

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 das selbe 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 Gesamfehler führen:

	-5ε	-4ε	-3ε	-2ε	$-\varepsilon$	0	$+\varepsilon$	$+2\varepsilon$	$+3\varepsilon$	$+4\varepsilon$	$+5\varepsilon$
$n = 0$						1					
$n = 1$			1		1						
$n = 2$			1	2	1						
$n = 3$		1	3	3	1						
$n = 4$	1	4	6	4	1						
$n = 5$	1	5	10	10	5	1					

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 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

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

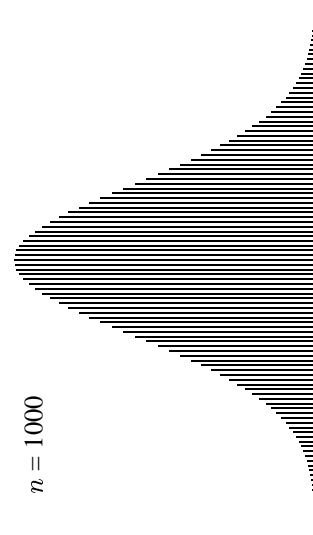
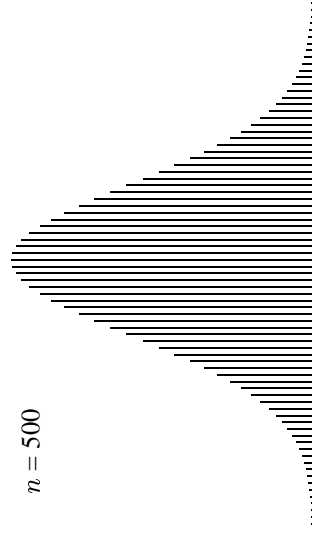
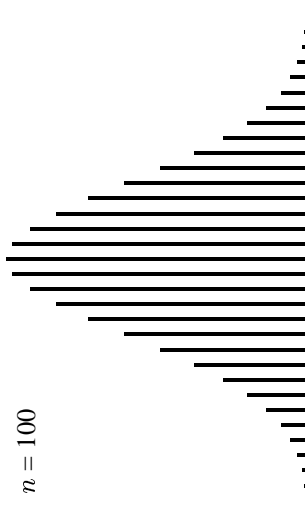
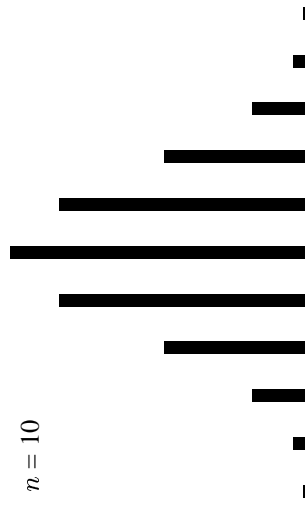
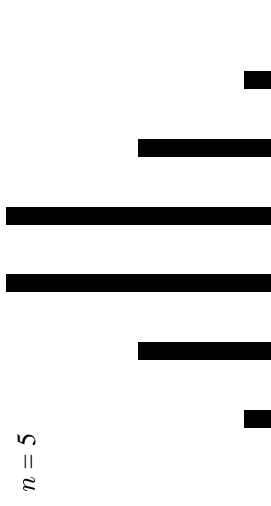
Wer diese Formel nicht kennt, kann sie leicht durch vollständige Induktion beweisen: Für $n = 1$ sowie allgemein für $i = 0$ oder $i = n$ ist alles klar; für $n > 1$ und $0 < i < n$ stehen über $\binom{n}{i}$ die beiden Zahlen $\binom{n-1}{i-1}$ und $\binom{n-1}{i}$, für die in der Tat gilt

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

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 $n = 5, 10, 50, 100, 500, 1000$ 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 $-\varepsilon$ 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 $-\varepsilon$ 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

$$\binom{n}{n-2k} = \binom{n}{k}$$



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

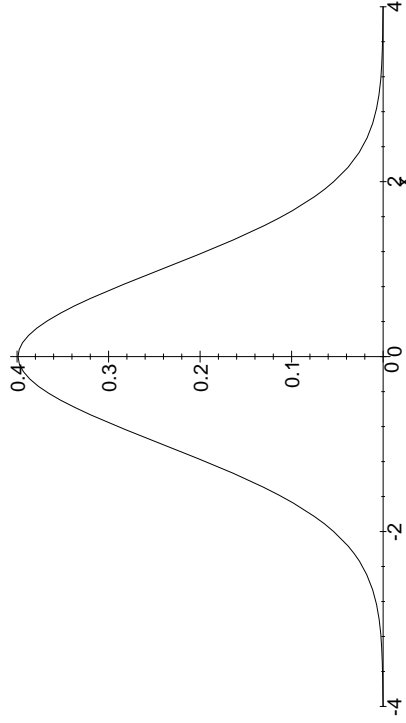


Abb. 65: Die „Glockenkurve“ $y = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

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 fundamentale 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

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

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 \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 Batragmaxima 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

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

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.

§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 geeignet zu wählende Konstante $\sigma > 0$, die Standardabweichung.

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

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

annehmen, was wir in der Form

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

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 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 $\{x\} = x - k$ für $x \in [k, k+1)$.

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

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

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

Für Binomialkoeffizienten folgt, daß

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

c) Die Stirlingsche Formel und die Normalverteilung

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

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

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.

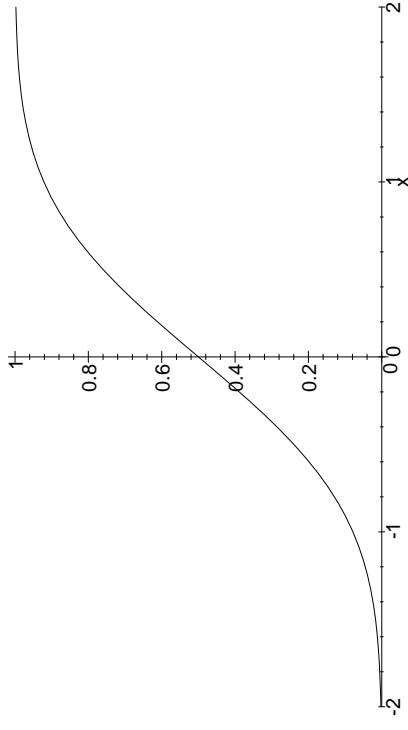


Abb. 68: Das Integral $F(z)$ über die „Glockenkurve“

Mißt man beispielsweise die Temperatur eines Wasserbads eine Viertelstunde lang jede Minute und erhält dabei 15 Meßwerte mit Mittelwert $20,1^\circ\text{C}$ 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 normalisierten 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.

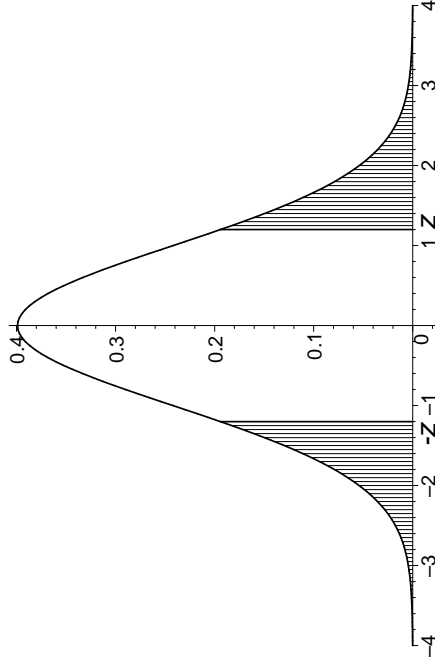


Abb. 69: Zusammenhang zwischen $F(z)$ und $F(-z)$

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

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

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 *Word probability* und das germanische *Word 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 viele gängige 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.

Aus dem letzten Semester kennen wir ein Maß für die gegenseitige Abhängigkeit von Daten: Als wir dort untersuchten, wie das Klausurergebnis eines Studenten von seiner Arbeit bei den wöchentlichen Übungen abhängt oder die Korruption eines Staats vom Bruttozialprodukt pro Einwohner, überprüften wir die Qualität unserer Modelle mit Hilfe des Korrelationskoeffizienten: Dieser lag bei ± 1 bei perfekter Übereinstimmung, und nahe Null, wenn das Modell keinen Zusammenhang zwischen den Daten lieferte.

Dieselbe Technik können wir auch anwenden, um Abhängigkeiten innerhalb einer Folge zu finden; bevor wir diese sogenannte Autokorrelation verstehen können, brauchen wir aber zunächst noch einige Vorbereitungen aus der Stochastik.

b) Zufallsvariablen und ihre statistischen Kenngrößen

Ein guter Komprimierungsalgorithmus muß auch für Bilder funktionieren, die wir erst in ein paar Jahren photographieren. 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* $p_i m$ -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_{ij}) 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

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

c) 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 Beispielbildern, 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 Beispielbildern: Bilder mit geringem Kontrast wie „Tiffany“ oder „Lenna“ haben deutlich geringere Werte als die kontrastreicheren Bilder „Peppers“ und „Sailboat“.

d) 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 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

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

sein muß. Für so ein Paar definieren wir