deren Definitionsbereich das Intervall $[1, n]$ umfaßt, ist

$$\sum_{k=1}^n f(k) = \int_1^n f(x) dx + \frac{f(1) + f(n)}{2} + \int_1^n (\{x\} - \tfrac{1}{2}) f'(x) dx. \quad \blacksquare$$

Für die Abschätzung der Binomialkoeffizienten und Fakultäten interessiert uns speziell der Fall $f(x) = \ln x$; hierfür wird die EULERSche

Summenformel zu

$$\begin{aligned} \ln n! &= \int_1^n \ln x dx + \frac{\ln n}{2} + \int_1^n \frac{\{x\} - \frac{1}{2}}{x} dx \\ &= x(\ln x - 1) \Big|_1^n + \frac{\ln n}{2} + \int_1^n \frac{\{x\} - \frac{1}{2}}{x} dx \\ &= n(\ln n - 1) + 1 + \frac{\ln n}{2} + \int_1^n \frac{\{x\} - \frac{1}{2}}{x} dx. \end{aligned}$$

In dieser Formel stört noch das rechte Integral; dieses können wir wie folgt abschätzen: Für eine natürliche Zahl k ist

$$\begin{aligned} \int_k^{k+1} \frac{\{x\} - \frac{1}{2}}{x} dx &= \int_{-\frac{1}{2}}^{\frac{1}{2}} \frac{x}{k + \frac{1}{2} + x} dx \\ &= \int_0^{\frac{1}{2}} \left(\frac{x}{k + \frac{1}{2} + x} - \frac{x}{k + \frac{1}{2} - x} \right) dx \\ &= \int_0^{\frac{1}{2}} \frac{-2x^2}{(k + \frac{1}{2})^2 - x^2} dx. \end{aligned}$$

Im Intervall von 0 bis $\frac{1}{2}$ ist der Integrand monoton fallend, d.h.

$$0 \geq \frac{-2x^2}{(k + \frac{1}{2})^2 - x^2} \geq \frac{-\frac{1}{2}}{(k + \frac{1}{2})^2 - \frac{1}{4}} = \frac{-2}{(2k + 1)^2 - 1} \geq -\frac{1}{4k^2},$$

und damit ist

$$0 \geq \int_k^{k+1} \frac{\{x\} - \frac{1}{2}}{x} dx = \int_0^{\frac{1}{2}} \frac{-2x^2}{(k + \frac{1}{2})^2 - x^2} dx \geq -\frac{1}{8k^2},$$

denn wir können das Integral abschätzen durch das Produkt aus der Länge des Integrationsintervalls und dem Minimum des Integranden. Summation von $k = 1$ bis $n - 1$ schließlich gibt die Abschätzung

$$0 \geq \int_1^n \frac{\{x\} - \frac{1}{2}}{x} dx \geq - \sum_{k=1}^{n-1} \frac{1}{4k^2}$$

für das störende Integral aus der obigen Formel.

Wie wohl jeder schon einmal in einer Analysis I Übungsaufgabe zeigen mußte, konvergiert die rechtsstehende Summe (egal ob mit oder ohne acht im Nenner) für $n \rightarrow \infty$; aus Kapitel III, §3f) wissen wir sogar, daß der Grenzwert $\pi^2/48$ ist. Auf jeden Fall können wir folgern, daß das uneigentliche Integral

$$\int_1^{\infty} \frac{\{x\} - \frac{1}{2}}{x} dx$$

konvergiert; den uns bislang noch unbekanntem Grenzwert wollen wir mit I bezeichnen. Damit ist

$$\ln n! = n(\ln n - 1) + \frac{\ln n}{2} + C + o(1) \quad \text{mit } C = I + 1$$

oder

$$n! \approx e^C \cdot n \cdot e^{-n} \sqrt{n}.$$

c) Die Stirlingsche Formel und die Normalverteilung

Mit dieser Formel können wir nun daran gehen, den Binomialkoeffizienten $\binom{n}{k}$ abzuschätzen:

$$\begin{aligned} \binom{n}{k} &= \frac{n!}{k!(n-k)!} \approx \frac{e^C n e^{-n} \sqrt{n}}{e^C k^k e^{-k} \sqrt{k} e^C (n-k)^{n-k} e^{-(n-k)} \sqrt{(n-k)}} \\ &= \frac{1}{e^C} \frac{n^n}{k^k (n-k)^{n-k}} \sqrt{\frac{n}{k \cdot (n-k)}}. \end{aligned}$$

Ausgedrückt durch $u = (n - 2k)\varepsilon = (n - 2k)\sigma/\sqrt{n}$ ist

$$k = \frac{n}{2} - \frac{u\sqrt{n}}{2\sigma},$$

und setzen wir zur Vereinfachung der Schreibweise

$$m = \frac{n}{2}, \quad v = \frac{u}{2\sigma} \quad \text{und} \quad \ell = \sqrt{2m} \cdot v,$$

so ergeben sich die Formeln

$$k = m - \sqrt{2m} \cdot v = m - \ell \quad \text{und} \quad n - k = m + \sqrt{2m} \cdot v = m + \ell.$$

Setzen wir dies alles in die Formel für die Wahrscheinlichkeitsdichte ein, erhalten wir

$$\begin{aligned} & \binom{n}{k} 2^{-n-1} \frac{\sqrt{n}}{\sigma} \\ & \approx \frac{1}{e^C} \frac{n^n}{k^k (n-k)^{n-k}} \sqrt{\frac{n}{k \cdot (n-k)}} 2^{-n-1} \frac{\sqrt{n}}{\sigma} \\ & = \frac{1}{e^C \sigma} \frac{n^n \cdot 2^{-n}}{k^k (n-k)^{n-k} \sqrt{k \cdot (n-k)}} \cdot 2^{-1} \\ & = \frac{1}{e^C \sigma} \frac{(2m)^{2m} \cdot 2^{-2m}}{(m-\ell)^{m-\ell} (m+\ell)^{m+\ell} \sqrt{(m-\ell)(m+\ell)}} \cdot 2m \cdot 2^{-1} \\ & = \frac{1}{e^C \sigma} \frac{m^{2m}}{(m-\ell)^m (m+\ell)^m} \left(\frac{m-\ell}{m+\ell}\right)^\ell \frac{m}{\sqrt{m^2 - \ell^2}} \\ & = \frac{1}{e^C \sigma} \frac{m^{2m}}{(m^2 - \ell^2)^m} \left(\frac{m-\ell}{m+\ell}\right)^\ell \frac{m}{\sqrt{m^2 - \ell^2}} \\ & = \frac{1}{e^C \sigma} \frac{1}{\left(1 - \frac{\ell^2}{m^2}\right)^m} \left(\frac{1-\ell/m}{1+\ell/m}\right)^\ell \frac{1}{\sqrt{1 - \ell^2/m^2}} \\ & = \frac{1}{e^C \sigma} \frac{1}{\left(1 - \frac{2v^2}{m}\right)^m} \left(\frac{1 - \sqrt{2/m} \cdot v}{1 + \sqrt{2/m} \cdot v}\right)^\ell \frac{1}{\sqrt{1 - 2v^2/m}}. \end{aligned}$$

Nun können wir langsam daran denken, n (und damit auch m) gegen unendlich gehen zu lassen; wir verwenden dazu die aus der Analysis I und wahrscheinlich auch aus der Schule bekannte Beziehung

$$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n.$$

Danach ist insbesondere

$$\lim_{m \rightarrow \infty} \left(1 - \frac{2v^2}{m} \right)^m = e^{-2v^2}$$

und

$$\begin{aligned} \lim_{m \rightarrow \infty} \left(1 \pm \sqrt{\frac{2}{m}} \cdot v \right)^{\sqrt{2m} \cdot v} &= \lim_{m \rightarrow \infty} \left(1 \pm \frac{\sqrt{2} \cdot v}{\sqrt{m}} \right)^{\sqrt{m} \cdot \sqrt{2} \cdot v} \\ &= \lim_{q \rightarrow \infty} \left(1 \pm \frac{\sqrt{2} \cdot v}{q} \right)^{q \cdot \sqrt{2} \cdot v} = \left(e^{\pm \sqrt{2} \cdot v} \right)^{\sqrt{2} \cdot v} = e^{\pm 2v^2}, \end{aligned}$$

denn es bleibt sich natürlich gleich, ob m oder $q = \sqrt{m}$ gegen unendlich geht. Da der Term v^2/m gegen null geht, erhalten wir somit als Grenzwert des gesamten obigen Ausdrucks

$$\frac{1}{e^{C\sigma}} \cdot \frac{1}{e^{-2v^2}} \cdot \frac{e^{-2v^2}}{e^{+2v^2}} \cdot 1 = \frac{1}{e^{C\sigma}} e^{-2v^2}.$$

Beachten wir nun noch, daß $v = u/2\sigma$ war, erhalten wir

$$\frac{1}{e^{C\sigma}} e^{-u^2/2\sigma^2}.$$

Damit sind wir fast am Ziel; das einzige, was noch fehlt, ist die Konstante C . Diese können wir bestimmen, indem wir ausnutzen, daß jeder Fehler mit Wahrscheinlichkeit eins zwischen $-\infty$ und ∞ liegt, d.h.

$$\frac{1}{e^{C\sigma}} \int_{-\infty}^{\infty} e^{-u^2/2\sigma^2} du = 1.$$

Aus [HM1], Kap. 2, §6c), wissen wir, daß

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$$

ist; mit der Substitution $x = u/\sqrt{2}\sigma$ folgt, daß dann

$$\int_{-\infty}^{\infty} e^{-u^2/2\sigma^2} du = \sqrt{2\pi}\sigma$$

ist und

$$e^C = \sqrt{2\pi} \quad \text{oder} \quad C = \frac{1}{2} \ln(2\pi).$$

Damit haben wir die Wahrscheinlichkeitsdichte endlich vollständig berechnet; das Endergebnis ist

$$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{u^2}{2\sigma^2}}.$$

Auch die Formel für $n!$ können wir nach der Bestimmung von C nun vollständig hinschreiben:

$$\ln n! = n(\ln n - 1) + \ln \sqrt{2\pi n} + o(1) \quad \text{oder} \quad n! \approx \left(\frac{n}{e}\right)^n \sqrt{2\pi n}.$$

Beides bezeichnet man als STIRLINGSche Formel.

Der schottische Mathematiker JAMES STIRLING (1692–1770) war Anhänger des gestürzten Königs Jakob II Stuart und hatte deshalb große politische Probleme bei seinem Studium; unter anderem wurde er deshalb von der Universität Oxford ausgeschlossen. 1717–1722 lebte er in Venedig und hatte auch gute Kontakte zu NICOLAUS BERNOULLI an der Universität von Padua; außerdem brachte er aus Venedig die Produktionsgeheimnisse der dortigen Glasbläser mit. Ab 1724 arbeitete er zehn Jahre lang als Mathematiklehrer in London, wo er viel mit NEWTON zusammentraf; 1735 wurde er Direktor einer schottischen Bergbaugesellschaft. In seine Londoner Zeit fällt die Veröffentlichung seines bedeutendsten Werks *Methodus Differentialis sive Tractatus de Summatione et Interpolatione Serierum Infinitarum* im Jahre 1730, das die obige Formel als Beispiel zwei zu Proposition 28 enthält. Ebenfalls ziemlich bekannt wurde seine 1735 veröffentlichte Arbeit über die Gestalt der Erde.

d) Eigenschaften der Normalverteilung

Oft interessiert nicht so sehr die Verteilung der Fehler, sondern die der Meßwerte selbst. Ist \hat{x} der korrekte Wert und x_i der i -te Meßwert dafür, der gemäß $x_i = \hat{x} + u_i$ mit dem Fehler u_i behaftet ist, so können wir mit $x = \hat{x} + u$ die obigen Wahrscheinlichkeitsdichte auch als

$$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\hat{x})^2}{2\sigma^2}}$$

schreiben.

Als *Normalverteilung mit Mittelwert a und Standardabweichung σ* bezeichnen wir daher die Verteilung mit Wahrscheinlichkeitsdichte

$$\varphi(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-a)^2}{2\sigma^2}}.$$

Diese Wahrscheinlichkeitsdichte hängt offensichtlich nur von der *normierten Variablen*

$$z = \frac{x-a}{\sigma}$$

ab; diese hat Mittelwert null und Standardabweichung eins. Daher gibt es für die Normalverteilung nicht – wie für viele andere statistische Verteilungen – je nach Parameterwerten verschiedene Tabellen, sondern man findet in allen Tabellenwerken nur die Normalverteilung mit Mittelwert null und Standardabweichung eins, man findet also die Wahrscheinlichkeitsdichte

$$\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

und deren Integral

$$F(z) = \int_{-\infty}^z f(u) du = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{u^2}{2}} du.$$

Dieses Integral läßt sich nicht weiter vereinfachen, da sich die Stammfunktion von $e^{-u^2/2}$ nicht durch elementare Funktionen ausdrücken läßt. Für die Bestimmung von $F(z)$ ist man daher auf Tabellen oder Computerprogramme angewiesen; eine graphische Darstellung von $F(z)$ ist in Abbildung 68 zu sehen. Mit dieser Funktion läßt sich die Wahrscheinlichkeit dafür, daß

$$c \leq z = \frac{x-a}{\sigma} \leq d$$

ist berechnen als $F(d) - F(c)$, und damit läßt sich auch leicht die Wahrscheinlichkeit berechnen, daß x selbst zwischen zwei gegebenen Schranken liegt.

Mißt man beispielsweise die Temperatur eines Wasserbads eine Viertelstunde lang jede Minute und erhält dabei 15 Meßwerte mit Mittel-

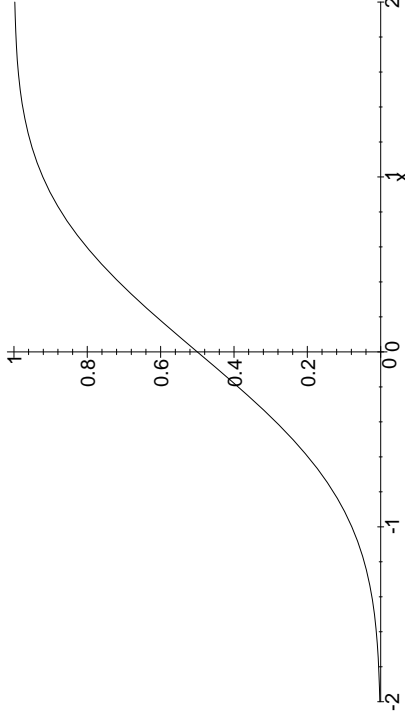


Abb. 68: Das Integral $F(z)$ über die „Glockenkurve“

wert $20,1^\circ\text{C}$ und Standardabweichung $0,2^\circ\text{C}$, so ist die Standardabweichung des Mittelwerts

$$\sigma_{\bar{y}} = \frac{0,2^\circ\text{C}}{\sqrt{14}} \approx 0,053^\circ\text{C}.$$

Wenn wir dann beispielsweise wissen wollen, mit welcher Wahrscheinlichkeit die „tatsächliche“ mittlere Temperatur zwischen $20,0^\circ\text{C}$ und $20,2^\circ\text{C}$ liegt, müssen wir dazu zunächst die normierten Werte berechnen:

$$z_1 = \frac{20,0 - 20,1}{0,053} \approx -1,89 \quad \text{und} \quad z_2 = \frac{20,2 - 20,1}{0,053} \approx 1,89.$$

Die Wahrscheinlichkeit ist also

$$F(1,89) - F(-1,89) \approx 0,94;$$

oder rund 94%.

Schaut man in einer Tabelle nach, wird man dort allerdings im allgemeinen nur den Wert $F(1,89)$ finden, nicht aber $F(-1,89)$. Der Grund dafür liegt in der Symmetrie des Graphen von F bezüglich des Punktes $(0, \frac{1}{2})$. Was dahinter steckt, sieht man am besten, wenn man die Dichtefunktion der Normalverteilung betrachtet, also die Glockenkurve: Für $z > 0$ ist $F(-z)$ die in Abbildung 69 links eingezeichnete schraffierte Fläche.

Diese Fläche ist wegen der Symmetrie der Glockenkurve zur senkrechten Achse gleich der rechts eingezeichneten schraffierten Fläche, und deren Komplement ist $F(z)$. Also ist

$$F(-z) = 1 - F(z),$$

und es reicht, wenn wir die Werte von F im positiven Bereich kennen.

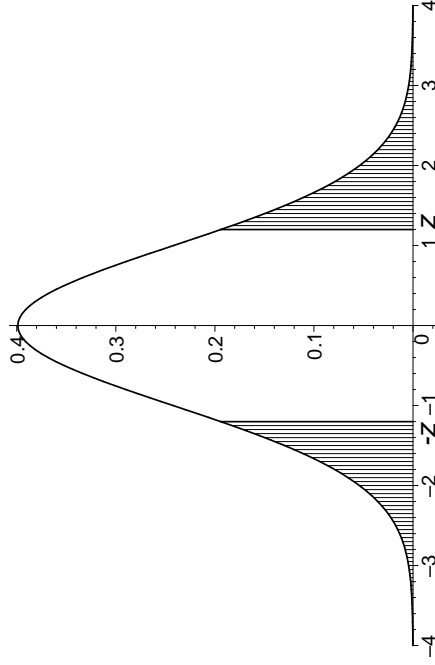


Abb. 69: Zusammenhang zwischen $F(z)$ und $F(-z)$

Oft interessiert auch die Wahrscheinlichkeit dafür, daß der Betrag des Fehlers unterhalb einer bestimmten Schranke liegt, etwa $z \cdot \sigma$; in Abbildung 69 wäre dies der nichtschraffierte Bereich unter der Glockenkurve.

Wie man sich anhand der Abbildung leicht klarmacht, ist diese Wahrscheinlichkeit gleich

$$F(z) - F(-z) = 2F(z) - 1;$$

die Wahrscheinlichkeit, daß wir im obigen Beispiel die mittlere Temperatur mit einem Fehler von höchstens $0,05^\circ$ gemessen haben, ist also

$$2F\left(\frac{0,05}{0,053}\right) \approx F(0,94) \approx 0,83.$$

Ein Wasserbad hat üblicherweise den Sinn, ein Experiment unter kontrollierten Temperaturbedingungen durchzuführen; daher interessiert vor allem, inwieweit es gelingt, die Temperatur innerhalb gewisser Schranken zu halten. Die Wahrscheinlichkeit dafür können wir mit denselben Methoden berechnen, allerdings müssen wir dazu mit der Standardabweichung der Meßreihe selbst arbeiten.

Wenn wir etwa wollen, daß die Temperatur immer zwischen $19,5$ und $20,5^\circ\text{C}$ liegt, so ist die Wahrscheinlichkeit, daß wir dies mit dem oben ausgemessenen Versuchsaufbau erreichen, gleich

$$F\left(\frac{20,5 - 20,1}{0,2}\right) - F\left(\frac{29,5 - 20,1}{0,2}\right) = F(2) - F(-3) \approx 0,976.$$

In knapp zweieinhalb Prozent aller Fälle, im Schnitt also alle vierzig Minuten, müssen wir also damit rechnen, daß die Toleranzgrenzen überschritten werden.

Wie Abbildung 68 zeigt, liegt $F(-2)$ sehr nahe bei null und $F(2)$ sehr nahe bei eins. In der Tat ist die Wahrscheinlichkeit dafür, daß ein Wert z Betrag größer z hat, nach obiger Diskussion gleich

$$1 - (2F(z) - 1) = 2F(z) - 2,$$

was für $z = 2$ zu $-0,0455$ wird; die Wahrscheinlichkeit ist also kleiner als 5%. Allgemein gilt für eine beliebige Normalverteilung, daß der Wert der Variablen mit folgenden Wahrscheinlichkeiten um höchstens $i\sigma$ vom Mittelwert abweicht:

$i =$	1	2	3	4
Wahrscheinlichkeit:	0,683	0,954	0,9973	0,99994

Damit liegen also etwa zwei Drittel aller Fehler zwischen $-\sigma$ und σ , 95% liegen zwischen -2σ und 2σ und 99,7% zwischen -3σ und 3σ ; die Wahrscheinlichkeit dafür, daß der Fehler größer als 3σ ist, beträgt nur etwa 0,27%. Da Ereignisse mit einer so geringen Wahrscheinlichkeit seltener als in einem von 300 Fällen auftreten, betrachtet man Fehler, die außerhalb des 3σ -Bereichs liegen, oft als „Ausreißer“, d.h. als grobe Meßfehler, die bei der Bestimmung des Ergebnisses nicht berücksichtigt

werden. Sehr vorsichtige Leute reden allerdings erst ab einer Abweichung von 4σ von Ausreißern; solche Fehler treten zufällig weniger als einmal pro 15 000 Messungen auf.

Für Leser, die ihren Computer selbst programmieren und keine spezielle Statistiksoftware haben, sei hier eine Näherungsformel für $F(z)$ angegeben: Mit einem Fehler von höchstens $7,5 \cdot 10^{-8}$ ist

$$F(z) = 1 - \varphi(z) \cdot (a_1 t + a_2 t^2 + a_3 a^3 + a_4 t^4 + a_5 t^5)$$

mit $t = \frac{1}{1 + pz}$ und $\varphi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$ sowie

$$\begin{aligned} a_1 &= 0,319\,381\,530 & a_2 &= -0,356\,563\,782 & a_3 &= 1,781\,477\,973 \\ a_4 &= -1,821\,255\,978 & a_5 &= 1,330\,274\,429 & p &= 0,231\,641\,9 \end{aligned}$$

Beim Rechnen mit dem Taschenrechner kann man sich auch mit einer vereinfachten Version begnügen, bei der $a_4 = a_5 = 0$ ist und

$$a_1 = 0,436\,1836 \quad a_2 = -0,120\,1676 \quad a_3 = 0,937\,2980 \quad p = 0,332\,67;$$

hier kann der Fehler bis zu 10^{-5} betragen.

e) Die Maximum Likelihood Methode

GAUSS gab im Laufe seines Lebens mehrere Begründungen für die Methode der kleinsten Quadrate (die er bei sowohl bei seinen astronomischen Arbeiten wie auch bei der von ihm geleiteten Vermessung des Königreichs Hannover zwischen 1818 und 1832 ständig benutzte); die unter dem Gesichtspunkt einer in sich geschlossenen Fehlertheorie interessanteste beruht auf dem LAPLACESchen Fehlermodell:

Danach sollte der Wert u_i für die korrekten Parameterwerte a, b, \dots aus einer Normalverteilung mit Mittelwert $f(a, b, \dots; t_i)$ kommen, deren Standardabweichung σ_i von der Genauigkeit abhängt, mit der u_i bestimmt werden kann. Die Wahrscheinlichkeit dafür, daß u_i zwischen zwei Werten a und b liegt, ist damit

$$\int_a^b e^{-\frac{(u - f(a, b, \dots; t_i))^2}{2\sigma_i^2}} \cdot$$

Von der Wahrscheinlichkeit, daß u_i gleich einem Wert c ist, können wir natürlich nicht reden, da diese nach obiger Formel ein Integral von c nach c wäre, also Null. Aber die Wahrscheinlichkeit dafür, daß u_i in einem kleinen Intervall der Länge ε_i um einen Wert c_i liegt, ist ungefähr proportional zum ε_i -fachen Wert des Integranden an der Stelle c_i , also zu

$$\varepsilon_i \cdot e^{-\frac{(c_i - f(a, b, \dots; t_i))^2}{2\sigma_i^2}}.$$

Entsprechend ist

$$\varepsilon_j \cdot e^{-\frac{(c_j - f(a, b, \dots; t_j))^2}{2\sigma_j^2}}$$

ungefähr gleich der Wahrscheinlichkeit dafür, daß u_j in einem Intervall der Breite ε_j um c_j liegt.

Wenn wir wie üblich davon ausgehen, daß keine systematischen Fehler auftreten, sind die Fehler von u_i und u_j voneinander unabhängig, die Wahrscheinlichkeit dafür, daß (u_i, u_j) in einem Rechteck mit Seiten ε_i und ε_j um (c_i, c_j) liegt, ist also proportional zum Produkt der beiden obigen Einzelwahrscheinlichkeiten, d.h. zu

$$\varepsilon_i \varepsilon_j e^{-\frac{(c_i - f(a, b, \dots; t_i))^2}{2\sigma_i^2} - \frac{(c_j - f(a, b, \dots; t_j))^2}{2\sigma_j^2}}.$$

Entsprechend kann auch die Wahrscheinlichkeit dafür berechnet werden, daß der Punkt (u_1, \dots, u_n) in einem kleinen gegebenen Quader mit Kantenlängen $\varepsilon_1, \dots, \varepsilon_n$ liegt; sie ergibt sich zu

$$L(a, b, \dots) \cdot \prod_{i=1}^n \varepsilon_i$$

mit

$$L(a, b, \dots) \stackrel{\text{def}}{=} e^{-\sum_{i=1}^n \frac{(u_i - f(a, b, \dots; t_i))^2}{2\sigma_i^2}}.$$

Diese Größe ist selbst keine Wahrscheinlichkeit, sondern der Quotient aus einer Wahrscheinlichkeit und einem Volumen; man spricht daher von einer *Wahrscheinlichkeitsdichte*.

Wenn wir diese Wahrscheinlichkeitsdichte als Funktion von a, b, \dots betrachten, macht sie eine Aussage über die Güte der Parameter: Schließlich wird man einem Modell, das dem beobachteten Ausgang eines Experiments eine hohe Wahrscheinlichkeit zuweist, eher glauben als einem alternativen Modell, das die beobachteten Daten zu Ausreißern erklärt. Aus diesem Grund kann die Funktion L auch als Maß dafür betrachtet werden, wie „wahrscheinlich“ in irgendeinem umgangssprachlichen (und schwer präzisierbaren) Sinne die Parameter a, b, \dots sind.

Im englischen gibt es zwei Wörter für Wahrscheinlichkeit: Das romanische Wort *probability* und das germanische Wort *likelihood*. Für den mathematisch exakten Wahrscheinlichkeitsbegriff verwendet man *probability*, für „Wahrscheinlichkeit“ im Sinne der Funktion L *likelihood*. Da es im deutschen kein zweites Wort für Wahrscheinlichkeit gibt, spricht man hier in Anlehnung an das Englische von einer *Likelihoodfunktion*.

Die Maximum Likelihood Methode besteht nun genau in dem, was ihr Name besagt: *Man wähle die Parameter a, b, \dots so, daß die Likelihoodfunktion maximal wird.*

Da $L(a, b, \dots)$ durch eine Exponentialfunktion beschrieben wird, wird die Likelihoodfunktion genau dann maximal, wenn ihr Exponent maximal wird. Dieser Exponent ist eine negative Zahl, wird also genau dann maximal, wenn sein Betrag *minimal* wird, das heißt, wenn die Quadratsumme

$$\sum_{i=1}^n \frac{(u_i - f(a, b, \dots; t_i))^2}{2\sigma_i^2}$$

minimal wird.

In vielen Fällen wird die Zuverlässigkeit der einzelnen Paare (t_i, u_i) miteinander vergleichbar sein, so daß alle σ_i gleich sind; in diesem Fall kann man die σ_i ignorieren und einfach die Quadratsumme

$$\sum_{i=1}^n (u_i - f(a, b, \dots; t_i))^2$$

minimieren, d.h. wir kommen wieder zur klassischen Methode der kleinsten Quadrate. Es gibt aber auch Anwendungen, wie etwa oben beim

überexponentiellen Bevölkerungswachstum, bei denen die Verschiedenheit des σ_i sehr wesentlich ist: Sicherlich wird man etwa der auf Volkszählungen beruhenden Weltbevölkerungszahl, die die Vereinten Nationen für 1995 veröffentlichten, mehr Vertrauen entgegenbringen als der Schätzung eines Historikers für die Weltbevölkerung des Jahres Null, und selbst bei ein und derselben Meßreihe im Labor kommt es gelegentlich vor, daß (beispielsweise aufgrund unterschiedlicher Genauigkeit eines Meßinstruments in verschiedenen Bereichen) manche Daten zuverlässiger sind als andere.

§6: Kompression von Bild- und Audiodaten

Zum Abschluß der Vorlesung wollen wir wenigstens kurz eine praktische Anwendung kennenlernen, in der mit Eigenwerten und Eigenvektoren symmetrischer Matrizen, FOURIER-Transformationen und Statistik gleich mehrere der Methoden aus diesem Semester gleichzeitig benötigt werden: die Komprimierung von Bild- und Audiodaten.

a) Datenkompression

Ziel der Datenkompression ist es, eine Datei für Zwecke der Speicherung oder Übertragung möglichst stark zu verkleinern, das aber in einer solchen Weise, daß sich die ursprüngliche Datei aus der verkleinerten wieder exakt rekonstruieren läßt.

Es ist klar, daß es keinen universellen Algorithmus zur Datenkompression geben kann: Gäbe es nämlich ein Verfahren, das für beliebige Dateien einen Kompressionsfaktor $\alpha < 1$ garantieren würde, so könnte man dieses Verfahren iterativ anwenden und nach n Anwendungen eine Kompressionsrate von α^n erreichen. Wenn man n nur hinreichend groß wählt, könnte man daher jede Datei auf weniger als ein Bit komprimieren, was natürlich absurd ist.

Ein Kompressionsverfahren kann also nur auf Dateien mit spezieller Struktur erfolgreich angewandt werden und muß die spezielle Redundanz in diesen Dateien ausnutzen. In Textdateien beispielsweise ist dies die Redundanz der Sprache, die schon bei bloßer Beachtung der

höchst unterschiedlichen Buchstabenhäufigkeiten Kompressionen von rund 50% gestattet.

Bilddaten werden typischerweise als Matrizen aus ganzen Zahlen zwischen 0 und 255 digitalisiert; bei Audiodaten nimmt man Vektoren von ganzen Zahlen zwischen 0 und $65\,535 = 2^{16} - 1$ oder $16\,777\,215 = 2^{24} - 1$. (Der Unterschied zwischen den Wertebereichen liegt darin begründet, daß unser Auge selbst bei gedruckten Bildern mit nur 64 Graustufen praktisch keine Artefakte mehr erkennen kann, wohingegen unser Gehör noch auf sehr feine Unterschiede reagiert.)

Bei einer Musik-CD etwa wird das Signal 44 100-mal pro Sekunde abgetastet (dies bedeutet nach dem Abtasttheorem von NYQUIST, daß ein auf den Bereich von 0 bis 22,05kHz bandbegrenztes Signal fehlerfrei rekonstruiert werden kann), und das Ergebnis wird dann so skaliert und quantisiert (d.h. gerundet), daß eine Zahl zwischen 0 und 65535 entsteht. Bei Bilddaten werden je nach Auflösung und Seitenverhältnis zwischen etwa 256×256 und 1024×1024 Bildpunkte abgetastet, für Schwarzweißbilder nur nach Helligkeit, für Farbbildern nach insgesamt drei Größen, die vom jeweiligen Farbmodell abhängen. Das Ergebnis dieser Abtastungen wird dann entsprechend skaliert und quantisiert.

Typische Komprimierungsverfahren arbeiten daher mit Vektoren oder Matrizen aus Zahlen zwischen 0 und einer geeigneten Zahl M , die aus praktischen Gründen meist von der Form $2^{8r} - 1$ ist, wobei die Zahl r der Empfindlichkeit unserer Sinne angepaßt zwischen eins und drei liegt.

Da man zur eindeutigen Festlegung von N beliebigen Zahlen zwischen 0 und 2^{8r} nicht mit weniger als den $8Nr$ Bit auskommen kann, die man zum Hinschreiben der Zahlen braucht, sehen wir auch hier wieder, daß kein Verfahren *alle* solchen Vektoren komprimieren kann; wir müssen also eine Teilmenge auszeichnen.

Die ideale solche Teilmenge wäre hier natürlich die Menge aller möglicher Bilder (oder Audiosequenzen), aber diese Menge dürfte mathematisch kaum definierbar sein: Schließlich hängt es sehr vom Betrachter ab, welches Pixelmuster er noch als „Bild“ gelten läßt und welches nicht. Sinnvoll läßt sich eine solche Menge daher höchstens definieren, wenn

von vornherein feststeht, welche Bilder berücksichtigt werden sollen – und dann ist wohl ein Verfahren, das statt vom Bildinhalt von einer Bildnummer ausgeht, unschlagbar.

Die meisten klassischen Verfahren, die beliebige, aber realistische Bilder komprimieren sollen, gehen aus von einem *statistischen Modell*, das zwar auch viele Matrizen produziert, die niemand als „Bilder“ anerkennen würde, das aber dennoch genügend viele Eigenschaften realer Bilder reproduziert, um eine große Anzahl von „Nichtbildern“ auszuschließen.

Ausgangspunkt ist die Beobachtung, daß es in einem Bild oder Musikstück nur wenige abrupte Übergänge gibt. Zwar gibt es natürlich immer wieder ein plötzliches *fortissimo*, das auf eine leise Stelle folgt, aber da das Signal 44 100-mal pro Sekunde abgetastet wird und solche Übergänge selbst bei der schrägsten Musik deutlich seltener als im Sekundenrhythmus erfolgen, sind diese Sprünge innerhalb des zu behandelnden Datenstroms in der Tat sehr seltene Ereignisse. Wir können daher davon ausgehen, daß sich die unmittelbaren Nachbarn eines Datums *im Mittel* nur wenig vom gegebenen Datum unterscheiden.

Dasselbe gilt auch für Bilddaten: Falls das Bild digital hinreichend fein dargestellt wird, so daß keine Rastereffekte erkennbar sind, kommen große Sprünge in den Helligkeitswerten nur selten vor.

Bei diesem engen Zusammenhang zwischen benachbarten Werten setzen die gängigen Komprimierungsalgorithmen an: Wenn zwei Größen typischerweise sehr ähnlich sind, wird bei der Übertragung oder Speicherung *beider* Werte ein großer Teil der Information doppelt betrachtet; die Informationsdichte kann also deutlich erhöht werden, wenn man nur Informationen betrachtet, die weitgehend unabhängig voneinander sind.

Bevor wir das in die Praxis umsetzen können, müssen wir zunächst mathematisch fassen, was die gegenseitige Abhängigkeit benachbarter Daten bedeutet. Aus dem letzten Paragraphen kennen wir schon ein Maß für die gegenseitige Abhängigkeit zweier Meßreihen, den Korrelationskoeffizienten; wir wollen etwas Entsprechendes nun für eine beliebige Anzahl von Meßreihen oder Datenvektoren definieren.

Zur Klärung der Begriffe beginnen wir mit einem Beispiel, das zwar nicht das geringste mit Bild- oder Audiodaten zu tun hat, das dafür aber einfach und übersichtlich ist:

b) Kovarianzen und Korrelationen bei der Europawahl 1999

Wir betrachten das Ergebnis der Europawahl vom 13. Juni 1999 in den verschiedenen Mannheimer Stadtteilen. Laut Mannheimer Morgen vom 15. Juni gab es folgende Prozentzahlen:

Bezirk	Bet.	CDU	SPD	Grüne	Rep.	FDP	PDS	Sonst.
(A)	37,7	40,7	31,8	13,0	2,4	5,0	3,4	4,2
(B)	30,9	33,1	39,8	12,9	3,9	2,1	4,1	4,1
(D)	36,1	39,5	36,1	10,8	2,8	3,1	3,7	4,0
(D)	44,9	46,0	26,6	12,8	1,9	6,6	2,5	3,6
(E)	46,0	44,2	30,0	12,0	2,4	5,6	2,1	3,7
(F)	48,6	51,1	23,7	10,8	2,6	6,3	2,5	3,0
(G)	36,8	43,3	42,0	4,6	3,1	2,2	1,3	3,5
(H)	28,3	39,2	46,0	4,0	4,1	1,6	1,5	3,6
(I)	37,9	40,1	42,7	5,7	3,7	2,8	1,7	3,3
(J)	37,9	44,9	37,3	6,4	3,3	2,6	1,8	3,7
(K)	40,9	44,8	40,9	3,8	3,2	2,5	1,6	3,2
(L)	49,1	46,3	35,9	7,7	2,2	3,7	1,2	3,0
(M)	47,4	47,0	30,2	11,9	1,6	4,2	1,8	3,3
(N)	42,5	48,8	32,7	9,1	1,9	3,5	1,4	2,6
(O)	40,9	44,8	38,1	7,4	1,8	3,2	1,8	2,9
(P)	42,5	47,8	30,3	10,2	2,2	3,8	2,1	3,5
(Q)	38,4	47,6	34,9	6,0	3,4	2,9	1,2	4,0
MA	39,8	44,2	34,6	9,0	2,7	3,7	2,1	3,7

Die Stadtteile sind, damit alles in eine Zeile paßt, durch Buchstaben bezeichnet; für Interessenten seien auch die richtigen Namen angegeben: (A) steht für Innenstadt/Jungbusch, (B) für Neckarstadt West, (C) für Neckarstadt Ost/Wohlgelegen, (D) ist die Oststadt/Schwezingenstadt, (E) der Lindenhof, (F) Neustadt/Neuhermsheim, (G) Sandhofen, (H) Schönau, (I) Waldhof, (J) Käfertal, (K) Vogelstang, (L) Wallstadt, (M) Feudenheim, (N) Seckenheim, (O) Friedrichsfeld, (P) Neckarau, (Q)

Rheinau und MA schließlich Mannheim insgesamt.

Wie man sieht, haben die einzelnen Parteien in den verschiedenen Stadtteilen recht unterschiedlich abgeschnitten; bei der FDP etwa ergibt sich mehr als ein Faktor vier zwischen den 1,6% in der Schönau und den 6,6% in Neustadt/Neuhermsheim. Es ist auch klar, daß die CDU dort unterdurchschnittlich viele Stimmen hat, wo die SPD überdurchschnittlich abgeschnitten hat; anders geht das bei den Stimmzahlen dieser beiden Parteien gar nicht. Wie steht es aber mit den kleinen Parteien?

Eine Partei hat in einem Stadtteil dann überdurchschnittlich gut abgeschnitten, wenn sie dort einen höheren Prozentsatz erreicht hat als für Mannheim insgesamt, wenn also die Differenz zwischen dem Prozentsatz dort und dem für ganz Mannheim positiv ist. Bei negativer Differenz hat sie unterdurchschnittlich abgeschnitten.

Wenn wir also vergleichen wollen, ob zwei Parteien ihre Hochburgen und ihre Schwächegebiete jeweils in denselben Stadtteilen haben, geben uns wie in §4 die Produkte ihrer Abweichungen von ihren Mittelwerten einen ersten Hinweis, d.h. also die Ausdrücke

$$\frac{1}{17} \sum_{i=1}^{17} (x_i - \bar{x})(y_i - \bar{y}),$$

wobei die x_i, y_i die Ergebnisse aus den einzelnen Stadtteilen sind und \bar{x}, \bar{y} die Ergebnisse für ganz Mannheim.

Da die verschiedenen Stadtteile verschieden groß sind, sind \bar{x} und \bar{y} nicht genau die Mittelwerte der x_i beziehungsweise y_i , sondern aussagekräftigere gewichtete Mittelwerte; da andererseits Wahlbezirke ungefähr gleich groß sein sollten, wollen wir diesen Unterschied nicht allzu ernst nehmen und können obige Summe im wesentlichen mit der Kovarianz

$$\text{cov}(\bar{x}, \bar{y}) = \frac{1}{17} \sum_{i=1}^{17} (x_i - \bar{x})(y_i - \bar{y})$$

identifizieren, wobei

$$\bar{x} = \frac{1}{17} \sum_{i=1}^{17} x_i \quad \text{und} \quad \bar{y} = \frac{1}{17} \sum_{i=1}^{17} y_i$$

die Mittelwerte der beiden Datensätze sind.

Bei m Datensätzen $\vec{x}_1, \dots, \vec{x}_m \in \mathbb{R}^N$ können wir die Kovarianzen $\text{cov}(\vec{x}_i, \vec{x}_j)$ zu einer $m \times m$ -Matrix zusammenfassen, der *Kovarianzmatrix*

$$\text{Cov}(\vec{x}_1, \dots, \vec{x}_m) = (c_{ij}) \quad \text{mit} \quad c_{ij} = \text{cov}(\vec{x}_i, \vec{x}_j).$$

Sie ist offensichtlich eine symmetrische Matrix.

Wenn wir die Wahlbeteiligung formal wie eine Partei namens W behandeln, erhalten wir hier die folgende Kovarianzmatrix:

	W	CDU	SPD	Grüne	Rep.	FDP	PDS	Sonst.
W	32,4	18,8	-24,3	5,66	-3,27	5,79	-1,32	-1,46
CDU	18,8	18,2	-15,7	-0,58	-1,96	3,00	-2,07	-1,08
SPD	-24,3	-15,7	34,9	-13,8	3,11	-7,66	-1,07	0,10
Grüne	5,66	-0,58	-13,8	10,2	-1,21	3,13	1,93	0,44
Rep.	-3,27	-1,96	3,11	-1,21	0,57	-0,69	0,05	0,14
FDP	5,79	3,00	-7,66	3,13	-0,69	2,04	0,26	-0,04
PDS	-1,32	-2,07	-1,07	1,93	0,05	0,26	0,73	0,22
Sonst.	-1,46	-1,08	0,10	0,44	0,14	-0,04	0,23	0,24

Übersichtlicher ist die *Korrelationsmatrix*, deren Einträge nicht die Kovarianzen sondern die Korrelationskoeffizienten sind; diese liegen bekanntlich stets zwischen -1 und 1. Im vorliegenden Fall ist das die Matrix

	W	CDU	SPD	Grüne	Rep.	FDP	PDS	Sonst.
W	1,00	0,77	-0,72	0,31	-0,76	0,71	-0,27	-0,52
CDU	0,77	1,00	-0,62	-0,04	-0,61	0,49	-0,57	-0,52
SPD	-0,72	-0,62	1,00	-0,73	0,70	-0,91	-0,21	0,03
Grüne	0,31	-0,04	-0,73	1,00	-0,50	0,69	0,71	0,28
Rep.	-0,76	-0,61	0,70	-0,50	1,00	-0,64	0,07	0,37
FDP	0,71	0,49	-0,91	0,69	-0,64	1,00	0,21	-0,05
PDS	-0,27	-0,56	-0,21	0,71	0,07	0,21	1,00	0,53
Sonst.	-0,52	-0,52	0,03	0,28	0,37	-0,05	0,53	1,00

Dieser Matrix können wir nun aller gewünschten Informationen entnehmen. Die erste Zeile beispielsweise sagt uns, daß die CDU, FDP und in geringerem Ausmaß auch die Grünen dort überdurchschnittlich gut abgeschnitten haben, wo auch die Wahlbeteiligung überdurchschnittlich war; diese drei Parteien konnten also anscheinend ihre Anhänger am besten mobilisieren.

Wo die CDU überdurchschnittlich gut abschnitt, waren alle anderen Parteien mit Ausnahme der FDP unterdurchschnittlich; wo die SPD überdurchschnittlich gut abschnitt, waren ebenfalls fast alle anderen unterdurchschnittlich; die Ausnahme hier sind die Republikaner. Zwischen SPD und FDP haben wir mit -0,91 den betragsgrößten Korrelationskoeffizienten; die Hochburgen dieser beiden Parteien sind also sehr unterschiedlich. Die Grünen schließlich sind positiv korreliert mit FDP, PDS und Sonstigen.

c) Zufallsvariablen und ihre statistischen Kenngrößen

Die Europawahl vom 13. Juni 1999 ist vorbei und ausgezählt, wir können daher im letzten Abschnitt mit wohlbekanntem Daten rechnen. Beim Umgang mit Bild- und Tondaten gilt das nur in eingeschränktem Maße: Zwar können Daten erst dann verarbeitet werden, wenn sie da sind, aber eine Digitalkamera beispielsweise muß doch schon heute auf die Bilder vorbereitet sein, die ihr Besitzer erst in ein paar Jahren aufnehmen wird.

Für den Grundalgorithmus zur Datenkompression müssen wir daher Daten zulassen, über die wir noch nichts konkretes wissen – abgesehen von gewissen vagen Gesetzmäßigkeiten, durch die sich echte Bilddaten von beliebigen Matrizen unterscheiden. Als Hilfsmittel dazu dient der Begriff der *Zufallsvariablen*.

Definition: Eine (diskrete) Zufallsvariablen ist ein Prozeß, der zufällig einen Wert aus einer vorgegebenen endlichen Menge

$$\{x_0, \dots, x_m\}$$

liefert.

Dieser „Zufall“ muß natürlich, falls er mathematisch faßbar sein soll, irgendwelchen Regeln genügen, und hier kommt der zweite fundamentale Begriff zum Einsatz: Wir nehmen an, daß für jeden der möglichen Werte x_i feststeht, mit welcher *Wahrscheinlichkeit* p_i er angenommen wird. Diese „Wahrscheinlichkeit“ definieren wir informell so, daß bei einer großen Anzahl m von Versuchen *ungefähr* $p_i \cdot m$ -mal der Wert x_i geliefert wird. Das „Gesetz der großen Zahlen“, das wir im nächsten Semester kennenlernen werden, sagt, daß diese Definition sinnvoll ist und man die Wahrscheinlichkeiten p_i damit in wohldefinierter Weise mit beliebiger Genauigkeit bestimmen kann.

Zwei formale Konsequenzen der Definition sind offensichtlich:

$$0 \leq p_i \leq 1 \quad \text{für alle } i \quad \text{und} \quad \sum_{i=0}^m p_i = 1.$$

Sollen beispielsweise alle x_i mit gleicher Wahrscheinlichkeit angenommen werden (was für Bilddaten nicht unbedingt ein sehr realistisches Modell ist), so sind alle $p_i = 1/(m+1)$.

In unserem Modell soll ein „Bild“ dann produziert werden durch eine Matrix $(X_{i,j})$ von geeigneten Zufallsvariablen; eine „Audiosequenz“ entsprechend durch einen Vektor (X_1, \dots, X_N) .

Welche Zufallsvariablen sind geeignet? In unserem einfachen Modell wollen wir uns nicht darauf festlegen, wie die p_i gewählt werden sollen, sondern stattdessen mit ihrer Hilfe Analoga zu den Kenngrößen aus dem vorigen Abschnitt definieren.

Definition: Der *Erwartungswert* $E(X)$ einer Zufallsvariablen X ist

$$E(X) = \sum_{i=0}^m p_i x_i.$$

Falls alle x_i mit gleicher Wahrscheinlichkeit angenommen werden, ist der Erwartungswert also einfach das arithmetische Mittel

$$\frac{1}{m+1} \sum_{i=0}^m x_i$$

der möglichen Werte; ansonsten ist er ein sogenanntes gewichtetes arithmetisches Mittel. Für einen Würfel etwa, der die Augenzahlen von $x_0 = 1$ bis $x_5 = 6$ mit jeweils gleicher Wahrscheinlichkeit produziert, ist

$$E(X) = \frac{1 + 2 + 3 + 4 + 5 + 6}{6} = 3\frac{1}{2}.$$

Als nächste wichtige Kenngröße hatten wir im letzten Abschnitt die mittlere quadratische Abweichung vom Mittelwert, die Varianz, betrachtet; in völliger Analogie zu dort definieren wir

Definition: Die Varianz einer Zufallsvariablen X mit Erwartungswert $E(X)$ ist

$$\sigma_X^2 = E((X - E(X))^2) = \sum_{i=0}^m p_i (x_i - E(X))^2;$$

ihre Standardabweichung ist $\sigma_X = \sqrt{\sigma_X^2}$.

Beim Würfel wäre also

$$\sigma_X^2 = \frac{(-2\frac{1}{2})^2 + (-1\frac{1}{2})^2 + (-\frac{1}{2})^2 + (\frac{1}{2})^2 + (1\frac{1}{2})^2 + (2\frac{1}{2})^2}{6} = \frac{35}{12}$$

und

$$\sigma_X = \sqrt{\frac{35}{12}} \approx 1,7078.$$

d) Beispiele aus der Bildverarbeitung

Um die Bedeutung dieser Kenngrößen in der Bildverarbeitung zu veranschaulichen, sind auf der nächsten Doppelseite sechs beliebige Testbilder zusammen mit den Werten dieser Kenngrößen abgedruckt. Die Werte sind entnommen aus

P.M. FARELLE: Recursive Block Coding for Image Data Compression, *Springer*, 1990 ;

sie beziehen sich natürlich auf die Originalbilder und nicht auf das, was der Druckvorgang hier im Skriptum daraus gemacht hat. Trotzdem sollte der Vergleich von Bildern und Daten einen einigermaßen korrekten

Eindruck zumindest der relativen Situation vermitteln, da hoffentlich alle hier abgedruckte Bilder in derselben Weise verunstaltet sind.

Die mittlere Helligkeit eines Bildes, dessen (viele) Pixel durch je eine Zufallsvariable mit Erwartungswert μ produziert werden, sollte ziemlich nahe bei μ liegen; der beste Schätzwert für den gemeinsamen Erwartungswert der Zufallsvariablen ist also die mittlere Helligkeit des Bildes. Typischerweise werden Helligkeiten durch Zahlen zwischen 0 und 255 kodiert, wobei schwarz der Zahl Null entspricht und weiß der 255. Dies sieht man gut an den Beispielen, wo das mit Abstand hellste Bild „Tiffany“ auch den mit Abstand größten Mittelwert μ hat; den kleinsten Wert hat das auch visuell dunkelste Bild „Lenna“.

Die nächste wichtige Kenngröße ist die Varianz, welche angibt, wie stark eine Zufallsvariable um ihren Erwartungswert streut. Auch hier wollen wir wieder davon ausgehen, daß alle Zufallsvariablen zu einem gegebenen Bild bzw. einer gegebenen Audiosequenz aus N dieselbe Varianz haben. Wir schätzen diese gemeinsame Varianz aufgrund der vorliegenden Daten als

$$\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \mu)^2,$$

wobei y_1, \dots, y_N die Helligkeits- bzw. Lautstärkewerte sind. (Wer sich wundert, daß vor dieser Summe mit N Summanden nur $N-1$ im Nenner steht, sollte zu §4d) zurückblättern.)

Was die Varianz und die Standardabweichung bedeuten, sieht man wieder deutlich an den Beispielen: Bilder mit geringem Kontrast wie „Tiffany“ oder „Lenna“ haben deutlich geringere Werte als die kontrastreicheren Bilder „Peppers“ und „Sailboat“.

e) Kovarianz und Korrelation von Zufallsvariablen

So, wie wir sie bislang definiert haben, ist jede Zufallsvariable ein eigenständiger Prozeß, und zwei verschiedene Zufallsvariablen haben nichts miteinander zu tun. Das ist natürlich nicht das, was wir für die Beschreibung von Bild- und Audiodaten brauchen; hier müssen wir davon ausgehen, daß ein einziger Prozeß gleichzeitig einen ganzen Vektor



Peppers

$$\begin{aligned} \mu &= 115,6 \\ \sigma^2 &= 5632 \\ \sigma &= 75,0 \\ \rho &= 0,98 \\ x_{\min} &= 0 \\ x_{\max} &= 237 \end{aligned}$$



Lenna

$$\begin{aligned} \mu &= 99,1 \\ \sigma^2 &= 2796 \\ \sigma &= 52,9 \\ \rho &= 0,97 \\ x_{\min} &= 3 \\ x_{\max} &= 248 \end{aligned}$$

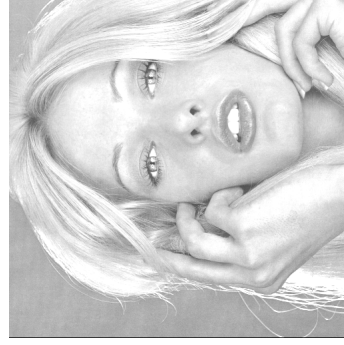


Sailboat

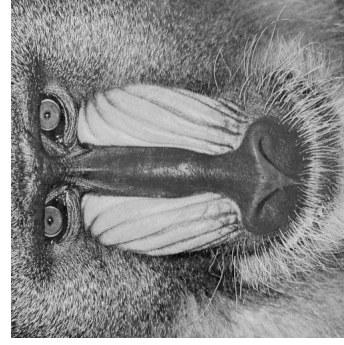
$$\begin{aligned} \mu &= 124,3 \\ \sigma^2 &= 6027 \\ \sigma &= 77,6 \\ \rho &= 0,97 \\ x_{\min} &= 0 \\ x_{\max} &= 249 \end{aligned}$$

**Stream**

$$\begin{aligned}\mu &= 113,8 \\ \sigma^2 &= 2996 \\ \sigma &= 54,7 \\ \rho &= 0,94 \\ x_{\min} &= 0 \\ x_{\max} &= 255\end{aligned}$$

**Tiffany**

$$\begin{aligned}\mu &= 208,6 \\ \sigma^2 &= 1126 \\ \sigma &= 33,6 \\ \rho &= 0,87 \\ x_{\min} &= 3 \\ x_{\max} &= 255\end{aligned}$$

**Baboon**

$$\begin{aligned}\mu &= 128,9 \\ \sigma^2 &= 2282 \\ \sigma &= 47,8 \\ \rho &= 0,86 \\ x_{\min} &= 0 \\ x_{\max} &= 236\end{aligned}$$

bzw. eine ganze Matrix von Zufallswerten erzeugt, wobei deren einzelne Komponenten dann sehr wohl voneinander abhängig sein können.

Für zwei solche Komponenten X und Y mit jeweiligen Wertebereichen $\{x_0, \dots, x_m\}$ und $\{y_0, \dots, y_n\}$ sowie Wahrscheinlichkeiten p_i für x_i und q_j für y_j ist die Wahrscheinlichkeit dafür, daß X den Wert x_i liefert und Y den Wert y_j dann nicht $p_i q_j$, wie das bei unabhängigen Variablen der Fall wäre, sondern irgendeine Wahrscheinlichkeit π_{ij} , von der wir nur wissen, daß aus offensichtlichen Gründen etwa

$$\sum_{j=0}^n \pi_{ij} = p_i \quad \text{und} \quad \sum_{i=0}^m \pi_{ij} = q_j$$

sein muß. Für so ein Paar definieren wir

Definition: a) Die Kovarianz eines solchen Paares (X, Y) ist

$$\begin{aligned}\text{cov}(X, Y) &= E((X - E(X))(Y - E(Y))) \\ &= \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} ((x_i - E(X))(y_j - E(Y))).\end{aligned}$$

b) Die Korrelation von (X, Y) ist $\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$.

Wir bezeichnen die beiden Zufallsvariablen X und Y entsprechend der üblichen Definition für Ereignisse als als voneinander unabhängig, falls für alle i, j gilt: $\pi_{ij} = p_i q_j$. Alsdann ist

$$\begin{aligned}\text{cov}(X, Y) &= \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} (x_i - E(X))(y_j - E(Y)) \\ &= \sum_{i=0}^m \sum_{j=0}^n [p_i (x_i - E(X))] [q_j (y_j - E(Y))] \\ &= \left(\sum_{i=0}^m p_i (x_i - E(X)) \right) \left(\sum_{j=0}^n q_j (y_j - E(Y)) \right) = 0,\end{aligned}$$

denn

$$\sum_{i=0}^m p_i (x_i - E(X)) = \sum_{i=0}^m p_i x_i - \sum_{i=0}^m p_i E(X) = \sum_{i=0}^m p_i x_i - E(X)$$

verschwindet nach Definition des Erwartungswerts.

Damit haben zwei voneinander unabhängige Zufallsvariablen also Kovarianz und Korrelation null; man sagt auch, sie seien unkorreliert. Bei Bilddaten wird das im allgemeinen nicht der Fall sein; hier wird man im Gegenteil davon ausgehen, daß die Zufallsvariablen zu benachbarten Pixeln sehr stark miteinander korrelieren. Wir können beispielsweise annehmen, daß

$$Y = \rho X + Z$$

ist mit einer von X unabhängigen Zufallsvariablen Z und einer positiven reellen Zahl $\rho < 1$. Dann ist

$$E(Y) = E(\rho X + Z) = \rho E(X) + E(Z);$$

falls wir annehmen, daß X und Y denselben Erwartungswert haben, ist daher

$$E(Z) = (1 - \rho)E(X),$$

was für ein ρ nahe eins deutlich kleiner ist als $E(Z)$.

In dieser Situation ist

$$\begin{aligned} \text{cov}(X, Y) &= \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} (x_i - E(X)) (y_j - E(Y)) \\ &= \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} (x_i - E(X)) (\rho x_i + z_j - \rho E(X) - E(Z)) \\ &= \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} (x_i - E(X)) [\rho(x_i - E(X)) + (z_j - E(Z))] \\ &= \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} \rho (x_i - E(X))^2 \\ &\quad + \sum_{i=0}^m \sum_{j=0}^n \pi_{ij} (x_i - E(X)) (z_j - E(Z)) \\ &= \rho \sum_{i=0}^m p_i (x_i - E(X))^2 + \text{cov}(X, Z) = \rho \sigma_X^2, \end{aligned}$$

da X und Z voneinander unabhängige Zufallsvariablen sind.

Wenn wir jetzt noch annehmen, daß $\sigma_X = \sigma_Y$ ist, folgt

$$\rho(X, Y) = \frac{\rho \sigma_X}{\sigma_X \sigma_Y} = \rho,$$

wir können auf diese Weise also für beliebiges $\rho \in [0, 1]$ ein Paar voneinander abhängiger Zufallsvariablen mit Korrelation ρ erzeugen.

Besser noch: Wann immer zwei Zufallsvariablen mit gleicher Standardabweichung Korrelation ρ haben, sind wir immer im obigen Fall, denn definieren wir eine neue Zufallsvariable Z durch $Z = Y - \rho X$, so ist

$E(Z) = E(Y) - \rho E(X)$, und das Paar (X, Z) ist unkorreliert, da

$$\begin{aligned} \text{cov}(X, Z) &= \sum_{i=0}^m \sum_{j=0}^n (x_i - E(X))(z_j - E(Z)) \\ &= \sum_{i=0}^m \sum_{j=0}^n (x_i - E(X)) \left[(y_j - E(Y)) - \rho(x_j - E(X)) \right] \\ &= \sum_{i=0}^m \sum_{j=0}^n (x_i - E(X))(y_j - E(Y)) \\ &\quad - \rho \sum_{i=0}^m \sum_{j=0}^n (x_i - E(X))^2 \\ &= \text{cov}(X, Y) - \rho \sigma_X^2 = \rho \sigma_X \sigma_Y - \rho \sigma_X^2 = 0. \end{aligned}$$

f) Modellierung von Bild- und Audiodaten durch Zufallsvariable

Unser Modell für Audio- und Bilddaten ist nun folgendes: Die Daten werden gegeben durch einen Vektor b bzw. eine Matrix voneinander abhängiger Zufallsvariablen mit gemeinsamem Erwartungswert μ und Standardabweichung σ ; zwei benachbarte Zufallsvariablen sollen jeweils Korrelation ρ haben. Diese Zahl ρ bezeichnen wir als die *Autokorrelation* des jeweiligen Signals. Wir schätzen sie anhand eines gegebenen Bilds oder Musikstücks als Korrelation zwischen benachbarten Daten.

Da der Begriff der Autokorrelation das wohl am schwersten verständliche der hier eingeführten statistischen Konzepte ist, sind die sechs Testbilder in Richtung fallender Autokorrelation geordnet: Die höchste Autokorrelation hat mit $\rho = 0,98$ das Bild „Peppers“, wo die recht homogenen Flächen der Paprikaschoten dafür sorgen, daß sich ein Pixel nur selten von seinen Nachbarn unterscheidet; auch „Lenna“ und „Sailboat“ werden von flächigen Strukturen dominiert. Bei „Stream“ kommen in stärkerem Maße feine Verästelungen von Bäumen und Büschen ins Spiel, so daß die Autokorrelation auf 0,94 absinkt, und bei „Tiffany“ und „Baboon“ schließlich sorgen die vielen Haare für feine Details, die die Autokorrelation auf 0,87 bzw. 0,86 herunterdrücken.

Um zu sehen, wie sich die Autokorrelation zur Komprimierung der Daten ausnutzen läßt, betrachten wir der Einfachheit halber zunächst nur eine Folge X_1, \dots, X_n von Zufallsvariablen; der gemeinsame Erwartungswert sei μ , die gemeinsame Standardabweichung σ , und die Korrelation zwischen X_i und X_{i+1} sei jeweils ρ .

Nach dem oben Gesagten gibt es dann für jede der Zufallsvariablen X_i mit $i < n$ eine davon unabhängige Zufallsvariable Z_i , so daß

$$X_{i+1} = \rho X_i + Z_i$$

ist; entsprechend ist für $i < n - 1$

$$X_{i+2} = \rho X_{i+1} + Z_{i+1} = \rho^2 X_i + \rho Z_i + Z_{i+1}$$

usw.; wenn wir zusätzlich annehmen, daß alle Z_i voneinander unabhängig sind, ist also $\rho(X_i, X_{i+2}) = \rho^2$ und allgemein

$$\rho(X_i, X_j) = \rho^{|i-j|}.$$

In der Signalverarbeitung spricht man bei einer solchen Folge von Zufallsvariablen von einem *autoregressive Prozeß*; der hier betrachtete allereinfachste Fall, bei dem alle Korrelationen nur von der Korrelation zwischen zwei benachbarten Zufallsvariablen abhängen, wird als AR(1)-Modell bezeichnet; in der Sprechweise der Wahrscheinlichkeitstheorie handelt es sich hier um spezielle sogenannte MARKOV-Ketten.



Der russische Mathematiker ANDREI ANDREEVICH MARKOV (1856–1922) studierte in Sankt Petersburg, wo er später auch Professor wurde. Er beschäftigte sich zunächst hauptsächlich mit Zahlentheorie und Analysis; erst später kommen die Wahrscheinlichkeitstheoretischen Arbeiten, für die er heute vor allem bekannt ist. MARKOV-Ketten sind Prozesse ohne Erinnerung, in denen das zukünftige Verhalten nur vom augenblicklichen Zustand abhängt, nicht aber von der Geschichte des Systems. Damit sind sie gerade hier bei Bilddaten nur eine unvollkommene Approximation an die Realität, aber dennoch sehr nützlich.

Falls wir bei einer solchen Folge von Zufallsvariablen die Werte von X_1, \dots, X_n nacheinander übertragen, übertragen wir zuerst den

Wert von X_1 , dann mit X_2 noch einmal zu $100 \times \rho$ % denselben Wert, mit X_3 dasselbe noch einmal zu $100 \times \rho^2$ %, usw.

Eine offensichtliche Alternative hierzu wäre, nur den Wert von X_1 zu übertragen und ansonsten nur die Werte der Z_i . Eine ähnliche Vorgehensweise wird tatsächlich gelegentlich angewandt, allerdings macht man es sich dann noch einfacher und überträgt nur die *Differenzen*, also

$$X_1, X_2 - X_1, \dots, X_n - X_{n-1}.$$

Der Nachteil dieses Verfahrens ist, daß sowohl diese Differenzen als auch die Z_i von Zeit zu Zeit sehr groß werden *müssen*, da es in fast jedem Bild oder Musikstück gelegentliche abrupte Veränderungen gibt.

g) Komprimierung durch Dekorrelation

Die Idee hinter allen Komprimierungsverfahren, die auf Transformationen beruhen, ist es, anstelle der Zufallsvariablen X_i geeignete Linearkombinationen

$$Y_i = \sum_{j=1}^n \alpha_{i,j} X_j$$

zu betrachten, wobei $(\alpha_{i,j})$ eine *invertierbare* $n \times n$ -Matrix ist, so daß sich auch umgekehrt die X_i wieder aus den Y_j rekonstruieren lassen. Diese Matrix wird so gewählt, daß die neuen Variablen möglichst unkorreliert sind und daß man die Größe der neuen Variablen möglichst gut abschätzen kann.

Letzterer Aspekt erfordert statistische Betrachtungen, auf die wir hier verzichten wollen; die Dekorrelation der Zufallsvariablen aber führt uns geradewegs zu Eigenvektoren symmetrischer Matrizen:

Wir definieren für eine Folge von Zufallsvariablen deren *Korrelationsmatrix*

$$\text{Kor}(X_1, \dots, X_n) \in \mathbb{R}^{n \times n}$$

dadurch, daß der Eintrag an der Stelle ij dieser Matrizen jeweils die Korrelation $\rho(X_i, X_j)$ sein soll.

Das Ideal, auf das wir hinarbeiten, sind Zufallsvariablen, deren Korrelationsmatrix eine Diagonalmatrix ist, denn dann sind je zwei verschiedene Variablen unkorreliert.

Da die Korrelationsmatrix eine symmetrische Matrix ist, gibt es eine Orthonormalbasis des \mathbb{R}^n aus reellen Eigenvektoren, bezüglich derer sie Diagonalgestalt hat; die Vektoren dieser Orthonormalbasis seien

$$\vec{b}_1 = \begin{pmatrix} \alpha_{11} \\ \vdots \\ \alpha_{1n} \end{pmatrix}, \dots, \vec{b}_n = \begin{pmatrix} \alpha_{n1} \\ \vdots \\ \alpha_{nn} \end{pmatrix}.$$

Wir definieren die neuen Zufallsvariablen durch

$$Y_i = \sum_{j=1}^n \alpha_{i,j} X_j;$$

dann ist

$$\begin{aligned} \rho(Y_i, Y_k) &= \vec{v}_{Y_i} \cdot \vec{v}_{Y_k} = \left(\sum_{j=1}^n \alpha_{i,j} \vec{v}_{X_j} \right) \cdot \left(\sum_{\ell=1}^n \alpha_{k,\ell} \vec{v}_{X_\ell} \right) \\ &= \sum_{j=1}^n \sum_{\ell=1}^n \alpha_{i,j} \vec{v}_{X_j} \cdot \vec{v}_{X_\ell} \alpha_{k,\ell} = \sum_{j=1}^n \sum_{\ell=1}^n \alpha_{i,j} \rho(X_j, X_\ell) \alpha_{k,\ell} \\ &= \sum_{\ell=1}^n \left(\sum_{j=1}^n \alpha_{i,j} \rho(X_j, X_\ell) \right) \alpha_{k,\ell}. \end{aligned}$$

Der Inhalt der großen Klammer ist offensichtlich der Eintrag an der Stelle $i\ell$ der Produktmatrix $A \cdot \text{Kor}(X_1, \dots, X_n)$, wobei $A = (\alpha_{i,j})$ die Matrix der Koeffizienten $\alpha_{i,j}$ ist, und die Summation über ℓ macht daraus den Eintrag an der Stelle ik des Produkts mit ${}^t A$. Insgesamt haben wir also gezeigt, daß

$$\text{Kor}(Y_1, \dots, Y_n) = A \cdot \text{Kor}(X_1, \dots, X_n) \cdot {}^t A$$

ist. Nun müssen wir nur noch beachten, daß die Spaltenvektoren der Matrix A als die Vektoren einer Orthonormalbasis des \mathbb{R}^n gewählt waren; der Eintrag an der Stelle ij der Matrix ${}^t A$ ist also das Standardskalarprodukt des i -ten und des j -ten Vektors aus einer Orthonormalbasis und

somit null für $i \neq j$ und eins für $i = j$. Daher ist $A \cdot {}^t A = E$, also ${}^t A^{-1}$ und somit auch

$$\text{Kor}(Y_1, \dots, Y_n) = A \cdot \text{Kor}(X_1, \dots, X_n) \cdot A^{-1}.$$

Damit ist $\text{Kor}(Y_1, \dots, Y_n)$ eine Diagonalmatrix, denn für jede Matrix $B \in \mathbb{R}^{n \times n}$ ist ABA^{-1} die Matrix B bezüglich der Basis aus den Spaltenvektoren von A . Diese Basis besteht hier aber aus lauter Eigenvektoren der Korrelationsmatrix, die transformierte Matrix ist also eine Diagonalmatrix.

Unter den Annahmen unseres statistischen Modells können wir also jede Folge von Zufallsvariablen durch eine lineare Transformation in eine Folge unkorrelierter Zufallsvariablen überführen. Diese Transformation bezeichnet man, obwohl sie zuerst von HOTELLING vorgeschlagen wurde, als KARHUNEN-LOÈVE-Transformation.



HAROLD HOTELLING (1895–1973) war ein amerikanischer Statistiker und Ökonom; er lehrte an der Columbia University und der University of North Carolina. In einer 1933 veröffentlichten Arbeit im *Journal of Educational Psychology* schlug er erstmalig diese Transformation vor, die von Statistikern heute in Anlehnung an den Titel seiner Arbeit meist als *Hauptkomponentenanalyse* bezeichnet wird. In Europa erschien die Transformation fast gleichzeitig um 1947 bzw. 1948 in wahrscheinlichkeits-theoretischen Arbeiten des Finnen KARI KARHUNEN (* 1915) und des Franzosen MICHEL LOÈVE (1907–1979), nach denen sie in der technischen Literatur benannt wird.

Die Matrix A der linearen Transformation hängt nur von ρ ab und kann daher für gängige Werte von ρ vorberechnet werden; die KARHUNEN-LOÈVE-Transformation ist also einfach die Multiplikation mit einer bekannten Matrix.

h) Die diskrete Cosinus-Transformation

Für die Multiplikation zweier $n \times n$ -Matrizen benötigt man allerdings n^3 Multiplikationen und noch einmal $n^2(n-1)$ Additionen; der Aufwand steigt mit großem n also sehr stark an. In der Praxis gibt man

sich daher mit einem Kompromiß zufrieden und zerlegt eine Folge von Zufallszahlen in kurze Teilsequenzen, die bei eindimensionalen Folgen typischerweise die Länge 8 haben; dies ist beispielsweise der Standard bei Musik-CDs.

Allerdings wird weder bei Musik-CDs noch sonstwo die KARHUNEN-LOÈVE-Transformation wirklich angewandt. Der Grund liegt an der Struktur der Eigenvektoren der Korrelationsmatrix: Betrachten wir etwa als typisches Beispiel den $n = 8$; dann haben wir die Matrix

$$\text{Cov}(X_1, \dots, X_8) = \begin{pmatrix} 1 & \rho & \rho^2 & \rho^3 & \rho^4 & \rho^5 & \rho^6 & \rho^7 \\ \rho & 1 & \rho & \rho^2 & \rho^3 & \rho^4 & \rho^5 & \rho^6 \\ \rho^2 & \rho & 1 & \rho & \rho^2 & \rho^3 && \rho^4 & \rho^5 \\ \rho^3 & \rho^2 & \rho & 1 & \rho & \rho^2 & \rho^3 & \rho^4 \\ \rho^4 & \rho^3 & \rho^2 & \rho & 1 & \rho & \rho^2 & \rho^3 \\ \rho^5 & \rho^4 & \rho^3 & \rho^2 & \rho & 1 & \rho & \rho^2 \\ \rho^6 & \rho^5 & \rho^4 & \rho^3 & \rho^2 & \rho & 1 & \rho \\ \rho^7 & \rho^6 & \rho^5 & \rho^4 & \rho^3 & \rho^2 & \rho & 1 \end{pmatrix}.$$

Ihre Eigenwerte für $\rho = 0,95$ können zumindest näherungsweise berechnet werden, und auch die Eigenvektoren lassen sich bestimmen. Diese sollen hier jedoch nicht numerisch angegeben werden: Eine Folge von acht reellen Zahlen ist schließlich im allgemeinen eher unanschaulich. Stattdessen sind in den Abbildungen 70 bis 77 die Eigenvektoren *graphisch* dargestellt, wobei einem Vektor

$$(a_1, \dots, a_8) \in \mathbb{R}^8$$

die acht Striche vom Punkt $(i, 0)$ bis (i, a_i) in der Ebenen entsprechen sollen. Zusätzlich ist in jedes dieser Diagramme noch eine der Kurven

$$y = \cos\left(\frac{(2x-1)(j-1)\pi}{16}\right)$$

für $j = 1, \dots, 8$ eingezeichnet; wie man sieht, lassen sich die Komponenten der Eigenvektoren sehr gut durch diese Cosinuswerte annähern. Dies gilt nicht nur für den speziellen Wert $\rho = 0,95$, sondern für jeden Wert von ρ , der hinreichend nahe bei eins liegt.

Aus diesem Grund arbeitet man in der Praxis lieber mit den Cosinuswerten; der Basiswechsel hin zur Basis aus den Cosinusvektoren bezeichnet

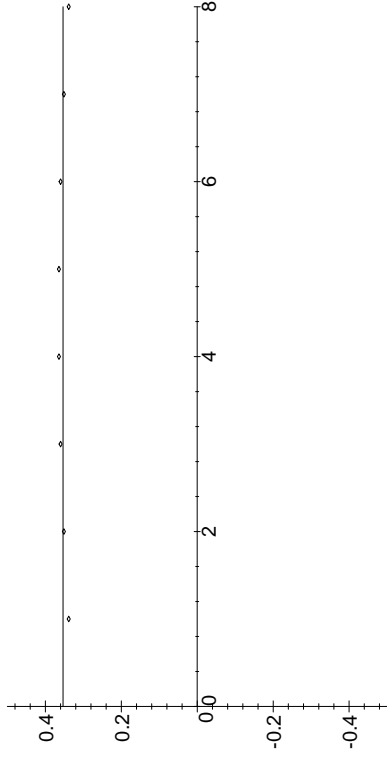


Abb. 70: Der erste Eigenvektor der Korrelationsmatrix

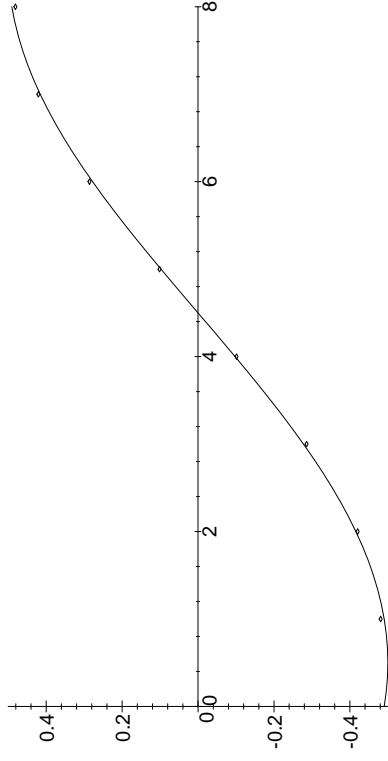


Abb. 71: Der zweite Eigenvektor der Korrelationsmatrix

man als *diskrete Cosinustransformation*. Ihr Hauptvorteil gegenüber der KARHUNEN-LOÈVE-Transformation ist, daß sie durch einen schnellen Algorithmus berechnet werden kann, der anstelle des Aufwands n^3 für eine Matrixmultiplikation nur den Aufwand $n^2 \log n$ hat. Für Einzelheiten sei auf die Vorlesung *Numerik I* verwiesen.

Die diskrete Cosinustransformation ist Teil fast aller gängiger Normen

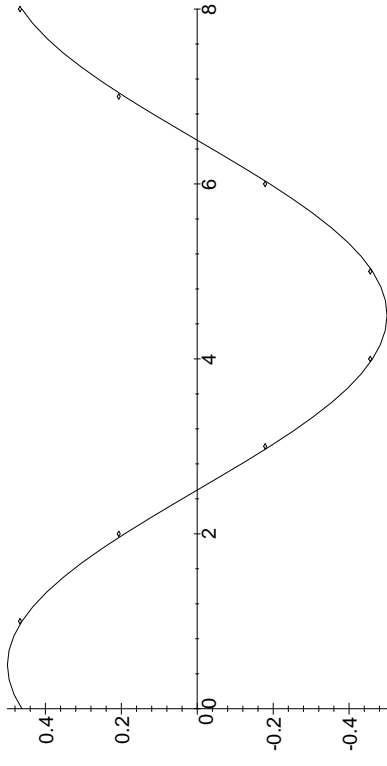


Abb. 72: Der dritte Eigenvektor der Korrelationsmatrix

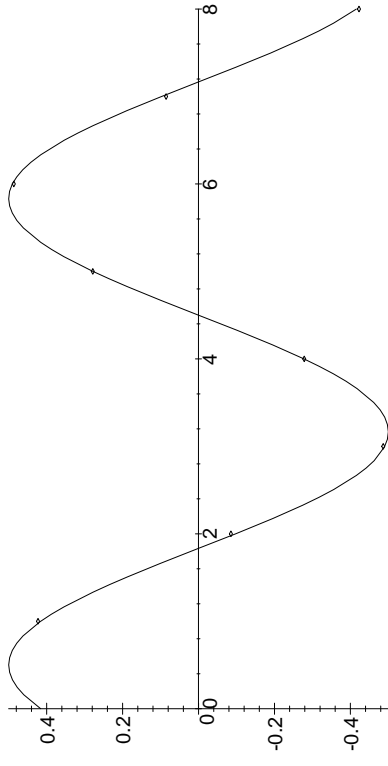


Abb. 73: Der vierte Eigenvektor der Korrelationsmatrix

zur Bildkomprimierung. Sowohl der JPEG-Standard für Photographien, die Standards MPEG 1 und 2 für digitale (Unterhaltungs-)Videos als auch der Standard CCITT H.261 für Videokonferenzen enthalten (neben anderen Bestandteilen) jeweils eine diskrete Cosinustransformation. Auch bei Audio-CDs ist sie ein Teil der Codierung.

Die Transformation allein ist natürlich noch keine Komprimierung:

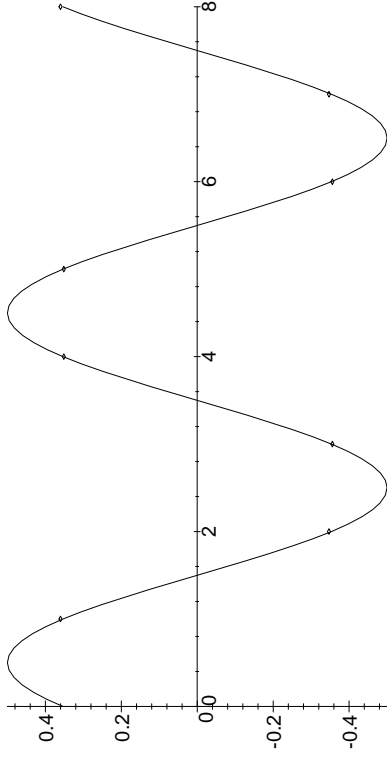


Abb. 74: Der fünfte Eigenvektor der Korrelationsmatrix

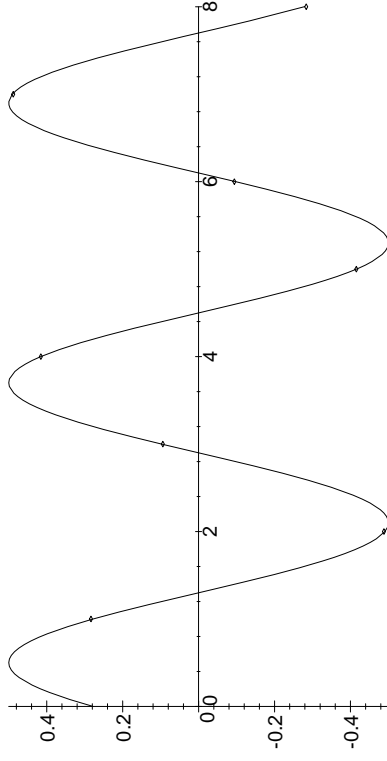


Abb. 75: Der sechste Eigenvektor der Korrelationsmatrix

Schließlich haben wir nur einen Vektor in einer anderen Basis hingeschrieben, und die Anzahl der reellen Zahlen, die man zur Beschreibung eines solchen Vektors benötigt, ist unabhängig von der Basis. Der wesentliche Vorteil der neuen Basis ist, daß man statistisch recht gute Aussagen über die Größe der Komponenten machen können. Hier wollen wir auf exakte statistische Berechnungen verzichten und stattdessen

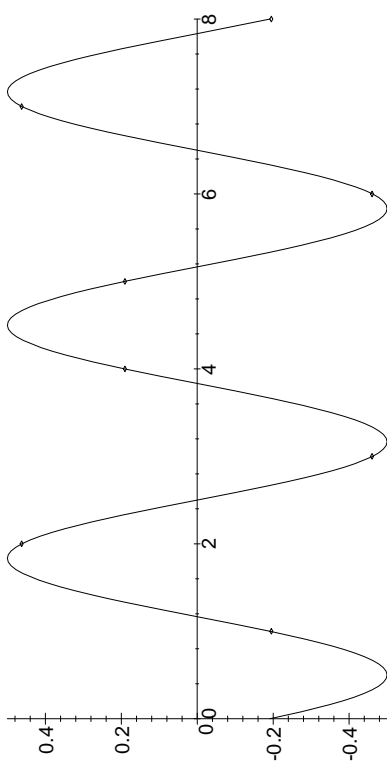


Abb. 76: Der siebte Eigenvektor der Korrelationsmatrix

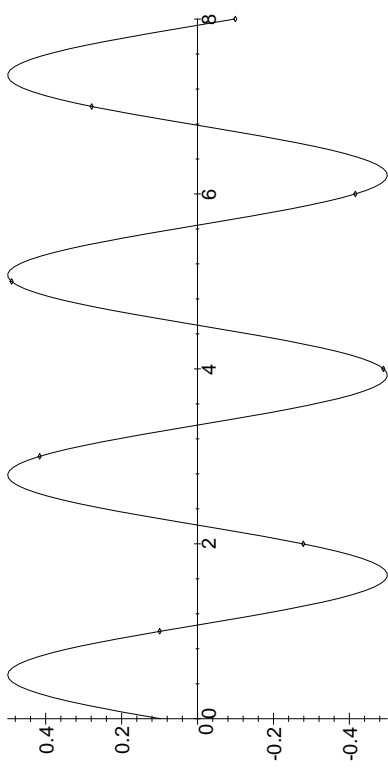


Abb. 77: Der achte Eigenvektor der Korrelationsmatrix

informell diskutieren, warum dies der Fall sein könnte.

Wie die Abbildungen der Basisvektoren zur KARHUNEN-LOÈVE-Transformation und die Formeln für die Basisvektoren zur diskreten Cosinustransformation zeigen, werden die Basisvektoren, wenn man sie in der hier angegebenen Reihenfolge betrachtet, immer hochfrequenter. In einem hinreichend fein abgetasteten Bild oder Audiosignal erwarten

wir, daß hochfrequente Schwankungen keine große Rolle spielen und somit die entsprechenden Basisvektoren nur kleine Koeffizienten haben oder in vielen Fällen sogar gleich gar nicht auftreten. Dementsprechend genügt es, für die Übertragung dieser Koeffizienten nur wenige Bits bereitzustellen; bei nur geringen Abstrichen an die Qualität kann man auf gewisse Koeffizienten sogar ganz verzichten.

Ein Kompressionsverfahren wird daher, je nach Anspruch an die Qualität, entweder alle Koeffizienten des Signals in der neuen Basis übertragen und durch eine geeignete Darstellung der Daten dafür sorgen, daß Folgen von Nullen nur wenig Platz benötigen, oder aber es wird nur eine Auswahl der Koeffizienten übertragen und auch für diese jeweils festlegen, wie viele Bit dafür in Anspruch genommen werden. Diese Anzahl wird umso geringer sein, je höher die Frequenz des jeweiligen Basisvektors ist; bei einigen Verfahren wie etwa JPEG können die Anzahlen auch variabel in Abhängigkeit von einer Qualitätszahl gewählt werden.

Zum Schluß sei noch ganz kurz erwähnt, daß die KARHUNEN-LOÈVE-Transformation und damit (mit ganz geringen Abstrichen) auch die diskrete Cosinustransformation zwar die Korrelationsmatrix in optimaler Weise diagonalisieren, daß aber daraus nicht folgt, daß sie auch optimale Kompressionsverfahren liefern: Ausßer der Kovarianz gibt es noch weitere Quellen für Redundanz eines Bildes.

Ein gewisser Nachteil der Cosinustransformation ist außerdem, daß man für abrupte Übergänge, wie sie etwa bei Kanten immer wieder einmal auftauchen, die hochfrequenten Basisvektoren braucht, die dann aber nicht nur die Kante selbst beeinflussen, sondern das gesamte Quadrat, auf das die Transformation angewandt wird.

Eine bessere Möglichkeit wäre es daher, wenn man anstelle von Cosinusfunktionen Funktionen verwenden könnte, die sowohl im Zeit- als auch im Frequenzbereich lokalisiert sind. Solche Funktionen gibt es in der Tat, etwa die sogenannten *Wavelets*. Hierbei handelt es sich um schnell abklingende Wellen, und neuere Arbeiten deuten darauf hin, daß diese für gewisse Bildmodelle (die im Gegensatz zum hier betrachteten nicht mit Wahrscheinlichkeiten arbeiten) nicht zu weit vom Optimum

entfernt sein sollten. Im Rahmen dieser Vorlesung ist es jedoch zeitlich weder möglich, auf diese Modelle einzugehen, noch ist an eine genauere Behandlung von Wavelets zu denken.

Einen allgemein verständlichen Überblick über Wavelets findet man etwa bei

BARBARA BURKE HUBBARD: *Wavelets: Die Mathematik der kleinen Wellen, Birkhäuser 1997*;

das zitierte Optimalitätsresultat ist beschrieben im Vortrag

STÉPHANE MALLAT: *Applied Mathematics meets signal processing*

auf dem Internationalen Mathematikerkongress 1998 in Berlin, nachzulesen in Band I der Proceedings, S. 319–338, oder unter <http://www.mathematik.uni-bielefeld.de/documenta/xvol-icm/00/Mallat.MAN.html>.

Eine für Technische Informatiker gut geeignete fundierte Einführung in diesen Themenkreis ist etwa

STÉPHANE MALLAT: *A wavelet tour of signal processing, Academic Press, 1998*.

$$\varepsilon \quad \mathcal{N} \quad \mathcal{D} \quad \varepsilon$$

$$S \quad C \quad \mathcal{H} \quad \ddot{O} \quad \mathcal{N} \quad \varepsilon \quad \mathcal{F} \quad \varepsilon \quad \mathcal{R} \quad \mathcal{I} \quad \varepsilon \quad \mathcal{N} \quad !$$