

Wolfgang K. Seiler

# Zahlentheorie

Vorlesung an der Universität Mannheim  
im Frühjahrssemester 2018

Dieses Skriptum entsteht parallel zur Vorlesung und soll mit möglichst geringer Verzögerung erscheinen. Es ist daher in seiner Qualität auf keinen Fall mit einem Lehrbuch zu vergleichen; insbesondere sind Fehler bei dieser Entstehungsweise nicht nur möglich, sondern **sicher**. Dabei handelt es sich wohl leider nicht immer nur um harmlose Tippfehler, sondern auch um Fehler bei den mathematischen Aussagen. Da mehrere Teile aus anderen Skripten für Hörerkreise der verschiedensten Niveaus übernommen sind, ist die Präsentation auch teilweise ziemlich inhomogen.

Das Skriptum sollte daher mit Sorgfalt und einem gewissen Mißtrauen gegen seinen Inhalt gelesen werden. Falls Sie Fehler finden, teilen Sie mir dies bitte persönlich oder per e-mail (seiler@math.uni-mannheim.de) mit. Auch wenn Sie Teile des Skriptums unverständlich finden, bin ich für entsprechende Hinweise dankbar.

Falls genügend viele Hinweise eingehen, werde ich von Zeit zu Zeit Listen mit Berichtigungen und Verbesserungen zusammenstellen. In der online Version werden natürlich alle bekannten Fehler korrigiert.

Biographische Angaben von Mathematikern beruhen größtenteils auf den entsprechenden Artikeln im *MacTutor History of Mathematics archive* ([www-history.mcs.st-andrews.ac.uk/history/](http://www-history.mcs.st-andrews.ac.uk/history/)), von wo auch die meisten abgedruckten Bilder stammen. Bei noch lebenden Mathematikern bezog ich mich, soweit möglich, auf deren eigenen Internetauftritt.

KAPITEL I: GANZE ZAHLEN UND IHRE PRIMZERLEGUNG .....	1
§1: Rationale und irrationale Zahlen .....	1
§2: Der Euklidische Algorithmus .....	5
§3: Der erweiterte Euklidische Algorithmus .....	8
§4: Der Aufwand des Euklidischen Algorithmus .....	13
§5: Die multiplikative Struktur der ganzen Zahlen .....	18
§6: Kongruenzenrechnung .....	21
§7: Der chinesische Restesatz .....	27
§8: Prime Restklassen .....	33
KAPITEL II: ANWENDUNGEN IN DER KRYPTOLOGIE .....	39
§1: New directions in cryptography .....	39
§2: Das RSA-Verfahren .....	47
§3: Weitere Anwendungen des RSA-Verfahrens .....	51
<i>a)</i> Identitätsnachweis .....	51
<i>b)</i> Elektronische Unterschriften .....	52
<i>c)</i> SSL und TLS .....	53
<i>d)</i> Blinde Unterschriften und elektronisches Bargeld .....	55
<i>e)</i> Bankkarten mit Chip .....	58
§4: Wie groß sollten die Primzahlen sein? .....	60
§5: Praktische Gesichtspunkte .....	63
§6: Verfahren mit diskreten Logarithmen .....	66
§6: DSA .....	69
§7: Ausblick .....	71

KAPITEL III: PRIMZAHLEN .....	73
§1: Die Verteilung der Primzahlen .....	73
§2: Das Sieb des Eratosthenes .....	92
§3: Fermat-Test und Fermat-Zahlen .....	94
§4: Der Test von Miller und Rabin .....	105
§5: Der Test von Agrawal, Kayal und Saxena .....	107
 KAPITEL IV: FAKTORISIERUNGSVERFAHREN .....	 119
§1: Die ersten Schritte .....	121
<i>a)</i> Test auf Primzahl .....	121
<i>b)</i> Abdividieren kleiner Primteiler .....	122
§2: Die Verfahren von Pollard und ihre Varianten .....	124
<i>a)</i> Die Monte-Carlo-Methode .....	125
<i>b)</i> Die $(p - 1)$ -Methode .....	130
<i>c)</i> Varianten .....	132
§3: Das Verfahren von Fermat und seine Varianten .....	134
 KAPITEL V: KETTENBRÜCHE .....	 145
§1: Der Kettenbruchalgorithmus .....	145
§2: Geometrische Formulierung .....	148
§3: Optimale Approximation .....	152
§4: Kettenbrüche und Kalender .....	162
§5: Eine kryptographische Anwendung .....	174
§6: Die Kettenbruchentwicklung der Eulerschen Zahl .....	177

KAPITEL VI: GAUSSSCHE ZAHLEN UND QUATERNIONEN .....	183
§1: Der Ring der Gaußschen Zahlen .....	183
§2: Euklidische Ringe .....	184
§3: Quaternionen .....	189
 KAPITEL VII: QUADRATISCHE FORMEN .....	 193
§1: Summen zweier Quadrate .....	193
§2: Anwendung auf die Berechnung von $\pi$ .....	199
§3: Der Satz von Lagrange .....	205
 KAPITEL VIII: QUADRATISCHE RESTE .....	 209
§1: Das Legendre-Symbol .....	209
§2: Das quadratische Reziprozitätsgesetz .....	211
§3: Das Jacobi-Symbol .....	215
§4: Anwendungen quadratischer Reste .....	219
a) Münzwurf per Telephon .....	219
b) Akustik von Konzerthallen .....	221
 KAPITEL IX: DIE FERMAT-VERMUTUNG FÜR ZAHLEN UND FÜR POLYNOME	227
§1: Zahlen und Funktionen .....	227
§2: Pythagoräische Tripel .....	229
§3: Der Satz von Mason .....	232
§4: Die abc-Vermutung .....	235
§5: Die Frey-Kurve .....	239



# Kapitel 1

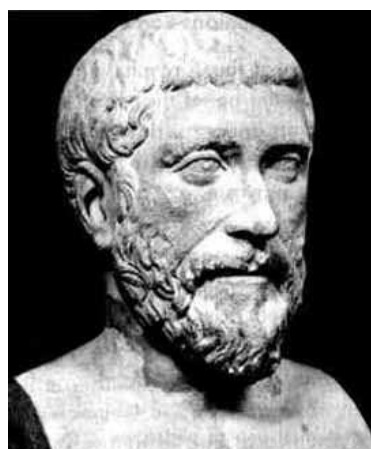
## Ganze Zahlen und ihre Primzerlegung

### § 1: Rationale und irrationale Zahlen

Die Zahlentheorie beschäftigt sich, wie schon ihr Name sagt, mit Zahlen. Nun würde allerdings eine Umfrage wohl ergeben, daß sich nach Ansicht eines Großteils der Bevölkerung die gesamte Mathematik mit Zahlen beschäftigt – auch wenn beispielsweise im neunbändigen Analysislehrbuch von JEAN DIEUDONNÉ abgesehen von den dreiteiligen Abschnittsnummern praktisch keine Zahlen außer 0, 1 und 2 vorkommen.

Die Zahlentheorie unterscheidet sich dadurch von anderen Teilen der Mathematik, daß es dort vor allem um *ganze* Zahlen geht, oft sogar einfach um die natürlichen Zahlen  $1, 2, 3, \dots$ . Die frühe griechische Philosophie der Pythagoräer beispielsweise stand unter dem Motto *Alles ist Zahl*. Sie konnten die musikalischen Harmonien auf einfache Verhältnisse natürlicher Zahlen zurückführen und waren überzeugt, daß dies auch für alle anderen Proportionen galt.

Umso größer war der Schock, als um 450 v.Chr. einer von ihnen, wahrscheinlich HIPASSOS VON METAPONT, herausfand, daß es in der Geometrie Längenverhältnisse gibt, die sich *nicht* so beschreiben lassen. Schlimmer noch: Ein Beispiel dafür bietet ausgerechnet das Wahrzeichen der Pythagoräer, das Pentagramm. HIPASSOS VON METAPONT nahm deshalb auch ein schlimmes Ende: Nach einigen Überlieferungen wurde er von den erzürnten Pythagoräern ertränkt, nach anderen ließen ihn die Götter als Strafe für seine Schandtat bei einem Schiffsuntergang ertrinken.



PYTHAGORAS VON SAMOS lebte etwa von 569 bis 475. Als ungefähr 18-Jähriger besuchte er THALES in Milot und ging auf dessen Rat 535 nach Ägypten, um mehr über Mathematik und Astronomie zu lernen. Im Tempel von Diospolis wurde er nach der dafür vorgesehen Ausbildung in die Priesterschaft aufgenommen. 525, bei der persischen Invasion Ägyptens, geriet er in Gefangenschaft und wurde nach Babylon gebracht, was er nutzte er, um die dortige Mathematik zu erlernen. 520 kam er frei und kehrte zurück nach Samos, zog aber schon bald weiter nach Croton in Süditalien, wo er eine religiöse und philosophische Schule gründete, die Pythagoräer.

Wir wollen uns hier mit einem einfacheren Fall als dem des Pentagramms beschäftigen, dem Verhältnis zwischen der Diagonale und der Seite eines Quadrats. In einem Quadrat der Seitenlänge  $a$  kann die Diagonale aufgefaßt werden als Hypotenuse eines gleichschenkligen rechtwinkligen Dreiecks mit Katheten der Länge  $a$ ; sie hat daher nach dem Satz des PYTHAGORAS (der tatsächlich wohl einige hundert Jahre älter als PYTHAGORAS sein dürfte) die Länge  $a\sqrt{2}$ , und das gesuchte Verhältnis ist  $\sqrt{2}$ .

Wäre  $\sqrt{2}$  als Verhältnis  $a/b$  zweier natürlicher Zahlen darstellbar, so könnten wir ohne Beschränkung der Allgemeinheit annehmen, daß mindestens eine der beiden Zahlen  $a$  und  $b$  ungerade ist: Andernfalls müßten wir einfach so lange durch zwei kürzen, bis dies der Fall ist.

Quadrieren wir beide Seiten der Gleichung  $a/b = \sqrt{2}$ , so erhalten wir die neue Gleichung  $a^2/b^2 = 2$ ; demnach müßte  $a^2 = 2b^2$  eine gerade Zahl sein. Damit müßte aber auch  $a$  gerade sein, denn wäre  $a = 2c + 1$  ungerade, so auch  $a^2 = 2 \cdot (2c^2 + c) + 1$ . Wenn aber  $a$  durch zwei teilbar ist, ist  $a^2 = 2b^2$  durch vier teilbar, also wäre auch  $b^2$  und damit auch  $b$  gerade, ein Widerspruch. Somit ist  $\sqrt{2}$  keine rationale Zahl.

Auch das Verhältnis zwischen Umfang und Durchmesser eines Kreises, für das wir heute die Bezeichnung  $\pi$  verwenden, ist nicht rational, allerdings erfordert der Beweis dafür etwas mehr Arbeit. Ich möchte ihn trotzdem hier vorstellen, denn er zeigt sehr schön, wie zahlentheoretische Aussagen teilweise nur auf Umwegen über andere Gebiete der Mathematik bewiesen werden können. Beim folgenden Beweis führt



dieser Umweg über die reelle Analysis:

Wir gehen zunächst aus von einem beliebigen Polynom  $P(x)$  mit reellen Koeffizienten von einem geraden Grad  $2n$ . Dazu definieren wir das Polynom

$$Q(x) \stackrel{\text{def}}{=} P(x) - P''(x) + P^{(4)}(x) - \dots + (-1)^n P^{(2n)}(x)$$

als die alternierende Summe der Ableitungen gerader Ordnung von  $P$ . Weiter betrachten wir die Funktion  $S(x) = Q'(x) \sin x - Q(x) \cos x$ ; ihre Ableitung ist

$$\begin{aligned} S'(x) &= Q''(x) \sin x + Q'(x) \cos x - Q'(x) \cos x + Q(x) \sin x \\ &= (Q''(x) + Q(x)) \sin x. \end{aligned}$$

In  $Q''(x) = P''(x) - P^{(4)}(x) + \dots + (-1)^{n-1} P^{(2n)}(x)$  kommen bis auf  $P(x)$  genau dieselben Terme vor wie in  $Q(x)$ , allerdings mit dem jeweils anderen Vorzeichen. Daher ist  $Q''(x) + Q(x) = P(x)$  und  $S(x) = Q'(x) \sin x - Q(x) \cos x$  ist eine Stammfunktion von  $P(x) \sin x$ . Damit folgt die

### Formel von Hermite:

$$\int_0^\pi P(x) \sin x \, dx = S(\pi) - S(0) = Q(\pi) + Q(0),$$

denn  $\sin 0 = \sin \pi = 0$ ,  $\cos 0 = 1$  und  $\cos \pi = -1$ .

Wir nehmen nun an,  $\pi = a/b$  sei eine rationale Zahl, und wenden die gerade bewiesene Formel an auf das spezielle Polynom

$$P_n(x) \stackrel{\text{def}}{=} \frac{x^n(a - bx)^n}{n!} = \frac{b^n}{n!} x^n (\pi - x)^n;$$

wir erhalten

$$I_n \stackrel{\text{def}}{=} \int_0^\pi P_n(x) \sin x \, dx = Q_n(\pi) + Q_n(0)$$

mit  $Q_n(x) \stackrel{\text{def}}{=} P_n(x) - P_n''(x) + P_n^{(4)}(x) - \dots + (-1)^n P_n^{(2n)}(x)$ .

$P_n(x)$  ist im Intervall  $(0, \pi)$  genau wie die Sinusfunktion überall positiv; daher ist  $I_n > 0$ . Außerdem ist  $P_n(x)$  symmetrisch zu  $\pi/2$ , denn wie die zweite Darstellung der Funktion zeigt, vertauschen sich einfach die beiden Faktoren  $x^n$  und  $(\pi - x)^n$ , wenn wir  $x$  durch  $\pi - x$  ersetzen.

Das Maximum von  $P_n$  in  $(0, \pi)$  wird an derselben Stelle angenommen wie das der Funktion  $f(x) = x(\pi - x)$ ; da  $f'(x) = \pi - 2x$  nur in Intervallmitte  $\pi/2 = a/2b$  verschwindet, liegt es dort und dort ist

$$P_n\left(\frac{\pi}{2}\right) = P_n\left(\frac{a}{2b}\right) = \frac{1}{n!} \left(\frac{a}{2b}\right)^n \left(\frac{ab}{2b}\right)^n = \frac{1}{n!} \frac{a^{2n}}{(4b)^n}.$$

Schätzen wir das Integral ab durch Intervalllänge mal Maximum des Integranden, erhalten wir daher die Ungleichung

$$I_n \leq \pi \cdot \frac{1}{n!} \frac{a^{2n}}{(4b)^n}$$

und sehen, daß  $I_n$  für  $n \rightarrow \infty$  gegen Null geht:  $n!$  wächst bekanntlich schneller als jede Potenz einer reellen Zahl. Somit ist

$$\lim_{n \rightarrow \infty} I_n = 0.$$

Andererseits ist aber  $I_n = Q_n(0) + Q_n(\pi)$ , was in diesem Fall wegen der Symmetrie von  $P_n$  einfach  $2Q_n(0)$  ist. Wir wollen uns überlegen, daß alle  $Q_n(0)$  ganze Zahlen sind. Dazu reicht es zu zeigen, daß alle Ableitungen von  $P_n(x)$  an der Stelle Null ganzzahlige Werte annehmen. Nach der binomischen Formel ist

$$P_n(x) = \frac{x^n(a - bx)^n}{n!} = \frac{1}{n!} \sum_{i=0}^n \binom{n}{i} a^{n-i} (-b)^i x^{n+i};$$

die  $k$ -te Ableitung verschwindet also für  $k < n$  an der Stelle Null. Für  $k = n + i \geq n$  ist

$$P_n^{(k)}(0) = \frac{1}{n!} \binom{n}{i} a^{n-i} (-b)^i (n+i)!$$

ebenfalls eine ganze Zahl, da der Nenner  $n!$  Teiler von  $(n+i)!$  ist und ansonsten nur ganze Zahlen dastehen.

Somit ist also  $I_n$  für jedes  $n \in \mathbb{N}$  eine positive ganze Zahl. Der Limes einer Folge positiver ganzer Zahlen kann aber unmöglich Null sein,

also führt die Annahme,  $\pi = a/b$  sei eine rationale Zahl, zu einem Widerspruch, der die Irrationalität von  $\pi$  zeigt.

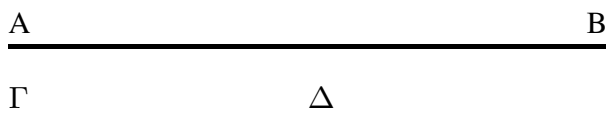
Wenn man schon dabei ist, kann man leicht auch noch viele andere wichtige Zahlen als irrational erkennen; auf dem ersten Übungsblatt ist ein Beweis für die Irrationalität der EULERSchen Zahl  $e$  skizziert, und mit nur wenig mehr Aufwand als im Fall der Quadratwurzel aus zwei läßt sich auch leicht zeigen, daß *jede* Quadratwurzel einer natürlichen Zahl entweder ganzzahlig oder irrational ist. Dasselbe gibt für höhere Wurzeln und sogar für Nullstellen aller Polynome mit ganzzahligen Koeffizienten und höchstem Koeffizient eins, allerdings brauchen wir zum Beweis einen Satz, den zwar fast jeder kennt, dessen Beweis man aber nur selten sieht: Die eindeutige Primzerlegung der natürlichen Zahlen. Den ersten vollständigen Beweis dieser Aussage fand GAUSS 1798; er erschien 1801 in seinen *Disquisitiones Arithmeticae*. Heute verwendet man beim Beweis meist eine Konstruktion, die wahrscheinlich bereits den Pythagoräern bekannt war, die wir aber als EUKLIDischen Algorithmus bezeichnen. Wie sich zeigen wird, ist er zusammen mit einer ganzen Reihe von Varianten ein nicht nur in der Zahlentheorie allgegenwärtiges Werkzeug; es lohnt sich also, ihn gleichs jetzt zum Beginn der Vorlesung etwas ausführlicher zu betrachten.

## §2: Der Euklidische Algorithmus

Bei EUKLID, in Proposition 2 des siebten Buchs seiner *Elemente*, wird er (in der Übersetzung von CLEMENS THAER in Oswalds Klassikern der exakten Wissenschaft) so beschrieben:

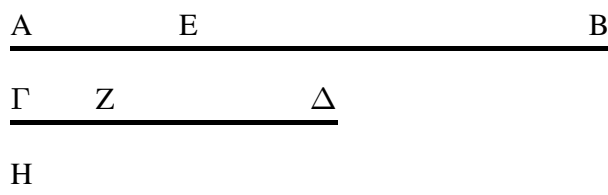
*Zu zwei gegebenen Zahlen, die nicht prim gegeneinander sind, ihr größtes gemeinsames Maß zu finden.*

Die zwei gegebenen Zahlen, die nicht prim, gegeneinander sind, seien  $AB, \Gamma\Delta$ . Man soll das größte gemeinsame Maß von  $AB, \Gamma\Delta$  finden.



Wenn  $\Gamma\Delta$  hier  $AB$  mißt – sich selbst mißt es auch – dann ist  $\Gamma\Delta$  gemeinsames Maß von  $\Gamma\Delta, AB$ . Und es ist klar, daß es auch das größte ist, denn keine Zahl größer  $\Gamma\Delta$  kann  $\Gamma\Delta$  messen.

Wenn  $\Gamma\Delta$  aber  $AB$  nicht mißt, und man nimmt bei  $AB$ ,  $\Gamma\Delta$  abwechselnd immer das kleinere vom größeren weg, dann muß (schließlich) eine Zahl übrig bleiben, die die vorangehende mißt. Die Einheit kann nämlich nicht übrig bleiben; sonst müßten  $AB, \Gamma\Delta$  gegeneinander prim sein, gegen die Voraussetzung. Also muß eine Zahl übrig bleiben, die die vorangehende mißt.  $\Gamma\Delta$  lasse, indem es  $BE$  mißt,  $EA$ , kleiner als sich selbst übrig; und  $EA$  lasse, indem es  $\Delta Z$  mißt,  $Z\Gamma$ , kleiner als sich selbst übrig; und  $\Gamma Z$  messe  $AE$ .



Da  $\Gamma Z$   $AE$  mißt und  $AE$   $\Delta Z$ , muß  $\Gamma Z$  auch  $\Delta Z$  messen; es mißt aber auch sich selbst, muß also auch das Ganze  $\Gamma\Delta$  messen.  $\Gamma\Delta$  mißt aber  $BE$ ; also mißt  $\Gamma Z$  auch  $BE$ ; es mißt aber auch  $EA$ , muß also auch das Ganze  $BA$  messen. Und es mißt auch  $\Gamma\Delta$ ;  $\Gamma Z$  mißt also  $AB$  und  $\Gamma\Delta$ ; also ist  $\Gamma Z$  gemeinsames Maß von  $AB, \Gamma\Delta$ . Ich behaupte, daß es auch das größte ist. Wäre nämlich  $\Gamma Z$  nicht das größte gemeinsame Maß von  $AB, \Gamma\Delta$ , so müßte irgendeine Zahl größer  $\Gamma Z$  die Zahlen  $AB$  und  $\Gamma\Delta$  messen. Dies geschehe; die Zahl sei  $H$ . Da  $H$  dann  $\Gamma\Delta$  mäßt und  $\Gamma\Delta$   $BE$  mißt, mäßt  $H$  auch  $BE$ ; es soll aber auch das Ganze  $BA$  messen, müßte also auch den Rest  $AE$  messen.  $AE$  mißt aber  $\Delta Z$ ; also müßte  $H$  auch  $\Delta Z$  messen; es soll aber auch das Ganze  $\Delta\Gamma$  messen, müßte also auch den Rest  $\Gamma Z$  messen, als größere Zahl die kleinere; dies ist unmöglich. Also kann keine Zahl größer  $\Gamma Z$  die Zahlen  $AB$  und  $\Gamma\Delta$  messen;  $\Gamma Z$  ist also das größte gemeinsame Maß von  $AB, \Gamma\Delta$ ; dies hatte man beweisen sollen.

Aus heutiger Sicht erscheint hier die Voraussetzung, daß die betrachteten Größen nicht teilerfremd sein dürfen, seltsam. Sie erklärt sich daraus, daß in der griechischen Philosophie und Mathematik die Einheit eine Sonderrolle einnahm und nicht als Zahl angesehen wurde: Die Zahlen begannen erst mit der Zwei. Dementsprechend führt EUKLID in Proposition 1 des siebten Buchs fast wörtlich dieselbe Konstruktion durch für den Fall von teilerfremden Größen. Schon wenig später wurde die Eins auch in Griechenland als Zahl anerkannt, und für uns heute ist die Unterscheidung ohnehin bedeutungslos. Wir können die Bedingung, daß der ggT ungleich eins sein soll, also einfach ignorieren.

Das dem EUKLIDischen Algorithmus zugrunde liegende Prinzip der *Wechselwegnahme* oder wechselseitigen Subtraktion war in der grie-

chischen Mathematik spätestens gegen Ende des fünften vorchristlichen Jahrhunderts bereits wohlbekannt unter dem Namen Antanairesis (ἀνταναιρέσις) oder auch Anthyphairesis (ἀνθυφαιρέσις), und auch der Algorithmus selbst geht mit ziemlicher Sicherheit, wie so vieles in den Elementen, *nicht* erst auf EUKLID zurück: Seine *Elemente* waren das wohl mindestens vierte Buchprojekt dieses Namens, und alles spricht dafür, daß er vieles von seinen Vorgängern übernommen hat. Seine Elemente waren dann aber mit Abstand die erfolgreichsten, so daß die anderen in Vergessenheit gerieten und verloren gingen; EUKLID wurde schließlich als *der* Stoichist bekannt nach dem griechischen Titel στοιχεῖα der Elemente.



Es ist nicht ganz sicher, ob EUKLID (Εὐκλείδης) wirklich gelebt hat; es ist möglich, wenn auch sehr unwahrscheinlich, daß EUKLID nur ein Pseudonym für eine Autorengruppe ist. (Das nebenstehende Bild aus dem 18. Jahrhundert ist reine Phantasie.) EUKLID ist vor allem bekannt als Autor der *Elemente*, in denen er die Geometrie seiner Zeit systematisch darstellte und (in gewisser Weise) auf wenige Definitionen sowie die berühmten fünf Postulate zurückführte; sie entstanden um 300 v. Chr. EUKLID arbeitete wohl am Museion in Alexandrien; außer den Elementen schrieb er ein Buch über Optik und weitere, teilweise verschollene Bücher.

Wenn wir nicht mit Zirkel und Lineal arbeiten, sondern rechnen, können wir die mehrfache „Wegnahme“ einer Strecke von einer anderen einfacher beschreiben durch eine Division mit Rest: Sind  $a$  und  $b$  die (als natürliche Zahlen vorausgesetzten) Längen der beiden Strecken und ist  $a : b = q$  Rest  $r$ , so kann man  $q$  mal die Strecke  $b$  von  $a$  wegnehmen; was übrig bleibt ist eine Strecke der Länge  $r < b$ .

EUKLIDS Konstruktion wird dann zu folgendem Algorithmus für zwei natürliche Zahlen  $a, b$ :

**Schritt 0:** Setze  $r_0 = a$  und  $r_1 = b$ .

**Schritt  $i, i \geq 1$ :** Falls  $r_i$  verschwindet, endet der Algorithmus mit  $\text{ggT}(a, b) = r_{i-1}$ ; andernfalls sei  $r_{i+1}$  der Rest bei der Division von  $r_{i-1}$  durch  $r_i$ .

EUKLID behauptet, daß dieser Algorithmus stets endet und daß das Ergebnis der größte gemeinsame Teiler der Ausgangszahlen  $a, b$  ist, d.h. die größte natürliche Zahl, die sowohl  $a$  als auch  $b$  teilt.

Da der Divisionsrest  $r_{i+1}$  stets echt kleiner ist als sein Vorgänger  $r_i$  und eine Folge immer kleiner werdender nichtnegativer ganzer Zahlen notwendigerweise nach endlich vielen Schritten die Null erreicht, muß der Algorithmus in der Tat stets enden. Daß er mit dem richtigen Ergebnis endet, ist ebenfalls leicht zu sehen, denn im  $i$ -ten Schritt ist

$$r_{i-1} = q_i r_i + r_{i+1} \quad \text{oder} \quad r_{i+1} = r_{i-1} - q_i r_i,$$

so daß jeder gemeinsame Teiler von  $r_i$  und  $r_{i+1}$  auch ein Teiler von  $r_{i-1}$  ist und umgekehrt jeder gemeinsame Teiler von  $r_{i-1}$  und  $r_i$  auch  $r_{i+1}$  teilt. Somit haben  $r_i$  und  $r_{i-1}$  dieselben gemeinsamen Teiler wie  $r_i$  und  $r_{i+1}$ , insbesondere haben sie denselben größten gemeinsamen Teiler. Durch Induktion folgt, daß in jedem Schritt  $\text{ggT}(r_i, r_{i-1}) = \text{ggT}(a, b)$  ist. Im letzten Schritt ist  $r_i = 0$ ; da jede natürliche Zahl Teiler der Null ist, ist dann  $r_{i-1} = \text{ggT}(r_i, r_{i-1}) = \text{ggT}(a, b)$ , wie behauptet.

### §3: Der erweiterte Euklidische Algorithmus

Mehr als zwei Tausend Jahre nach der Entdeckung von Anthyphaire-sis und EUKLIDISCHEM Algorithmus, 1624 in Bourg-en-Bresse, stellte BACHET DE MÉZIRIAC in der zweiten Auflage seines Buchs *Problèmes plaisants et délectables qui se font par les nombres* Aufgaben wie die folgende:

*Il y a 41 personnes en un banquet tant hommes que femmes et enfants qui en tout dépensent 40 sous, mais chaque homme paye 4 sous, chaque femme 3 sous, chaque enfant 4 deniers. Je demande combien il y a d'hommes, combien de femmes, combien d'enfants.*

(Bei einem Bankett sind 41 Personen, Männer, Frauen und Kinder, die zusammen vierzig Sous ausgeben, aber jeder Mann zahlt vier Sous, jede Frau drei Sous und jedes Kind 4 Deniers. Ich frage, wie viele Männer, wie viele Frauen und wie viele Kinder es sind.)

Sobald man weiß, daß zwölf Deniers ein Sou sind (und zwanzig Sous ein Pfund), kann man dies in ein lineares Gleichungssystem übersetzen:

Ist  $x$  die Zahl der Männer,  $y$  die der Frauen und  $z$  die der Kinder, so muß gelten  $x + y + z = 41$  und  $4x + 3y + \frac{1}{3}z = 40$ .

Im Gegensatz zum Fall der in Schule und Linearer Algebra betrachteten Gleichungssystemen kommen hier natürlich nur natürliche Zahlen als Lösungen in Frage.



CLAUDE GASPAR BACHET SIEUR DE MÉZIRIAC (1581-1638) verbrachte den größten Teil seines Lebens in seinem Geburtsort Bourg-en-Bresse. Er studierte bei den Jesuiten in Lyon und Milano und trat 1601 in den Orden ein, trat aber bereits 1602 wegen Krankheit wieder aus und kehrte nach Bourg zurück. Sein Buch erschien 1612; 1959 brachte der Verlag Blanchard eine vereinfachte Ausgabe heraus. Am bekanntesten ist BACHET für seine lateinische Übersetzung der *Arithmetika* von DIOPHANTOS. In einem Exemplar davon schrieb FERMAT seine Vermutung an den Rand. Auch Gedichte von BACHET sind erhalten. 1635 wurde er Mitglied der französischen Akademie der Wissenschaften.

Zur Lösung kann man zunächst die erste Gleichung nach  $z$  auflösen und in die zweite Gleichung einsetzen; dies führt auf die Gleichung

$$\frac{11}{3}x + \frac{8}{3}y = \frac{79}{3} \quad \text{oder} \quad 11x + 8y = 79.$$

Bei einer solchen Gleichung ist *a priori* nicht klar, ob es überhaupt Lösungen gibt: Die Gleichung  $10x + 8y = 79$  beispielsweise kann keine haben, denn für ganze Zahlen  $x, y$  ist  $10x + 8y$  stets gerade. Allgemein kann  $ax + by = c$  höchstens dann ganzzahlige Lösungen haben, wenn der ggT von  $a$  und  $b$  Teiler von  $c$  ist.

BACHET DE MÉZIRIAC hat bewiesen, daß sie in diesem Fall auch stets Lösungen hat; das Kernstück dazu ist seine Proposition XVIII, wo er zu zwei teilerfremden Zahlen  $a, b$  ganze Zahlen  $x, y$  konstruiert, für die  $ax - by = 1$  ist: *Deux nombres premiers entre eux estant donnéz, treuver le moindre multiple de chascun d'iceux, surpassant de l'unité un multiple de l'autre*. Die Methode ist eine einfache Erweiterung des EUKLIDischen Algorithmus, und genau wie letzterer nach EUKLID benannt ist, da ihn dieser rund 150 Jahre nach seiner Entdeckung in seinem Lehrbuch darstellte, heißt auch BACHETS Satz heute *Identität von*

BÉZOUT, weil dieser ihn 142 Jahre später, im Jahre 1766, in seinem Lehrbuch beschrieb (und auf Polynome verallgemeinerte).



ETIENNE BÉZOUT (1730-1783) wurde in Nemours in der Ile-de-France geboren, wo seine Vorfahren Magistrate waren. Er ging stattdessen an die Akademie der Wissenschaften; seine Hauptbeschäftigung war die Zusammenstellung von Lehrbüchern für die Militärausbildung. Im 1766 erschienenen dritten Band (von vier) seines *Cours de Mathématiques à l'usage des Gardes du Pavillon et de la Marine* ist die Identität von BÉZOUT dargestellt. Seine Bücher waren so erfolgreich, daß sie auch ins Englische übersetzt und als Lehrbücher z.B. in Harvard benutzt wurden. Heute ist er vor allem auch bekannt durch seinen Beweis, daß sich zwei Kurven der Grade  $n$  und  $m$  in höchstens  $nm$  Punkten schneiden können.

Zur Lösung von Problemen wie dem von BACHET wollen wir gleich allgemein den größten gemeinsamen Teiler zweier Zahlen als Linearkombination dieser Zahlen darstellen. Dazu ist nur eine kleine Erweiterung des EUKLIDISCHEN Algorithmus notwendig, so daß man oft auch einfach vom erweiterten EUKLIDISCHEN Algorithmus spricht.

Die Gleichung

$$r_{i-1} = q_i r_i + r_{i+1}$$

läßt sich umschreiben als

$$r_{i+1} = r_{i-1} - q_i r_i,$$

so daß  $r_{i+1}$  eine ganzzahlige Linearkombination von  $r_i$  und  $r_{i-1}$  ist. Da entsprechend auch  $r_i$  Linearkombination von  $r_{i-1}$  und  $r_{i-2}$  ist, folgt induktiv, daß der ggT von  $a$  und  $b$  als ganzzahlige Linearkombination von  $a$  und  $b$  dargestellt werden kann.

Algorithmisch sieht dies folgendermaßen aus:

**Schritt 0:** Setze  $r_0 = a$ ,  $r_1 = b$ ,  $\alpha_0 = \beta_1 = 1$  und  $\alpha_1 = \beta_0 = 0$ . Für  $i = 1$  ist dann

$$r_{i-1} = \alpha_{i-1}a + \beta_{i-1}b \quad \text{und} \quad r_i = \alpha_i a + \beta_i b.$$

Im  $i$ -ten Schritt werden neue Zahlen berechnet derart, daß diese Gleichungen auch für  $i + 1$  gelten:



**Schritt  $i$ ,  $i \geq 1$ :** Falls  $r_i$  verschwindet, endet der Algorithmus mit

$$\text{ggT}(a, b) = r_{i-1} = \alpha_{i-1}a + \beta_{i-1}b.$$

Andernfalls dividiere man  $r_{i-1}$  durch  $r_i$ ; der Divisionsrest sei  $r_{i+1}$ . Dann ist

$$\begin{aligned} r_{i+1} &= r_{i-1} - q_i r_i = (\alpha_{i-1}a + \beta_{i-1}b) - q_i(\alpha_i a + \beta_i b) \\ &= (\alpha_{i-1} - q_i \alpha_i)a + (\beta_{i-1} - q_i \beta_i)b; \end{aligned}$$

die gewünschten Gleichungen gelten also für

$$\alpha_{i+1} = \alpha_{i-1} - q_i \alpha_i \quad \text{und} \quad \beta_{i+1} = \beta_{i-1} - q_i \beta_i.$$

Genau wie oben folgt, daß der Algorithmus für alle natürlichen Zahlen  $a$  und  $b$  endet und daß am Ende der richtige ggT berechnet wird; außerdem sind die  $\alpha_i$  und  $\beta_i$  so definiert, daß in jedem Schritt  $r_i = \alpha_i a + \beta_i b$  ist, insbesondere wird also im letzten Schritt der ggT als Linearkombination der Ausgangszahlen dargestellt.

Als Beispiel wollen wir den ggT von 200 und 148 als Linearkombination darstellen. Im nullten Schritt haben wir 200 und 148 als die trivialen Linearkombinationen

$$200 = 1 \cdot 200 + 0 \cdot 148 \quad \text{und} \quad 148 = 0 \cdot 200 + 1 \cdot 148.$$

Im ersten Schritt dividieren wir, da 148 nicht verschwindet, 200 mit Rest durch 148:

$$200 = 1 \cdot 148 + 52 \quad \text{und} \quad 52 = 1 \cdot 200 - 1 \cdot 148.$$

Da auch  $52 \neq 0$ , dividieren wir im zweiten Schritt 148 durch 52:

$$148 = 2 \cdot 52 + 44 \quad \text{und} \quad 44 = 148 - 2 \cdot (1 \cdot 200 - 1 \cdot 148) = 3 \cdot 148 - 2 \cdot 200.$$

Auch  $44 \neq 0$ ; wir machen also weiter:  $52 = 1 \cdot 44 + 8$  und

$$8 = 52 - 44 = (1 \cdot 200 - 1 \cdot 148) - (3 \cdot 148 - 2 \cdot 200) = 3 \cdot 200 - 4 \cdot 148.$$

Im nächsten Schritt erhalten wir  $44 = 5 \cdot 8 + 4$  und

$$4 = 44 - 5 \cdot 8 = (3 \cdot 148 - 2 \cdot 200) - 5 \cdot (3 \cdot 200 - 4 \cdot 148) = 23 \cdot 148 - 17 \cdot 200.$$

Bei der Division von acht durch vier schließlich ist der Divisionsrest null; damit ist  $4 = 23 \cdot 148 - 17 \cdot 200$  der ggT von 148 und 200.

Zur Lösung des Problems von BACHET müssen wir die Gleichung  $11x + 8y = 79$  betrachten. Dazu stellen wir zunächst den ggT von 11 und 8 als Linearkombination dieser Zahlen dar.

Elf durch acht ist eins Rest drei, also ist  $3 = 1 \cdot 11 - 1 \cdot 8$ .

Im nächsten Schritt dividieren wir acht durch drei mit dem Ergebnis zwei Rest zwei, also ist  $2 = 1 \cdot 8 - 2 \cdot 3 = 1 \cdot 8 - 2 \cdot (1 \cdot 11 - 1 \cdot 8) = -2 \cdot 11 + 3 \cdot 8$ .

Im letzten Schritt wird daher drei durch zwei dividiert und wir sehen erstens, daß der ggT gleich eins ist (was hier keine Überraschung ist), und zweitens, daß gilt  $1 = 3 - 2 = (1 \cdot 11 - 1 \cdot 8) - (-2 \cdot 11 + 3 \cdot 8) = 3 \cdot 11 - 4 \cdot 8$ .

Damit haben wir auch eine Darstellung von 79 als Linearkombination von elf und acht:

$$79 = 79 \cdot (3 \cdot 11 - 4 \cdot 8) = 237 \cdot 11 - 316 \cdot 8.$$

Dies ist allerdings nicht die gesuchte Lösung: BACHET dachte sicherlich nicht an 237 Männer,  $-316$  Frauen und 119 Kinder.

Nun ist aber die obige Gleichung  $1 = 3 \cdot 11 - 4 \cdot 8$  nicht die einzige Möglichkeit zur Darstellung der Eins als Linearkombination von acht und elf: Da  $8 \cdot 11 - 11 \cdot 8$  verschwindet, können wir ein beliebiges Vielfaches dieser Gleichung dazuaddieren und bekommen die allgemeinere Lösung

$$(3 + 8k) \cdot 11 - (4 + 11k) \cdot 8 = 1.$$

Entsprechend können wir auch ein beliebiges Vielfaches dieser Gleichung zur Darstellung von 79 addieren:

$$79 = (237 + 8k) \cdot 11 - (316 + 11k) \cdot 8.$$

Wir müssen  $k$  so wählen, daß sowohl die Anzahl  $237 + 8k$  der Männer als auch die Anzahl  $-(316 + 11k)$  der Frauen positiv oder zumindest nicht negativ wird, d.h.  $-\frac{237}{8} \leq k \leq -\frac{316}{11}$ . Da  $k$  ganzzahlig sein muß, kommt nur  $k = -29$  in Frage; es waren also fünf Männer, drei Frauen und dazu noch  $41 - 5 - 3 = 33$  Kinder. Ihre Gesamtausgaben belaufen sich in der Tat auf  $5 \cdot 4 + 3 \cdot 3 + 33 \cdot \frac{1}{3} = 40$  Sous.

Entsprechend kann der erweiterte EUKLIDISCHE Algorithmus zur Lösung anderer diophantischer Gleichungen verwendet werden, von Gleichungen also, bei denen nur ganzzahlige Lösungen interessieren. Wir betrachten hier nur die lineare Gleichung  $ax + by = c$  mit  $a, b, c \in \mathbb{Z}$  für zwei Unbekannte  $x, y \in \mathbb{Z}$ .

Der größte gemeinsame Teiler  $d = \text{ggT}(a, b)$  von  $a$  und  $b$  teilt offensichtlich jeden Ausdruck der Form  $ax + by$  mit  $x, y \in \mathbb{Z}$ ; falls  $d$  kein Teiler von  $c$  ist, kann es also keine ganzzahlige Lösung geben.

Ist aber  $c = rd$  ein Vielfaches von  $d$  und ist  $d = \alpha a + \beta b$  die lineare Darstellung des ggT nach dem erweiterten EUKLIDISCHEN Algorithmus, so haben wir mit  $x = r\alpha$  und  $y = r\beta$  offensichtlich eine Lösung gefunden.

Ist  $(x', y')$  eine weitere Lösung, so ist

$$a(x - x') + b(y - y') = c - c = 0 \quad \text{oder} \quad a(x - x') = b(y' - y).$$

$v = a(x - x') = b(y' - y)$  ist also ein gemeinsames Vielfaches von  $a$  und  $b$  und damit auch ein Vielfaches des kleinsten gemeinsamen Vielfachen von  $a$  und  $b$ . Dieses kleinste gemeinsame Vielfache ist  $ab/d$ , es muß also eine ganze Zahl  $m$  geben mit

$$x - x' = m \cdot \frac{b}{d} \quad \text{und} \quad y' - y = m \cdot \frac{a}{d}.$$

Die allgemeine Lösung der obigen Gleichung ist somit

$$x = r\alpha - m \cdot \frac{b}{d} \quad \text{und} \quad y = r\beta + m \cdot \frac{a}{d} \quad \text{mit} \quad m \in \mathbb{Z}.$$

#### §4: Der Aufwand des Euklidischen Algorithmus

Im Beweis, daß der EUKLIDISCHE Algorithmus stets nach endlich vielen Schritten abbricht, hatten wir argumentiert, daß der Divisionsrest stets kleiner ist als der Divisor, so daß er irgendwann einmal null werden muß; dann endet der Algorithmus.

Damit haben wir auch eine obere Schranke für den Rechenaufwand zur Berechnung von  $\text{ggT}(a, b)$ : Wir müssen höchstens  $b$  Divisionen durchführen.

Das erscheint zwar auf den ersten Blick als ein recht gutes Ergebnis; wenn man aber bedenkt, daß der EUKLIDISCHE Algorithmus heute in der Kryptographie auf über 600-stellige Zahlen angewendet wird, verliert diese Schranke schnell ihre Nützlichkeit: Da unser Universum ein geschätztes Alter von zehn Milliarden Jahren, also ungefähr  $3 \cdot 10^{18}$  Sekunden hat, ist klar, daß auch der schnellste heutige Computer, der zu Beginn des Universum zu rechnen begann, bis heute nur einen verschwindend kleinen Bruchteil von  $10^{600}$  Divisionen ausgeführt hätte. Wäre  $10^{600}$  eine realistische Aufwandsabschätzung, könnten wir an eine Anwendung des EUKLIDISCHEN Algorithmus auf 600-stellige Zahlen nicht einmal denken.

Tatsächlich ist  $10^{600}$  aber natürlich nur eine obere Schranke, von der wir bislang noch nicht wissen, wie realistisch sie ist. Um dies zu entscheiden, suchen wir die kleinsten natürlichen Zahlen  $a, b$ , für die  $n$  Divisionen notwendig sind; dies wird uns auf ein bekanntes Problem aus dem 13. Jahrhundert führen.

Im Falle  $n = 1$  sind offensichtlich  $a = b = 1$  die kleinstmöglichen Zahlen; wenn  $a = b$  ist, kommt man immer mit genau einer Division aus.

Dies ist allerdings ein eher untypischer Fall, der sich insbesondere nicht rekursiv verallgemeinern läßt, denn ab dem zweiten Schritt des EUKLIDISCHEN Algorithmus ist der Divisor stets kleiner als der Dividend: Ersterer ist schließlich der Rest bei der vorangegangenen Division und letzterer der Divisor. Die kleinsten natürlichen Zahlen  $a \neq b$ , für die man mit nur einer Division auskommt, sind offensichtlich  $a = 2$  und  $b = 1$ .

Als nächstes suchen wir die kleinsten Zahlen  $a, b$  für die zwei Divisionen notwendig sind. Ist  $r$  der Rest bei der ersten Division, so ist  $b : r$  die zweite Division. Für diese muß  $r \geq 1$  und  $b \geq 2$  sein, und  $a = qb + r$ , wobei  $q$  der Quotient bei der ersten Division ist. Dieser ist mindestens eins, die kleinstmöglichen Werte sind damit

$$r = 1, \quad b = 2 \quad \text{und} \quad a = b + r = 3.$$

Allgemeiner seien  $a_n$  und  $b_n$  die kleinsten Zahlen, für die  $n$  Divisionen notwendig sind, und  $r$  sei der Rest bei der ersten Division. Für die zweite

Division  $b : r$  ist dann  $b_n \geq a_{n-1}$  und  $r \geq b_{n-1}$ ; die kleinstmöglichen Werte sind damit

$$r = b_{n-1}, \quad b_n = a_{n-1} \quad \text{und} \quad a_n = b_n + r = a_{n-1} + b_{n-1} = a_{n-1} + a_{n-2}.$$

Da wir  $a_1 = 2$  und  $b_1 = 1$  kennen, können wir somit alle  $a_n$  und  $b_n$  berechnen; was wir erhalten, sind die sogenannten FIBONACCI-Zahlen.

Sie sind durch folgende Rekursionsformel definiert:

$$F_0 = 0, \quad F_1 = 1 \quad \text{und} \quad F_n = F_{n-1} + F_{n-2} \quad \text{für } n \geq 2.$$

FIBONACCI führte sie ein, um die Vermehrung einer Karnickelpopulation durch ein einfaches Modell zu berechnen. In seinem 1202 erschienenen Buch *Liber abaci* schreibt er:

*Ein Mann bringt ein Paar Karnickel auf einen Platz, der von allen Seiten durch eine Mauer umgeben ist. Wie viele Paare können von diesem Paar innerhalb eines Jahres produziert werden, wenn man annimmt, daß jedes Paar jeden Monat ein neues Paar liefert, das vom zweiten Monat nach seiner Geburt an produktiv ist?*



LEONARDO PISANO (1170–1250) ist heute vor allem unter seinem Spitznamen FIBONACCI bekannt; gelegentlich nannte er sich auch BIGOLLO, auf Deutsch *Tunichtgut* oder *Reisender*. Er ging in Nordafrika zur Schule, kam aber 1202 zurück nach Pisa. Seine Bücher waren mit die ersten, die die indisch-arabischen Ziffern in Europa einführten. Er behandelt darin nicht nur Rechenaufgaben für Kaufleute, sondern auch zahlentheoretische Fragen, beispielsweise daß man die Quadratzahlen durch Aufaddieren der ungeraden Zahlen erhält. Auch betrachtet er nichtlineare Gleichungen, die er approximativ löst, und erinnert an viele in Vergessenheit geratene Ergebnisse der antiken Mathematik.

Wie wir gerade gesehen haben, kann man mit den FIBONACCI-Zahlen nicht nur Karnickelpopulationen beschreiben, sondern – wie GABRIEL LAMÉ 1844 entdeckte – auch eine Obergrenze für den Aufwand beim EUKLIDischen Algorithmus angeben:

**Satz von Lamé (1844):** Die kleinsten natürlichen Zahlen  $a, b$ , für die beim EUKLIDISCHEN Algorithmus  $n \geq 2$  Divisionen benötigt werden, sind  $a = F_{n+2}$  und  $b = F_{n+1}$ . ■

(Für  $n = 1$  gilt der Satz nur, wenn wir zusätzlich voraussetzen, daß  $a \neq b$  ist; für  $n \geq 2$  ist dies automatisch erfüllt.)



GABRIEL LAMÉ (1795–1870) studierte von 1813 bis 1817 Mathematik an der Ecole Polytechnique, danach bis 1820 Ingenieurwissenschaften an der Ecole des Mines. Auf Einladung Alexanders I. kam er 1820 nach Rußland, wo er in St. Petersburg als Professor und Ingenieur unter anderem Vorlesungen über Analysis, Physik, Chemie und Ingenieurwissenschaften hielt. 1832 erhielt er einen Lehrstuhl für Physik an der Ecole Polytechnique in Paris, 1852 einen für mathematische Physik und Wahrscheinlichkeitstheorie an der Sorbonne. 1836/37 war er wesentlich am Bau der Eisenbahnlinien Paris-Versailles und Paris-S<sup>t</sup>. Germain beteiligt.

Um die Zahlen  $F_n$  durch eine geschlossene Formel darzustellen, können wir (genau wie man es auch für die rekursive Berechnung per Computer tun würde) die Definitionsgleichung der FIBONACCI-Zahlen als

$$\begin{pmatrix} F_{n+1} \\ F_n \end{pmatrix} = A \begin{pmatrix} F_n \\ F_{n-1} \end{pmatrix} \quad \text{mit} \quad A = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$$

schreiben; dann ist

$$\begin{pmatrix} F_n \\ F_{n-1} \end{pmatrix} = A^{n-1} \begin{pmatrix} F_1 \\ F_0 \end{pmatrix} = A^{n-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Das charakteristische Polynom von  $A$  ist

$$\det(A - \lambda E) = (1 - \lambda)(-\lambda) - 1 = \lambda^2 - \lambda - 1;$$

die Eigenwerte von  $A$  sind daher  $\lambda_{1/2} = \frac{1}{2} \pm \frac{1}{2}\sqrt{5}$ . Bezeichnet  $B$  die Matrix, deren Spalten aus den zugehörigen Eigenvektoren besteht, so ist

also  $A = B^{-1} \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} B$  und

$$\begin{pmatrix} F_n \\ F_{n-1} \end{pmatrix} = A^{n-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = B^{-1} \begin{pmatrix} \lambda_1^{n-1} & 0 \\ 0 & \lambda_2^{n-1} \end{pmatrix} B \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Auch ohne die Matrix  $B$  zu berechnen, wissen wir somit, daß sich  $F_n$  in der Form  $F_n = u'\lambda_1^{n-1} + v'\lambda_2^{n-1}$  darstellen läßt. Indem wir  $u'$  durch  $\lambda_1$  und  $v'$  durch  $\lambda_2$  dividieren, erhalten wir auch eine Darstellung der Form  $F_n = u\lambda_1^n + v\lambda_2^n$ , deren Koeffizienten wir durch Einsetzen von  $n = 0$  und  $n = 1$  bestimmen können:

$$F_0 = 0 = v\lambda_1^0 + u\lambda_2^0 = u + v \quad \text{und} \quad F_1 = 1 = u\lambda_1 + v\lambda_2.$$

Damit ist  $v = -u$ , und die zweite Gleichung wird zu

$$u(\lambda_1 - \lambda_2) = u\sqrt{5} = 1 \implies u = \frac{1}{\sqrt{5}}, \quad v = -\frac{1}{\sqrt{5}} \quad \text{und} \quad F_n = \frac{\lambda_1^n - \lambda_2^n}{\sqrt{5}}.$$

Numerisch ist

$$\lambda_1 = \frac{1 + \sqrt{5}}{2} \approx 1,618034, \quad \lambda_2 = 1 - \lambda_1 = \frac{1 - \sqrt{5}}{2} \approx -0,618034$$

und  $\sqrt{5} \approx 2,236068$ ; der Quotient  $\lambda_2^n/\sqrt{5}$  ist also für jedes  $n$  betragsmäßig kleiner als  $1/2$ . Daher können wir  $F_n$  auch einfacher berechnen als nächste ganze Zahl zu  $\lambda_1^n/\sqrt{5}$ . Insbesondere folgt, daß  $F_n$  exponentiell mit  $n$  wächst.

Die Gleichung  $\lambda^2 - \lambda - 1 = 0$  läßt sich umschreiben als  $\lambda(\lambda - 1) = 1$  oder  $\lambda : 1 = 1 : (\lambda - 1)$ . Diese Gleichung charakterisiert den *goldenen Schnitt*: Stehen zwei Strecken  $x$  und  $y$  in diesem Verhältnis, so auch die beiden Strecken  $y$  und  $x - y$ . Die positive Lösung  $\lambda_1$  wird traditionell mit dem Buchstaben  $\phi$  bezeichnet;  $F_n$  ist also der zur nächsten ganzen Zahl gerundete Wert von  $\phi^n/\sqrt{5}$ .

Die beiden kleinsten Zahlen, für die wir  $n$  Divisionen brauchen, sind nach LAMÉ  $a = F_{n+2}$  und  $b = F_{n+1}$ . Aus der geschlossenen Formel für die FIBONACCI-Zahlen folgt

$$\begin{aligned} n &\approx \log_\phi \sqrt{5} b - 1 = \log_\phi b + \log_\phi \sqrt{5} - 1 = \frac{\ln b}{\ln \phi} + \frac{\ln \sqrt{5}}{\ln \phi} - 1 \\ &\approx 2,078 \ln b + 0,672. \end{aligned}$$

Für beliebige Zahlen  $a > b$  können nicht mehr Divisionen notwendig sein als für die auf  $b$  folgenden nächstgrößeren FIBONACCI-Zahlen, also gibt obige Formel für jedes  $b$  eine obere Grenze. Die Anzahl der Divisionen wächst daher nicht (wie oben bei der naiven Abschätzung)

wie  $b$ , sondern höchstens wie  $\log b$ . Für sechshundertstellige Zahlen  $a, b$  müssen wir daher nicht mit  $10^{600}$  Divisionen rechnen, sondern mit weniger als drei Tausend, was auch mit weniger leistungsfähigen Computern problemlos und schnell möglich ist.

Tatsächlich gibt natürlich auch die hier berechnete Schranke nur selten den tatsächlichen Aufwand wieder; fast immer werden wir mit erheblich weniger auskommen. Im übrigen ist auch alles andere als klar, ob wir den ggT auf andere Weise nicht möglicherweise schneller berechnen können. Da wir aber für Zahlen der Größenordnung, die in heutigen Anwendungen interessieren, selbst mit der Schranke für den schlimmsten Fall ganz gut leben können, sei hier auf diese Fragen nicht weiter eingegangen. Interessenten finden mehr dazu z.B. in den Abschnitten 4.5.2+3 des Buchs

DONALD E. KNUTH: *The Art of Computer Programming, vol. 2: Seminumerical Algorithms, Addison-Wesley, 2* 1981

Eine deutsche Übersetzung des hier relevanten vierten Kapitels erschien 2001 unter dem Titel *Arithmetik* bei Springer.

## §5: Die multiplikative Struktur der ganzen Zahlen

Eine Primzahl ist bekanntlich eine natürliche Zahl  $p$ , die genau zwei Teiler hat, nämlich die Eins und sich selbst. Der erweiterte EUKLIDISCHE Algorithmus liefert eine wichtige Folgerung aus dieser Definition:

**Lemma:** Wenn eine Primzahl das Produkt  $ab$  zweier natürlicher Zahlen teilt, teilt sie mindestens einen der Faktoren.

*Beweis:* Angenommen, die Primzahl  $p$  sei kein Teiler von  $a$ , teile aber  $ab$ . Da der ggT von  $a$  und  $p$  Teiler von  $p$  und ungleich  $p$  ist, muß er notgedrungen gleich eins sein; es gibt also eine Darstellung

$$1 = \alpha a + \beta p \quad \text{mit} \quad \alpha, \beta \in \mathbb{Z}.$$

Dann ist  $b = \alpha ab + \beta pb$  durch  $p$  teilbar, denn sowohl  $ab$  also auch  $pb$  sind Vielfache von  $p$ . ■

Daraus folgt induktiv



**Korollar:** Wenn eine Primzahl ein Produkt  $a_1 \cdot \dots \cdot a_r$  natürlicher Zahlen teilt, teilt sie mindestens einen der Faktoren. ■

**Hauptsatz der elementaren Zahlentheorie:** Jede natürliche Zahl läßt sich bis auf Reihenfolge eindeutig als ein Produkt von Primzahlen schreiben.

*Beweis:* Die Eins läßt sich nach der üblichen Konvention über leere Produkte als das leere Produkt von Primzahlen darstellen, und es gibt offensichtlich keine andere Darstellung.

Nun sei  $n > 1$ , und wir nehmen an, der Satz gelte für alle Zahlen kleiner  $n$ . Falls  $n$  eine Primzahl ist, haben wir mit  $n = n$  eine Darstellung der gewünschten Art, und diese ist eindeutig, weil  $n$  keine Teiler außer der Eins und sich selbst hat. Andernfalls ist der kleinste von eins verschiedene Teiler  $p$  von  $n$  echt kleiner als  $n$  und muß prim sein, denn ansonsten hätte  $p$  einen Teiler  $q$  mit  $1 < q < p$ , der auch Teiler von  $n$ , aber kleiner als  $p$  wäre. Sei  $n = p \cdot a$ . Dann ist  $a < n$ , läßt sich also eindeutig als Produkt von Primzahlen schreiben. Erweitert man dieses Produkt noch um einen Faktor  $p$ , hat man auch  $n$  als Produkt von Primzahlen dargestellt. Angenommen,  $n$  hat noch eine weitere solche Darstellung. Dann kann in dieser Darstellung kein  $p$  vorkommen, denn sonst hätte auch  $a$  mehrere Produktdarstellungen. Ist andererseits  $n = q_1 \cdot \dots \cdot q_r$  irgendeine Darstellung von  $n$  als Produkt von Primzahlen, so muß  $p$  als Teiler von  $n$  nach obigem Korollar mindestens eine der Primzahlen  $q_i$  teilen, was nur möglich ist, wenn  $q_i = p$  ist. Also gibt es keine weitere Darstellung von  $n$ . ■

Die Verwendung des obigen Lemmas macht diesen Beweis recht kurz und durchsichtig; andererseits geht dadurch der erweiterte EUKLIDISCHE Algorithmus in den Beweis ein, und für diesen mußten wir doch einige Zeit investieren. Mit einer kleinen Modifikation, die auf ERNST ZERMELO (1871–1953) zurückgeht, läßt sich der Gebrauch des Korollars im obigen Beweis vermeiden:

Wir müssen zeigen, daß  $n = pa$  nicht dargestellt werden kann als ein Produkt von Primzahlen, die allesamt von  $p$  verschieden sind. Dazu sei

$q \neq p$  irgendeine Primzahl aus einer Primzerlegung von  $n$  und  $n = qb$ . Dann läßt sich  $d = n - pb$  faktorisieren als

$$d = pa - pb = p(a - b) \quad \text{und} \quad d = qb - pb = (q - p)b.$$

Da  $p$  der kleinste von eins verschiedene Teiler von  $n$  ist, muß  $q > p$  und damit  $b < a$  sein; daher sind die Zahlen  $d, q - p, a - b$  allesamt positiv, und natürlich sind sie auch alle kleiner als  $n$ , haben also nach Induktionsannahme eine eindeutige Primzerlegung. Die von  $d = p(a - b)$  enthält  $p$ . Andererseits ist  $d = (q - p)b$ , die Primzerlegung von  $d$  kann also auch als Produkt der Primzerlegungen von  $q - p$  und  $b$  interpretiert werden. In der Primzerlegung von  $q - p$  kann kein  $p$  vorkommen, denn sonst wäre  $p$  ein Teiler von  $q$ , was für zwei verschiedene Primzahlen nicht möglich ist. Daher kommt  $p$  in der Primzerlegung von  $b$  vor, also auch in der von  $n = qb$ . Somit muß  $p$  in jeder Primzerlegung von  $n$  vorkommen, und damit gibt es bis auf Reihenfolge nur eine Darstellung von  $n$  als Produkt von Primzahlen.

Aus dem Hauptsatz der elementaren Zahlentheorie läßt sich auch das Lemma zu Beginn dieses Paragraphen neu beweisen, und zwar im wesentlichen so, wie es EUKLID in Satz 30 des siebten Buch seiner Elemente bewiesen hat: Angenommen, die Primzahl  $p$  teilt das Produkt zweier Zahlen  $a, b$ , ist aber kein Teiler von  $a$ . Dann läßt sich  $ab$  schreiben als  $pc$  mit einer natürlichen Zahl  $c$ , und  $a/p = c/b$ . Da  $p$  prim ist und  $a$  nicht teilt, ist  $a/p$  ein gekürzter Bruch, der den gleichen Wert hat wie  $c/b$ , also (hier kommt der Hauptsatz ins Spiel, auch wenn es bei EUKLID auch ohne geht) entsteht  $c/b$  aus  $a/p$  durch Erweiterung. Somit ist  $b$  ein Vielfaches von  $p$ .

Als erste Anwendung dieses Hauptsatzes der elementaren Zahlentheorie können wir das bereits am Ende von §1 erwähnte Resultat über die Irrationalität von Nullstellen ganzzahliger Polynome beweisen:

**Satz:** Die reelle Zahl  $x$  erfülle die Gleichung

$$x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0 = 0 \quad \text{mit} \quad a_i \in \mathbb{Z}.$$

Dann ist  $x$  entweder ganzzahlig oder irrational.

*Beweis:* Jede rationale Zahl  $x$  kann als Quotient  $x = p/q$  zweier zueinander teilerfremder ganzer Zahlen  $p$  und  $q$  geschrieben werden. Multiplizieren wir die Gleichung

$$\left(\frac{p}{q}\right)^n + a_{n-1} \left(\frac{p}{q}\right)^{n-1} + \cdots + a_1 \left(\frac{p}{q}\right) + a_0 = 0$$

mit  $q^n$ , erhalten wir die nennerlose Gleichung

$$p^n + a_{n-1}p^{n-1}q + \cdots + a_1pq^{n-1} + a_0q^n = 0.$$

Auflösen nach  $p^n$  führt auf

$$\begin{aligned} p^n &= -a_{n-1}p^{n-1}q - \cdots - a_1pq^{n-1} - a_0q^n \\ &= q(-a_{n-1}p^{n-1} - \cdots - a_1pq^{n-2} - a_0q^{n-1}), \end{aligned}$$

d.h.  $q$  muß ein Teiler von  $p^n$  sein, was wegen der Eindeutigkeit der Primfaktorzerlegung von  $p^n$  sowie der vorausgesetzten Teilerfremdheit von  $p$  und  $q$  nur für  $q = \pm 1$  der Fall sein kann. Somit ist  $x$  eine ganze Zahl, wie behauptet. ■

## §6: Kongruenzenrechnung

Zwei ganze Zahlen lassen sich im allgemeinen nicht durcheinander dividieren. Trotzdem – oder gerade deshalb – spielen Teilbarkeitsfragen in der Zahlentheorie eine große Rolle. Das technische Werkzeug zu ihrer Behandlung ist die Kongruenzenrechnung.

**Definition:** Wir sagen, zwei ganze Zahlen  $x, y \in \mathbb{Z}$  seien kongruent modulo  $m$  für eine natürliche Zahl  $m$ , in Zeichen

$$x \equiv y \pmod{m},$$

wenn  $x - y$  durch  $m$  teilbar ist.

Die Kongruenz modulo  $m$  definiert offensichtlich eine Äquivalenzrelation auf  $\mathbb{Z}$ : Jede ganze Zahl ist kongruent zu sich selbst, denn  $x - x = 0$  ist durch jede natürliche Zahl teilbar; wenn  $x - y$  durch  $m$  teilbar ist, so auch  $y - x = -(x - y)$ , und ist schließlich  $x \equiv y \pmod{m}$  und

$y \equiv z \pmod{m}$ , so sind  $x - y$  und  $y - z$  durch  $m$  teilbar, also auch ihre Summe  $x - z$ , und damit ist auch  $x \equiv z \pmod{m}$ .

Zwei Zahlen  $x, y \in \mathbb{Z}$  liegen genau dann in derselben Äquivalenzklasse, wenn sie bei der Division durch  $m$  denselben Divisionsrest haben; es gibt somit  $m$  Äquivalenzklassen, die den  $m$  möglichen Divisionsresten  $0, 1, \dots, m - 1$  entsprechen.

**Lemma:** Ist  $x \equiv x' \pmod{m}$  und  $y \equiv y' \pmod{m}$ , so ist auch

$$x \pm y \equiv x' \pm y' \pmod{m} \quad \text{und} \quad x'y' \equiv xy \pmod{m}.$$

*Beweis:* Sind  $x - x'$  und  $y - y'$  durch  $m$  teilbar, so auch

$$\begin{aligned} (x \pm y) - (x' \pm y') &= (x - x') \pm (y - y') && \text{und} \\ xy - x'y' &= x(y - y') + y'(x - x') \end{aligned} \quad \blacksquare$$

Im folgenden wollen wir das Symbol „mod“ nicht nur in Kongruenzen wie  $x \equiv y \pmod{m}$  benutzen, sondern auch – wie in vielen Programmiersprachen üblich – als Rechenoperation:

**Definition:** Für eine ganze Zahl  $x$  und eine natürliche Zahl  $m$  bezeichnet  $x \bmod m$  jene ganze Zahl  $0 \leq r < m$  mit  $x \equiv r \pmod{m}$ .

$x \bmod m$  ist also einfach der Divisionsrest bei der Division von  $x$  durch  $m$ .

Da nach dem gerade bewiesenen Lemma die Addition, Subtraktion und Multiplikation mit Kongruenzen vertauschbar sind, können wir auf der Menge aller Äquivalenzklassen Rechenoperationen einführen. Übersichtlicher wird das, wenn wir statt dessen die Menge

$$\mathbb{Z}/m \stackrel{\text{def}}{=} \{0, 1, \dots, m - 1\}$$

betrachten. Wir definieren eine Addition durch

$$x \oplus y = (x + y) \bmod m = \begin{cases} x + y & \text{falls } x + y < m \\ x + y - m & \text{sonst} \end{cases}$$

und entsprechend eine Multiplikation gemäß

$$x \odot y = (xy) \bmod m.$$

Für  $m = 4$  haben wir also folgende Operationen:

$\oplus$	0	1	2	3		$\odot$	0	1	2	3
0	0	1	2	3		0	0	0	0	0
1	1	2	3	0	und	1	0	1	2	3
2	2	3	0	1		2	0	2	0	2
3	3	0	1	2		3	0	3	2	1

Um diese Tabellen zu interpretieren, sollten wir uns an einige Grundbegriffe aus der Algebra erinnern:

**Definition:** a) Eine *Gruppe* ist eine Menge  $G$  zusammen mit einer Verknüpfung  $\times: G \times G \rightarrow G$ , für die gilt

- 1.)  $(x \times y) \times z = x \times (y \times z)$  für alle  $x, y, z \in G$ .
- 2.) Es gibt ein Element  $e \in G$ , so daß  $e \times x = x \times e = x$  für alle  $x \in G$ .
- 3.) Zu jedem  $x \in G$  gibt es ein  $x' \in G$ , so daß  $x \times x' = x' \times x = e$  ist.

Die Gruppe heißt *kommutativ* oder *abelsch*, wenn zusätzlich noch gilt

- 4.)  $x \times y = y \times x$  für alle  $x, y \in G$ .

b) Eine Abbildung  $\varphi: G \rightarrow H$  zwischen zwei Gruppen  $G$  und  $H$  mit Verknüpfungen  $\times$  und  $\otimes$  heißt (*Gruppen- $\varphi$* )*Homomorphismus*, falls für alle  $x, y \in G$  gilt:  $\varphi(x \times y) = \varphi(x) \otimes \varphi(y)$ . Ist  $\varphi$  zusätzlich bijektiv, reden wir von einem *Isomorphismus*. Die Gruppen  $G$  und  $H$  heißen *isomorph*, in Zeichen  $G \cong H$ , wenn es einen Isomorphismus  $\varphi: G \rightarrow H$  gibt.

c) Ein *Ring* ist eine Menge  $R$  zusammen mit zwei Verknüpfungen  $+, \cdot: R \times R \rightarrow R$ , so daß gilt

- 1.) Bezüglich  $+$  ist  $R$  eine abelsche Gruppe.
- 2.)  $(x \cdot y) \cdot z = x \cdot (y \cdot z)$  für alle  $x, y, z \in R$ .
- 3.) Es gibt ein Element  $1 \in R$ , so daß  $1 \cdot x = x \cdot 1 = x$  für alle  $x \in R$ .
- 4.)  $x \cdot (y + z) = x \cdot y + x \cdot z$  und  $(x + y) \cdot z = x \cdot z + y \cdot z$  für alle  $x, y, z \in R$ .

Der Ring heißt *kommutativ*, wenn zusätzlich noch gilt

- 5.)  $x \cdot y = y \cdot x$  für alle  $x, y \in R$ .

d) Eine Abbildung  $\varphi: R \rightarrow S$  zwischen zwei Ringen heißt (Ring-) *Homomorphismus*, wenn für alle  $x, y \in R$  gilt

$$\varphi(r + s) = \varphi(r) + \varphi(s) \quad \text{und} \quad \varphi(r \cdot s) = \varphi(r) \cdot \varphi(s),$$

wobei  $+$  und  $\cdot$  auf der linken Seite jeweils die Operationen von  $R$  bezeichnen und rechts die von  $S$ . Auch hier reden wir von einem *Isomorphismus*, wenn  $\varphi$  bijektiv ist, und bezeichnen  $R \cong S$  als *isomorph*, wenn es einen Isomorphismus  $\varphi: R \rightarrow S$  gibt.

**Lemma:** Für jedes  $m \in \mathbb{N}$  ist  $\mathbb{Z}/m$  mit den Operationen  $\oplus$  und  $\odot$  ein Ring.

*Beweis:* Wir betrachten die Abbildung

$$\varphi: \begin{cases} \mathbb{Z} \rightarrow \mathbb{Z}/m \\ x \mapsto x \bmod m \end{cases}.$$

Nach dem obigen Lemma ist

$$\varphi(x + y) = \varphi(x) \oplus \varphi(y) \quad \text{und} \quad \varphi(xy) = \varphi(x) \odot \varphi(y).$$

Da  $\mathbb{Z}$  bezüglich  $+$  eine abelsche Gruppe ist, gilt somit dasselbe für  $\mathbb{Z}/m$  bezüglich  $\oplus$ : Wenn zwei ganze Zahlen gleich sind, sind schließlich auch ihre Divisionsreste modulo  $m$  gleich. Das Neutralelement ist  $\varphi(0) = 0$ , und das additive Inverse ist  $\varphi(-x) = m - \varphi(x)$ . Auch die Eigenschaften von  $\odot$  folgen sofort aus den entsprechenden Eigenschaften der Multiplikation ganzer Zahlen; das Neutralelement ist  $\varphi(1) = 1$ . ■

Man beachte, daß  $\mathbb{Z}/m$  im allgemeinen kein Körper ist: In  $\mathbb{Z}/4$  beispielsweise ist  $2 \odot 2 = 0$ , und in einem Körper kann ein Produkt nur verschwinden, wenn mindestens einer der beiden Faktoren verschwindet.

Im folgenden werden wir die Rechenoperationen in  $\mathbb{Z}/m$  einfach mit  $+$  und  $\cdot$  bezeichnen, wobei jedesmal aus dem Zusammenhang klar sein sollte, ob wir von der Addition und Multiplikation in  $\mathbb{Z}/m$  oder der in  $\mathbb{Z}$  reden. Der Malpunkt wird dabei, wie üblich, oft weggelassen.

Restklassen werden oft benutzt für Prüfziffern. Bei der internationalen Standardbuchnummer ISBN etwa werden nur die ersten neun Ziffern zur

Identifikation des Buchs benötigt, die zehnte ist eine daraus berechnete Prüfziffer, Sind  $a_1, \dots, a_9$  die ersten zehn Ziffern, so ist

$$a_{10} = \sum_{i=1}^9 i a_i \pmod{11}.$$

Falls  $a_{10} = 10$  ist, wird es durch die römische Ziffer X dargestellt.

Die zehnstellige ISBN wird heute zunehmend durch eine dreizehnstellige ersetzt, die mit der *Global Trade Item Number* GTIN (früher *European Article Number* EAN) kompatibel ist. Hier wird die dreizehnte Ziffer  $a_{13}$  so aus den ersten zwölf Ziffern berechnet, daß

$$\sum_{i=1}^{13} g_i a_i \equiv 0 \pmod{10} \quad \text{ist mit} \quad g_i = \begin{cases} 1 & \text{falls } i \equiv 1 \pmod{2} \\ 3 & \text{falls } i \equiv 0 \pmod{2} \end{cases}.$$

Üblicherweise geben die ersten Ziffern der GTIN das Land an, in dem der Produzent der Ware sitzt; für Bücher sind die ersten drei Ziffern bislang stets 978, darauf folgen die ersten neun Ziffern der zehnstelligen ISBN und die nach GTIN-Regeln berechnete Prüfziffer. Da die Nummern langsam knapp werden, können auch dreizehnstellige ISBN-Nummern vergeben werden, die mit 979 beginnen; ihnen entspricht dann keine zehnstellige ISBN mehr. Die GTIN ist üblicherweise als Barcode auf Waren zu finden, die dadurch von Scannerkassen identifiziert werden.

Auch die *International Bank Account Number* IBAN arbeitet mit Restklassen als Prüfziffern: Die IBAN beginnt mit einem Länderkennzeichen aus zwei Buchstaben, darauf folgen zwei Prüfziffern und bis zu dreißig sonstige Zeichen, die das Konto identifizieren. Deutschland (Länderkennzeichen DE) verwendet nur achtzehn dieser Stellen; die ersten acht enthalten die Bankleitzahl, die restlichen zehn die gegebenenfalls links mit führenden Nullen aufgefüllte Kontonummer, so daß eine deutsche IBAN aus insgesamt 22 Zeichen besteht. In anderen Ländern gelten andere Regeln; beispielsweise hat eine norwegische IBAN nur 15 Stellen, eine maltesische dagegen 31 und eine aus St. Lucia sogar 32.

Zur Berechnung der Prüfziffern wird das Länderkennzeichen hinten an die IBAN angehängt, dahinter folgen an Stelle der noch unbekannt

Prüfziffern zwei Nullen. In dieser nun mit dem fünften Zeichen beginnenden Zeichenkette werden alle Buchstaben durch Zahlen ersetzt nach dem Schema  $A = 10$  bis  $Z = 35$ ; die entstehende Zahl wird modulo 97 reduziert. Die Prüfziffern sind 98 minus dieser Restklasse, geschrieben als zweiziffrige Zahl.

Die Semesterbeiträge der Universität Mannheim müssen beispielsweise einbezahlt werden auf das Konto 1 379 273 bei der LBBW, BLZ 600 501 01. Hier muß die Kontonummer durch drei führende Nullen erweitert werden, so daß man startet mit der Zeichenkette

$$6005\ 0101\ 0001\ 3792\ 73DE\ 00.$$

Einsetzen von  $D = 13$  und  $E = 14$  macht daraus

$$6005\ 0101\ 0001\ 3792\ 7313\ 1400 \equiv 75 \pmod{97}.$$

$98 - 75 = 23$ ; die IBAN ist also

$$DE23\ 6005\ 0101\ 0001\ 3792\ 73.$$

Für die Berechnung des Divisionsrest darf man die obige 24-stellige Zahl 6005 0101 0001 3792 7313 1400 natürlich nicht einfach so in den Taschenrechner eingeben; da die meisten Rechner intern mit höchstens dreizehn geltenden Ziffern rechnen, würde sie dabei mit Sicherheit verfälscht. Stattdessen muß man gemäß der in der Schule gelernten Methode dividieren, wobei man als Basis natürlich auch eine größere Zehnerpotenz nehmen kann. Beispielsweise ist

$$6005\ 0101 \equiv 20 \pmod{97}$$

$$\implies 6005\ 0101\ 0001\ 3792\ 7313\ 1400 \equiv 20\ 0001\ 3792\ 7313\ 1400 \pmod{97}$$

$$2000\ 0137 \equiv 95 \pmod{97}$$

$$\implies 20\ 0001\ 3792\ 7313\ 1400 \equiv 9592\ 7313\ 1400 \pmod{97}$$

$$9592\ 7313 \equiv 36 \pmod{97} \implies 9592\ 7313\ 1400 \equiv 36\ 1400 \equiv 75 \pmod{97},$$

so daß man auch durch Rechnen mit maximal achtstelligen Zahlen in wenigen Schritten das Ergebnis erhält.

Eine weitere Anwendung von Kongruenzen sind sogenannte Modulararithmetiken: Beim konkreten Rechnen mit großen Zahlen verwendet man in der Numerik meist Gleitkommaarithmetiken, hat dann aber



das Problem von Rundungsfehlern und eventuell auch numerischer Instabilität. Beim Rechnen modulo einer natürlichen Zahl  $N$  hat man stets exakte Ergebnisse, allerdings ist das Ergebnis nur modulo  $N$  bestimmt. Falls man weiß, daß es beispielsweise zwischen  $-(N - 1)/2$  und  $(N - 1)/2$  liegt, reicht das. Für große  $N$  ist allerdings auch das Rechnen modulo  $N$  schon recht aufwendig; tatsächlich verwenden daher beispielsweise Computeragebrasysteme bei auf Modulararithmetik beruhenden Algorithmen mehrere paarweise teilerfremde Moduln  $N_i$  und rekonstruieren dann das Endergebnis modulo dem Produkt aller  $N_i$ . Wie dies geschieht, zeigt der nächste Paragraph.

## §7: Der chinesische Restesatz

Der Legende nach zählten chinesische Generäle ihre Truppen, indem sie diese mehrfach antreten ließen in Reihen verschiedener Breiten  $m_1, \dots, m_r$  und jedesmal nur die Anzahl  $a_r$  der Soldaten in der letzten Reihe zählten. Aus den  $r$  Relationen

$$x \equiv a_1 \pmod{m_1}, \quad \dots, \quad x \equiv a_r \pmod{m_r}$$

bestimmten sie dann die Gesamtzahl  $x$  der Soldaten.

Ob es im alten China wirklich Generäle gab, die soviel Mathematik konnten, sei dahingestellt; Beispiele zu diesem Satz finden sich jedenfalls bereits 1247 in den chinesischen *Mathematischen Abhandlungen in neun Bänden* von CH'IN CHIU-SHAO (1202–1261), allerdings geht es dort nicht um Soldaten, sondern um Reis.

CH'IN CHIU-SHAO oder QIN JIUSHAO wurde 1202 in der Provinz Sichuan geboren. Auf eine wilde Jugend mit vielen Affären folgte ein wildes und alles andere als gesetztreues Berufsleben in Armee, öffentlicher Verwaltung und illegalem Salzhandel. Als Jugendlicher studierte er an der Akademie von Hang-chou Astronomie, später brachte ihm ein unbekannter Lehrer Mathematik bei. Insbesondere studierte er die in vorchristlicher Zeit entstandenen *Neun Bücher der Rechenkunst*, nach deren Vorbild er 1247 seine deutlich anspruchsvolleren *Mathematischen Abhandlungen in neun Bänden* publizierte, die ihn als einen der bedeutendsten Mathematiker nicht nur Chinas ausweisen. Zum chinesischen Restesatz schreibt er, daß er ihn von den Kalendermachern gelernt habe, diese ihn jedoch nur rein mechanisch anwendeten ohne ihn zu verstehen. CH'IN CHIU-SHAO starb 1261 in Meixian, wohin er nach einer seiner vielen Entlassungen wegen krimineller Machenschaften geschickt worden war.

Wir wollen uns zunächst überlegen, unter welchen Bedingungen ein solches Verfahren überhaupt funktionieren kann. Offensichtlich können die obigen  $r$  Relationen eine natürliche Zahl nicht eindeutig festlegen, denn ist  $x$  eine Lösung und  $M$  irgendein gemeinsames Vielfaches der sämtlichen  $m_i$ , so ist  $x + M$  auch eine –  $M$  ist schließlich modulo aller  $m_i$  kongruent zur Null.

Außerdem gibt es Relationen obiger Form, die unlösbar sind, beispielsweise das System

$$x \equiv 2 \pmod{4} \quad \text{und} \quad x \equiv 3 \pmod{6},$$

dessen erste Gleichung nur gerade Lösungen hat, während die zweite nur ungerade hat. Das Problem hier besteht darin, daß zwei ein gemeinsamer Teiler von vier und sechs ist, so daß jede der beiden Kongruenzen auch etwas über  $x \pmod{2}$  aussagt: Nach der ersten ist  $x$  gerade, nach der zweiten aber ungerade.

Dieses Problem können wir dadurch umgehen, daß wir nur Moduln  $m_i$  zulassen, die paarweise teilerfremd sind. Dies hat auch den Vorteil, daß jedes gemeinsame Vielfache der  $m_i$  Vielfaches des Produkts aller  $m_i$  sein muß, so daß wir  $x$  modulo einer vergleichsweise großen Zahl kennen.

**Chinesischer Restesatz:** Das System von Kongruenzen

$$x \equiv a_1 \pmod{m_1}, \quad \dots, \quad x \equiv a_r \pmod{m_r}$$

hat für paarweise teilerfremde Moduln  $m_i$  genau eine Lösung  $x$  mit  $0 \leq x < m_1 \cdots m_r$ . Jede andere Lösung  $y \in \mathbb{Z}$  läßt sich in der Form  $x + km_1 \cdots m_r$  schreiben mit  $k \in \mathbb{Z}$ .

Mit den Begriffen aus dem vorigen Paragraphen läßt sich dies auch anders formulieren: Die Zahl  $x \pmod{m_i}$  können wir auffassen als Element von  $\mathbb{Z}/m_i$ , das  $r$ -Tupel aus allen diesen Zahlen also als Element von  $\mathbb{Z}/m_1 \times \cdots \times \mathbb{Z}/m_r$ . Man überlegt sich leicht, daß das kartesische Produkt von zwei oder mehr Gruppen wieder eine Gruppe ist: Die Verknüpfung wird einfach komponentenweise definiert, und das Neutralelement ist dasjenige Tupel, dessen sämtliche Komponenten Neutralelemente der jeweiligen Faktoren sind. Genauso folgt, daß das kartesische Produkt von zwei oder mehr Ringen wieder ein Ring ist.

In algebraischer Formulierung haben wir dann die folgende Verschärfung des obigen Satzes:

**Chinesischer Restesatz (Algebraische Form):** Die Abbildung

$$\varphi: \begin{cases} \mathbb{Z}/m_1 \cdots m_r \rightarrow \mathbb{Z}/m_1 \times \cdots \times \mathbb{Z}/m_r \\ x \mapsto (x \bmod m_1, \dots, x \bmod m_r) \end{cases}$$

ist ein Isomorphismus von Ringen.

Wir *beweisen* den Satz in dieser algebraischen Form:

Zunächst ist

$$\psi: \begin{cases} \mathbb{Z} \rightarrow \mathbb{Z}/m_1 \times \cdots \times \mathbb{Z}/m_r \\ x \mapsto (x \bmod m_1, \dots, x \bmod m_r) \end{cases}$$

ein Ringhomomorphismus, denn nach dem Lemma aus dem vorigen Paragraphen ist der Übergang zu den Restklassen modulo jeder der Zahlen  $m_i$  mit Addition und Multiplikation vertauschbar. Da  $\psi(x)$  nur von  $x \bmod m_1 \cdots m_r$  abhängt, folgt daraus, daß auch  $\varphi$  ein Ringhomomorphismus ist.

$\varphi$  ist injektiv, denn ist  $\varphi(x) = \varphi(y)$ , so ist  $x \bmod m_i = y \bmod m_i$  für alle  $i$ ; da die  $m_i$  paarweise teilerfremd sind, ist  $x - y$  somit durch das Produkt der  $m_i$  teilbar, was für  $x, y \in \mathbb{Z}/m_1 \cdots m_r$  nur im Fall  $x = y$  möglich ist.

Nun ist  $\varphi$  aber eine Abbildung zwischen endlichen Mengen, die beide aus je  $m_1 \cdots m_r$  Elementen bestehen. Jede injektive Abbildung zwischen zwei gleichmächtigen endlichen Mengen ist zwangsläufig auch surjektiv, also bijektiv, und somit ist  $\varphi$  ein Isomorphismus. ■

Aus Sicht der chinesischen Generäle ist dieser Beweis enttäuschend: Angenommen, ein General weiß, daß höchstens Tausend Soldaten vor ihm stehen. Er läßt sie in Zehnerreihen antreten; in der letzten Reihe stehen fünf Soldaten. Bei Elferreihen sind es neun, bei Dreizehnerreihen sechs. Da  $10 \cdot 11 \cdot 13 = 1430$  größer ist als Tausend, weiß er, daß dies die Anzahl eindeutig festlegt. Er weiß aber nicht, wie viele Soldaten nun tatsächlich vor ihm stehen. Als General hat er natürlich die

Möglichkeit, einige Soldaten abzukommandieren, die für jede Zahl bis Tausend die Divisionsreste modulo 9, 10 und 13 berechnen müssen, bis sie auf das Tripel (5, 9, 6) stoßen. Wir als Mathematiker sollten jedoch eine effizientere Methode finden.

Dazu verhilft uns der erweiterte EUKLIDISCHE Algorithmus:

Wir beginnen mit dem Fall zweier Kongruenzen

$$x \equiv a \pmod{m} \quad \text{und} \quad x \equiv b \pmod{n}$$

mit zueinander teilerfremden Zahlen  $m$  und  $n$ . Ihr ggT eins läßt sich nach dem erweiterten EUKLIDISCHEN Algorithmus als  $1 = \alpha m + \beta n$  schreiben. Somit ist

$$1 - \alpha m = \beta n \equiv \begin{cases} 1 & \pmod{m} \\ 0 & \pmod{n} \end{cases} \quad \text{und} \quad 1 - \beta n = \alpha m \equiv \begin{cases} 0 & \pmod{m} \\ 1 & \pmod{n} \end{cases},$$

also löst

$$x = \beta n a + \alpha m b \equiv \begin{cases} a & \pmod{m} \\ b & \pmod{n} \end{cases}$$

das Problem.

$x$  ist natürlich nicht die einzige Lösung; wenn wir ein gemeinsames Vielfaches von  $m$  und  $n$  addieren, ändert sich nichts an den Kongruenzen. Da wir von teilerfremden Zahlen ausgegangen sind, ist das Produkt das kleinste gemeinsame Vielfache; die allgemeine Lösung ist daher

$$x + (\beta n + \lambda b)a + (\alpha m - \lambda a)b.$$

Insbesondere ist die Lösung eindeutig modulo  $nm$ .

Als Beispiel betrachten wir die beiden Kongruenzen

$$x \equiv 1 \pmod{17} \quad \text{und} \quad x \equiv 5 \pmod{19}.$$

Wir müssen als erstes den erweiterten EUKLIDISCHEN Algorithmus auf die beiden Moduln 17 und 19 anwenden:

$$19 : 17 = 1 \text{ Rest } 2 \implies 2 = 19 - 17$$

$$17 : 2 = 8 \text{ Rest } 1 \implies 1 = 17 - 8 \cdot 2 = 9 \cdot 17 - 8 \cdot 19$$

Also ist  $9 \cdot 17 = 153 \equiv 0 \pmod{17}$  und  $\equiv 1 \pmod{19}$ ; außerdem ist  $-8 \cdot 19 = -152$  durch 19 teilbar und  $\equiv 1 \pmod{17}$ . Die Zahl

$$x = -152 \cdot 1 + 153 \cdot 5 = 613$$

löst somit das Problem. Da  $613$  größer ist als  $17 \cdot 19 = 323$ , ist allerdings nicht  $613$  die kleinste positive Lösung, sondern  $613 - 323 = 290$ .

Bei mehr als zwei Kongruenzen gehen wir rekursiv vor: Wir lösen die ersten beiden Kongruenzen  $x \equiv a_1 \pmod{m_1}$  und  $x \equiv a_2 \pmod{m_2}$  wie gerade besprochen; das Ergebnis ist eindeutig modulo  $m_1 m_2$ . Ist  $c_2$  eine feste Lösung, so läßt sich die Lösung schreiben als Kongruenz

$$x \equiv c_2 \pmod{m_1 m_2},$$

und da die  $m_i$  paarweise teilerfremd sind, ist auch  $m_1 m_2$  teilerfremd zu  $m_3$ . Mit EUKLID können wir daher das System

$$x \equiv c_2 \pmod{m_1 m_2} \quad \text{und} \quad x \equiv a_3 \pmod{m_3}$$

lösen und die Lösung schreiben als

$$x \equiv c_3 \pmod{m_1 m_2 m_3}$$

und so weiter, bis wir schließlich  $x$  modulo dem Produkt aller  $m_i$  kennen und somit das Problem gelöst haben.

Im Beispiel des oben angesprochenen Systems

$$x \equiv 5 \pmod{10}, \quad x \equiv 9 \pmod{11}, \quad x \equiv 6 \pmod{13}$$

lösen wir also zunächst nur das System

$$x \equiv 5 \pmod{10} \quad \text{und} \quad x \equiv 9 \pmod{11}.$$

Da  $1 = 11 - 10$ , ist  $11 \equiv 0 \pmod{11}$  und  $11 \equiv 1 \pmod{10}$ ; entsprechend ist  $-10 \equiv 0 \pmod{10}$  und  $-10 \equiv 1 \pmod{11}$ . Also ist

$$x = 5 \cdot 11 - 9 \cdot 10 = -35$$

eine Lösung; die allgemeine Lösung ist  $-35 + 110k$  mit  $k \in \mathbb{Z}$ . Die kleinste positive Lösung ist  $-35 + 110 = 75$ .

Unser Ausgangssystem ist somit äquivalent zu den beiden Kongruenzen

$$x \equiv 75 \pmod{110} \quad \text{und} \quad x \equiv 6 \pmod{13}.$$

Um es zu lösen, müssen wir zunächst die Eins als Linearkombination von  $110$  und  $13$  darstellen. Hier bietet sich keine offensichtliche Lösung an, also verwenden wir den erweiterten EUKLIDischen Algorithmus:

$$110 : 13 = 8 \text{ Rest } 6 \implies 6 = 110 - 8 \cdot 13$$

$$13 : 6 = 2 \text{ Rest } 1 \implies 1 = 13 - 2 \cdot 6 = 17 \cdot 13 - 2 \cdot 110$$

Also ist  $17 \cdot 13 = 221 \equiv 1 \pmod{110}$  und  $\equiv 0 \pmod{13}$ ; genauso ist  $-2 \cdot 110 = 220 \equiv 1 \pmod{13}$  und  $\equiv 9 \pmod{110}$ . Eine ganzzahlige Lösung unseres Problems ist somit

$$75 \cdot 221 - 6 \cdot 220 = 15\,255.$$

Die allgemeine Lösung ist

$$15\,255 + k \cdot 110 \cdot 13 = 15\,255 + 1\,430k \quad \text{mit } k \in \mathbb{Z}.$$

Da  $15\,255 : 1\,430 = 10$  Rest  $955$  ist, hatte der General also  $955$  Soldaten vor sich stehen.

Alternativ läßt sich die Lösung eines Systems aus  $r$  Kongruenzen auch in einer geschlossenen Form darstellen allerdings um den Preis einer  $n$ -maligen statt  $(n - 1)$ -maligen Anwendung des EUKLIDischen Algorithmus und größeren Zahlen schon von Beginn an: Um das System

$$x \equiv a_i \pmod{m_i} \quad \text{für } i = 1, \dots, r$$

zu lösen, berechnen wir zunächst für jedes  $i$  das Produkt

$$\widehat{m}_i = \prod_{j \neq i} m_j$$

der sämtlichen übrigen  $m_j$  und bestimmen dazu ganze Zahlen  $\alpha_i, \beta_i$ , für die gilt  $\alpha_i m_i + \beta_i \widehat{m}_i = 1$  Dann ist

$$x = \sum_{j=1}^n \beta_j \widehat{m}_j a_j \equiv \beta_i \widehat{m}_i a_i = (1 - \alpha_i m_i) a_i \equiv a_i \pmod{m_i}.$$

Natürlich wird  $x$  hier – wie auch bei der obigen Formel – oft größer sein als das Produkt der  $m_i$ ; um die kleinste Lösung zu finden, müssen wir also noch modulo diesem Produkt reduzieren.

Im obigen Beispiel wäre

$$\begin{aligned} m_1 = 10 & \quad \widehat{m}_1 = 11 \cdot 13 = 143 & \quad 1 = 43 \cdot 10 - 3 \cdot 143 \\ m_2 = 11 & \quad \widehat{m}_2 = 10 \cdot 13 = 130 & \quad 1 = -59 \cdot 11 + 5 \cdot 130 \\ m_3 = 13 & \quad \widehat{m}_3 = 10 \cdot 11 = 110 & \quad 1 = 17 \cdot 13 - 2 \cdot 110, \end{aligned}$$

also

$$x = -3 \cdot 143 \cdot 5 + 5 \cdot 130 \cdot 9 - 2 \cdot 110 \cdot 6 = -2145 + 5850 - 1320 = 2385.$$

Modulo  $10 \cdot 11 \cdot 13$  erhalten wir natürlich auch hier wieder 955.

Damit kennen wir nun auch zwei konstruktive Beweise des chinesischen Restesatzes und wissen, wie man Systeme von Kongruenzen mit Hilfe des erweiterten EUKLIDischen Algorithmus lösen kann.

## §8: Prime Restklassen

Wie wir gesehen haben, können wir auch in  $\mathbb{Z}/m$  im allgemeinen nicht dividieren. Allerdings ist Division doch sehr viel häufiger möglich als in den ganzen Zahlen. Dies wollen wir als nächstes genauer untersuchen:

**Lemma:** Zu zwei gegebenen natürlichen Zahlen  $a, m$  gibt es genau dann ein  $x \in \mathbb{N}$  mit  $ax \equiv 1 \pmod{m}$ , wenn  $\text{ggT}(a, m) = 1$  ist.

*Beweis:* Wenn es ein solches  $x$  gibt, gibt es dazu ein  $y \in \mathbb{N}$ , so daß  $ax = 1 + my$ , d.h.  $1 = ax - my$ . Damit muß jeder gemeinsame Teiler von  $a$  und  $m$  Teiler der Eins sein,  $a$  und  $m$  sind also teilerfremd.

Sind umgekehrt  $a$  und  $m$  teilerfremd, so gibt es nach dem erweiterten EUKLIDischen Algorithmus  $x, y \in \mathbb{Z}$  mit  $ax + my = 1$ . Durch (gegebenenfalls mehrfache) Addition der Gleichung  $am - ma = 0$  läßt sich nötigenfalls erreichen, daß  $x$  positiv wird, und offensichtlich ist  $ax \equiv 1 \pmod{m}$ . ■

**Definition:** Ein Element  $a \in \mathbb{Z}/m$  heißt prime Restklasse, wenn  $\text{ggT}(a, m) = 1$  ist.

Nach dem gerade bewiesenen Lemma gibt es somit zu jeder primen Restklasse  $a$  ein  $x \in \mathbb{Z}/m$ , so daß dort  $ax = 1$  ist. Damit ist das folgende Lemma nicht verwunderlich:

**Lemma:** Die primen Restklassen aus  $\mathbb{Z}/m$  bilden bezüglich der Multiplikation eine Gruppe.

*Beweis:* Wir müssen uns zunächst überlegen, daß das Produkt zweier primen Restklassen wieder eine prime Restklasse ist. Sind  $a, b \in \mathbb{Z}/m$  beide teilerfremd zu  $m$ , so auch  $ab$ , denn wäre  $p$  ein gemeinsamer

Primteiler von  $ab$  und  $m$ , so wäre  $p$  als Primzahl auch Teiler von  $a$  oder  $b$ , also gemeinsamer Teiler von  $a$  und  $m$  oder von  $b$  und  $m$ . Die Eins ist natürlich eine prime Restklasse, und auch die Existenz von Inversen ist kein Problem: Nach dem vorigen Lemma gibt es ein  $x \in \mathbb{Z}$ , so daß  $ax \equiv 1 \pmod{m}$  ist, und die andere Richtung dieses Lemmas zeigt, daß auch  $x \pmod{m}$  eine prime Restklasse ist. Das Assoziativgesetz der Multiplikation gilt für alle Elemente von  $\mathbb{Z}/m$ , erst recht also für die primen Restklassen. ■

**Definition:** Die Gruppe  $(\mathbb{Z}/m)^\times$  der primen Restklassen heißt *prime Restklassengruppe*, ihre Ordnung wird mit  $\varphi(m)$  bezeichnet.  $\varphi : \mathbb{N} \rightarrow \mathbb{N}$  heißt EULERSche  $\varphi$ -Funktion.



LEONHARD EULER (1707–1783) wurde in Basel geboren. Er ging auch dort zur Schule und, im Alter von 14 Jahren, zur Universität. Zwei Jahre später legte er die Magisterprüfung in Philosophie ab und begann mit dem Studium der Theologie; daneben hatte er sich seit Beginn seines Studium unter Anleitung von JOHANN BERNOULLI mit Mathematik beschäftigt. 1726 beendete er sein Studium in Basel und bekam eine Stelle an der Petersburger Akademie der Wissenschaften, die er 1727 antrat. Auf Einladung FRIEDRICHS DES GROSSEN wechselte er 1741 an die preußische Akademie der Wissenschaften. Nachdem sich das Verhältnis zwischen den beiden dramatisch verschlechtert hatte, kehrte er 1766 nach St. Petersburg zurück. Im gleichen Jahr erblindete er vollständig; trotzdem schrieb er rund die Hälfte seiner zahlreichen Arbeiten (Seine gesammelten Abhandlungen umfassen 73 Bände) danach. Sie enthalten bedeutende Beiträge zu zahlreichen Teilgebieten der Mathematik, Physik, Astronomie und Kartographie.

beiden dramatisch verschlechtert hatte, kehrte er 1766 nach St. Petersburg zurück. Im gleichen Jahr erblindete er vollständig; trotzdem schrieb er rund die Hälfte seiner zahlreichen Arbeiten (Seine gesammelten Abhandlungen umfassen 73 Bände) danach. Sie enthalten bedeutende Beiträge zu zahlreichen Teilgebieten der Mathematik, Physik, Astronomie und Kartographie.

**Lemma:** a) Für zwei zueinander teilerfremde Zahlen  $n, m \in \mathbb{N}$  ist  $\varphi(nm) = \varphi(n)\varphi(m)$ .

b) Für  $m = \prod_{i=1}^r p_i^{e_i}$  ist  $\varphi(m) = \prod_{i=1}^r (p_i^{e_i-1}(p_i - 1))$ .

*Beweis:* a) Eine Zahl  $a$  ist genau dann teilerfremd zum Produkt  $nm$ , wenn  $a \pmod{n}$  teilerfremd zu  $n$  und  $a \pmod{m}$  teilerfremd zu  $m$  ist. Da nach dem chinesischen Restesatz  $\mathbb{Z}/nm \cong \mathbb{Z}/n \times \mathbb{Z}/m$  ist, ist daher auch  $(\mathbb{Z}/nm)^\times \cong (\mathbb{Z}/n)^\times \times (\mathbb{Z}/m)^\times$ .



b) Wegen a) genügt es, dies für Primzahlpotenzen  $p^e$  zu beweisen. Eine Zahl  $a$  ist genau dann teilerfremd zu  $p^e$ , wenn sie kein Vielfaches von  $p$  ist. Unter den Zahlen von 1 bis  $p^e$  gibt es genau  $p^{e-1}$  Vielfache von  $p$ , also ist  $\varphi(p^e) = p^e - p^{e-1} = p^{e-1}(p - 1)$ . ■

**Korollar:**  $\mathbb{Z}/m$  ist genau dann ein Körper, wenn  $m$  eine Primzahl ist.

*Beweis:* Das einzige, was  $\mathbb{Z}/m$  zu einem Körper eventuell fehlt, ist die Existenz von multiplikativen Inversen für alle von null verschiedenen Elemente. Dies ist offenbar äquivalent zur Formel  $\varphi(m) = m - 1$ , und die gilt nach dem Lemma genau dann, wenn  $m$  prim ist. ■

Der Körper  $\mathbb{Z}/p$  mit  $p$  Elementen wird üblicherweise mit  $\mathbb{F}_p$  bezeichnet; die zugehörige prime Restklassengruppe  $(\mathbb{Z}/p)^\times = \mathbb{F}_p \setminus \{0\}$  entsprechend als  $\mathbb{F}_p^\times$ . Dabei steht das „ $\mathbb{F}$ “ für *finit*. Im Englischen werden endliche Körper gelegentlich auch als *Galois fields* bezeichnet, so daß man hier auch die Abkürzung  $\text{GF}(p)$  sieht. *Field* ist das englische Wort für Körper; das gelegentlich in Informatikbüchern zu lesende Wort *Galoisfeld* ist also ein Übersetzungsfehler.

Wir wollen uns als nächstes überlegen, daß die multiplikative Gruppe dieses Körpers aus den Potenzen eines einzigen Elements besteht. Dazu brauchen wir zunächst noch ein Lemma aus der Gruppentheorie:

**Definition:** Die Ordnung eines Elements  $a$  einer (multiplikativ geschriebenen) Gruppe  $G$  ist die kleinste natürliche Zahl  $r$ , für die  $a^r$  gleich dem Einselement ist. Falls es keine solche Zahl gibt, sagen wir,  $a$  habe unendliche Ordnung.

**Lemma (LAGRANGE):** In einer endlichen Gruppe teilt die Ordnung eines jeden Elements die Gruppenordnung.

*Beweis:* Die Potenzen des Elements  $a$  bilden zusammen mit der Eins eine Untergruppe  $H$  von  $G$ , deren Elementanzahl gerade die Ordnung  $r$  von  $H$  ist. Wir führen auf  $G$  eine Äquivalenzrelation ein durch die Vorschrift  $g \sim h$ , falls  $gh^{-1}$  in  $H$  liegt. Offensichtlich besteht die

Äquivalenzklasse eines jeden Elements  $g \in G$  aus genau  $r$  Elementen, nämlich  $g, ga, \dots, ga^{r-1}$ . Da  $G$  die Vereinigung aller Äquivalenzklassen ist, muß die Gruppenordnung somit ein Vielfaches von  $r$  sein. ■



JOSEPH-LOUIS LAGRANGE (1736–1813) wurde als GIUSEPPE LODOVICO LAGRANGIA in Turin geboren und studierte dort zunächst Latein. Erst eine alte Arbeit von HALLEY über algebraische Methoden in der Optik weckte sein Interesse an der Mathematik, woraus ein ausgedehnter Briefwechsel mit EULER entstand. In einem Brief vom 12. August 1755 berichtete er diesem unter anderem über seine Methode zur Berechnung von Maxima und Minima; 1756 wurde er, auf EULERS Vorschlag, Mitglied der Berliner Akademie; zehn Jahre später zog er nach Berlin und wurde dort EULERS Nachfolger als mathematischer Direktor der

Akademie. 1787 wechselte er an die Pariser Académie des Sciences, wo er bis zu seinem Tod blieb und unter anderem an der Einführung des metrischen Systems beteiligt war. Seine Arbeiten umspannen weite Teile der Analysis, Algebra und Geometrie.

**Korollar:** Für zwei zueinander teilerfremde Zahlen  $a, m$  ist

$$a^{\varphi(m)} \equiv 1 \pmod{m}.$$

*Beweis:* Klar, denn  $\varphi(m)$  ist die Ordnung der primen Restklassengruppe modulo  $m$ . ■

Für eine Primzahl  $N = p$  bezeichnet man diese Aussage auch als den *kleinen Satz von FERMAT*:

**Satz (FERMAT):** Für jede nicht durch die Primzahl  $p$  teilbare ganze Zahl  $a$  ist  $a^{p-1} \equiv 1 \pmod{p}$ . Für alle  $a \in \mathbb{Z}$  ist  $a^p \equiv a \pmod{p}$ .

*Beweis:* Die erste Aussage ist klar, da  $\varphi(p) = p - 1$  ist. Für die zweite müssen wir nur noch beachten, daß für durch  $p$  teilbare Zahlen  $a$  sowohl  $a^p$  als auch  $a$  kongruent null modulo  $p$  sind. ■



Der französische Mathematiker PIERRE DE FERMAT (1601–1665) wurde in Beaumont-de-Lomagne im Département Tarn et Garonne geboren. Bekannt ist er heutzutage vor allem für seine 1637 von ANDREW WILES bewiesene Vermutung, wonach die Gleichung  $x^n + y^n = z^n$  für  $n \geq 3$  keine ganzzahlige Lösung mit  $xyz \neq 0$  hat. Dieser „große“ Satz von FERMAT, von dem FERMAT lediglich in einer Randnotiz behauptete, daß er ihn beweisen könne, erklärt den Namen der obigen Aussage. Obwohl FERMAT sich sein Leben lang sehr mit Mathematik beschäftigte und wesentliche Beiträge zur Zahlentheorie, Wahrscheinlichkeitstheorie und Analysis lieferte, war er hauptberuflich Jurist.

**Satz:** Die multiplikative Gruppe eines endlichen Körpers ist zyklisch.

*Beweis:* Da die multiplikative Gruppe eines Körpers mit  $q$  Elementen aus allen Körperelementen außer der Null besteht, hat sie die Ordnung  $q - 1$ , d.h. nach LAGRANGE ist die Ordnung eines jeden Elements ein Teiler von  $q - 1$ . Wir müssen zeigen, daß es mindestens ein Element gibt, dessen Ordnung *genau*  $q - 1$  ist.

Für jeden Primteiler  $p_i$  von  $q - 1$  hat die Polynomgleichung

$$x^{(q-1)/p_i} = 1$$

höchstens  $(q - 1)/p_i$  Lösungen im Körper; es gibt also zu jedem  $p_i$  ein Körperelement  $a_i$  mit  $a_i^{(q-1)/p_i} \neq 1$ .

$q_i$  sei die größte Potenz von  $p_i$ , die  $q - 1$  teilt, und  $g_i = a_i^{(q-1)/q_i}$  die  $(q - 1)/q_i$ -te Potenz von  $a_i$ . Dann ist

$$g_i^{q_i} = a_i^{q-1} = 1 \quad \text{und} \quad g_i^{\frac{q_i}{p_i}} = a_i^{\frac{q-1}{p_i}} \neq 1;$$

$g_i$  hat also die Ordnung  $q_i$ . Da die verschiedenen  $q_i$  Potenzen verschiedener Primzahlen  $p_i$  sind, hat daher das Produkt  $g$  aller  $g_i$  das Produkt aller  $q_i$  als Ordnung, also  $q - 1$ . Damit ist die multiplikative Gruppe des Körpers zyklisch. ■

**Definition:** Ein Element  $g$  eines endlichen Körpers  $k$  heißt *primitive Wurzel*, wenn es die zyklische Gruppe  $k^\times$  erzeugt.

Selbst im Fall der Körper  $\mathbb{F}_p$  gibt es keine Formel, mit der man eine solche primitive Wurzel explizit in Abhängigkeit von  $p$  angeben kann. Üblicherweise wählt man zufällig ein Element aus und testet, ob es die Ordnung  $p - 1$  hat. Die Wahrscheinlichkeit dafür ist offenbar  $\varphi(p - 1) : (p - 1)$ , was für die meisten Werte von  $p$  recht gut ist. Der Test, ob die Ordnung gleich  $p - 1$  ist, läßt sich allerdings nur dann effizient durchführen, wenn die Primteiler  $p_i$  von  $p - 1$  bekannt sind, denn dann kann man einfach testen, ob alle Potenzen mit den Exponenten  $(p - 1)/p_i$  von eins verschieden sind. Für große Werte von  $p$ , wie sie in der Kryptographie benötigt werden, kann dies ein Problem sein, so daß man hier im allgemeinen von faktorisierten Zahlen  $r$  ausgeht und dann testet, ob  $r + 1$  prim ist. Im Kapitel über Primzahlen werden wir geeignete Tests kennenlernen.

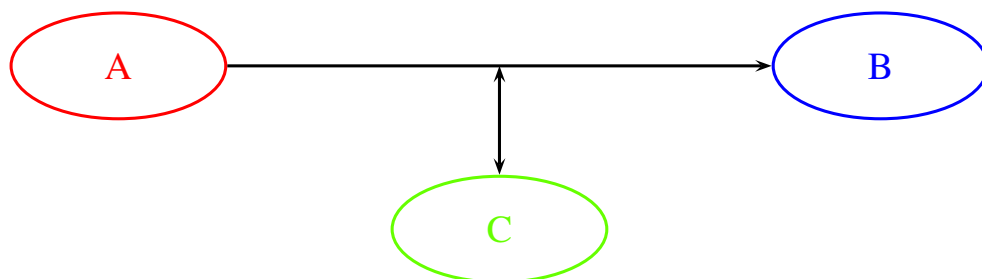
## Kapitel 2

# Anwendungen in der Kryptologie

### §1: New directions in cryptography

Kryptologie ist zusammengesetzt aus den beiden griechischen Wörtern κρυπτός = verborgen, versteckt und λόγος = Rede, Darlegung, Vernunft; sie ist also die Wissenschaft vom Geheimen. Sie besteht aus der Kryptographie (von γράφή = Das Schreiben), die Geheimschriften entwickelt, und der Kryptanalyse (von ἀναλύειν = auflösen, zerlegen), die versucht, letztere zu analysieren mit dem Ziel, sie zu knacken.

Die Grundsituation ist also die folgende:



A möchte eine Nachricht  $x$  an B übermitteln, jedoch besteht die Gefahr, daß alles, was er an B schickt, auf dem Weg dorthin von C gelesen und vielleicht auch verändert wird; außerdem könnte C eventuell versuchen, sich gegenüber B als A ausgeben oder umgekehrt.

Die Kryptographie versucht, dies zu verhindern, indem A anstelle von  $x$  einen Chiffretext  $c$  schickt, aus der zwar B, nicht aber C die Nachricht  $x$  und gegebenenfalls weitere Informationen rekonstruieren kann. Natürlich bietet diese Verschlüsselung für sich allein

noch keinen Schutz, denn C könnte zum Beispiel auch den Computer von A oder B angreifen, um so den Text vor der Verschlüsselung oder nach der Entschlüsselung abzugreifen. Auch könnte er das Verschlüsselungsprogramm so manipulieren, daß die Verschlüsselung für ihn keine Hürde mehr ist, er könnte elektromagnetische Strahlung von Computer und/oder Monitor auffangen und auswerten, und so weiter. Die Enthüllungen über NSA und GCHQ zeigen, daß es praktisch nichts gibt, was ein hinreichend entschlossener Gegner nicht anwendet. Trotzdem macht es gute Kryptographie zumindest für manche Gegner schwierig oder sogar unmöglich, die Nachricht zu lesen.

In der klassischen Kryptographie verläuft die Entschlüsselung entweder genauso oder zumindest sehr ähnlich wie die Verschlüsselung; insbesondere kann jeder, der eine Nachricht verschlüsseln kann, jede andere entsprechend verschlüsselte Nachricht auch entschlüsseln. Man bezeichnet diese Verfahren daher als *symmetrisch*.

Der Nachteil eines solchen Verfahrens besteht darin, daß in einem Netzwerk jeder Teilnehmer mit jedem anderen einen Schlüssel vereinbaren muß. In militärischen Netzen war dies üblicherweise so geregelt, daß das gesamte Netz denselben Schlüssel benutzte, der in einem Codebuch für jeden Tag im voraus festgelegt war. In kommerziellen Netzen wie beispielsweise einem Mobilfunknetz ist so etwas natürlich unmöglich.

1976 publizierten MARTIN HELLMAN, damals Assistenzprofessor in Stanford, und sein Forschungsassistent WHITFIELD DIFFIE eine Arbeit mit dem Titel *New directions in cryptography* (IEEE Trans. Inform. Theory **22**, 644–654), in der sie vorschlugen, den Vorgang der Verschlüsselung und den der Entschlüsselung völlig voneinander zu trennen: Es sei schließlich nicht notwendig, daß der Sender einer verschlüsselten Nachricht auch in der Lage sei, diese zu *entschlüsseln*.

Der Vorteil eines solchen *asymmetrischen* Verfahrens wäre, daß jeder potentielle Empfänger nur einen einzigen Schlüssel bräuchte und dennoch sicher sein könnte, daß nur er selbst seine Post entschlüsseln kann. Der Schlüssel für die Verschlüsselung müßte nicht einmal geheimgehalten werden, da es ja (meistens) nichts schadet, wenn jedermann Nachrichten verschlüsseln kann. In einem Netzwerk mit  $n$  Teilnehmern bräuchte

man also nur  $n$  Schlüssel, um es jedem Teilnehmer zu ermöglichen, mit jedem anderen sicher zu kommunizieren. Die Schlüssel könnten sogar in einem öffentlichen Verzeichnis stehen. Bei einem symmetrischen Kryptosystem wäre der gleiche Zweck nur erreichbar mit  $\frac{1}{2}n(n - 1)$  Schlüsseln, die auf einem sicheren Weg wie etwa bei einem persönlichen Treffen oder durch vertrauenswürdige Boten ausgetauscht werden müßten.



BAILEY WHITFIELD DIFFIE wurde 1944 geboren. Erst im Alter von zehn Jahren lernte er lesen; im gleichen Jahr hielt eine Lehrerin an seiner New Yorker Grundschule einen Vortrag über Chiffren. Er ließ sich von seinem Vater alle verfügbare Literatur darüber besorgen, entschied sich dann 1961 aber doch für ein Mathematikstudium am MIT. Um einer Einberufung zu entgehen, arbeitete er nach seinem Bachelor bei Mitre; später, nachdem sein Interesse an der Kryptographie wieder erwacht war, kam er zu MARTIN HELLMAN nach Stanford, der ihn als Forschungsassistent einstellte. Ab 1991 arbeitete er als *chief security officer* bei Sun Microsystems, von 2010 bis 2012 war er bei ICANN für Sicherheit zuständig.



MARTIN HELLMAN wurde 1945 in New York geboren. Er studierte Elektrotechnik zunächst bis zum Bachelor an der dortigen Universität; für Master und Promotion studierte er in Stanford. Nach kurzen Zwischenaufenthalten am Watson Research Center der IBM und am MIT wurde er 1971 Professor an der Stanford University. Seit 1996 ist er emeritiert, gibt aber immer noch Kurse, mit denen er Schüler für mathematische Probleme interessieren will. Seine home page findet man unter <http://www-ee.stanford.edu/~hellman/> .

DIFFIE und HELLMAN machten nur sehr vage Andeutungen, wie so ein System mit öffentlichen Schlüsseln aussehen könnte. Es ist zunächst einmal klar, daß ein solches System keinerlei Sicherheit gegen einen Gegner mit unbeschränkter Rechenkraft (In der Kryptographie spricht man von einem BAYESSchen Gegner) bieten kann, denn für jeden gegebenen Chiffretext ist die Länge der in Frage kommenden Quelltexte beschränkt, und damit gibt es für den Quelltext auch nur endlich viele Möglichkei-

ten. Ein Gegner mit unbeschränkter Rechenkraft muß daher einfach auf jeden dieser potentiellen Quelltexte die bekannte Verschlüsselungsfunktion anwenden bis er den zu entschlüsselnden Chiffretext erhält.

Wer im Gegensatz zum BAYESSchen Gegner nur über begrenzte Ressourcen verfügt, kann diesen Algorithmus natürlich auch anwenden; zu Programmieren ist er einfach genug. Bei einem guten Kryptosystem, das kompetent angewandt wird, kann er aber ziemlich sicher sein, daß seine Computer auch nach Jahrzehnten noch nicht den gesuchten Quelltext identifiziert haben.

Was ein in diesem Sinne „gutes“ Kryptosystem ist, hängt natürlich davon ab, über welche Möglichkeiten die anzunehmenden Gegner verfügen. Bei ernstzunehmenden Gegner kann man sicher sein, daß diese deutlich mehr als nur einige PCs einsetzen können: Größere Organisationen können beispielsweise ihr ganzes Computernetz so programmieren, daß jeder Rechner, der gerade zu nichts anderem gebraucht wird, sich mit einem Teil der Entschlüsselung beschäftigt. Wenn sie häufiger an solchen Entschlüsselungen interessiert sind, lohnt sich auch der Einsatz von Spezialhardware wie FPGAs oder ASICs, und spätestens bei Regierungsstellen und insbesondere Geheimdiensten können wir davon ausgehen, daß diese auch Technologien einsetzen, die auf dem offenen Markt nicht verfügbar und eventuell sogar unbekannt sind. Außerdem könnten Gegner aller Art eventuell über Angriffsmöglichkeiten verfügen, die effizienter sind als systematisches Probieren; sie könnten beispielsweise eine bislang unbekannte Schwachstelle des Verschlüsselungsverfahrens finden und diese ausnutzen.

Bei der Auswahl eines Verfahrens tut man also gut daran, die Fähigkeiten der zu erwartenden Gegner deutlich zu überschätzen, um so noch einen Sicherheitsspielraum zu haben. Da die verfügbare Hardware von Jahr zu Jahr leistungsfähiger wird, ist auch klar, daß kein solches Verschlüsselungsverfahren für alle Ewigkeit sicher ist, sondern höchstens für eine einigermaßen vorhersehbare Zukunft von relativ wenigen Jahren.

Zur Abschätzung der Sicherheit eines Verfahrens verwendet man den Begriff der „ $n$ -Bit-Sicherheit“. Er soll besagen, daß ein Gegner nur dann



eine nicht vernachlässigbare Chance auf eine unbefugte Entschlüsselung hat, wenn er mindestens  $2^n$  Rechenoperationen durchführen kann. Das Bundesamt für Sicherheit in der Informationstechnik hielt bislang ein Sicherheitsniveau von 100 Bit für ausreichend, strebt aber für die nächsten Jahre, spätestens ab Anfang 2023, eines von 120 Bit an. Auf europäischer Ebene gibt es eine SOG-IS (Senior Officials Group Information Security) Crypto Working Group, in der beispielsweise Deutschland durch Beamte des Bundesamts für Sicherheit in der Informationstechnik vertreten ist. Diese Arbeitsgruppe unterscheidet zwischen *legacy mechanisms* und *recommended mechanisms*. Das englische Wort *legacy* bedeutet in diesem Zusammenhang *überliefert*, *hergebracht* oder gar *Altlast*; solche Verfahren müssen eine Sicherheit von mindestens hundert Bit haben und sind akzeptabel bis Ende 2020. Danach sollten nur noch *recommended mechanisms* verwendet werden; deren Sicherheit liegt bei mindestens 125 Bit. Kryptologen im akademischen Bereich empfehlen schon seit mehreren Jahren 128-Bit-Sicherheit. Die meisten deutschen Bankkarten benutzen zur Verschlüsselung der PIN noch den sogenannten Triple-DES; seine Sicherheit liegt bei 112 Bit.

Um ein Gefühl für diese verschiedenen Sicherheitsniveaus zu bekommen, wollen wir uns überlegen, was  $n$ -Bit-Sicherheit bei heutiger Technologie bedeutet. Angenommen, wir haben einen Prozessor, der  $10^{10}$  Rechenoperationen pro Sekunde ausführen kann. (Man beachte, daß ein mit 10 GHz getakteter Prozessor dazu nicht in der Lage ist, denn die meisten Operationen benötigen deutlich mehr als einen Takt.) Pro Jahr kann dieser Prozessor dann ungefähr

$$10^{10} \times 60 \times 60 \times 365,25 = 3,15576 \times 10^{17}$$

Rechenoperationen ausführen. Für  $2^n$  Rechenoperationen benötigte er

für  $n = 100$  mit  $2^n \approx 1,27 \times 10^{30}$  ungefähr  $4,0 \times 10^{13}$  Jahre

für  $n = 112$  mit  $2^n \approx 5,19 \times 10^{33}$  ungefähr  $1,6 \times 10^{16}$  Jahre

für  $n = 120$  mit  $2^n \approx 1,33 \times 10^{36}$  ungefähr  $4,2 \times 10^{18}$  Jahre

für  $n = 125$  mit  $2^n \approx 4,25 \times 10^{37}$  ungefähr  $1,3 \times 10^{20}$  Jahre

für  $n = 128$  mit  $2^n \approx 3,40 \times 10^{38}$  ungefähr  $1,1 \times 10^{21}$  Jahre

Angesichts dieser Zahlen würden wir natürlich nicht einen Prozessor allein rechnen lassen, sondern vielleicht eine Million Prozessoren pa-

rallel; dadurch erniedrigen sich die Exponenten in der letzten Spalte um sechs.

Die besten heutigen Supercomputer haben einen Durchsatz im Bereich von mehreren Tera-Flops, d.h. sie können mehrere Billionen Gleitkommaoperationen pro Sekunde (*floating points operations per second*) durchführen. Gleitkommaarithmetik spielt bei den heute gebräuchlichen Kryptoverfahren keine nennenswerte Rolle; mit hinreichend viel Geld kann man aber wohl ähnliche Maschinen bauen, die eine entsprechende Anzahl von für die Kryptanalyse relevanten Operationen meistern. Wenn wir unseren Chip durch solche Maschinen ersetzen, können wir die Exponenten nochmals um drei erniedrigen. Falls ein Gegner bislang unbekannte neue Technologien einsetzen kann, beispielsweise Quanteneffekte, reduziert sich der Zeitaufwand noch weiter, genauso wenn er für ein spezielles Verfahren einen Weg findet, um deutlich weniger Berechnungen durchzuführen. Solche Effekte lassen sich naturgemäß nicht abschätzen und müssen daher durch einen möglichst großzügigen Sicherheitszuschlag kompensiert werden.

DIFFIE und HELLMAN bezeichnen eine Funktion, deren Umkehrfunktion nicht mit vertretbarem Aufwand berechnet werden kann, als *Einwegfunktion* und wollen solche Funktionen zur Verschlüsselung verwenden. Das allein führt allerdings noch nicht zu einem praktikablen Kryptosystem, denn bei einer echten Einwegfunktion ist es auch für den legitimen Empfänger nicht möglich, seinen Posteingang zu entschlüsseln. DIFFIE und HELLMAN schlagen deshalb eine Einwegfunktion mit *Falltür* vor, wobei der legitime Empfänger zusätzlich zu seinem öffentlichen Schlüssel noch über einen geheimen Schlüssel verfügt, mit dem er (und nur er) diese Falltür öffnen kann.

Natürlich hängt alles davon ab, ob es solche Einwegfunktionen mit Falltür wirklich gibt. DIFFIE und HELLMAN gaben keine an, und es gab damals durchaus Experten, die nicht an die Existenz solcher Funktionen glaubten.

Tatsächlich gab es wohl bereits damals Systeme, die auf solchen Funktionen beruhten, auch wenn sie nicht in der offenen Literatur dokumentiert waren: Die britische *Communications-Electronics Security Group*

(CESG) hatte bereits Ende der sechziger Jahre begonnen, nach entsprechenden Verfahren zu suchen, um die Probleme des Militärs mit dem Schlüsselmanagement zu lösen, aufbauend auf (impraktikablen) Ansätzen von AT&T zur Sprachverschlüsselung während des zweiten Weltkriegs. Die CESG sprach nicht von Kryptographie mit öffentlichen Schlüsseln, sondern von *nichtgeheimer Verschlüsselung*, aber das Prinzip war das gleiche.

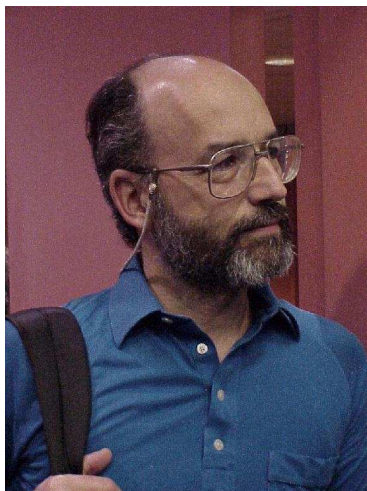
Erste Ideen dazu sind in einer auf Januar 1970 datierten Arbeit von JAMES H. ELLIS zu finden, ein praktikables System in einer auf den 20. November 1973 datierten Arbeit von CLIFF C. COCKS. Wie im Milieu üblich, gelangte nichts über diese Arbeiten an die Öffentlichkeit; erst 1997 veröffentlichten die *Government Communications Headquarters* (GCHQ), zu denen CESG gehört, einige Arbeiten aus der damaligen Zeit. Eine Zeitlang waren sie auch auf dem Server <http://www.cesg.gov.uk/> zu finden, wo sie allerdings inzwischen anscheinend wieder verschwunden sind.

Im akademischen Bereich gab es ein Jahr nach Erscheinen der Arbeit von DIFFIE und HELLMAN das erste Kryptosystem mit öffentlichen Schlüsseln: Drei Wissenschaftler am Massachusetts Institute of Technology fanden nach rund vierzig erfolglosen Ansätzen 1977 schließlich jenes System, das heute nach ihren Anfangsbuchstaben mit RSA bezeichnet wird: RON RIVEST, ADI SHAMIR und LEN ADLEMAN.

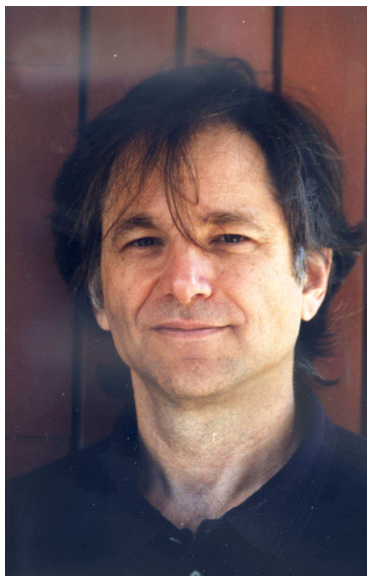
RIVEST, SHAMIR und ADLEMAN gründeten eine Firma namens RSA Computer Security Inc., die 1983 das RSA-Verfahren patentieren ließ und auch nach Auslaufen dieses Patents im September 2000 weiterhin erfolgreich im Kryptobereich tätig ist. 2002 erhielten RIVEST, SHAMIR und ADLEMAN für die Entdeckung des RSA-Systems den TURING-Preis der *Association for Computing Machinery* ACM, ein jährlich vergebener Preis, der als eine der höchsten Auszeichnungen der Informatik gilt. Die Firma RSA Computer Security Inc. wurde 2006 von einer ebenfalls nach den Initialen ihrer Gründer benannten Firma übernommen, der 1979 von RICHARD EGAN, ROGER MARINE und JOHN CURLY gegründeten Firma EMC, die vor allem Speichermedien entwickelt hatte, u.a. das erste 64 KB-Board. 2016 schließlich wurde EMC (einschließlich RSA) von Dell übernommen.



RONALD LINN RIVEST wurde 1947 in Schenectady im US-Bundesstaat New York geboren. Er studierte zunächst Mathematik an der Yale University, wo er 1969 seinen Bachelor bekam; danach studierte er in Stanford Informatik. Nach seiner Promotion 1974 wurde er Assistenzprofessor am Massachusetts Institute of Technology, wo er heute einen Lehrstuhl hat. Er arbeitet immer noch auf dem Gebiet der Kryptographie und entwickelte eine ganze Reihe weiterer Verfahren, auch symmetrische Verschlüsselungsalgorithmen und Hashverfahren. Er ist Koautor eines Lehrbuchs über Algorithmen. Seine home page ist <http://people.csail.mit.edu/rivest/>.



ADI SHAMIR wurde 1952 in Tel Aviv geboren. Er studierte zunächst Mathematik an der dortigen Universität; nach seinem Bachelor wechselte er ans Weizmann Institut, wo er 1975 seinen Master und 1977 die Promotion in Informatik erhielt. Nach einem Jahr als Postdoc an der Universität Warwick und drei Jahren am MIT kehrte er ans Weizmann Institut zurück, wo er bis heute Professor ist. Außer für RSA ist er bekannt sowohl für die Entwicklung weiterer Kryptoverfahren als auch für erfolgreiche Angriffe gegen Kryptoverfahren. Er schlug auch einen optischen Spezialrechner zur Faktorisierung großer Zahlen vor. Seine home page ist erreichbar unter <http://www.wisdom.weizmann.ac.il/profile/scientists/shamir-profile.html>

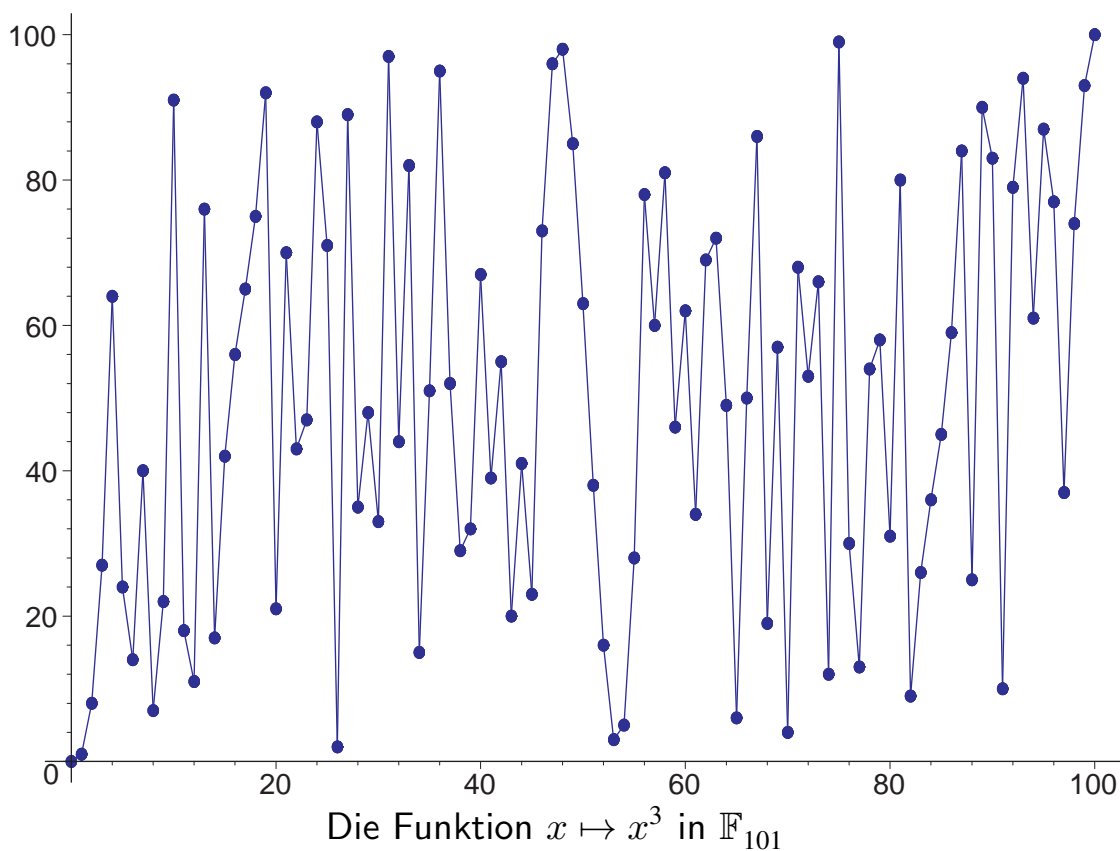


LEONARD ADLEMAN wurde 1945 in San Francisco geboren. Er studierte in Berkeley, wo er 1968 einen BS in Mathematik und 1976 einen PhD in Informatik erhielt. Thema seiner Dissertation waren zahlentheoretische Algorithmen und ihre Komplexität. Von 1976 bis 1980 war er an der mathematischen Fakultät des MIT; seit 1980 arbeitet er an der University of Southern California in Los Angeles. Seine Arbeiten beschäftigen sich mit Zahlentheorie, Kryptographie und Molekularbiologie. Er führte nicht nur 1994 die erste Berechnung mit einem „DNS-Computer“ durch, sondern arbeitete auch auf dem Gebiet der Aidsforschung. Heute hat er einen Lehrstuhl für Informatik und Molekularbiologie. <http://www.usc.edu/directory/profile/?lname=Adleman&fname=Leonard>

Das RSA-Verfahren ist im wesentlichen identisch mit dem von der CESG entwickelten System, so daß man auch Zweifel an den Behauptungen der GCHQ haben kann. Die Beschreibung durch RIVEST, SHAMIR und ADLEMAN erschien 1978 unter dem Titel *A method for obtaining digital signatures and public-key cryptosystems* in *Comm. ACM* **21**, 120–126.

## §2: Das RSA-Verfahren

Für eine natürliche Zahl  $e$  ist die reelle Funktion  $x \mapsto x^e$  für positive  $x$  monoton ansteigend und bijektiv; ihre Umkehrfunktion  $x \mapsto \sqrt[e]{x}$  läßt sich mit etwa demselben Aufwand berechnen wie die Funktion selbst. Betrachten wir  $x \mapsto x^e$  allerdings als Funktion von  $\mathbb{Z}/N \rightarrow \mathbb{Z}/N$ , so erhalten wir einen sehr chaotisch aussehenden Graphen und können uns daher Hoffnungen machen, daß diese Funktion vielleicht als Grundlage einer kryptographischen Verschlüsselung brauchbar sein könnte.



Dazu muß sie natürlich zunächst einmal injektiv sein. Da die Ordnung

eines Elements von  $(\mathbb{Z}/N\mathbb{Z})^\times$  Teiler von  $\varphi(N)$  ist, muß insbesondere  $e$  teilerfremd zu  $\varphi(N)$  sein. Dann lassen sich mit dem erweiterten EUKLIDischen Algorithmus Zahlen  $d, k \in \mathbb{N}$  finden, so daß  $de - k\varphi(N) = 1$  ist, d.h. für jedes zu  $N$  teilerfremde  $x$  ist

$$(x^e)^d = x^{ed} = x^{1+k\varphi(N)} \equiv x \pmod{N}.$$

Somit sind die Funktionen

$$\left\{ \begin{array}{l} (\mathbb{Z}/N\mathbb{Z})^\times \rightarrow (\mathbb{Z}/N\mathbb{Z})^\times \\ x \mapsto x^e \end{array} \right. \quad \text{und} \quad \left\{ \begin{array}{l} (\mathbb{Z}/N\mathbb{Z})^\times \rightarrow (\mathbb{Z}/N\mathbb{Z})^\times \\ x \mapsto x^d \end{array} \right.$$

zueinander invers.

Die Beschränkung auf prime Restklassen ist für kryptographische Anwendungen ungünstig: Am einfachsten wäre es, wenn wir jede Bitfolge, deren Länge kleiner ist als die der Binärdarstellung von  $N$ , als Zahl zwischen 0 und  $N - 1$  auffassen, verschlüsseln und übertragen könnten. Der Empfänger könnte dann die Zahl entschlüsseln, als Bitfolge hinschreiben und daraus die Nachricht rekonstruieren. Zum Glück ist das zumindest für *quadratifreie* Zahlen  $N$ , d.h. Zahlen, die durch kein Primzahlquadrat teilbar sind, auch möglich:

**Satz:** Für eine quadratifreie natürliche Zahl  $N$  sind die beiden Funktionen

$$\left\{ \begin{array}{l} \mathbb{Z}/N\mathbb{Z} \rightarrow \mathbb{Z}/N\mathbb{Z} \\ x \mapsto x^e \end{array} \right. \quad \text{und} \quad \left\{ \begin{array}{l} \mathbb{Z}/N\mathbb{Z} \rightarrow \mathbb{Z}/N\mathbb{Z} \\ x \mapsto x^d \end{array} \right.$$

bijektiv und invers zueinander.

*Beweis:* Als quadratifreie Zahl ist  $N$  ein Produkt verschiedener Primzahlen  $p_i$ , und  $\varphi(N)$  ist das Produkt der zugehörigen  $\varphi(p_i) = p_i - 1$ . Somit ist auch  $ed \equiv 1 \pmod{\varphi(p_i)}$  für alle  $i$ , und nach dem chinesischen Restesatz genügt es, wenn wir den Satz für die einzelnen  $p_i$  beweisen. Ist  $N = p$  prim, so ist die Null das einzige Element von  $\mathbb{Z}/p$ , das nicht in  $(\mathbb{Z}/p)^*$  liegt, und sie wird von beiden Funktionen auf sich selbst abgebildet. ■

Jeder, der  $e$  und  $N$  kennt, kann auch  $d$  berechnen, allerdings muß er dazu als erstes  $\varphi(N)$  bestimmen. Nach der Formel aus Kapitel 1, §7 ist das

möglich, wenn er die Primfaktorzerlegung von  $N$  kennt. Einfachere alternative Verfahren sind nicht bekannt, und wie wir im Kapitel über Faktorisierungsalgorithmen sehen werden, liegt der in der offenen Literatur veröffentlichte Rekord für die Faktorisierung eines Produkts  $N = pq$  zweier gut gewählter Primzahlen bei einem  $N$  mit etwas mehr als 250 Dezimalstellen. Natürlich wird diese Schranke im Laufe der Jahrzehnte ansteigen, und wahrscheinlich können einzelne Geheimdienste schon heute etwas mehr als der Rest der Welt. Es ist aber unwahrscheinlich, daß bei einem schon seit Jahrhunderten untersuchten Problem wie der Faktorisierung ganzer Zahlen ausgerechnet einem Geheimdienst ein Durchbruch gelingen sollte, von dem der Rest der Welt nichts bemerkt. Die Effekte einer leistungsfähigeren Hardware lassen sich durch großzügige Sicherheitszuschläge kompensieren.

Für eine Primzahl  $N = p$  kann natürlich jeder ganz einfach  $\varphi(p) = p - 1$  berechnen; ist  $N = pq$  dagegen das Produkt zweier Primzahlen, so ist die Bestimmung von

$$\varphi(N) = (p - 1)(q - 1) = N - (p + q) + 1$$

äquivalent zur Kenntnis der Faktorisierung: Kennt man nämlich das Produkt  $N = pq$  sowie die Summe  $N + 1 - \varphi(N) = p + q$  der beiden Primzahlen, so kann man sie einfach berechnen als Lösungen der quadratischen Gleichung

$$(x - p)(x - q) = x^2 - (N + 1 - \varphi(N))x + N = 0.$$

Auch bei Produkten von mehr als drei Primzahlen ist keine Methode zur Berechnung der EULERSchen  $\varphi$ -Funktion bekannt, die effizienter wäre als der Umweg über die Faktorisierung, allerdings wird die Faktorisierung bei konstanter Größenordnung von  $N$  tendenziell einfacher, wenn die Anzahl der Faktoren steigt, da wir es dann zumindest teilweise mit kleineren Faktoren zu tun haben.

Zur praktischen Durchführung des RSA-Verfahrens wählt sich daher jeder Teilnehmer zwei verschiedene Primzahlen  $p, q$ , die unbedingt geheim gehalten werden müssen, und eine natürliche Zahl  $e$ , die keinen gemeinsamen Teiler mit  $(p - 1)(q - 1)$  hat. Die Zahlen  $N = pq$  und  $e$  sind sein öffentlicher Schlüssel, der beispielsweise in einem Verzeichnis publiziert werden kann.

Des weiteren berechnet er ein gemeinsames Vielfaches  $\lambda$  von  $p - 1$  und  $q - 1$ , zum Beispiel das kleinste gemeinsame Vielfache oder aber einfach  $\lambda = \varphi(N) = (p - 1)(q - 1)$ , und dazu nach dem erweiterten EUKLIDischen Algorithmus natürliche Zahlen  $d$  und  $k$  so daß  $de - k\lambda = 1$  ist. Dabei kann erreicht werden, daß  $d < \lambda$  (und  $k < e$ ) ist; ein kleineres  $\lambda$  führt also im Allgemeinen auch zu einem kleineren  $d$ . Die so bestimmte Zahl  $d$  ist sein geheimer Schlüssel; da

$$(a^e)^d = a^{ed} = a^{1+k\lambda} = a \cdot a^{k\lambda} \equiv a \pmod{N}$$

ist für alle  $a$ , läßt sich die Entschlüsselung rückgängig machen durch Potenzieren mit  $d$ .

Jeder, der den öffentlichen Schlüssel  $(N, e)$  kennt, kann damit Nachrichten verschlüsseln: Er bricht die Nachricht auf in Blöcke, die durch ganze Zahlen zwischen 0 und  $N - 1$  dargestellt werden können, berechnet für jeden so dargestellten Block  $a$  den Chiffretext  $b = a^e \pmod{N}$  und schickt diesen an den Inhaber des geheimen Schlüssels. Dieser berechnet  $b^d \pmod{N} = a^{ed} \pmod{N} = a$ , und da er dazu seinen geheimen Schlüssel braucht, kann dies niemand außer ihm auf diese Weise berechnen.

In der hier vorgestellten Reinform ist das Verfahren natürlich noch nicht praktikabel: Ein Gegner, der weiß, daß nur wenige Klartexte in Frage kommen, kann diese leicht verschlüsseln und so sehen, welcher auf den gegebenen Chiffretext führt. Falls die Nachricht  $a$  kleiner ist als  $\sqrt[3]{N}$  und ist  $e = 3$ , so ist der Chiffretext  $b = a^3$  kleiner als  $N$ , und damit kann  $a$  einfach als die gewöhnliche Kubikwurzel aus  $b$  berechnet werden. Um Attacken dieser Art zu verhindern, muß das Verfahren daher noch etwas modifiziert werden. Üblicherweise geschieht das in der Weise, daß die mögliche Länge des Nachrichtenblocks nicht ganz ausgenutzt wird und stattdessen ganz links zunächst ein oder zwei vorgegebene *magic bytes* stehen und dahinter eine gewisse Anzahl von Zufallsbits. Ein gegebener Nachrichtenblock kann somit je nach Wahl der Zufallsbits auf  $2^n$  verschiedene Weisen verschlüsselt sein, und wenn  $n$  dem vorgegebenen Sicherheitsniveau entspricht, kann der Gegner nicht alle diese Möglichkeiten durchprobieren. Außerdem gibt es bei dieser Vorgehensweise keine „kleinen“ Nachrichtenblöcke mehr.



### §3: Weitere Anwendungen des RSA-Verfahrens

Im Gegensatz zu symmetrischen Kryptoverfahren bietet das RSA-Verfahrens nicht nur die Möglichkeit einer Verschlüsselung, sondern erlaubt noch eine ganze Reihe weiterer Anwendungen:

#### a) Identitätsnachweis

Hier geht es darum, in Zugangskontrollsystemen, vor Geldautomaten oder bei einer Bestellung im Internet die Identität einer Person zu beweisen: Mit RSA ist das beispielsweise dadurch möglich, daß nur der Inhaber des geheimen Schlüssels  $d$  zu einer gegebenen Zahl  $a$  eine Zahl  $b$  berechnen kann, für die  $b^e \equiv a \pmod{N}$  ist. Letzteres wiederum kann jeder überprüfen, der den öffentlichen Schlüssel  $(N, e)$  kennt.

Falls also der jeweilige Gegenüber eine Zufallszahl  $a$  erzeugt und als Antwort das zugehörige  $b$  verlangt, kann er anhand eines öffentlichen Schlüsselverzeichnisses die Richtigkeit von  $b$  überprüfen und sich so von der Identität seines Partners überzeugen. Im Gegensatz zu Kreditkarteninformation oder Paßwörtern ist dieses Verfahren auch immun gegen Abhören: Falls jedesmal ein neues zufälliges  $a$  erzeugt wird, nützt ein einmal abgehörtes  $b$  nichts.

Grundsätzlich bräuchte man hier kein Kryptosystem mit öffentlichen Schlüsseln; in der Tat funktionierten die ersten Freund-/Feinderkennungssysteme für Flugzeuge zur Zeit des zweiten Weltkriegs nach diesem Prinzip, aber damals natürlich mit einem klassischen symmetrischen Kryptosystem, wobei alle Teilnehmer mit demselben Schlüssel arbeiteten. Der Vorteil eines asymmetrischen Systems besteht darin, daß sich keiner der Teilnehmer für einen anderen ausgeben kann, was beispielsweise wichtig ist, wenn man sich gegenüber weniger vertrauenswürdigen Personen identifizieren muß.

Trotzdem ist das Verfahren in dieser Form nicht als Ersatz zur Übertragung von rechtlich bindender Information geeignet, da der Gegenüber anhand des öffentlichen Schlüssels jederzeit zu einer willkürlich gewählten Zahl  $b$  die Zahl  $a = b^e \pmod{N}$  erzeugen kann um dann zu behaupten, er habe  $b$  als Antwort darauf empfangen. Daher kann der Inhaber des geheimen Schlüssels zwar seine Identität beweisen, aber sein

Gegenüber kann später nicht beispielsweise vor Gericht beweisen, daß er dies (zum Beispiel bei einer Geldabhebung oder Bestellung) getan hat. Falls dies eventuell nötig werden könnte, ist das hier vorgestellte Verfahren also ungeeignet; es funktioniert nur zwischen Personen, die einander vertrauen können.

Eine mögliche Modifikation bestünde darin, daß man beispielsweise noch zusätzlich verlangt, daß die Zahl  $a$  eine spezielle Form hat, etwa daß die vordere Hälfte der Ziffernfolge identisch mit der hinteren Hälfte ist. Ohne Kenntnis von  $d$  hat man praktisch keine Chancen eine Zahl  $b$  zu finden, für die  $b^e \bmod N$  eine solche Gestalt hat: Bei Zahlen mit  $2r$  Ziffern liegt die Wahrscheinlichkeit dafür bei  $10^{-r}$ .

## b) Elektronische Unterschriften

Praktische Bedeutung hat vor allem eine andere Variante: die elektronische Unterschrift. Hier geht es darum, daß der Empfänger erstens davon überzeugt wird, daß eine Nachricht tatsächlich vom behaupteten Absender stammt, und daß er dies zweitens auch einem Dritten gegenüber beweisen kann. (In Deutschland sind solche elektronischen Unterschriften, sofern gewisse formale Voraussetzungen erfüllt sind, rechtsverbindlich.)

Um einen Nachrichtenblock  $a$  mit  $0 \leq a < N$  zu unterschreiben, berechnet der Inhaber des öffentlichen Schlüssels  $(N, e)$  mit seinem geheimen Schlüssel  $d$  die Zahl

$$b = a^d \bmod N$$

und sendet das Paar  $(a, b)$  an den Empfänger. Dieser überprüft, ob

$$b^e \equiv a \bmod N ;$$

falls ja, akzeptiert er dies als unterschriebene Nachricht  $a$ . Da er ohne Kenntnis des geheimen Schlüssels  $d$  nicht in der Lage ist, den Block  $(a, b)$  zu erzeugen, kann er auch gegenüber einem Dritten beweisen, daß der Absender selbst die Nachricht  $a$  unterschrieben hat.

Für kurze Nachrichten ist dieses Verfahren in der vorgestellten Form praktikabel; in vielen Fällen kann man sogar auf die Übermittlung von  $a$

verzichten, da  $b^e \bmod N$  für ein falsch berechnetes  $b$  mit an Sicherheit grenzender Wahrscheinlichkeit keine sinnvolle Nachricht ergibt.

Falls die übermittelte Nachricht geheimgehalten werden soll, müssen  $a$  und  $b$  natürlich noch vor der Übertragung mit dem öffentlichen Schlüssel des Empfängers oder nach irgendeinem anderen Kryptoverfahren verschlüsselt werden.

Bei langen Nachrichten ist die Verdoppelung der Nachrichtenlänge nicht mehr akzeptabel, und selbst, wenn man auf die Übertragung von  $a$  verzichten kann, ist das Unterschreiben jedes einzelnen Blocks sehr aufwendig. Deshalb unterschreibt man meist nicht die Nachricht selbst, sondern einen daraus extrahierten Hashwert. Dieser Wert muß natürlich erstens von der gesamten Nachricht abhängen, und zweitens muß es für den Empfänger (praktisch) unmöglich sein, zwei Nachrichten zu erzeugen, die zum gleichen Hashwert führen. Letzteres bedeutet wegen des sogenannten *Geburtstagsparadoxons*, daß für  $n$ -Bit Sicherheit Hashwerte der Länge  $2n$  erforderlich sind. Die immer noch recht verbreiteten Hashalgorithmen, die 160 Bit liefern, haben somit nur eine Sicherheit von etwa 80 Bit, was heute nicht mehr als wirklich sicher gelten kann. Die heute gebräuchlichen Hashalgorithmen liefern Werte mit 224 oder 256 Bit, was einer 112- oder 128-Bit Sicherheit entspricht. Die Algorithmen funktionieren ähnlich wie symmetrische Kryptoverfahren; sie versuchen durch Konfusion und Diffusion ein Ergebnis zu berechnen, dessen sämtliche Bits in einer nicht offensichtlichen Weise von jedem einzelnen Nachrichtenbit abhängen.

### c) SSL und TLS

Eine wichtige Anwendung elektronischer Unterschriften ist die Veröffentlichung von RSA-Schlüsseln: Falls es einem Angreifer gelingt, einem Teilnehmer  $A$  einen falschen öffentlichen Schlüssel von Teilnehmer  $B$  unterzuschreiben, kann (nur) der Angreifer die Nachrichten von  $A$  an  $B$  lesen, und er kann sich gegenüber  $A$  mittels elektronischer Unterschrift als  $B$  ausgeben. Daher sind öffentliche Schlüssel meist unterschrieben von einer Zertifizierungsstelle. Auch deren Unterschrift muß natürlich gegen Manipulationen gesichert sein, beispielsweise indem sie von der nächsthöheren Zertifizierungsstelle unterschrieben ist.

An der Spitze der Zertifizierungshierarchie stehen Stellen, deren elektronische Unterschrift jeder Teilnehmer kennen sollte, weil es sich entweder um staatliche Stellen handelt, deren elektronische Unterschriften auf leicht zugänglichen Webseiten verifiziert werden können, oder aber – in der Praxis häufiger – weil die Unterschriften dieser Stellen in Mail- und Browserprogramme eingebaut sind. Letzteres bietet selbstverständlich keine Sicherheit gegen manipulierte Browserprogramme aus dubiosen Quellen, die möglicherweise auch die Mafia als Zertifizierungsstelle anerkennen.

Zertifizierte Unterschriften werden insbesondere angewandt bei den Standards SSL und TLS für sichere Internetverbindungen.

SSL steht für *secure socket layer*, TLS für *transport layer security*; Zweck ist jeweils der Aufbau einer sicheren Internetverbindung. Wie im Internet üblich, können dazu die verschiedensten Verfahren benutzt werden; die auf Grundlage von RSA zählen derzeit zu den populärsten.

Natürlich ist RSA zu aufwendig, um damit eine längere Kommunikation wie beispielsweise eine *secure shell* Sitzung zu verschlüsseln; tatsächlich dient RSA daher nur zur Vereinbarung eines Schlüssels für ein konventionelles Kryptoverfahren wie AES oder teilweise auch noch Triple-DES oder gar noch Schlimmeres, auf das sich die Beteiligten unter SSL/TLS ebenfalls einigen müssen.

Am einfachsten wäre es, wenn der Client einen Schlüssel für ein solches Verfahren wählt und dann diesen mit dem RSA-Schlüssel des Servers verschlüsselt an diesen schickt – vorausgesetzt, er kennt diesen RSA-Schlüssel. Letzteres ist im Allgemeinen nicht der Fall; daher muß zunächst der Server dem Client seinen Schlüssel mitteilen.

Da der Client nicht sicher sein kann, mit dem richtigen Server verbunden zu sein, schickt der Server diesen Schlüssel meist zusammen mit einem Zertifikat, das sowohl seine Identität als auch seinen RSA-Schlüssel enthält und von einer Zertifizierungsstelle unterschrieben ist.

Die öffentlichen Schlüssel der gängigen Zertifizierungsstellen sind, wie bereits erwähnt, in die Browserprogramme eingebaut; bei weniger bekannten Zertifizierungsstellen wie etwa dem Rechenzentrum der Univer-

sität Mannheim fragt der Browser den Benutzer, ob er das Zertifikat anerkennen will oder nicht. Bei *secure shell* schließlich, wo die Gegenseite typischerweise keinerlei Zertifikat vorweisen kann, fragt das Programm beim ersten Verbindungsaufbau zu einem Server, ob dessen Schlüssel anerkannt werden soll und speichert dann einen sogenannten *fingerprint* davon; dieser wird bei späteren Verbindungen zur Identitätsfeststellung benutzt.

Tatsächlich funktioniert das Verfahren nicht so einfach; da erfolgreiche Kryptographie nur mit maximaler Paranoia möglich ist, wird der tatsächliche Schlüssel für das symmetrische Verfahren in einem komplizierten Verfahren aus Eingaben beider Seiten berechnet.

#### **d) Blinde Unterschriften und elektronisches Bargeld**

Einer der erfolgversprechendsten Ansätze zum Aushebeln eines Kryptosystems besteht darin, sich auf die Dummheit seiner Mitmenschen zu verlassen.

So sollte es durch gutes Zureden nicht schwer sein, jemanden zu Demonstrationzwecken zum Unterschreiben einer sinnlosen Nachricht zu bewegen: Eine Folge von Nullen und Einsen ohne sinnvolle Interpretation hat schließlich keine rechtliche Wirkung.

Nun muß eine sinnlose Nachricht aber nicht unbedingt eine Zufallszahl sein: Sie kann sorgfältig präpariert sein. Sei dazu etwa  $m$  eine Nachricht, die ein Zahlungsverprechen enthält,  $(N, e)$  der öffentliche Schlüssel des Opfers und  $r$  eine Zufallszahl zwischen 2 und  $N - 2$ . Dann wird

$$x = m \cdot r^e \pmod{N}$$

wie eine Zufallsfolge aussehen, für die man eine Unterschrift

$$u = x^d \pmod{N} = (mr^e)^d \pmod{N} = m^d r \pmod{N}$$

bekommt. Multiplikation mit  $r^{-1}$  macht daraus eine Unterschrift unter die Zahlungsverpflichtung  $m$ .

Das angegebene Verfahren kann nicht nur von Trickbetrügern benutzt werden; blinde Unterschriften sind auch die Grundlage von *digitalem Bargeld*.

Zahlungen im Internet erfolgen meist über Kreditkarten; die Kreditkartengesellschaften haben also einen recht guten Überblick über die Ausgaben ihrer Kunden und machen teilweise auch recht gute Geschäfte mit Kundenprofilen.

Digitales Bargeld will die Anonymität von Geldscheinen mit elektronischer Übertragbarkeit kombinieren und so ein anonymes Zahlungssystem z.B. für das Internet bieten.

Es wird ausgegeben von einer Bank, die für jede angebotene Stückelung einen öffentlichen Schlüssel  $(N, e)$  bekanntgibt. Eine Banknote ist eine mit dem zugehörigen geheimen Schlüssel unterschriebene Seriennummer.

Die Seriennummer kann natürlich nicht einfach *jede* Zahl sein; sonst wäre jede Zahl kleiner  $N$  eine Banknote. Andererseits dürfen die Seriennummern aber auch nicht von der Bank vergeben werden, denn sonst wüßte diese, welcher Kunde Scheine mit welchen Seriennummern hat. Als Ausweg wählt man Seriennummern einer sehr speziellen Form: Ist  $N > 10^{150}$ , kann man etwa als Seriennummer eine 150-stellige Zahl wählen, deren Ziffern spiegelsymmetrisch zur Mitte sind, d.h. ab der 76. Ziffer werden die vorherigen Ziffern rückwärts wiederholt. Die Wahrscheinlichkeit, daß eine zufällige Zahl  $x$  nach Anwendung des öffentlichen Exponenten auf so eine Zahl führt, ist  $10^{-75}$  und damit vernachlässigbar.

Seriennummern werden von den Kunden zufällig erzeugt. Für jede solche Seriennummer  $m$  erzeugt der Kunde eine Zufallszahl  $r$ , schickt  $mr^e \bmod N$  an die Bank und erhält (nach Belastung seines Kontos) eine Unterschrift  $u$  für diese Nachricht zurück. Wie oben berechnet er daraus durch Multiplikation mit  $r^{-1}$  die Unterschrift  $v = m^d \bmod N$  für die Seriennummer  $N$ , und mit diesem Block kann er bezahlen.

Der Zahlungsempfänger berechnet  $v^e \bmod N$ ; falls dies die Form einer gültigen Seriennummer hat, kann er sicher sein, einen von der Bank unterschriebenen Geldschein vor sich zu haben. Er kann allerdings noch nicht sicher sein, daß dieser Geldschein nicht schon einmal ausgegeben wurde.

Deshalb muß er die Seriennummer an die Bank melden, die mit ihrer Datenbank bereits ausbezahlter Seriennummern vergleicht. Falls sie darin noch nicht vorkommt, wird sie eingetragen und der Händler bekommt sein Geld; andernfalls verweigert sie die Zahlung.

Bei  $10^{75}$  möglichen Nummern liegt die Wahrscheinlichkeit dafür, daß zwei Kunden, die eine (wirklich) zufällige Zahl wählen, dieselbe Nummer erzeugen, bei etwa  $10^{-37,5}$ . Die Wahrscheinlichkeit, mit jeweils einem Spielschein fünf Wochen lang hintereinander sechs Richtige im Lotto zu haben, liegt dagegen bei  $\binom{49}{6}^{-5} \approx 5 \cdot 10^{-35}$ , also etwa um den Faktor sechzig höher. Zwei gleiche Seriennummern sind also praktisch auszuschließen, wenn auch theoretisch möglich.

Falls wirklich einmal zufälligerweise zwei gleiche Seriennummern erzeugt worden sein sollten, kann das System nur funktionieren, wenn der zweite Geldschein mit derselben Seriennummer nicht anerkannt wird, so daß der zweite Kunde sein Geld verliert. Dies muß als eine zusätzliche Gebühr gesehen werden, die mit an Sicherheit grenzender Wahrscheinlichkeit nie fällig wird, aber trotzdem nicht ausgeschlossen werden kann.

Da digitales Bargeld nur in kleinen Stückelungen sinnvoll ist (Geldscheine im Millionenwert wären auf Grund ihrer Seltenheit nicht wirklich anonym und würden wegen der damit verbundenen Möglichkeiten zur Geldwäsche auch in keinem seriösen Wirtschaftssystem akzeptiert), wäre der theoretisch mögliche Verlust ohnehin nicht sehr groß.

Digitales Bargeld der gerade beschriebenen Form wurde 1982 von DAVID CHAUM vorgestellt; 1990 gründete er eine Firma namens DigiCash, die es kommerziell vermarkten sollte. Zu deren Kunden gehörte beispielsweise auch die Deutsche Bank, die allerdings nur 27 Kunden fand, die Zahlungen damit akzeptierten. DigiCash wurde 1998 zahlungsunfähig und später trotz allem von der Deutschen Bank übernommen, aber ziemlich bald eingestellt. Derzeit gibt es meines Wissens kein entsprechendes Zahlungssystem. Die heutigen digitalen Währungen wie Bitcoin beruhen zwar auch auf kryptographischen Verfahren, den sogenannten Blockchains, aber diese haben nichts mit RSA zu tun und auch nichts mit Zahlentheorie..

### e) Bankkarten mit Chip

Bankkarten speichern ihre Information sowohl in einem Chip, als auch, unabhängig davon, auf einem einen Magnetstreifen. Dort stehen Informationen wie Kontenname und -nummer, Bankleitzahl, Gültigkeitsdauer *usw.*; dazu kommt verschlüsselte Information, die unter anderem die Geheimzahl enthält, die aber auch von den obengenannten Daten abhängt. Zur Verschlüsselung verwendet man hier ein konventionelles, d.h. symmetrisches Kryptoverfahren; derzeit noch meist Triple-DES.

Der Schlüssel, mit dem dieses arbeitet, muß natürlich streng geheimgehalten werden: Wer ihn kennt, kann problemlos die Geheimzahlen fremder Karten ermitteln und eigene Karten zu beliebigen Konten erzeugen.

Um eine Karte nur anhand der Magnetstreifeninformation zu überprüfen, muß daher eine Verbindung zu einem Zentralrechner aufgebaut werden, an den sowohl der Inhalt des Magnetstreifens als auch die vom Kunden eingetippte Zahl übertragen werden; dieser wendet Triple-DES mit dem Systemschlüssel an und meldet dann, wie die Prüfung ausgefallen ist.

Der Chip enthält ebenfalls die Kontendaten; zusätzlich ist dort auch noch in einem auslesesicheren Register Information über die Geheimzahl gespeichert. Daher muß die eingegebene PIN nicht an einen Zentralrechner übertragen werden, sondern wird vom Lesegerät an den Chip weitergegeben, der dann entscheidet, ob die Eingabe akzeptiert wird oder nicht.

Da frei programmierbare Chipkarten relativ billig sind, muß dafür Sorge getragen werden, daß ein solches System nicht durch einen *Yes-Chip* unterlaufen werden kann, der ebenfalls die Konteninformationen enthält, ansonsten aber ein Programm, das ihn *jede* Geheimzahl akzeptieren läßt. Das Terminal muß also, bevor es überhaupt eine Geheimzahl anfordert, zunächst einmal den Chip authentisieren, d.h. sich davon überzeugen, daß es sich um einen vom Bankenkonsortium ausgegebenen Chip handelt.

Aus diesem Grund sind die Kontendaten auf dem Chip mit dem privaten RSA-Schlüssel des Konsortiums unterschrieben. Die Terminals



kennen den öffentlichen Schlüssel dazu und können so die Unterschrift überprüfen.

Solche Chipkarten wurden hier in Mitteleuropa als erstes in Frankreich ausgegeben; Einzelheiten über die verwendeten Algorithmen und deren technische Implementierung wurden vom Bankenkonsortium streng geheimgehalten. Trotzdem machte sich 1997 ein elsässischer Ingenieur namens SERGE HUMPICH daran, den Chip genauer zu untersuchen. Er verschaffte sich dazu ein (im freien Verkauf erhältliches) Terminal und untersuchte sowohl die Kommunikation zwischen Chip und Terminal als auch die Vorgänge innerhalb des Terminals mit Hilfe eines Logikanalysators. Damit gelang es ihm nach und nach, die Funktionsweise des Terminals zu entschlüsseln und in ein äquivalentes PC-Programm zu übersetzen. Durch dessen Analyse konnte er die Authentisierungsprozedur und die Prüflöge entschlüsseln und insbesondere auch feststellen, daß hier mit RSA gearbeitet wurde.

Blieb noch das Problem, den Modul zu faktorisieren. Dazu besorgte er sich ein japanisches Programm aus dem Internet, das zwar eigentlich für kleinere Zahlen gedacht war, aber eine Anpassung der Wortlänge ist natürlich auch für jemanden, der den Algorithmus hinter dem Programm nicht versteht, kein Problem. Nach sechs Wochen Laufzeit hatte sein PC damit den Modul faktorisiert:

$$\begin{aligned} & 213598703592091008239502270499962879705109534182 \backslash \\ & 6417406442524165008583957746445088405009430865999 \\ = & 1113954325148827987925490175477024844070922844843 \\ \times & 1917481702524504439375786268230862180696934189293 \end{aligned}$$

Als er seine Ergebnisse über einen Anwalt dem Bankenkonsortium mitteilte, zeigte sich, was dieses sich unter Sicherheitsstandards vorstellt: Es erreichte, daß HUMPICH wegen des Eindringens in ein DV-System zu zehn Monaten Haft auf Bewährung sowie einem Franc Schadenersatz plus Zinsen verurteilt wurde; dazu kamen 12 000 F Geldstrafe. Einzelheiten findet man in seinem Buch

SERGE HUMPICH: *Le cerveau bleu*, Xo, 2001.

Ab November 1999 hatten neu ausgegebene Bankkarten noch ein zusätzliches Feld mit einer Unterschrift, die im Gegensatz zum obigen 320-Bit-Modul einen 768-Bit-Modul verwendet. Natürlich können damit erzeugte Unterschriften nur von neueren Terminals überprüft werden, so daß viele Transaktionen weiterhin nur über den 320-Bit-Modul mit inzwischen wohlbekannter Faktorisierung „geschützt“ waren. Die heutigen Standards behandelt der nächste Paragraph.

#### §4: Wie groß sollten die Primzahlen sein?

Das Beispiel der ersten französischen Bankkarten zeigt, daß RSA auch für recht groß aussehende Primzahlen ziemlich unsicher sein kann. Wir müssen uns daher die Frage stellen, wie groß die Primzahlen nach heutigen Standards sein müssen, um einen Angriff mit hinreichender Sicherheit auszuschließen. In §1 lernten wir das Konzept der  $n$ -Bit-Sicherheit kennen; dieses sollten wir auf RSA anwenden.

Der offensichtlichste Angriff auf RSA besteht darin, den Modul  $N$  zu faktorisieren; wir brauchen also Aussagen der Art *Für die Faktorisierung eines Produkts zweier Primzahlen der Längen  $x$  und  $y$  sind mindestens  $2^n$  Rechenoperationen notwendig*. Leider gibt es aber keine Aussagen über den minimalen Aufwand für die Faktorisierung einer gegebenen Zahl; wir können nur abschätzen, welchen Aufwand die besten *bekannt* Algorithmen benötigen. Einige dieser Algorithmen werden wir im Kapitel über Faktorisierung kennen lernen. Wie die historische Entwicklung zeigt, wurden immer wieder neue Algorithmen gefunden, die (zumindest für hinreichend große Zahlen) besser waren als alle bekannten, und wir können natürlich nicht darauf vertrauen, daß diese Entwicklung nun abgeschlossen ist. Ein neues Verfahren muß nicht in der offenen Literatur vorgestellt werden; sein Entdecker kann es auch geheim halten in der Hoffnung, damit Nachrichten entschlüsseln und Unterschriften fälschen zu können. Für diese und andere Möglichkeiten müssen großzügige Sicherheitszuschläge einkalkuliert werden, so daß die Wahl einer mit hoher Wahrscheinlichkeit sicheren Primzahlgröße alles andere als einfach ist.

Bis 2017 hielt es daher der deutsche Staat für seine Pflicht, die Bürger bei

einer derart wichtigen Frage nicht allein lassen: Zwar gibt es in Deutschland keine oberste Bundesbehörde für Primzahlen, aber das Bundesamt für Sicherheit in der Informationstechnik (BSI) und die Bundesnetzagentur für Elektrizität, Gas, Telekommunikation, Post und Eisenbahnen publizierten jedes Jahr ein Dokument mit dem Titel *Geeignete Kryptoalgorithmen zur Erfüllung der Anforderungen nach §17 Abs. 1 bis 3 SigG vom 22. Mai 2001 in Verbindung mit Anlage 1 Abschnitt I Nr. 2 SigV vom 22. November 2001*.

SigV steht für die aufgrund des Signaturgesetzes SigG erlassene Signaturverordnung; beide gemeinsam legen fest, daß elektronische Unterschriften in Deutschland grundsätzlich zulässig und in vielen Fällen rechtlich gleichwertig zu einer klassischen Unterschrift sind, sofern sie gewisse Bedingungen erfüllen. Zu diesen Bedingungen gehörte unter anderem, daß das Verfahren und die Schlüssellänge gemeinsam ein „geeigneter Kryptoalgorithmus“ im Sinne der jeweils gültigen Veröffentlichung der Bundesnetzagentur waren.

Da Rechner immer schneller und leistungsfähiger werden und auch auf der mathematisch-algorithmischen Seite fast jedes Jahr kleinere oder größere Fortschritte zu verzeichnen sind, galten die jeweiligen Empfehlungen nur für etwa sechs Jahre. Im Augenblick ist die Lage relativ stabil; daher waren die Empfehlungen der letzten Jahre ausnahmsweise sieben Jahre gültig. Für Dokumente, die länger gültig sein sollen, sind elektronische Unterschriften somit nicht vorgesehen.

Die letzten Richtlinien stammen vom 7. Dezember 2016 und wurden am 30. Dezember 2016 im Bundesanzeiger veröffentlicht. Dieses „Amtsblatt“ des Bundes erscheint inzwischen nur noch in einer elektronischen Version unter [www.bundesanzeiger.de](http://www.bundesanzeiger.de); die Richtlinien wurden veröffentlicht unter BAnz AT 30.12.2016 B5. Bis Ende 2022 empfehlen sie 2048 Bit; wirklich verbindlich sind allerdings nur 1976. (Der minimale Unterschied hängt mit Implementierungsproblemen beim Chipkarten-Betriebssystem SECCOS zusammen.) Danach, bis Ende 2023, sind mindestens drei Tausend Bit vorgeschrieben.

Die beiden Primfaktoren  $p, q$  sollen zufällig und unabhängig voneinan-

der erzeugt werden und aus einem Bereich stammen, in dem

$$\varepsilon_1 < |\log_2 p - \log_2 q| < \varepsilon_2$$

gilt. Als *Anhaltspunkte* werden dabei die Werte  $\varepsilon_1 \approx 0,1$  und  $\varepsilon_2 \approx 30$  vorgeschlagen; ist  $p$  die kleinere der beiden Primzahlen, soll also  $2^{-10}p < q < 2^{30}p \approx 10^9p$  gelten. Die beiden Primzahlen sollten somit zwar ungefähr dieselbe Größenordnung haben, aber nicht *zu nahe* beieinander liegen. Der Grund dafür ist ein auf FERMAT zurück gehendes Faktorisierungsverfahren auf Grundlage der dritten binomischen Formel: Falls für eine Zahl  $N$  und eine natürliche Zahl  $y$  die Zahl  $N + y^2$  eine Quadratzahl  $x^2$  ist, ist  $N = x^2 - y^2 = (x + y)(x - y)$ , womit zwei Faktoren gefunden sind. Probiert man alle kleinen natürlichen Zahlen  $y$  systematisch durch, führt dieses Verfahren offensichtlich umso schneller zum Erfolg, je näher die beiden Faktoren von  $N$  beieinander liegen. Wir werden uns in Kapitel über Faktorisierung noch genauer damit befassen.

Nachdem die Primzahlen gefunden sind, muß als nächstes der öffentliche Exponent  $e$  gewählt werden. Für diesen soll  $2^{16} + 1 \leq e < 2^{256}$  sein, was aber erst ab 2021 verpflichtend werden sollte. (Heute dürfte noch oft  $e = 3$  verwendet werden.) Danach wird der private Exponent  $d$  so gewählt, daß  $ed \equiv 1 \pmod{\text{kgV}(p - 1, q - 1)}$  ist. Auch wenn das Verfahren primär für Unterschriften verwendet werden soll, darf man also nicht vom privaten Exponenten ausgehen, denn wie wir im Kapitel über Kettenbrüche sehen werden, läßt sich ein kleiner privater Exponent aus  $e$  und  $N = pq$  mit recht geringem Aufwand bestimmen.

Im Sommer 2017 wurde das Signaturgesetz außer Kraft gesetzt, womit auch die Rechtsgrundlage für die „Geeigneten Algorithmen“ entfallen war. Die Bundesnetzagentur lehnte es ab, stattdessen unverbindliche Empfehlungen zu erarbeiten. Stattdessen veröffentlichte das BSI am 22. Januar 2018 eine Technische Richtlinie (TR-02102-1). Danach soll  $N$  für den Einsatz bis 2022 mindestens zwei Tausend Bit haben, danach mindestens drei Tausend. Außerdem soll  $2^{16} + 1 \leq e < 2^{1824}$  gelten.

Die meisten im Signaturgesetz geregelten Fragen gingen 2017 über in die Zuständigkeit der Europäischen Union; die bereits zu Beginn

des Kapitels erwähnte SOG-IS Crypto Working Group hatte im Mai 2016 ein Dokument mit dem Titel *SOG-IS Crypto Evaluation Scheme – Agreed Cryptographic Mechanisms* veröffentlicht. Für *RSA legacy mechanisms* werden Moduln mit mehr als zwei Tausend Bit gefordert, für *recommended mechanisms* drei Tausend. Der Exponent  $e$  muß in beiden Fällen mehr als sechzehn Bit haben; die Primzahlen  $p$  und  $q$  sollen zufällig erzeugt und gleich lang sein; bei einem  $n$ -Bit Modul muß außerdem  $|p - q| \geq 2^{n/2-100}$  sein.

## §5: Praktische Gesichtspunkte

Wenn  $N = pq$  um die zwei oder drei Tausend Bit hat, wird im allgemeinen auch das kgV von  $p - 1$  und  $q - 1$  nicht viel kleiner sein, und bei der obigen Vorgehensweise können wir erwarten, daß dann auch zumindest der private Exponent  $d$  ebenfalls in dieser Größenordnung ist. Damit ist klar, daß  $x^d \bmod N$  nicht einfach durch sukzessive Multiplikation mit  $x$  berechnet werden kann: Unser Sicherheitsstandard beruht schließlich auf der Annahme, daß niemand  $2^{128}$  oder gar noch mehr Rechenoperationen ausführen kann. Hinzu kommt, daß  $x^d$  so groß ist, daß kein heutiger Computer diese Zahl speichern könnte. Um die Länge der Zwischenergebnisse in Grenzen zu halten, muß nach jeder Multiplikation sofort modulo  $N$  reduziert werden.

Auch das Problem der vielen Multiplikationen läßt sich in den Griff bekommen: Um beispielsweise  $x^{32}$  zu berechnen brauchen wir keine 31 Multiplikationen, sondern erhalten das Ergebnis über die Formel

$$x^{32} = \left( \left( \left( \left( (x^2)^2 \right)^2 \right)^2 \right)^2 \right)^2$$

mit nur fünf Multiplikationen (genauer: Quadrierungen).

Entsprechend können wir für jede gerade Zahl  $n = 2m$  die Potenz  $x^n$  als Quadrat von  $x^m$  berechnen. Für einen ungeraden Exponenten  $e$  ist  $e - 1$  gerade, wenn wir also  $x^e$  als Produkt von  $x$  und  $x^{e-1}$  berechnen, können wir zumindest im nächsten Schritt wieder die Formel für gerade Exponenten verwenden. Somit reichen pro Binärziffer des Exponenten ein bis zwei Multiplikationen; der Aufwand wächst also nur

proportional zur Stellenzahl von  $e$ . Für den ebenfalls recht populären Verschlüsselungsexponenten  $e = 2^{16} + 1 = 65537$  beispielsweise braucht man nur 17 Multiplikationen, nicht 65536.

Hinreichend große Primzahlen sind eine notwendige Bedingung für die Sicherheit des RSA-Verfahrens, aber keine hinreichende: Der Gegner, vor dem eine Nachricht geschützt werden soll, ist schließlich frei in der Wahl seiner Mittel und kann auch anders als mit einem Faktorisierungsversuch angreifen. Soweit entsprechende Strategien bekannt sind, muß man sich daher auch dagegen schützen.

In der bislang dargestellten sogenannten „Lehrbuchversion“ bietet RSA eine ganze Reihe alternativer Angriffsmöglichkeiten, beispielsweise bei kleinen öffentlichen oder privaten Schlüsseln.

Da der Aufwand einer Exponentiation proportional zur Stellenzahl des Exponenten wächst, bevorzugen viele Anwender kleine Exponenten; insbesondere wird in der Praxis sehr oft der kleinstmögliche öffentliche Exponent  $e = 3$  verwendet (was natürlich voraussetzt, daß die verwendeten Primzahlen kongruent zwei modulo drei sind), so daß zumindest die Verschlüsselung recht schnell geht. Außerdem läßt sich dann der private Exponent  $d$  sehr einfach bestimmen: Wie wir gesehen haben, gibt es Zahlen  $d < \lambda = \text{kgV}(p-1, q-1)$  und  $k < e$ , so daß  $de - k\lambda = 1$  ist. Im Falle  $e = 3$  kommen nur  $k = 1$  und  $k = 2$  in Frage, und wir müssen nur testen, welche der beiden Zahlen  $(1 + \varphi(N))/3$  und  $(1 + 2\varphi(N))/3$  ganzzahlig ist.

Bei so vielen Vorteilen muß es auch Nachteile geben, darunter einen ganz offensichtlichen: Ist nämlich die Nachricht  $x$  kleiner als die dritte Wurzel aus  $N$ , so ist  $x^3 < N$ , d.h. der Chiffretext ist einfach  $x^3$ . Die Kubikwurzel aus dieser ganzen Zahl kann natürlich leicht gezogen werden.

Kurze Nachrichten sind allerdings auch für größere Exponenten  $e$  problematisch. Ein Grund liegt darin, daß die Verschlüsselungsfunktion mit der Multiplikation verträglich ist: Auch im Ring  $\mathbb{Z}/N$  ist  $(yz)^e = y^e \cdot z^e$ .

Nehmen wir an, unsere Nachricht  $x$  habe höchstens  $2\ell$  Bit, sei also kleiner als  $2^{2\ell}$ . Dann gibt es eine nicht vernachlässigbare Chance, daß sich

$x = yz$  als Produkt zweier Zahlen darstellen läßt, die beide kleiner als  $2^\ell$  (oder als eine etwas größere Schranke  $M$ ) sind. Für viele Nachrichten wird das zwar nicht der Fall sein, aber wir haben schon ein ernstes Sicherheitsproblem, wenn ein Angriff nur für einen nicht vernachlässigbaren Bruchteil aller Nachrichten funktioniert, und das ist hier definitiv der Fall.

Für den Angriff berechnen wir für alle Zahlen  $y$  zwischen Null bis  $M$  deren Verschlüsselung  $y^e \bmod N$  und notieren die Ergebnisse in einer Tabelle. Um nun vom Chiffretext  $c = x^e \bmod N$  auf  $x$  zurückzuschließen, berechnen wir für jedes dieser Ergebnisse in  $\mathbb{Z}/N$  den Quotienten  $c/y^e$ . (Falls sich dieser Quotient nicht bilden läßt, ist  $y^e$  nicht teilerfremd zu  $N = pq$ , und wir erhalten sogar eine Faktorisierung von  $N$ ; das ist aber für gut gewählte, große  $N$  extrem unwahrscheinlich.) Falls einer dieser Quotienten als Eintrag  $z^e \bmod N$  in unserer Tabelle auftaucht, haben wir eine Relation der Form  $c \equiv y^e \cdot z^e = (yz)^e \bmod N$  gefunden, und damit kennen wir  $x = yz$ . Um diese Attacke zu verhindern, muß bei  $n$ -Bit-Sicherheit  $\ell \geq n$  sein, die übermittelten Nachrichten müssen also mindestens  $2n$  Bit lang sein. Bei der am Ende von §2 skizzierten Methode des Auffüllens mit Zufallsbits ist das automatisch gewährleistet.

Falls eine Nachricht an mehrere Empfänger geschickt wird, müssen – vor allem bei kleinen Verschlüsselungsexponenten wie  $e = 3$  – die Zufallsbits für jeden Empfänger neu erzeugt werden, denn wenn jedes Mal derselbe Block  $x$  verschlüsselt wird und dabei – wie dies häufig in der Praxis der Fall ist – stets mit drei potenziert wird, kennt ein Gegner anschließend  $x^3 \bmod N_i$  für die Moduln  $N_i$  der sämtlichen Empfänger, kann also nach dem chinesischen Restesatz  $x^3$  modulo dem Produkt der  $N_i$  berechnen, und da dieses Produkt bei mindestens drei Empfängern größer ist als  $x^3$ , kennt er  $x^3$  und damit auch  $x$ .

Bei der Wahl der Schlüsseldaten geht man stets aus vom öffentlichen Exponenten  $e$  und berechnet dann dazu nach dem EUKLIDischen Algorithmus den privaten Exponenten  $d$ . Dadurch ist praktisch sichergestellt, daß dieser in der Größenordnung von  $N$  liegen wird, und das muß auch so sein: Im Kapitel über Kettenbrüche werden wir sehen, daß  $N$  leicht faktorisiert werden kann, wenn  $d$  zu klein ist.

## §6: Verfahren mit diskreten Logarithmen

Kurz nach der Veröffentlichung des RSA-Algorithmus fanden auch DIFFIE und HELLMAN ein Verfahren, das im Gegensatz zu RSA sogar ganz ohne vorvereinbarte Schlüssel auskommt: Zwei Personen vereinbaren über eine unsichere Leitung einen Schlüssel, den anschließend nur sie kennen.

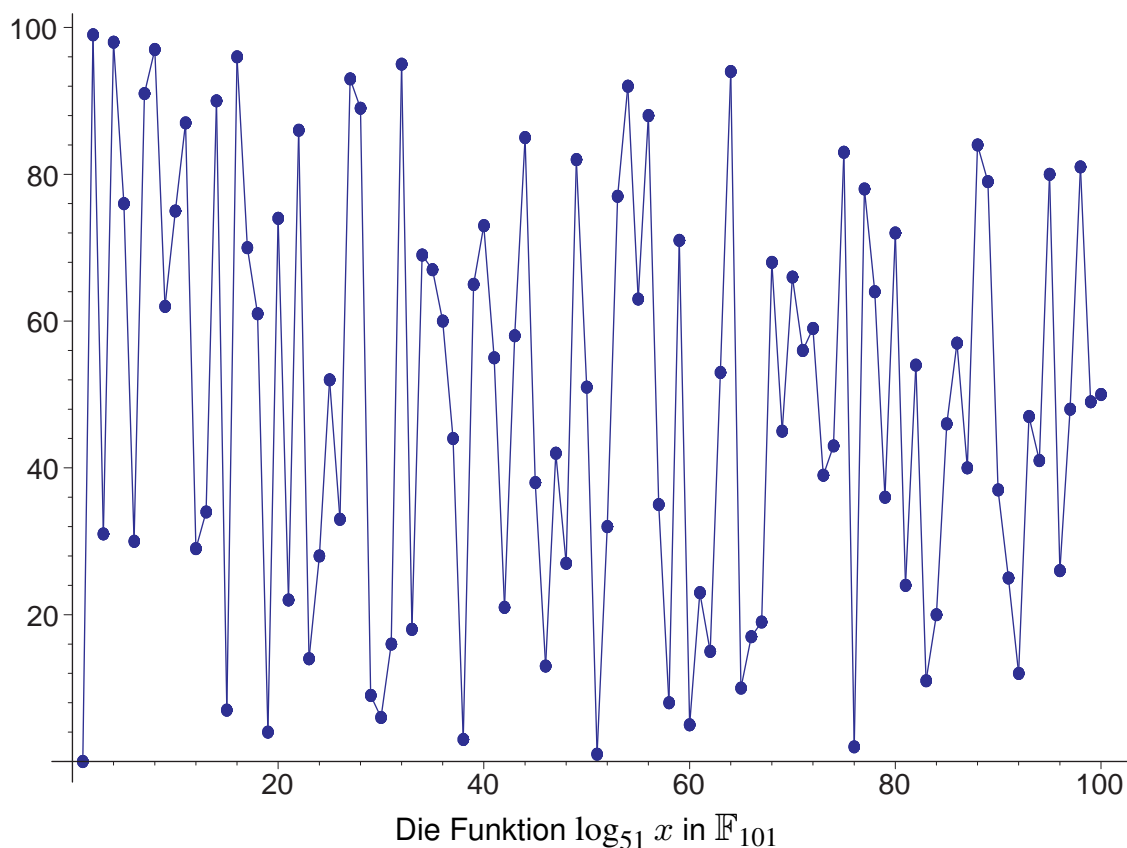
Ausgangspunkt ist wieder das Potenzieren im Körper  $\mathbb{F}_p$ ; hier betrachten wir aber die Exponentialfunktion  $x \mapsto a^x$  zu einer geeigneten Basis  $a$ . Ihre Umkehrfunktion bezeichnet man als *Index* oder *diskreten Logarithmus* zur Basis  $a$ :

$$y = a^x \implies x = \log_a y.$$

Trotz dieser formalen Übereinstimmung gibt es es allerdings große Unterschiede zwischen reellen Logarithmen und ihren Analoga in endlichen Körpern: Während reelle Logarithmen sanft ansteigende stetige Funktionen sind, die man leicht mit beliebig guter Genauigkeit annähern kann, sieht der diskrete Logarithmus typischerweise so aus, wie es in der Abbildung zu sehen ist. Auch ist im Reellen der Logarithmus zur Basis  $a > 1$  für jede positive Zahl definiert; in endlichen Körpern ist es viel schwerer zu entscheiden, ob ein bestimmter Logarithmus existiert: Modulo sieben etwa sind 2, 4 und 1 die einzigen Zweierpotenzen, so daß 3, 5 und 6 keine Zweierlogarithmen haben. Wie wir am Ende von Kapitel I gesehen haben, ist aber die multiplikative Gruppe eines Körpers zyklisch, so daß es stets Elemente  $a$  gibt, für die  $a^x$  jeden Wert außer der Null annimmt, die sogenannten primitiven Wurzeln. In  $\mathbb{F}_7$  wären dies etwa drei und fünf.

Die Berechnung der Potenzfunktion durch sukzessives Quadrieren und Multiplizieren ist auch in endlichen Körpern einfach, für ihre Umkehrfunktion, den diskreten Logarithmus gibt es aber derzeit nur deutlich schlechtere Verfahren. Die derzeit besten Verfahren zur Berechnung von diskreten Logarithmen in Körpern mit  $N$  Elementen erfordern etwa denselben Aufwand wie die Faktorisierung eines RSA-Moduls der Größenordnung  $N$ . Diese Diskrepanz zwischen Potenzfunktion und Logarithmen kann kryptologisch ausgenutzt werden.





Als Körper verwendet man entweder Körper von Zweipotenzordnung, die wir in dieser Vorlesung nicht betrachten werden, oder Körper von Primzahlordnung. Da es für viele interessante Körper von Zweipotenzordnung bereits Chips gibt, die dort diskrete Logarithmen berechnen, dürften Körper von Primzahlordnung bei ungefähr gleicher Elementanzahl wohl etwas sicherer sein: Es gibt einfach viel mehr Primzahlen als Zweierpotenzen, und jeder Fall erfordert einen neuen Hardwareentwurf. Falls man die Primzahlen hinreichend häufig wechselt, dürfte sich dieser Aufwand für kaum einen Gegner lohnen. Außerdem ist das Rechnen modulo einer Primzahl einfacher als das Rechnen in einem Körper von Zweipotenzordnung.

Beim DIFFIE-HELLMAN-Verfahren, dem ältesten auf der Grundlage diskreter Logarithmen, geht es wie gesagt darum, daß zwei Teilnehmer, die weder über gemeinsame Schlüsselinformation noch über eine sichere Leitung verfügen, einen Schlüssel vereinbaren wollen.

Dazu einigen sie sich zunächst (über die unsichere Leitung) auf eine

Primzahl  $p$  und eine natürliche Zahl  $a$  derart, daß die Potenzfunktion  $x \mapsto a^x$  möglichst viele Werte annimmt. Als nächstes wählt Teilnehmer A eine Zufallszahl  $x < p$  und B entsprechend  $y < p$ ; A schickt  $u = a^x \bmod p$  an B und erhält dafür  $v = a^y \bmod p$ .

Sodann berechnet A die Zahl  $v^x \bmod p = (a^y)^x \bmod p = a^{xy} \bmod p$  und B entsprechend  $u^y \bmod p = (a^x)^y \bmod p = a^{xy} \bmod p$ . Beide haben also auf verschiedene Weise dieselbe Zahl berechnet, die sie nun als Schlüssel in einem klassischen Kryptosystem verwenden können, wobei sie sich wohl meist auf einen Teil der Bits beschränken müssen, da solche Schlüssel typischerweise eine Länge von 128 bis 256 Bit haben, während die Primzahl  $p$  erheblich größer sein muß.

Ein Gegner, der den Datenaustausch abgehört hat, kennt die Zahlen  $p, a, u$  und  $v$ ; um  $a^{xy} \bmod p$  zu finden, muß er den diskreten Logarithmus von  $u$  oder  $v$  berechnen.

Mit den besten heute bekannten Algorithmen ist die möglich, wenn  $p$  eine Primzahl von bis zu etwa 200 Dezimalstellen ist; dies entspricht etwa 665 Bit. Auch in diesem Fall dauert die Berechnung allerdings selbst bei massiver Parallelisierung über das Internet mehrere Monate.

Natürlich gibt es keine Garantie, daß kein Gegner mit einem besseren als den bislang bekannten Verfahren diskrete Logarithmen oder Faktorisierungen auch in weitaus größeren Körpern berechnen kann. Dazu bräuchte er allerdings einen Durchbruch entweder auf der mathematischen oder auf der technischen Seite, für den weit und breit keine Grundlage zu sehen ist.

Trotzdem gibt es einen verhältnismäßig einfachen Angriff auf den Schlüsselaustausch nach DIFFIE und HELLMAN, die sogenannte *man in the middle attack*. Dabei unterbricht der Angreifer die Verbindung zwischen A und B und gibt sich gegenüber A als B aus und umgekehrt. So kann er mit beiden Teilnehmern je einen Schlüssel vereinbaren, und die damit verschlüsselte Kommunikation kann (nur) von ihm gelesen und gegebenenfalls manipuliert werden. In der vorgestellten Form funktioniert das Verfahren also nur, wenn man sicher sein weiß, mit wem man kommuniziert.

## §6: DSA

DSA steht für *Digital Signature Algorithm*, ein Algorithmus der im *Digital Signature Standard* DSS der USA festgelegt ist und neben RSA auch zu den von der Bundesnetzagentur empfohlenen „Geeigneten Algorithmen“ zählt. Seine Sicherheit beruht auf diskreten Logarithmen, allerdings wird das klassische Verfahren dadurch modifiziert, daß die Sicherheit zwar auf dem diskreten Logarithmenproblem in einem großen Körper beruht, die Rechenoperationen bei der Anwendung des Algorithmus aber nur eine deutlich kleinere Untergruppe verwenden.

Für diese kleine Untergruppe wählt man eine Primzahl  $q$  mit einer Länge von mindestens 224 Bit. (Das entspricht der Länge der Hashwerte, die in der Praxis anstelle des Texts unterzeichnet werden.) Zu dieser Primzahl  $q$  sucht man eine Primzahl  $p \equiv 1 \pmod{q}$ , für deren Länge die Bundesnetzagentur mindestens 2048 Bit vorschreibt.

Daß die mit der empfohlenen RSA-Modullänge übereinstimmt, ist kein Zufall: Auch wenn kein direkter Zusammenhang zwischen Faktorisierung und der Berechnung diskreter Logarithmen bekannt ist, hat bislang doch jede neue Idee für einen Faktorisierungsalgorithmus auch zu einem Algorithmus zur Berechnung diskreter Logarithmen geführt, und auch die Laufzeiten dieser Algorithmen sind bei gleicher Zahlenlänge ungefähr gleich.

Als nächstes muß ein Element  $g$  gefunden werden, dessen Potenzen im Körper  $\mathbb{F}_p$  eine Gruppe der Ordnung  $q$  bilden. Das ist einfach: Man starte mit irgendeinem Element  $g_0 \in \mathbb{F}_p \setminus \{0\}$  und berechne seine  $(p-1)/q$ -te Potenz. Falls diese ungleich eins ist, muß sie wegen  $g_0^{p-1} = 1$  die Ordnung  $q$  haben; andernfalls muß ein neues  $g_0$  betrachtet werden.

Die so bestimmten Zahlen  $q, p$  und  $g$  werden veröffentlicht und können auch in einem ganzen Netzwerk global eingesetzt werden. Geheimer Schlüssel jedes Teilnehmers ist eine Zahl  $x$  zwischen eins und  $q-1$ ; der zugehörige öffentliche Schlüssel ist  $y = g^x \pmod{p}$ .

Unterschreiben lassen sich mit diesem Verfahren Nachrichtenblöcke  $m$  mit  $0 \leq m < q$ , insbesondere also 224 Bit lange Hashwerte. Dazu wählt

man für jede Nachricht eine Zufallszahl  $k$  mit  $0 < k < q$  und berechnet

$$r = (g^k \bmod p) \bmod q .$$

Da  $q$  eine Primzahl ist, hat  $k$  ein multiplikatives Inverses modulo  $q$ ; man kann also durch  $k$  dividieren und erhält eine Zahl  $s$ , für die

$$sk \equiv m + xr \pmod{q}$$

ist; die Unterschrift unter die Nachricht  $m$  besteht dann aus den beiden Zahlen  $r$  und  $s$  modulo  $q$ . Sie kann nur berechnet werden von jemandem, der den geheimen Schlüssel  $x$  kennt.

Überprüfen kann die Unterschrift allerdings jeder: Ist  $t$  das multiplikative Inverse zu  $s$  modulo  $q$ , so ist  $k \equiv tsk \equiv tm + xtr \pmod{q}$ , also, da  $g$  die Ordnung  $q$  hat,

$$r \equiv g^k \equiv g^{tm} g^{xtr} \equiv g^{tm} y^{tr} \pmod{p} .$$

In dieser Gleichung sind die linke wie auch die rechte Seite *modulo*  $q$  öffentlich bekannt, die Gleichung kann also modulo  $q$  überprüft werden. Die Unterschrift wird anerkannt, wenn beide Seiten modulo  $q$  gleich sind. Ein Angreifer müßte sich  $x$  aus  $y$  verschaffen, müßte also ein diskretes Logarithmenproblem modulo der großen Primzahl  $p$  lösen.

Die Berechnung diskreter Logarithmen modulo einer Primzahl  $p$  mit  $n$  Bit ist etwa genauso aufwendig, wie die Faktorisierung eines RSA-Moduls mit  $n$  Bit. Die Empfehlungen bezüglich der Größe von  $p$  entsprechen daher denen für die Größe von RSA-Moduln. Die Primzahl  $q$  kann deutlich kleiner gewählt werden; die SOG-IS Empfehlungen etwa geben für *legacy mechanisms* eine Mindestgröße von 200 Bit vor, für *recommended mechanisms* 250.

Tatsächlich richtet sich die Größe von  $q$  nach der des zu unterschreibenden Hashwerts; da dieser für  $n$ -Bit-Sicherheit wegen des Geburtstagsparadoxons mindestens die Länge  $2n$  haben muß und es gängige Verfahren für die Bitlängen 160, 225, 256, 384 und 512 gibt, muß  $q$  für 100 und 112-Bit-Sicherheit mehr als 225 Bit haben, für 120–128-Bit-Sicherheit mehr als 256.

## §7: Ausblick

Dieses kurze Kapitel konnte selbstverständlich keine umfassende Übersicht über die Kryptographie oder auch nur die asymmetrische Kryptographie geben: Auch das RSA-Verfahren kann mit anderen Methoden angegriffen werden als der direkten Faktorisierung des Moduls. Eine dieser Methoden werden wir im Kapitel über Kettenbrüche kennenlernen, und auch sonst werden im weiteren Verlauf der Vorlesungen noch gelegentlich Themen aus der Kryptologie angeschnitten werden.

Mit Ausnahme von Verfahren wie dem sogenannten *one time pad* gibt es für keines der heute benutzten Kryptoverfahren einen Sicherheitsbeweis, nicht einmal in dem Sinn, daß man den Aufwand eines Gegners zum Knacken des Verfahrens in irgendeiner realistischen Weise nach unten abschätzen könnte. Seriöse Kryptographie außerhalb des Höchstsicherheitsbereichs muß sich daher damit begnügen, daß die Verantwortlichen für den Einsatz eines Verfahrens und der Wahl seiner Parameter (wie den Primzahlen bei RSA) darauf achten, auf dem neuesten Stand der Forschung zu bleiben und ihre Wahl so treffen, daß nicht nur die bekannten Angriffsmethoden versagen, sondern daß auch noch ein recht beträchtlicher Sicherheitszuschlag für künftige Entwicklungen und für nicht publizierte Entwicklungen bleibt.

Auf ewige Sicherheit kann man mit Verfahren wie RSA ohnehin nicht hoffen: Als RSA 1977 von MARTIN GARDNER im *Scientific American* vorgestellt wurde, bekam er von RIVEST, SHAMIR und ADLEMAN die 129-stellige Zahl

11438162575788886766923577997614661201021829672124236256256184293\  
5706935245733897830597123563958705058989075147599290026879543541

(seither bekannt als RSA-129) und eine damit verschlüsselte Nachricht, für deren Entschlüsselung die drei einen Preis von hundert Dollar ausgesetzt hatten. Sie schätzten, daß eine solche Entschlüsselung etwa vierzig Quadrillionen ( $4 \cdot 10^{25}$ ) Jahre dauern würde. (Heute sagt RIVEST, daß dies auf einem Rechenfehler beruhte.) Tatsächlich wurde der Modul 1994 faktorisiert in einer gemeinsamen Anstrengung von 600 Freiwilligen, deren Computer immer dann, wenn sie nichts besseres zu tun hatten,

daran arbeiteten. Nach acht Monaten war die Faktorisierung gefunden:  
Die obige Zahl ist gleich

$$490529510847650949147849619903898133417764638493387843990820577 \\ \times 32769132993266709549961988190834461413177642967992942539798288533 .$$

Mit dem Schema  $A = 01$  bis  $Z = 26$  und Zwischenraum gleich 00 ergab sich die Nachricht *The Magic Words are Squeamish Ossifrage*.

Auch bei den heute als sicher geltenden symmetrischen Kryptoverfahren rechnet niemand ernsthaft damit, daß sie noch in hundert Jahren sicher sind: Diese Verfahren werden üblicherweise so gewählt, daß man auf eine Sicherheit für etwa dreißig Jahren hoffen kann – garantieren kann aber auch das niemand.

Falls sich sogenannte *Quantencomputer* realisieren lassen, werden alle heute bekannten Verfahren der Kryptographie mit öffentlichen Schlüsseln, egal ob mit diskreten Logarithmen, RSA oder elliptischen Kurven, unsicher sein. Bisläng können Quantencomputer kaum mit acht Bit rechnen, und nicht alle Experten sind davon überzeugt, daß es je welche geben wird, die mit mehreren Tausend Bit rechnen können.

Wer mehr über Kryptographie wissen will, findet einen ersten Überblick beispielsweise bei

BUCHMANN: Einführung in die Kryptographie, *Springer*, <sup>6</sup>2016

oder natürlich auch im Skriptum zur hier immer wieder angebotenen Kryptologie-Vorlesung.

Mehr über die Geschichte der Kryptographie mit öffentlichen Schlüsseln ist (mathematikfrei) zu finden in

STEVEN LEVY: **crypto**: how the rebels beat the government – saving privacy in the digital age, *Penguin Books*, 2002.

Eine geschichtliche Darstellung der Kryptographie bis etwa zur Zeit des zweiten Weltkriegs mit ausführlicher Darstellung einiger der verwendeten Verfahren und Angriffen findet man bei

DAVID KAHN: The codebreakers, *Scribner*, <sup>2</sup>1996.

## Kapitel 3

### Primzahlen

Wie wir aus dem ersten Kapitel wissen, sind Primzahlen die Grundbausteine für die multiplikative Struktur der ganzen Zahlen, und aus Kapitel zwei wissen wir, daß sie auch wichtige Anwendungen außerhalb der Zahlentheorie haben. Es lohnt sich also auf jeden Fall, sie etwas genauer zu untersuchen.

#### § 1: Die Verteilung der Primzahlen

Als erstes stellt sich die Frage, wie viele Primzahlen es gibt. Die Antwort finden wir schon in Satz 20 des neunten Buchs von EUKLIDS Elementen:

*Es gibt mehr Primzahlen als jede vorgelegte Anzahl von Primzahlen.*

In heutiger Sprache ausgedrückt: Zu jeder endlichen Menge  $M$  von Primzahlen (die auch leer sein darf) gibt es eine Primzahl, die nicht in  $M$  liegt.

EUKLIDS Argument dafür dürfte immer noch der einfachste Beweis für die Unendlichkeit der Menge aller Primzahlen sein: Er bildet das Produkt  $P$  über alle Primzahlen aus  $M$ ; im Falle der leeren Menge ist dieses Produkt gleich eins. Zum Ergebnis addiert er eins. Die Zahl  $P + 1$  kann durch keine der Primzahlen aus der Menge teilbar sein, denn diese sind Teiler von  $P$  und wären daher als Teiler von  $P + 1$  auch Teiler der Eins. Falls  $P + 1$  prim ist, haben wir eine Primzahl gefunden, die schon wegen ihrer Größe nicht in der vorgegebenen Menge liegen kann. Andernfalls sei  $p$  der kleinste von eins und  $P + 1$  verschiedene Teiler von  $P + 1$ . Dieser ist eine Primzahl, die von allen Elementen der Menge verschieden sein muß, da diese im Gegensatz zu  $p$  keine Teiler von  $P + 1$  sind.

Um nicht ganz auf dem Stand von vor rund zweieinhalb Jahrtausenden stehen zu bleiben, wollen wir uns noch einen zweiten, auf EULER zurückgehenden Beweis ansehen.

Dazu betrachten wir für eine reelle Zahl  $s > 1$  die unendliche Reihe

$$\zeta(s) = \sum_{n=1}^{\infty} \frac{1}{n^s}.$$

Als erstes müssen wir uns überlegen, daß diese Reihe konvergiert. Da alle Summanden positiv sind, müssen wir dafür nur zeigen, daß es eine gemeinsame obere Schranke für alle Teilsummen gibt. Da die Funktion  $x \mapsto 1/x^s$  für  $x > 0$  monoton fallend ist, haben wir für  $n-1 \leq x \leq n$  die Abschätzung  $1/n^s \leq 1/x^s$ , d.h.

$$\begin{aligned} \sum_{n=1}^N \frac{1}{n^s} &= 1 + \sum_{n=2}^N \frac{1}{n^s} \leq 1 + \int_1^N \frac{dx}{x^s} \\ &< 1 + \int_1^{\infty} \frac{dx}{x^s} = 1 + \frac{1}{s-1} = \frac{s}{s-1}. \end{aligned}$$

Somit ist  $\zeta(s)$  für alle  $s > 1$  wohldefiniert.

Einen Zusammenhang mit Primzahlen liefert die Darstellung von  $\zeta(s)$  als sogenanntes EULER-Produkt im folgenden

**Satz:** a) Für  $s > 1$  ist  $\zeta(s) = \prod_{p \text{ prim}} \frac{1}{1 - \frac{1}{p^s}}$ .

b) Für alle  $N \in \mathbb{N}$  und alle reellen  $s > 0$  ist  $\sum_{n=1}^N \frac{1}{n^s} \leq \prod_{\substack{p \leq N \\ p \text{ prim}}} \frac{1}{1 - \frac{1}{p^s}}$ .

*Beweis:* Wir beginnen mit b). Für  $N = 1$  steht hier die triviale Formel  $1 \leq 1$ ; sei also  $N \geq 2$ , und seien  $p_1, \dots, p_r$  die sämtlichen Primzahlen kleiner oder gleich  $N$ . Nach der Summenformel für die geometrische Reihe ist

$$\frac{1}{1 - \frac{1}{p_k^s}} = \sum_{\ell=0}^{\infty} \frac{1}{p_k^{\ell s}},$$



und das Produkt der rechtsstehenden Reihen über  $k = 1$  bis  $r$  ist wegen der Eindeutigkeit der Primzerlegung die Summe über alle jene  $1/n^s$ , für die  $n$  keinen Primteiler größer  $N$  hat. Darunter sind insbesondere alle  $n \leq N$ , womit  $b)$  bewiesen wäre.

Die Differenz zwischen  $\zeta(s)$  und dem Produkt auf der rechten Seite von  $b)$  ist gleich der Summe über alle  $1/n^s$ , für die  $n$  mindestens einen Primteiler größer  $N$  hat. Diese Summe ist natürlich höchstens gleich der Summe aller  $1/n^s$  mit  $n > N$ , und die geht wegen der Konvergenz von  $\zeta(s)$  gegen null für  $N \rightarrow \infty$ . Damit ist auch  $a)$  bewiesen. ■

Auch daraus folgt, daß es unendlich viele Primzahlen gibt: Gäbe es nämlich nur endlich viele, so stünde auf der rechten Seite von  $b)$  für jedes hinreichend große  $N$  das Produkt über die *sämtlichen* Primzahlen. Da es nur endlich viele Faktoren hat, wäre es auch für  $s = 1$  endlich, und damit müßte

$$\sum_{n=1}^{\infty} \frac{1}{n}$$

kleiner oder gleich dieser Zahl sein, im Widerspruch zur Divergenz der harmonischen Reihe. ■

Verglichen mit dem Beweis aus EUKLIDS Elementen ist EULERS Methode erheblich komplizierter. Um trotzdem ihre Existenzberechtigung zu haben, sollte sie uns daher auch mehr Informationen liefern. In welchem Maße sie dies tatsächlich leistet, geht wahrscheinlich sogar noch deutlich über alles hinaus, was EULER seinerzeit träumen konnte.

Zunächst einmal können wir Teil  $b)$  für  $s = 1$  zu einer quantitativen Abschätzung bezüglich der Anzahl  $\pi(N)$  der Primzahlen kleiner oder gleich  $N$  umformulieren: Wie oben im Konvergenzbeweis für  $\zeta(s)$  können wir aus der Monotonie der Funktion  $x \mapsto 1/x$  folgern, daß für alle  $N \in \mathbb{N}$  gilt

$$\log(N+1) = \int_1^{N+1} \frac{dx}{x} < \sum_{n=1}^N \frac{1}{n} < 1 + \int_1^N \frac{dx}{x} = 1 + \log N.$$

Zur Abschätzung der linken Seite beachten wir einfach, daß der Faktor  $1/(1 - 1/p)$  für  $p = 2$  gleich zwei ist, ansonsten aber kleiner. Somit ist

$$\log(N + 1) < \sum_{n=1}^N \frac{1}{n} \leq \prod_{\substack{p \leq N \\ p \text{ prim}}} \frac{1}{1 - \frac{1}{p}} \leq 2^{\pi(N)}$$

und damit

$$\pi(N) \geq \frac{\log \log(N + 1)}{\log 2}.$$

Wie wir bald sehen werden, ist das allerdings eine sehr schwache Abschätzung.

EULERS Methode erlaubt uns auch, die Dichte der Primzahlen zu vergleichen mit der Dichte beispielsweise der Quadratzahlen: Wie wir oben gesehen haben, konvergiert  $\zeta(s)$  für alle  $s > 1$ , insbesondere also konvergiert die Summe  $\zeta(2)$  der inversen Quadratzahlen. EULER konnte mit seiner Methode zeigen, daß die Summe der inversen Primzahlen *divergiert*, so daß die Primzahlen zumindest in diesem Sinne dichter liegen als die Quadratzahlen und alle anderen Potenzen mit (reellem) Exponenten  $x > 1$ .

Zum Beweis fehlt uns nur noch eine Analysis I Übungsaufgabe: Wir wollen uns überlegen, daß für alle  $0 \leq x \leq \frac{1}{2}$  gilt  $(1 - x) \geq 4^{-x}$ . An den Intervallenden stimmen beide Funktionen überein, und  $1 - x$  ist eine lineare Funktion. Es reicht daher, wenn wir zeigen, daß  $4^{-x}$  eine konvexe Funktion ist, daß also ihre zweite Ableitung überall im Intervall positiv ist. Das ist aber klar, denn die ist einfach  $\log(4)^2 \cdot 4^{-x}$ . Für jede Primzahl  $p$  ist daher

$$1 - \frac{1}{p} \geq 4^{-1/p} \quad \text{und} \quad \frac{1}{1 - \frac{1}{p}} \leq 4^{1/p}.$$

Zusammen mit der vorigen Abschätzung folgt

$$\log(N + 1) < \prod_{\substack{p \leq N \\ p \text{ prim}}} \frac{1}{1 - \frac{1}{p}} \leq \prod_{\substack{p \leq N \\ p \text{ prim}}} 4^{1/p} = 4^{\sum \frac{1}{p}},$$

wobei die Summe im Exponenten über alle Primzahlen  $p \leq N$  geht. Da  $\log(N + 1)$  für  $N \rightarrow \infty$  gegen unendlich geht, muß auch die Summe der inversen Primzahlen divergieren.

Mit diesen Bemerkungen fängt allerdings die Nützlichkeit der Funktion  $\zeta(s)$  für das Verständnis der Funktion  $\pi(N)$  gerade erst an: Ein Jahrhundert nach EULER erkannte RIEMANN, daß die Funktion  $\zeta(s)$  ihre wahre Nützlichkeit für das Studium von  $\pi(N)$  erst zeigt, wenn man sie auch für komplexe Argumente  $s$  betrachtet. Jeder, der sich ein bißchen mit Funktionen einer komplexer Veränderlichen auskennt, kann leicht zeigen, daß  $\zeta(s)$  auch für komplexe Zahlen mit Realteil größer ein konvergiert: Der Imaginärteil des Exponenten führt schließlich nur zu einem Faktor vom Betrag eins.



GEORG FRIEDRICH BERNHARD RIEMANN (1826-1866) war Sohn eines lutherischen Pastors und schrieb sich 1846 auf Anraten seines Vaters an der Universität Göttingen für das Studium der Theologie ein. Schon bald wechselte an die Philosophische Fakultät, um dort unter anderem bei GAUSS Mathematikvorlesungen zu hören. Nach Promotion 1851 und Habilitation 1854 erhielt er dort 1857 einen Lehrstuhl. Trotz seines frühen Todes initiierte er grundlegende auch noch heute fundamentale Entwicklungen in der Geometrie, der Zahlentheorie und über abelsche Funktionen. Wie sein Nachlaß zeigte, stützte er seine 1859 aufgestellte Vermutung über die Nullstellen der  $\zeta$ -Funktion auf umfangreiche Rechnungen.

RIEMANNNS wesentliche Erkenntnis war, daß sich  $\zeta(s)$  fortsetzen läßt zu einer analytischen Funktion auf der gesamten Menge der komplexen Zahlen mit Ausnahme der Eins (wo die  $\zeta$ -Funktion wegen der Divergenz der harmonischen Reihe keinen endlichen Wert haben kann).

Für Leser, die nicht mit dem Konzept der analytischen Fortsetzung vertraut sind, möchte ich ausdrücklich darauf hinweisen, daß dies selbstverständlich nicht bedeutet, daß die definierende Summe der  $\zeta$ -Funktion für reelle Zahlen kleiner eins oder komplexe Zahlen mit Realteil kleiner oder gleich eins konvergiert: Analytische Fortsetzung besteht darin, daß eine differenzierbare Funktion (die im Komplexen automatisch beliebig oft differenzierbar ist und um jeden Punkt in eine TAYLOR-Reihe entwickelt werden kann) via TAYLOR-Reihen über ihren eigentlichen Definitionsbereich hinweg ausgedehnt wird. Man kann beispielsweise zeigen, daß  $\zeta(-1) = -\frac{1}{12}$  ist. Setzt man  $s = -1$  in die für

$s > 1$  gültige Reihe ein, erhält man die Summe aller natürlicher Zahlen, die selbstverständlich nicht gleich  $-\frac{1}{12}$  ist, sondern divergiert. Entsprechend hat  $\zeta(s)$  Nullstellen bei allen geraden negativen Zahlen, obwohl auch hier die entsprechenden Reihen divergieren. Diese Nullstellen bezeichnet man als die sogenannten *trivialen* Nullstellen der  $\zeta$ -Funktion, da sie sich sofort aus einer bei der Konstruktion der analytischen Fortsetzung zu beweisenden Funktionalgleichung ablesen lassen. Für die Primzahlverteilung spielen vor allem die übrigen, die sogenannten nicht-trivialen Nullstellen, eine große Rolle.

Wie wir gerade gesehen haben, liegen die Primzahlen zumindest in einem gewissen Sinne dichter als die Quadratzahlen. Zur Einstimmung auf das Problem der Primzahlverteilung wollen wir uns kurz mit der (deutlich einfacheren) Verteilung der Quadratzahlen beschäftigen.

Die Folge der Abstände zwischen zwei aufeinanderfolgenden Quadratzahlen ist einfach die Folge der ungeraden Zahlen, denn

$$(n + 1)^2 - n^2 = 2n + 1 .$$

Zwei aufeinanderfolgende Quadratzahlen  $Q < Q'$  haben daher die Differenz  $Q' - Q = 2\sqrt{Q} + 1$ .

Bei den Primzahlen ist die Situation leider sehr viel unübersichtlicher: EULER meinte sogar, die Verteilung der Primzahlen sei ein Geheimnis, das der menschliche Verstand nie erfassen werde. Der kleinstmögliche Abstand zwischen zwei verschiedenen Primzahlen ist offensichtlich eins, der Abstand zwischen zwei und drei. Er kommt nur an dieser einen Stelle vor, denn mit Ausnahme der Zwei sind alle Primzahlen ungerade.

Der Abstand zwei ist schon deutlich häufiger: Zwei ist beispielsweise der Abstand zwischen drei und fünf, aber auch der zwischen den Primzahlen  $10^{100} + 35737$  und  $10^{100} + 35739$ . Seit langer Zeit wird vermutet, daß es unendlich viele solcher *Primzahlzwillinge* gibt; experimentelle Untersuchungen deuten sogar darauf hin, daß ihre Dichte für Zahlen der Größenordnung  $n$  bei ungefähr  $1 : (\log n)^2$  liegen sollte, aber bislang konnte noch niemand auch nur beweisen, daß es unendlich viele gibt. Das derzeit größte bekannte Paar von Primzahlzwillingen ist

$2\,996\,863\,034\,895 \cdot 2^{1\,290\,000} \pm 1$ ; die beiden Primzahlen haben jeweils 388 342 Dezimalstellen.

Eine obere Grenze für den Abstand zwischen zwei aufeinanderfolgenden Primzahlen gibt es genauso wenig wie bei den Quadratzahlen: Ist  $n \geq 2$  und  $2 \leq i \leq n$ , so ist die Zahl  $n! + i$  durch  $i$  teilbar und somit keine Primzahl. Der Abstand zwischen der größten Primzahl kleiner oder gleich  $n! + 1$  und ihrem Nachfolger ist somit mindestens  $n - 1$ .

Da die Summe der Kehrwerte der Primzahlen genau wie die harmonische Reihe divergiert, können wir uns als nächstes fragen, ob nicht trotz aller Schwankungen der *mittlere* Abstand zweier Primzahlen (asymptotisch gesehen) einen festen Wert  $r$  hat. Dann wäre die  $n$ -te Primzahl ungefähr in der Größenordnung von  $n/r$ , und die Primzahlen kleiner oder gleich  $rm$  für  $m \in \mathbb{N}$  lägen ungefähr bei  $r, 2r, \dots, mr$ ; ihr Produkt hätte also die Größenordnung von  $m!r^m$ . Tatsächlich ist es aber nach dem folgenden Lemma viel kleiner:

**Lemma:** Für jede reelle Zahl  $x \geq 2$  ist

$$\prod_{\substack{p \leq x \\ p \text{ prim}}} p < 4^x.$$

*Beweis:* Wir wollen uns zunächst überlegen, daß es genügt, dieses Lemma für alle ungeraden Zahlen  $n \geq 3$  zu beweisen: Ist  $2 \leq x < 3$ , so ist zwei die einzige Primzahl kleiner oder gleich  $x$ , und  $4^x \geq 16$ , die Behauptung ist also richtig. Für  $x \geq 3$  sei  $n$  die größte ungerade Zahl kleiner oder gleich  $x$ . Da  $n + 1$  gerade und damit zusammengesetzt ist, muß jede Primzahl kleiner oder gleich  $x$  kleiner oder gleich  $n$  sein; falls die Behauptung für  $n$  gilt, ist also

$$\prod_{\substack{p \leq x \\ p \text{ prim}}} p = \prod_{\substack{p \leq n \\ p \text{ prim}}} p < 4^n \leq 4^x,$$

wie gewünscht.

Für ungerade Zahlen beweisen wir das Lemma durch vollständige Induktion. Für  $n = 3$  haben wir links das Produkt  $2 \cdot 3 = 6$ , und rechts steht  $4^3 = 64$ , die Ungleichung ist also erfüllt.

Für  $n \geq 5$  sei  $k = \frac{1}{2}(n \pm 1)$ , wobei das Zeichen so gewählt wird, daß wir eine ungerade Zahl  $k$  erhalten. Für  $n = 5$  ist  $k = \frac{1}{2}(5 + 1) = 3$ , für  $n = 7$  ist  $k = \frac{1}{2}(7 - 1) = 3$ , und offensichtlich ist  $k \geq 3$  für alle ungeraden Zahlen  $n \geq 5$ . Wegen  $n = 2k \pm 1$  ist

$$n - k = 2k \pm 1 - k = k \pm 1 \leq k + 1;$$

jede Primzahl  $p \leq n - k$  muß daher kleiner oder gleich  $k$  sein. Im Binomialkoeffizienten

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

sind daher alle Primzahlen  $p$  mit  $k < p \leq n$  zwar Teiler des Zählers, nicht aber des Nenners, und damit Teiler von  $\binom{n}{k}$ . Somit ist

$$\prod_{\substack{k < p \leq n \\ p \text{ prim}}} p \leq \binom{n}{k} = \binom{n}{n-k}.$$

Weiter ist

$$2 \binom{n}{k} = \binom{n}{k} + \binom{n}{n-k} < \sum_{i=0}^n \binom{n}{i} = (1+1)^n = 2^n;$$

somit ist  $\binom{n}{k} < 2^{n-1}$ . Das Produkt aller Primzahlen bis höchstens  $k$  ist nach Induktionsannahme kleiner als  $4^k$ , also ist

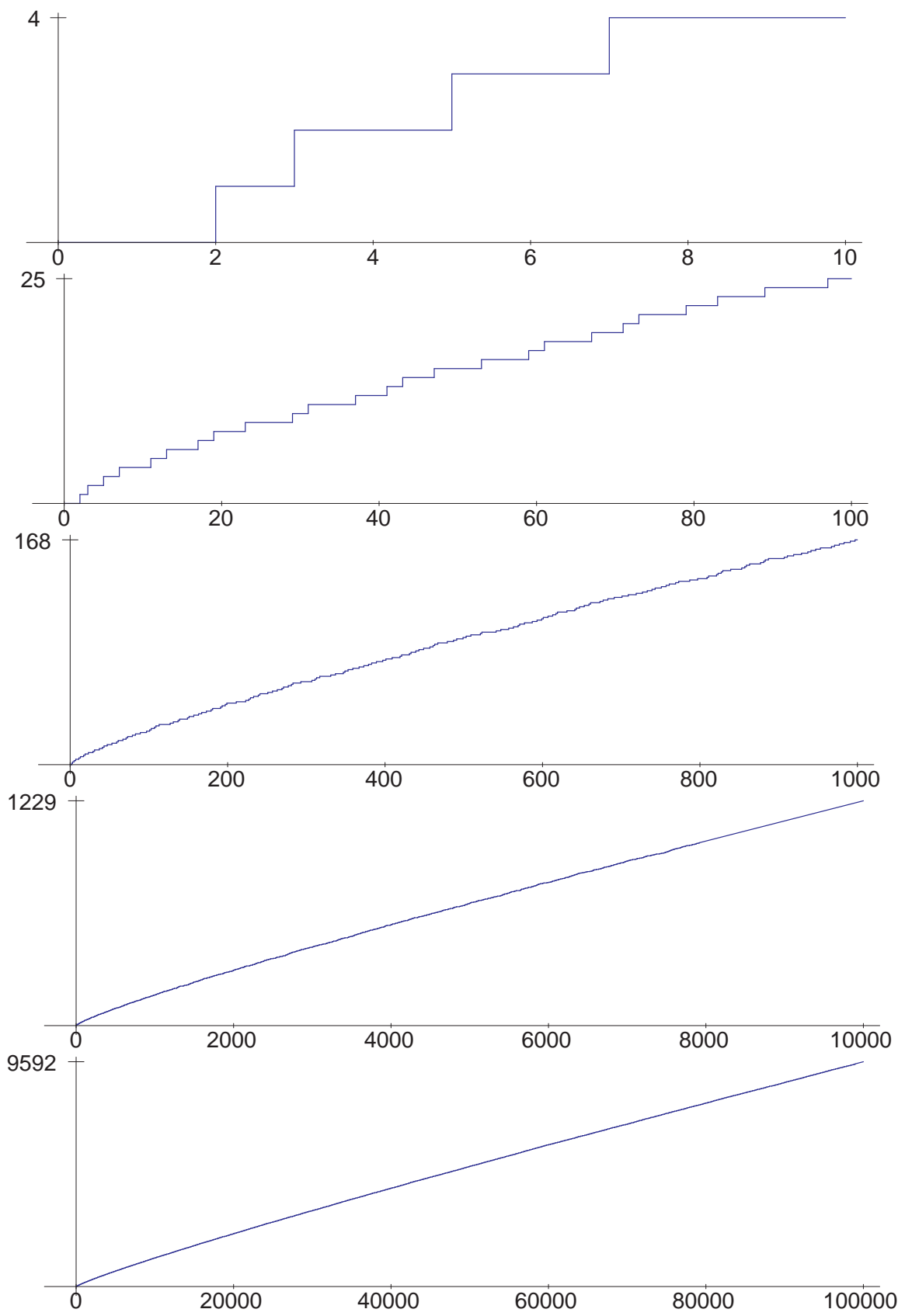
$$\prod_{\substack{p \leq n \\ p \text{ prim}}} p = \prod_{\substack{p \leq k \\ p \text{ prim}}} p \cdot \prod_{\substack{k < p \leq n \\ p \text{ prim}}} p < 4^k \cdot 2^{n-1} = 2^{n \pm 1} \cdot 2^{n-1} \leq 2^{2n} = 4^n,$$

wie behauptet. ■

Um einen ersten Eindruck von der Verteilung der Primzahlen zu bekommen, betrachten wir den Graphen der Funktion

$$\pi: \begin{cases} \mathbb{R}_{>0} \rightarrow \mathbb{N}_0 \\ x \mapsto \text{Anzahl der Primzahlen} \leq x \end{cases}.$$

Die Abbildungen auf der nächsten Seite zeigen ihn für die Intervalle von null bis  $10^i$  für  $i = 1, \dots, 5$ . Wie man sieht, werden die Graphen immer glatter, und bei den beiden letzten Bildern könnte man glauben, es handle



sich um den Graphen einer differenzierbaren Funktion; daher auch die Schreibweise  $\pi(x)$  statt – wie bisher –  $\pi(n)$ .

Auf den ersten Blick sieht diese Funktion fast linear aus; sieht man sich allerdings die Zahlenwerte genauer an, so sieht man schnell, daß  $\pi(x)$  etwas langsamer wächst als eine lineare Funktion; die Funktion  $x/\log x$  ist eine deutlich bessere Approximation. In der Tat können wir auch mit unseren sehr elementaren Mitteln eine entsprechende Aussage beweisen; sie ist eine vereinfachte Version eines Resultats von PAFNUTI L'VOVIČ ČEBYŠEV.



PAFNUTI LWOWITSCH TSCHEBYSCHJEFF (1821–1894) ist die in Deutschland übliche Transkription von Пафну-ты Львович Чебышев; im Englischen schreibt man heute meist PAFNUTY LVOVICH CHEBYSHEV. Er studierte Mathematik in Moskau und publizierte bereits während seines Studiums in deutschen und französischen Fachzeitschriften. 1847 bekam er eine Stelle an der Universität von St. Petersburg, wo er bis 1882 lehrte, ab 1860 als Professor. Seine Ergebnisse spielen noch heute eine große Rolle in der Wahrscheinlichkeitstheorie, der Zahlentheorie (insbesondere Primzahlverteilung), der Approximationstheorie und in anderen Gebieten der Mathematik.

**Satz:** Es gibt Konstanten  $c_1, c_2 > 0$ , so daß gilt:

$$c_1 \frac{x}{\log x} < \pi(x) < c_2 \frac{x}{\log x}.$$

*Beweis:* Ausgangspunkt für erste Abschätzungen von  $\pi(n)$  sind die Primzahlpotenzen, die in  $n!$  stecken. Von den Zahlen zwischen eins und  $n$  ist jede  $p$ -te durch die Primzahl  $p$  teilbar; das sind  $\left[\frac{n}{p}\right]$  Stück. Die Anzahl der sogar durch  $p^2$  teilbaren ist  $\left[\frac{n}{p^2}\right]$ , und so weiter. Die Anzahl der Faktoren  $p$ , die in  $n!$  stecken, ist daher die Summe der Zahlen  $\left[\frac{n}{p^k}\right]$ , wobei  $k$  die natürlichen Zahlen durchläuft. (Für alle  $k$  mit  $p^k > n$  erhält man natürlich einen Summanden Null.) Somit ist

$$n! = \prod_{p \leq n} p^{e_p} \quad \text{mit} \quad e_p = \sum_{k \geq 1} \left[ \frac{n}{p^k} \right],$$



wobei Produkte (und Summen) mit einem Laufindex  $p$  stets nur über Primzahlen gehen sollen. Der Binomialkoeffizient  $\binom{2n}{n} = \frac{(2n)!}{(n!)^2}$  enthält eine Primzahl  $p$  daher zur Potenz

$$m_p = \sum_{k \geq 1} \left[ \frac{2n}{p^k} \right] - 2 \sum_{k \geq 1} \left[ \frac{n}{p^k} \right] = \sum_{k \geq 1} \left( \left[ \frac{2n}{p^k} \right] - 2 \left[ \frac{n}{p^k} \right] \right).$$

Für jeden einzelnen Summanden ist

$$\left[ \frac{2n}{p^k} \right] - 2 \left[ \frac{n}{p^k} \right] < \frac{2n}{p^k} - 2 \left( \frac{n}{p^k} - 1 \right) = 2;$$

da die Differenz ganzzahlig ist, ist also jeder Summand höchstens gleich eins. Ist  $p^{k_p}$  die größte  $p$ -Potenz kleiner oder gleich  $2n$ , so verschwinden  $\left[ \frac{2n}{p^k} \right]$  und  $\left[ \frac{n}{p^k} \right]$  für alle  $k > k_p$ ; somit ist  $m_p \leq k_p$ . Somit ist

$$\binom{2n}{n} = \prod_{p \leq 2n} p^{m_p} \leq \prod_{p \leq 2n} p^{k_p} \leq \prod_{p \leq 2n} (2n) = (2n)^{\pi(2n)}$$

und wir erhalten die Abschätzung

$$\pi(2n) \geq \frac{\log \binom{2n}{n}}{\log(2n)}.$$

Eine Primzahl  $p$  mit  $n < p \leq 2n$  ist zwar Teiler von  $(2n)!$ , nicht aber von  $n!$ , also muß sie  $\binom{2n}{n}$  teilen. Somit ist

$$n^{\pi(2n) - \pi(n)} < \prod_{n < p \leq 2n} p \leq \binom{2n}{n},$$

also

$$n^{\pi(2n) - \pi(n)} < \binom{2n}{n} \leq (2n)^{\pi(2n)}.$$

Für den Binomialkoeffizienten in der Mitte können wir eine weitere, ganz grobe Abschätzung finden: In

$$\binom{2n}{n} = \frac{(n+1)(n+2) \cdots (n+n)}{n!} = \prod_{k=1}^n \frac{n+k}{k} = \prod_{k=1}^n \left( 1 + \frac{n}{k} \right)$$

ist jeder Faktor mindestens gleich zwei, also ist  $\binom{2n}{n} \geq 2^n$ . Andererseits kommt  $\binom{2n}{n}$  in der Summe

$$2^{2n} = (1 + 1)^{2n} = \sum_{k=0}^{2n} \binom{2n}{k}$$

vor, also ist

$$2^n \leq \binom{2n}{n} \leq 2^{2n}.$$

Zusammen mit dem obigen Resultat erhalten wir die beiden Ungleichungen

$$n^{\pi(2n) - \pi(n)} < \binom{2n}{n} \leq 2^{2n} \quad \text{und} \quad 2^n \leq \binom{2n}{n} \leq (2n)^{\pi(2n)}.$$

Logarithmieren führt auf die beiden Ungleichungen

$$(\pi(2n) - \pi(n)) \log n < 2n \log 2 \quad \text{und} \quad n \log 2 \leq \pi(2n) \log(2n)$$

und somit

$$\pi(2n) - \pi(n) > \frac{2n \log 2}{\log n} \quad \text{und} \quad \pi(2n) \geq \frac{n \log 2}{\log(2n)}.$$

Für jede reelle Zahl  $x \geq 2$  gibt es genau eine natürliche Zahl  $n$  mit  $2n \leq x < 2n + 2$ ; für diese ist dann

$$\pi(x) \geq \pi(2n) \geq \frac{n \log 2}{\log(2n)}.$$

Für jede natürliche Zahl  $n$  ist  $2n \geq n + 1$ , also

$$n \geq \frac{n + 1}{2} = \frac{2n + 2}{4};$$

damit folgt

$$\pi(x) \geq \frac{(2n + 2) \log 2}{4 \log x} > \frac{x \log 2}{4 \log x} = \frac{\log 2}{4} \frac{x}{\log x}.$$

Mit  $c_1 = \frac{1}{4} \log 2 \approx 0,1732867951$  gilt also

$$c_1 \frac{x}{\log x} < \pi(x).$$

Für eine Abschätzung nach oben betrachten wir zunächst nur den Fall, daß  $2n = 2^s$  eine Zweierpotenz ist mit  $s \geq 3$ . Nach den obigen Ungleichungen ist dann

$$\pi(2n) - \pi(n) = \pi(2^s) - \pi(2^{s-1}) < \frac{2^s \log 2}{(s-1) \log 2} = \frac{2^s}{s-1}$$

für alle  $s \geq 3$ . Summieren wir die Differenzen  $\pi(2^s) - \pi(2^{s-1})$  für  $s = 3, \dots, 2t$  erhalten wir

$$\sum_{s=3}^{2t} (\pi(2^s) - \pi(2^{s-1})) = \pi(2^{2t}) - \pi(2^2).$$

Da  $\pi(4) = 2 < 4 = 2^2/(2-1)$ , erhalten wir durch Addition der rechten Seiten ab  $s = 2$  die Ungleichung

$$\pi(2^{2t}) = \pi(4) + \sum_{s=3}^{2t} (\pi(2^s) - \pi(2^{s-1})) < \sum_{s=2}^{2t} \frac{2^s}{s-1}.$$

Die Summe rechts teilen wir auf in die Summe bis  $s = t$  und die ab  $s = t+1$ . Die erste schätzen wir ab durch die geometrische Reihe

$$\sum_{s=2}^t \frac{2^s}{s-1} < \sum_{s=2}^t 2^s = 2^{t+1} - 4 < 2^{t+1},$$

die zweite durch

$$\sum_{s=t+1}^{2t} \frac{2^s}{s-1} < \sum_{s=t+1}^{2t} \frac{2^s}{t} = \frac{2^{2t+1} - 2^{t+1}}{t} < \frac{2^{2t+1}}{t}.$$

Damit ist

$$\pi(2^{2t}) < 2^{t+1} + \frac{2^{2t+1}}{t}.$$

Wegen  $t < 2^t$  ist

$$2^{t+1} = 2 \cdot 2^t < 2 \cdot 2^t \cdot \frac{2^t}{t} = \frac{2^{2t+1}}{t},$$

also

$$\pi(2^{2t}) < 2 \cdot \frac{2^{2t+1}}{t} = \frac{4}{t} \cdot 2^{2t}.$$

Somit ist

$$\frac{\pi(2^{2t})}{2^{2t}} < \frac{4}{t} \quad \text{für alle } t \geq 2.$$

Tatsächlich gilt diese Ungleichung sogar für alle  $t \geq 1$ , denn für  $t = 1$  erhalten wir

$$\frac{\pi(4)}{4} = \frac{2}{4} = \frac{1}{2} < \frac{4}{1} = 4.$$

Für jede reelle Zahl  $x \geq 2$  gibt es eine natürliche Zahl  $t \geq 1$ , so daß  $2^{2t-1} < x \leq 2^{2t}$  ist. Mit diesem  $t$  gilt dann

$$\frac{\pi(x)}{x} \leq \frac{\pi(2^{2t})}{2^{2t-2}} = 4 \cdot \frac{\pi(2^{2t})}{2^{2t}} < \frac{16}{t}.$$

Wegen  $x \leq 2^{2t}$  ist  $\log x \leq 2t \log 2$ , also ist

$$\frac{1}{2t} < \frac{\log 2}{\log x} \quad \text{und} \quad \frac{16}{t} < \frac{32 \log 2}{\log x}.$$

Damit folgt

$$\frac{\pi(x)}{x} < \frac{32 \log 2}{\log x} \quad \text{und} \quad \pi(x) < (32 \log 2) \frac{x}{\log x}.$$

Mit  $c_2 = 32 \log 2 \approx 22,180709778$  und  $c_1$  wie oben gilt somit

$$c_1 \frac{x}{\log x} < \pi(x) < c_2 \frac{x}{\log x},$$

womit der Satz vollständig bewiesen ist. ■

Der bewiesene Satz ist nur ein schwacher Abglanz dessen, was über die Funktion  $\pi(x)$  bekannt ist. Zum Abschluß des Kapitels seien kurz einige der wichtigsten bekannten und vermuteten Eigenschaften von  $\pi(x)$  zusammengestellt. Diese knappe Übersicht folgt im wesentlichen dem Artikel *Primzahlsatz* aus

DAVID WELLS: Prime Numbers – The Most Mysterious Figures in Math, Wiley, 2005,

einer Zusammenstellung im Lexikonformat von interessanten Tatsachen und auch bloßen Kuriosa aus dem Umkreis der Primzahlen.

GAUSS kam 1792, im Alter von 15 Jahren also, durch seine Experimente zur Vermutung, daß  $\pi(x)$  ungefähr gleich dem sogenannten *Integrallogarithmus* von  $x$  sein sollte:

$$\pi(x) \approx \text{Li}(x) \stackrel{\text{def}}{=} \int_2^x \frac{d\xi}{\log \xi}.$$

Auch LEGENDRE versuchte,  $\pi(x)$  anhand experimenteller Daten anzunähern. Er stellte dazu eine Liste aller Primzahlen bis 400 000 zusammen, das sind immerhin 33 860 Stück, und suchte eine glatte Kurve, die den Graphen von  $\pi$  möglichst gut annähert. In seinem 1798 erschienenen Buch *Essai sur la théorie des nombres* gab er sein Ergebnis an als

$$\pi(x) \approx \frac{x}{\log x - 1,08366}.$$

Über ein halbes Jahrhundert später gab es den ersten Beweis einer Aussage: TSCHEBYTSCHEFF zeigte 1851: *Falls*

$$\lim_{x \rightarrow \infty} \frac{\pi(x)}{x / \log x}$$

existiert, dann muß er den Wert eins haben. Ein Jahr später zeigte er, daß

$$c_1 \cdot \frac{x}{\log x} < \pi(x) < c_2 \cdot \frac{x}{\log x} \quad \text{mit} \quad c_1 \approx 0,92 \quad \text{und} \quad c_2 \approx 1,105$$

für *hinreichend große* Werte von  $x$ .

1896 schließlich bewiesen der französische Mathematiker JACQUES SALOMON HADAMARD (1865–1963) und sein belgischer Kollege CHARLES JEAN GUSTAVE NICOLAS BARON DE LA VALLÉE POUSSIN (1866–1962) unabhängig voneinander die Aussage, die heute als **Primzahlsatz** bekannt ist:

$$\pi(x) \sim \frac{x}{\log x}.$$

Der Beweis benutzt Methoden aus der Funktionentheorie, d.h. der Analysis für Funktionen komplexer Zahlen, insbesondere die analytische Fortsetzung der  $\zeta$ -Funktion; er kann daher in einer einführenden Vorlesung, die keine Funktionentheorie voraussetzt, nicht behandelt werden.

Der Beweis des Primzahlsatzes bedeutet nun freilich nicht, daß damit die Formeln von GAUSS oder LEGENDRE überflüssig geworden wären: Die Tatsache, daß der Quotient zweier Funktionen asymptotisch gleich eins ist, erlaubt schließlich immer noch beträchtliche Unterschiede zwischen den beiden Funktionen: Nur der *relative* Fehler muß gegen Null gehen.

Offensichtlich ist für jedes  $a \in \mathbb{R}$

$$\lim_{x \rightarrow \infty} \frac{x/\log x}{x/(\log x - a)} = \lim_{x \rightarrow \infty} \frac{\log x - a}{\log x} = 1 - \lim_{x \rightarrow \infty} \frac{a}{\log x} = 1,$$

und es ist auch nicht schwer zu zeigen, daß

$$\lim_{x \rightarrow \infty} \frac{x/\log x}{\text{Li}(x)} = 1$$

ist. Nach dem Primzahlsatz ist daher auch für jedes  $a \in \mathbb{R}$

$$\pi(x) \sim \frac{x}{\log x - a} \quad \text{und} \quad \pi(x) \sim \text{Li}(x).$$

Wie DE LA VALLÉE POUSSIN zeigte, liefert der Wert  $a = 1$  unter allen reellen Zahlen  $a$  die beste Approximation an  $\pi(x)$ , aber  $\text{Li}(x)$  liefert eine noch bessere Approximation. Für kleine Werte von  $x$  sieht man das auch in der folgenden Tabelle, in der alle reellen Zahlen zur nächsten ganzen Zahl gerundet sind. Wie kaum anders zu erwarten, liefert LEGENDRES Formel für  $10^4$  und  $10^5$  die besten Werte:

$n$	$\pi(n)$	$\frac{n}{\log n}$	$\frac{n}{\log n - 1}$	$\frac{n}{\log n - 1,08366}$	$\text{Li}(n)$
$10^3$	168	145	169	172	178
$10^4$	1 229	1 086	1 218	1 231	1 246
$10^5$	9 592	8 686	9 512	9 588	9 630
$10^6$	78 489	72 382	78 030	78 534	78 628
$10^7$	664 579	620 420	661 459	665 138	664 918
$10^8$	5 761 455	5 428 681	5 740 304	5 769 341	5 762 209
$10^9$	50 847 478	48 254 942	50 701 542	50 917 519	50 849 235

Wenn wir genaue Aussagen über  $\pi(x)$  machen wollen, sollten wir also etwas über die Differenz  $\text{Li}(x) - \pi(x)$  wissen. Hier kommen wir in das Reich der offenen Fragen, und nach derzeitigem Verständnis hängt alles ab von der oben erwähnten RIEMANNschen Zetafunktion. Nach einer

berühmten Vermutung von RIEMANN haben alle nichttrivialen Nullstellen von  $\zeta(s)$  den Realteil ein halb. Falls dies stimmt, ist

$$\pi(x) = \text{Li}(x) + O(\sqrt{x} \log x).$$

Die RIEMANNsche Vermutung ist eines der wichtigsten ungelösten Probleme der heutigen Mathematik; sie war 1900 eines der HILBERTschen Probleme und ist auch eines der sieben *Millennium problems*, für deren Lösung das CLAY Mathematics Institute in Cambridge, Mass. 2000 einen Preis von jeweils einer Million Dollar ausgesetzt hat; für Einzelheiten siehe <http://www.claymath.org/millennium/>.

Die hier verwendeten elementaren Methoden für den Beweis einer groben Abschätzung für  $\pi(x)$  durch  $x/\log x$  reichen auch uns, um zumindest eine grobe Obergrenze für den Abstand zur nächsten Primzahl zu beweisen. Der Satz von BERTRAND besagt, daß es zwischen  $n$  und  $2n$  stets eine Primzahl gibt. Oft bezeichnet man diese Aussage auch als BERTRANDs Postulat, denn BERTRAND konnte seine Behauptung 1845 nur für  $n \leq 3\,000\,000$  beweisen. Erst fünf Jahre später fand TSCHEBYSCHEFF einen Beweis der allgemeinen Aussage.

**Satz von Bertrand:** Für jede natürliche Zahl  $n \geq 2$  gibt es mindestens eine Primzahl  $p$  mit  $n < p < 2n$ .

*Beweis:* Für kleine Werte von  $n$  läßt sich  $p$  leicht explizit angeben: Für  $n = 2$  ist  $p = 3$ , für  $n = 3, 4$  können wir  $p = 5$  nehmen, für  $n = 5, 6$  entsprechend  $p = 7$ , für  $7 \leq n \leq 12$  ist  $p = 13$  eine Möglichkeit, für  $13 \leq n \leq 22$  geht  $p = 23$ ,  $p = 43$  geht für  $23 \leq n \leq 42$ , für  $43 \leq n \leq 82$  etwa  $p = 83$ , und für  $83 \leq n \leq 162$  schließlich  $p = 163$ .

Wir nehmen an, es gäbe ein  $n \geq 2$ , so daß zwischen  $n$  und  $2n$  keine Primzahl liegt. Wir betrachten wieder den Binomialkoeffizienten

$$\binom{2n}{n} = \prod_{\substack{p \leq 2n \\ p \text{ prim}}} p^{m_p} \quad \text{mit} \quad m_p = \left[ \frac{2n}{p} \right] - 2 \left[ \frac{n}{p} \right].$$

Aus dem obigen Beweis wissen wir, daß  $m_p \leq k_p$  ist, wobei  $k_p$  der größte Exponent ist für den  $p^{k_p}$  ein Teiler von  $2n$  ist.

Da es keine Primzahl zwischen  $n$  und  $2n$  gibt, reicht es beim obigen Produkt, den Index nur bis  $n$  laufen zu lassen.

Für  $\frac{2}{3}n < p \leq n$  ist  $\frac{1}{n} \leq \frac{1}{p} < \frac{3}{2n}$ ; Multiplikation mit  $n$  bzw.  $2n$  macht daraus die Ungleichungen

$$1 \leq \frac{n}{p} < \frac{3}{2} \quad \text{und} \quad 2 \leq \frac{2n}{p} < 3.$$

Somit ist in diesem Fall  $\left[\frac{n}{p}\right] = 1$ , und  $\left[\frac{2n}{p}\right] = 2$ , also  $m_p = 0$ , so daß diese Primzahlen in der Faktorisierung von  $\binom{2n}{n}$  ebenfalls nicht vorkommen.

Im Fall  $\sqrt{2n} < p \leq \frac{2}{3}n$  ist  $p^2 > 2n$ , also  $k_p = 1$  und damit  $m_p \leq 1$ . Für  $p \leq \sqrt{2n}$  schließlich haben wir die Abschätzung  $p^{m_p} \leq p^{k_p} \leq 2n$ . Somit ist

$$\binom{2n}{n} = \prod_{\substack{p \leq 2n \\ p \text{ prim}}} p^{m_p} \leq \prod_{\substack{p \leq \sqrt{2n} \\ p \text{ prim}}} (2n) \prod_{\substack{\sqrt{2n} < p \leq 2n/3 \\ p \text{ prim}}} p.$$

Das erste Produkt auf der rechten Seite ist  $(2n)^{\pi(\sqrt{2n})}$ . Für eine grobe Abschätzung des Exponenten verwenden wir, daß für jede reelle Zahl  $x \geq 2$  die Anzahl der Primzahlen kleiner oder gleich  $x$  höchstens gleich der Anzahl der ungeraden Zahlen kleiner oder gleich  $x$  ist: Es gibt zwar auch die gerade Primzahl zwei, dafür ist aber die ungerade Eins keine Primzahl. Somit ist  $\pi(x) \leq \frac{1}{2}(x + 1)$ . Das reicht uns nicht ganz: Falls  $x \geq 16$  gibt es mit 9 und 15 mindestens zwei zusammengesetzte ungerade Zahlen kleiner oder gleich  $x$ , so daß  $\pi(x) \leq \frac{1}{2}(x + 1) - 2 < \frac{x}{2} - 1$  ist.

Das zweite Produkt geht über alle Primzahlen kleiner oder gleich  $2n/3$ ; wie wir im Lemma zu Beginn dieses Paragraphen gesehen haben, ist das Produkt aller Primzahlen kleiner oder gleich  $x$  kleiner als  $4^x$ . Insgesamt ist daher

$$\binom{2n}{n} < (2n)^{\sqrt{n/2}-1} \cdot 4^{2n/3}.$$

Wenden wir den binomischen Lehrsatz an auf  $(1 + 1)^{2n}$  erhalten wir  $2^{2n}$  als eine Summe der  $2n + 1$  Binomialkoeffizienten  $\binom{2n}{i}$  für  $0 \leq i \leq 2n$ .



Der größte dieser Binomialkoeffizienten ist  $\binom{2n}{n}$ ; somit ist

$$2^{2n} < (2n) \binom{2n}{n} < (2n)^{\sqrt{2n}} \cdot 4^{2n/3} = (2n)^{\sqrt{2n}} \cdot 2^{4n/3}.$$

Division durch  $2^{4n/3}$  führt zur Ungleichung

$$2^{2n/3} < (2n)^{\sqrt{n/2}},$$

die wir zum Widerspruch führen wollen.

Logarithmieren beider Seiten ergibt

$$\frac{2n}{3} \log 2 < \sqrt{\frac{n}{2}} \log(2n) \implies \sqrt{8n} \log 2 < 3 \log(2n)$$

durch Multiplikation mit  $3\sqrt{2}/\sqrt{n}$ . Ausgedrückt durch die reelle Funktion

$$F(x) = \sqrt{8x} \log 2 - 3 \log(2x) \quad \text{für } x > 0$$

heißt dies, daß  $F(n)$  negativ ist.

Wie wir wissen, ist  $n \geq 128$ , tatsächlich sogar mindestens 163.

$$F(128) = \sqrt{1024} \log 2 - 3 \log 256 = 32 \log 2 - 3 \cdot 8 \log 2 = 8 \cdot \log 2$$

ist positiv und

$$F'(x) = \frac{\sqrt{8}}{2\sqrt{x}} \log 2 - 3 \cdot 2 \cdot \frac{1}{2x} = \frac{\sqrt{2x} \log 2 - 3}{x} \cdot \log 2.$$

Für  $x \geq 128$  ist der Zähler des Bruchs mindestens

$$\sqrt{256} \log 2 - 3 = 8 \log 2 - 3 = \log 256 - 3 = \log 256 - \log e^3.$$

Dies ist positiv, denn wegen  $e < 3$  ist  $e^3 < 27 < 256$ . Der Nenner ist ohnehin positiv, also ist  $F'(x) > 0$  für alle  $x \geq 128$ , d.h.  $F(x)$  ist in diesem Bereich monoton steigend. Damit muß  $F(x)$  für alle  $x \geq 128$  positiv sein, es gibt also kein  $n$  mit  $F(n) < 0$ . Somit führt die Annahme, es gäbe ein  $n \geq 2$ , für das es keine Primzahl  $p$  mit  $n < p < 2n$  gibt, zum Widerspruch, und der Satz von BERTRAND ist bewiesen. ■



JOSEPH BERTRAND (1822–1900) wurde in Paris geboren und verbrachte dort auch sein gesamtes Leben. Schon mit elf Jahren besuchte er als Gasthörer Vorlesungen an der Ecole polytechnique; im Alter von siebzehn hatte er zwei Bachelorgrade, eine *licence* (vergleichbar mit unserem Staatsexamen für Lehrer) und hatte mit einer Dissertation über die mathematische Theorie der Elektrizität promoviert. Er arbeitete als Lehrer an verschiedenen Gymnasien sowie in verschiedenen Funktionen an der Ecole polytechnique, der Ecole normale supérieure und dem Collège de France. Seit 1856 war er Mitglied der Academie des Sciences, seit 1884 auch der

Académie française. Seine zahlreichen Bücher behandeln nicht nur mathematische, sondern auch physikalische Themen; außerdem schrieb er mehrere Bücher über Wissenschaftsgeschichte.

## §2: Das Sieb des Eratosthenes

Das klassische Verfahren zur Bestimmung aller Primzahlen unterhalb einer bestimmten Schranke geht zurück auf ERATOSTHENES im dritten vorchristlichen Jahrhundert. Es funktioniert folgendermaßen:

Um alle Primzahlen kleiner oder gleich einer Zahl  $N$  zu finden, schreibe man zunächst die Zahlen von eins bis  $N$  in eine Reihe.

Eins ist nach Definition keine Primzahl – für griechische Mathematiker wie EUKLID war die Eins nicht einmal eine Zahl. Also streichen wir die Eins durch. Die Zwei ist prim, aber ihre echten Vielfachen sind natürlich keine Primzahlen, werden also durchgestrichen. Dazu müssen wir nicht von jeder Zahl nachprüfen, ob sie durch zwei teilbar ist, sondern wir streichen einfach nach der Zwei jede zweite Zahl aus der Liste durch.

Die erste nichtdurchgestrichene Zahl der Liste ist dann die Drei. Sie muß eine Primzahl sein, denn hätte sie einen von eins verschiedenen kleineren Teiler, könnte das nur die Zwei sein, und alle Vielfachen von zwei (außer der Zwei selbst) sind bereits durchgestrichen.

Auch die echten Vielfachen der Drei sind keine Primzahlen, werden also durchgestrichen. Auch dazu streichen wir wieder einfach jede dritte Zahl aus der Liste durch, unabhängig davon, ob sie bereits durchgestrichen ist

oder nicht. (Alle durch sechs teilbaren Zahlen sind offensichtlich schon durchgestrichen.)

Genauso geht es weiter mit der Fünf *usw.*; nach jedem Durchgang durch die Liste muß offenbar die erste noch nicht durchgestrichene Zahl eine Primzahl sein, denn alle Vielfache von kleineren Primzahlen sind bereits durchgestrichen, und wenn eine Zahl überhaupt einen echten Teiler hat, dann ist sie natürlich auch durch eine echt kleinere Primzahl teilbar.

Wie lange müssen wir dieses Verfahren durchführen? Wenn eine Zahl  $n$  Produkt zweier echt kleinerer Faktoren  $u, v$  ist, können  $u$  und  $v$  nicht beide größer sein als  $\sqrt{n}$ : Sonst wäre schließlich  $n = uv$  größer als  $n$ . Also ist einer der beiden Teiler  $u, v$  kleiner oder gleich  $\sqrt{n}$ , so daß  $n$  mindestens einen Teiler hat, dessen Quadrat kleiner oder gleich  $n$  ist. Damit ist eine zusammengesetzte Zahl  $n$  durch mindestens eine Primzahl  $p$  teilbar mit  $p^2 \leq n$ .



ERATOSTHENES (Ερατοσθένης) wurde 276 v.Chr. in Cyrene im heutigen Libyen geboren, wo er zunächst von Schülern des Stoikers ZENO ausgebildet wurde. Danach studierte er noch einige Jahre in Athen, bis ihn 245 der Pharao PTOLEMAIOS III als Tutor seines Sohns nach Alexandrien holte. 240 wurde er dort Bibliothekar der berühmten Bibliothek im Museion.

Heute ist er außer durch sein Sieb vor allem durch seine Bestimmung des Erdumfangs bekannt. Er berechnete aber auch die Abstände der Erde von Sonne und Mond und entwickelte einen Kalender, der Schaltjahre enthielt. 194 starb er in Alexandrien, nach einigen Überlieferungen, indem er sich, nachdem er blind geworden war, zu Tode hungerte.

Für das Sieb des ERATOSTHENES, angewandt auf die Zahlen von eins bis  $N$  heißt das, daß wir aufhören können, sobald die erste nicht-durchgestrichene Zahl  $p$  ein Quadrat  $p^2 > N$  hat; dann können wir sicher sein, daß jede zusammengesetzte Zahl  $n \leq N$  bereits einen kleineren Primteiler als  $p$  hat und somit bereits durchgestrichen ist. Die noch nicht durchgestrichenen Zahlen in der Liste sind also Primzahlen.

Damit lassen sich leicht von Hand alle Primzahlen bis hundert finden, mit etwas Fleiß auch die bis Tausend, aber sicher nicht die hundertstelligen.

Trotzdem kann uns ERATOSTHENES helfen, zumindest zu zeigen, daß gewisse Zahlen nicht prim sind: Wenn wir Primzahlen in einem Intervall  $[a, b]$  suchen, d.h. also Primzahlen  $p$  mit

$$a \leq p \leq b,$$

so können wir ERATOSTHENES auf dieses Intervall fast genauso anwenden wie gerade eben auf das Intervall  $[1, N]$ :

Wir gehen aus von einer Liste  $p_1, \dots, p_r$  der ersten Primzahlen; dabei wählen wir  $r$  so, daß die Chancen auf nicht durch  $p_r$  teilbare Zahlen im Intervall  $[a, b]$  noch einigermaßen realistisch sind, d.h. wir gehen bis zu einer Primzahl  $p_r$ , die ungefähr in der Größenordnung der Intervalllänge  $b - a$  liegt.

Nun können wir mit jeder der Primzahlen  $p_i$  sieben wie im klassischen Fall; wir müssen nur wissen, wo wir anfangen sollen.

Dazu berechnen wir für jedes  $p_i$  den Divisionsrest  $r_i = a \bmod p_i$ . Dann ist  $a - r_i$  durch  $p_i$  teilbar, liegt allerdings nicht im Intervall  $[a, b]$ . Die erste Zahl, die wir streichen müssen, ist also  $a - r_i + p_i$ , und von da an streichen wir einfach, ohne noch einmal dividieren zu müssen, wie gehabt jede  $p_i$ -te Zahl durch.

Was nach  $r$  Durchgängen noch übrig bleibt, sind genau die Zahlen aus  $[a, b]$ , die durch keine der Primzahlen  $p_i$  teilbar sind. Sie können zwar noch größere Primteiler haben, aber wichtig ist, daß wir mit minimalem Aufwand für den Großteil aller Zahlen aus  $[a, b]$  gesehen haben, daß sie keine Primzahlen sind. Für den Rest brauchen wir andere Verfahren, aber die sind allesamt erheblich aufwendiger als ERATOSTHENES, so daß sich diese erste Reduktion auf jeden Fall lohnt.

### §3: Fermat-Test und Fermat-Zahlen

Nach dem kleinen Satz von FERMAT gilt für jede Primzahl  $p$  und jede nicht durch  $p$  teilbare Zahl  $a$  die Formel  $a^{p-1} \equiv 1 \pmod{p}$ . Im Umkehrschluß folgt sofort:

*Falls für eine natürliche Zahl  $1 \leq a \leq p - 1$  gilt  $a^{p-1} \not\equiv 1 \pmod{p}$ , kann  $p$  keine Primzahl sein.*

Beispiel: Ist  $p = 129$  eine Primzahl? Falls ja, ist nach dem kleinen Satz von FERMAT  $2^{128} \equiv 1 \pmod{129}$ . Tatsächlich ist aber

$$2^7 = 128 \equiv -1 \pmod{129},$$

also hat die Zwei in  $(\mathbb{Z}/129)^\times$  die Ordnung 14. Da 14 kein Teiler von 128 ist, kann  $2^{128} \pmod{129}$  nicht eins sein. (Wegen  $128 \equiv 2 \pmod{14}$  ist  $2^{128} \equiv 2^2 = 4 \pmod{129}$ .) Somit ist 129 keine Primzahl.

Dieses Ergebnis hätten wir natürlich auch durch Probedivisionen leicht gefunden: Da 129 die Quersumme 12 hat, ist die Zahl durch drei teilbar; ihre Primzerlegung ist  $129 = 3 \cdot 43$ .

Keine Kopfrechenaufgabe ist die Frage, ob  $F_{20} = 2^{2^{20}} + 1$  eine Primzahl ist. Falls ja, wäre nach dem kleinen Satz von FERMAT insbesondere

$$3^{F_{20}-1} = 1 \pmod{F_{20}}, \quad \text{also} \quad 3^{(F_{20}-1)/2} \equiv \pm 1 \pmod{F_{20}}.$$

Nachrechnen zeigt, daß dies nicht der Fall ist, allerdings ist das „Nachrechnen“ bei dieser 315 653-stelligen Zahl natürlich keine Übungsaufgabe für Taschenrechner: 1988 brauchte eine Cray X-MP dazu 82 Stunden, der damals schnellste Supercomputer Cray-2 immerhin noch zehn; siehe

JEFF YOUNG, DUNCAN A. BUELL: The Twentieth Fermat Number is Composite, *Math. Comp.* **50** (1988), 261–263.

Damit war gezeigt, daß  $F_{20}$  keine Primzahl ist. (Die anscheinend etwas weltabgewandt lebenden Autoren meinten, dies sei die aufwendigste bis dahin produzierte 1-Bit-Information.)

Umgekehrt können wir leider nicht folgern, daß  $p$  eine Primzahl ist, wenn für ein  $a \in \mathbb{N}$  mit  $1 < a < p - 1$  gilt  $a^{p-1} \equiv 1 \pmod{p}$ . So ist beispielsweise  $18^{322} \equiv 1 \pmod{323}$ , aber  $323 = 17 \cdot 19$  ist zusammengesetzt. Immerhin gibt es nicht viele  $a \leq 323$  mit  $a^{322} \equiv 1 \pmod{323}$ : Die einzigen Möglichkeiten sind  $a = \pm 1$  und  $a = \pm 18$ .

Es kann nicht vorkommen, daß für eine zusammengesetzte Zahl  $n$  und alle  $1 \leq a \leq n$  gilt  $a^{n-1} \equiv 1 \pmod{n}$ , denn ist  $p$  ein Primteiler von  $n$ , so ist für jedes Vielfache  $a$  von  $p$  natürlich auch  $a^{n-1}$  durch  $p$  teilbar, kann also nicht kongruent eins modulo des Vielfachen  $n$  von  $p$  sein.

Zumindest für die  $a$  mit  $\text{ggT}(a, n) > 1$  kann die Gleichung also nicht erfüllt sein. Bei großen Zahlen  $n$  mit nur wenigen Primfaktoren ist aber die Chance, ein solches  $a$  zu erwischen, recht klein; wenn dies die einzigen Gegenbeispiele sind, wird uns der FERMAT-Test daher fast immer in die Irre führen.

**Definition:** Eine natürliche Zahl  $n$  heißt CARMICHAEL-Zahl, wenn sie keine Primzahl ist, aber trotzdem für jede natürliche Zahl  $a$  mit  $\text{ggT}(a, n) = 1$  gilt:  $a^{n-1} \equiv 1 \pmod{n}$ .

ROBERT DANIEL CARMICHAEL (1879–1967) war ein amerikanischer Mathematiker, der unter anderem Bücher über die Relativitätstheorie, über Zahlentheorie, über Analysis und über Gruppentheorie veröffentlichte. Ab 1915 lehrte er an der University of Illinois. Er zeigte 1910, daß 561 die gerade definierte Eigenschaft hat und publizierte auch später noch eine Reihe von Arbeiten über solche Zahlen.

**Satz:** Eine natürliche Zahl  $n$  ist genau dann eine CARMICHAEL-Zahl, wenn sie das Produkt von mindestens drei paarweise verschiedenen ungeraden Primzahlen ist, wobei für jeden Primfaktor  $p$  auch  $p - 1$  Teiler von  $n - 1$  ist.

*Beweis:* Sei zunächst  $n = \prod p_i$  ein Produkt paarweise verschiedener Primzahlen, für die  $p_i - 1$  Teiler von  $n - 1$  ist. Nach dem chinesischen Restesatz ist dann  $(\mathbb{Z}/N)^\times \cong \prod (\mathbb{Z}/p_i)^\times$ . In der Gruppe  $(\mathbb{Z}/p_i)^\times$  ist die Ordnung eines jeden Elements ein Teiler von  $p_i - 1$  und damit von  $n - 1$ ; also ist auch in  $(\mathbb{Z}/n)^\times$  die Ordnung eines jeden Elements Teiler von  $n - 1$ . Damit gilt für jedes zu  $n$  teilerfremde  $a \in \mathbb{Z}$  die Kongruenz  $a^{n-1} \equiv 1 \pmod{n}$ , d.h.  $n$  ist eine CARMICHAEL-Zahl.

Umgekehrt sei  $n$  eine CARMICHAEL-Zahl. Dann ist  $n$  ungerade, denn für gerade Zahlen  $n$  ist  $(n - 1)^{n-1} \equiv (-1)^{n-1} = -1 \pmod{n}$ .

Als nächstes wollen wir uns überlegen, daß  $n$  Produkt verschiedener Primzahlen sein muß: Angenommen, in der Primzerlegung von  $n$  tritt eine Primzahl  $p$  mehrfach auf, d.h.  $n = p^e q$  mit einer zu  $p$  teilerfremden Zahl  $q$ . Nach dem binomischen Lehrsatz gilt

$$(p + 1)^{p^{e-1}} = \sum_{k=0}^{p^{e-1}} \binom{p^{e-1}}{k} p^k,$$

und für alle  $k \neq 0$  ist

$$\begin{aligned} \binom{p^{e-1}}{k} p^k &= \frac{p^{e-1}(p^{e-1}-1)\cdots(p^{e-1}-k+1)}{k!} p^k \\ &= p^{e-1} \cdot \frac{p^{e-1}}{1} \cdot \frac{p^{e-1}-2}{2} \cdots \frac{p^{e-1}-(k-1)}{k-1} \cdot \frac{p^k}{k}. \end{aligned}$$

In jedem der Brüche  $(p^{e-1}-\ell)/\ell$  kommt  $p$  in Zähler und Nenner mit der gleichen Potenz vor, denn  $\ell < p^{e-1}$  und  $p^{e-1}-\ell \equiv -\ell \pmod{p^{e-1}}$ . Im letzten Bruch  $p^k/k$  steht  $p$  im Zähler offensichtlich mit einer höheren Potenz als im Nenner; insgesamt ist der Ausdruck also mindestens durch  $p^e$  teilbar. Somit ist  $(1+p)^{p^{e-1}} \equiv 1 \pmod{p^e}$ ; in  $(\mathbb{Z}/p^e)^\times$  gibt es daher Elemente, deren Ordnung ein Vielfaches von  $p$  ist. Da nach dem chinesischen Restesatz  $(\mathbb{Z}/n)^\times \cong (\mathbb{Z}/p^e)^\times \times (\mathbb{Z}/q)^\times$  ist, gibt es dann auch in  $(\mathbb{Z}/n)^\times$  ein solches Element  $a$ . Da  $n-1$  nicht durch  $p$  teilbar ist, ist  $n$  kein Vielfaches dieser Ordnung, so daß  $a^{n-1}$  modulo  $n$  nicht eins sein kann. Damit ist  $n$  keine CARMICHAEL-Zahl; eine CARMICHAEL-Zahl muß also Produkt verschiedener Primzahlen sein.

Für jeden Primteiler  $p$  von  $n$  muß  $p-1$  ein Teiler von  $n-1$  sein, denn nach dem chinesischen Restesatz ist  $(\mathbb{Z}/n)^\times$  das Produkt der Gruppen  $(\mathbb{Z}/p)^\times$ , es gibt also in  $(\mathbb{Z}/n)^\times$  eine primitive Wurzel  $a$  modulo  $p$ . Da  $a^{n-1} \equiv 1 \pmod{n}$  und damit insbesondere modulo  $p$  ist, muß  $n-1$  ein Vielfaches der Ordnung  $p-1$  von  $a$  modulo  $p$  sein.

Schließlich müssen wir uns noch überlegen, daß  $n$  ein Produkt von mindestens drei Primzahlen ist: Da  $n$  nach Definition keine Primzahl ist, wäre  $n = pq$  sonst das Produkt zweier Primzahlen. Wie wir gerade gesehen haben, müßte

$$n-1 = pq-1 = (p-1)q + (q-1) = p(q-1) + (p-1)$$

sowohl durch  $p-1$  als auch durch  $q-1$  teilbar sein, also müßten  $p-1$  und  $q-1$  durcheinander teilbar sein, d.h.  $p=q$ , was wir bereits ausgeschlossen haben. ■

Als Beispiel können wir ein Produkt  $n = (6t+1)(12t+1)(18t+1)$  mit drei primen Faktoren betrachten, z.B.

$$1729 = 7 \times 13 \times 19 \quad \text{für } t = 1 \quad \text{oder} \quad 294409 = 37 \times 73 \times 109$$

für  $t = 6$ . Hier ist  $n - 1 = 1296t^3 + 396t^2 + 36t = 36t \cdot (36t^2 + 11t + 1)$  offensichtlich durch  $6t$ ,  $12t$  und  $18t$  teilbar,  $n$  ist also eine CARMICHAEL-Zahl.

Natürlich muß nicht jede CARMICHAEL-Zahl von dieser Form sein; die kleinste CARMICHAEL-Zahl beispielsweise ist  $561 = 3 \cdot 11 \cdot 17$  und hat keinen einzigen Primfaktor kongruent eins modulo sechs.

Eine größte CARMICHAEL-Zahl gibt es nicht, denn nach

W.R. ALFORD, ANDREW GRANVILLE, CARL POMERANCE: There are infinitely many Carmichael numbers, *Ann. Math.* **140** (1994), 703–722

gibt es unendlich viele. Konkret zeigen sie, daß die Anzahl der CARMICHAEL-Zahlen kleiner oder gleich  $x$  für große Zahlen  $x$  mindestens gleich  $x^{2/7}$  ist. Die tatsächliche Anzahl dürfte wohl deutlich größer sein, ist aber immer noch sehr viel kleiner als die der Primzahlen.

Für große Zahlen  $p$  wird es zunehmend unwahrscheinlich, daß sie auch nur für ein  $a$  den FERMAT-Test bestehen, ohne Primzahl zu sein. Rechnungen von

SU HEE KIM, CARL POMERANCE: The probability that a Random Probable Prime is Composite, *Math. Comp.* **53** (1989), 721–741

geben folgende obere Schranke für die Wahrscheinlichkeit  $\varepsilon$ , daß eine zufällig gewählte Zahl  $p$  der angegebenen Größenordnung den FERMAT-Test mit einem vorgegebenen  $a$  besteht und trotzdem keine Primzahl ist:

$p \approx 10^{60}$	$10^{70}$	$10^{80}$	$10^{90}$	$10^{100}$
$\varepsilon \leq 7,16 \cdot 10^{-2}$	$2,87 \cdot 10^{-3}$	$8,46 \cdot 10^{-5}$	$1,70 \cdot 10^{-6}$	$2,77 \cdot 10^{-8}$
$p \approx 10^{120}$	$10^{140}$	$10^{160}$	$10^{180}$	$10^{200}$
$\varepsilon \leq 5,28 \cdot 10^{-12}$	$1,08 \cdot 10^{-15}$	$1,81 \cdot 10^{-19}$	$2,76 \cdot 10^{-23}$	$3,85 \cdot 10^{-27}$
$p \approx 10^{300}$	$10^{400}$	$10^{500}$	$10^{600}$	$10^{700}$
$\varepsilon \leq 5,8 \cdot 10^{-29}$	$5,7 \cdot 10^{-42}$	$2,3 \cdot 10^{-55}$	$1,7 \cdot 10^{-68}$	$1,8 \cdot 10^{-82}$
$p \approx 10^{800}$	$10^{900}$	$10^{1000}$	$10^{2000}$	$10^{3000}$
$\varepsilon \leq 5,4 \cdot 10^{-96}$	$1,0 \cdot 10^{-109}$	$1,2 \cdot 10^{-123}$	$8,6 \cdot 10^{-262}$	$3,8 \cdot 10^{-397}$



(Sie geben natürlich auch eine allgemeine Formel an, jedoch ist diese zu grausam zum Abtippen.)

Wenn wir RSA-Moduln von 2048 Bit konstruieren wollen, brauchen wir etwa dreihundertstellige Primzahlen; hier liegt die Irrtumswahrscheinlichkeit bei einem einzigen FERMAT-Test also bei höchstens  $5,8 \cdot 10^{-29}$ . Wenn das zu hoch ist, kann man mit mehreren zufällig gewählten Basen testen und dadurch die Fehlerwahrscheinlichkeit deutlich verringern – auch wenn es wohl gewagt wäre, zwei solche Tests als unabhängig anzunehmen.

Die Bundesnetzagentur empfiehlt, bei probabilistischen Primzahltests für die Erzeugung von RSA-Moduln eine Irrtumswahrscheinlichkeit von höchstens  $2^{-100} \approx 7,89 \cdot 10^{-31}$  zuzulassen. Da die Wahrscheinlichkeiten in obiger Tabelle obere Schranken sind, könnte das vielleicht schon mit einem Test erreicht sein; besser sind auf jeden Fall mehrere oder, noch besser, ein Test, der wirklich *beweisen* kann, daß eine Zahl prim ist.

Einige Leute reden bei Zahlen, die einen FERMAT-Test bestanden haben, von „wahrscheinlichen Primzahlen“. Das ist natürlich Unsinn: Eine Zahl ist entweder *sicher* prim oder *sicher* zusammengesetzt; für Wahrscheinlichkeiten gibt es hier keinen Spielraum. Besser ist der ebenfalls gelegentlich zu hörende Ausdruck „industrial grade primes“, also „Industriepriimzahlen“, der ausdrücken soll, daß wir zwar nicht *bewiesen* haben, daß die Zahl wirklich prim ist, daß sie aber für (manche) „industrielle Anwendungen“ gut genug ist.

Zumindest grundsätzlich läßt sich der FERMAT-Test auch ausbauen zu einem echten Primzahltest; die einfachste Art ist die folgende sehr schwache Version eines Satzes von POCKLINGTON:

**Satz:** Ist für zwei natürliche Zahlen  $p, a$  zwar  $a^{p-1} \equiv 1 \pmod{p}$ , aber für jeden Primteiler  $q$  von  $p-1$  gilt  $a^{(p-1)/q} \not\equiv 1 \pmod{p}$ , so ist  $p$  eine Primzahl.

*Beweis:* Offensichtlich muß dann die Ordnung von  $a$  in  $(\mathbb{Z}/p\mathbb{Z})^\times$  gleich  $p-1$  sein. Wie wir aus Kapitel 1, §8 wissen, hat  $(\mathbb{Z}/p\mathbb{Z})^\times$  die Ordnung  $\varphi(p)$ , und für jede zusammengesetzte Zahl folgt aus der dort

angegebenen Formel leicht, daß  $\varphi(p) < p - 1$  ist. Also muß  $p$  prim sein. ■

HENRY CABOURN POCKLINGTON (1870–1952) wurde im englischen Exeter geboren; er studierte an dem College, aus 1904 die University of Leeds wurde. Damals wurden Studenten dort auf Examen in London vorbereitet; POCKLINGTON erhielt 1889 Bachelorgrade sowohl in Experimentalphysik als auch in Mathematik und erhielt anschließend verschiedene Stipendien des St. John's College in Cambridge, zunächst als Student, dann als *fellow*. 1896 erhielt er den Doktorgrad der Universität London. Als seine Stipendien 1900 ausliefen, wurde er Physiklehrer an einer Schule in Leeds und blieb dies auch bis zu seiner Pensionierung 1926, obwohl er mehrfach Angebote von Universitäten bekam und 1907 sogar *fellow* der *Royal Society* wurde. Zwischen 1895 und 1940 publizierte er rund vierzig Arbeiten, zunächst hauptsächlich aus dem Gebiet der Physik und Astronomie, ab 1910 aber vor allem aus der Mathematik,

Der Nachteil des gerade bewiesenen Satzes ist, daß wir alle Primteiler von  $p - 1$  kennen müssen; wenn wir Primzahlen mit mehreren hundert Dezimalstellen suchen, ist das meist keine sehr realistische Annahme. Für Zahlen spezieller Bauart kann dieser Test jedoch sehr nützlich sein: Wenn wir von einer Zahl  $n$  mit wenigen, uns bekannten Primteilern ausgehen, läßt sich so testen, ob  $n + 1$  eine Primzahl ist.

Die einfachsten Kandidaten für  $n$  sind Primzahlpotenzen  $n = p^r$ ; für eine ungerade Primzahl  $p$  ist allerdings  $n + 1$  gerade und somit keine Primzahl. Auf den Fall  $p = 2$  werden wir gleich zurückkommen.

Bei kleinen geraden Zahlen  $b$  mit wenigen Primfaktoren gibt es gelegentlich Chancen, daß  $b^r + 1$  eine Primzahl ist, allerdings kommt auch das nur selten vor. Ein Beispiel wäre etwa

$$p = 24^4 + 1 = 331\,777.$$

Hier hat  $p - 1 = 24^4$  nur 2 und 3 als Primteiler, wir müssen also eine Zahl  $a$  finden, so daß

$$a^{p-1} \equiv 1 \pmod{p}, \quad a^{(p-1)/2} \not\equiv 1 \pmod{p} \quad \text{und} \quad a^{(p-1)/3} \not\equiv 1 \pmod{p}$$

ist. Dazu berechnen wir am besten zunächst  $x = a^{(p-1)/6} \pmod{p}$ , sodann  $y = a^{(p-1)/3} \pmod{p} = x^2 \pmod{p}$  und  $z = a^{(p-1)/2} \pmod{p} = x^3 \pmod{p}$ . Wenn  $p$  eine Primzahl ist, muß offensichtlich  $z \equiv -1 \pmod{p}$  sein, und wenn dies gilt, ist auch  $a^{(p-1)} \equiv 1 \pmod{p}$ .

Für  $a = 2$  erhalten wir  $x = 92553, y = 239223$  und  $z = 1$ ; das beweist nichts. Für  $a = 3$  ist sogar bereits  $x = 1$ , aber für  $a = 5$  wird  $x = 92554, y = 92553$  und  $z = 331776 \equiv -1 \pmod{p}$ . Dies *beweist*, daß  $p$  eine Primzahl ist.

Als nächstes wollen wir uns überlegen, wann  $n = 2^r + 1$  prim sein kann. Für ungerade  $r$  ist  $n$  durch drei teilbar, denn  $2^r \equiv (-1)^r \equiv -1 \pmod{3}$ ; für  $r > 1$  kann  $n$  dann also nicht prim sein. Auch wenn  $r$  nur durch eine ungerade Zahl  $u > 1$  teilbar ist, kann  $2^r + 1$  nicht prim sein, denn ist  $r = uv$ , so ist  $2^r = (2^v)^u \equiv (-1)^u \equiv -1 \pmod{2^v + 1}$  so daß  $2^v + 1$  ein nichttrivialer Teiler ist. Die einzigen Kandidaten für Primzahlen sind daher die Zahlen

$$F_n = 2^{2^n} + 1,$$

bei denen der Exponent  $r$  eine Zweierpotenz ist. Sie heißen FERMAT-Zahlen, weil FERMAT zwischen 1630 und 1640 in mehreren Briefen, unter anderem an PASCAL und an MERSENNE, die Vermutung äußerte, diese Zahlen seien allesamt prim.

Für  $F_0 = 3, F_1 = 5, F_2 = 17$  sieht man das mit bloßem Auge und mit auch  $F_3 = 257$  gibt es keine nennenswerten Schwierigkeiten. Für größere Werte von  $n$  können wir den obigen Test anwenden; da 2 der einzige Primteiler von  $F_n - 1$  ist, wird er hier einfach einfach zur folgenden Aussage:

**Lemma:**  $F_n$  ist genau dann eine Primzahl, wenn es ein  $a$  gibt mit

$$a^{(F_n - 1)/2} \equiv -1 \pmod{F_n}.$$

■

Für  $F_4 = 2^{16} + 1 = 65\,537$  ist  $(F_n - 1)/2 = 2^{15}$ , wir müssen also ein  $a$  finden, so daß  $a^{2^{15}} \equiv -1 \pmod{F_4}$  ist. Für  $a = 2$  bekommen wir nach 15 Quadrierungen das nutzlose Ergebnis eins, für  $a = 3$  aber die gewünschte  $-1$ . Damit ist  $F_4$  als Primzahl erkannt. Für so kleine Zahlen lohnt sich dieser Test freilich noch nicht wirklich; FERMAT hat sich wohl klassisch mit Probedivisionen davon überzeugt, daß  $F_4$  prim ist.

Für  $F_5 = 2^{32} + 1 = 429\,4967\,297$  müssen wir  $a^{2^{31}} \pmod{F_5}$  berechnen, brauchen also schon 31 Quadrierungen. Für  $a = 2$  erhalten wir wieder

die nutzlose Eins, was übrigens kein Wunder ist: Wie Aufgabe 1b) des siebten Übungsblatts zeigt, gilt dies für *jedes*  $F_n$ . Für  $a = 3$  erhalten wir aber das Ergebnis 10 324 303, das sich offenbar auch modulo  $F_5$  deutlich von  $\pm 1$  unterscheidet. Somit ist  $F_5$  *keine* Primzahl und FERMATs obige Vermutung ist widerlegt. Sie ist übrigens die einzige seiner vielen Vermutungen, die sich als falsch erwies.

Der erste, der erkannte, daß  $F_5$  zusammengesetzt ist, war EULER, den GOLDBACH in Sankt Petersburg auf FERMATs Vermutung aufmerksam gemacht hatte. Nachdem EULERS Beweisversuche erfolglos blieben, fand er 1732 schließlich die Faktorisierung  $F_5 = 641 \cdot 6\,700\,417$ . Obwohl EULER viel rechnete, fand er diese Faktorisierung natürlich nicht durch systematisches Ausprobieren aller Primzahlen bis zur Quadratwurzel von  $F_5$ : dafür hätte er im schlimmsten Fall immerhin durch 6542 Primzahlen dividieren müssen!

Stattdessen benutzte EULER den kleinen Satz von FERMAT, um zunächst Aussagen über mögliche Teiler von FERMAT-Zahlen zu bekommen. Er fand

**Lemma:** Jeder Primteiler  $p$  von  $F_n$  ist kongruent eins modulo  $2^{n+1}$ .

*Beweis:* Wegen  $2^{2^n} \equiv -1 \pmod{F_n}$  ist erst recht  $2^{2^n} \equiv -1 \pmod{p}$ , also  $2^{2^{n+1}} \equiv 1 \pmod{p}$ . Die Ordnung der Zwei in  $\mathbb{F}_p^\times$  ist somit ein Teiler von  $2^{n+1}$ , also eine Zweierpotenz. Wäre diese kleiner als  $2^{n+1}$ , wäre sie ein Teiler von  $2^n$ , so daß  $2^{2^n} \equiv +1 \pmod{p}$  sein müßte. Somit ist  $2^{n+1}$  die genaue Ordnung. Diese muß aber nach dem Satz von LAGRANGE ein Teiler der Gruppenordnung sein, d.h.  $2^{n+1}$  teilt  $p-1$ , was die Behauptung beweist. ■

Somit wußte EULER, daß jeder Primteiler von  $F_5$  die Form  $p = 64k + 1$  haben muß; Primzahlen dieser Form gibt es nur 209 unterhalb von  $2^{15}$ , was das Problem schon viel handhabbarer erscheinen läßt. Dazu kam, daß er Glück hatte: Schon für  $k = 10$  bekam er einen Teiler. Nach 193, 257, 449 und 577 ist 641 bereits die fünfte Primzahl dieser Form!

Tatsächlich aber machte er sich trotzdem zuviel Arbeit, denn wie der französische Mathematiker EDOUARD LUCAS (1842–1891) fand, gilt

sogar

**Lemma:** Für  $n \geq 3$  ist jeder Primteiler  $p$  von  $F_n$  kongruent eins modulo  $2^{n+2}$ .

*Beweis:* Wir gehen genauso vor wie EULER, suchen aber ein Element von  $\mathbb{F}_p^\times$ , dessen Ordnung  $2^{n+2}$  ist. Ein solches Element haben wir gefunden, wenn wir in  $\mathbb{F}_p^\times$  ein  $x$  finden mit  $x^2 = 2$ . Wegen der Beziehung

$$(2^{2^{n-1}} + 1)^2 = 2^{2^n} + 1 + 2^{2^{n-1}+1} \equiv 2^{2^{n-1}+1} \pmod{F_n}$$

haben wir zunächst eine Zahl  $y$  gefunden mit  $y^2 \equiv 2^u \pmod{F_n}$ , wobei  $u = 2^{n-1} + 1$  eine ungerade Zahl ist. Mit dem erweiterten EUKLIDischen Algorithmus können wir daher ganze Zahlen  $v, w$  finden mit  $uv + 2w = 1$ . Damit ist auch

$$2 = 2^{uv+2w} = (2^u)^v \cdot (2^w)^2 = y^{2v} \cdot 2^{2w} = (y^v \cdot 2^w)^2$$

ein Quadrat in  $\mathbb{F}_p^\times$ . ■

Mit dieser Verschärfung seines Lemmas hätte sich EULER auf Primzahlen der Form  $128k + 1$  beschränken können und hätte bereits im zweiten Anlauf (nach einem vergeblichen Versuch mit 257) seinen Faktor gefunden.

Auch wenn er nie etwas darüber publiziert hat, hätte er auch beweisen können und bewies vielleicht auch, daß sein Kofaktor  $q = 6\,700\,417$  ebenfalls eine Primzahl ist: Da jeder Primteiler  $p$  von  $q$  insbesondere auch  $F_5$  teilt, muß  $p \equiv 1 \pmod{128}$  sein, und wenn es einen echten Teiler gibt, gibt es auch einen der kleiner ist als  $\sqrt{q} \approx 2588,5$ . Es gibt nur zwanzig Zahlen der Form  $128k + 1$  unterhalb von  $\sqrt{q}$ , und nur fünf davon sind prim: Neben den bereits bekannten Kandidaten 257 und 641 sind das noch 769, 1153 und 1409. Fünf einfache Divisionen mit Rest zeigen, daß  $q$  durch keine dieser Zahlen teilbar ist; somit haben wir bewiesen, daß auch  $q$  prim sein muß.

Die Suche nach FERMAT-Primzahlen sowie nach Faktoren von FERMAT-Zahlen beschäftigt auch heute noch eine ganze Reihe von Mathematikern; die jeweils neuesten Ergebnisse sind zum Beispiel auf den Webseiten von [www.fermatsearch.org](http://www.fermatsearch.org) zu finden. Unter anderem ist bewiesen,

daß  $F_n$  für  $5 \leq n \leq 32$  sowie viele andere Werte von  $n$  zusammengesetzt ist; oft sind auch zumindest einige Faktoren bekannt. FERMATsche Primzahlen mit  $n > 4$  wurden bislang keine gefunden, allerdings ist auch noch nicht bewiesen, daß es unendlich viele FERMAT-Zahlen gibt, die *keine* Primzahlen sind.

Der oben angegebene Beweis von LUCAS zeigt übrigens auch, warum wir beim Primzahltest für  $F_4$  mit der Basis zwei keinen Erfolg hatten: Da 2 modulo  $F_4$  ein Quadrat ist, *muß* die Potenz mit Exponent  $(F_4 - 1)/2$  gleich eins sein, denn sie ist schließlich die  $(F_4 - 1)$ -te Potenz der Wurzel aus 2. Aus diesem Grund wurde bei der Überprüfung von  $F_{20}$  nicht mit  $a = 2$  gearbeitet, obwohl dies rechnerisch sehr viel effizienter gewesen wäre. Wie wir im Kapitel über quadratische Reste sehen werden, konnten sich die Autoren *vor* Beginn ihrer Rechnung leicht davon überzeugen, daß drei modulo  $F_{20}$  *keine* Quadratzahl ist; der Aufwand für die Entscheidung, ob eine Zahl  $a$  Quadrat modulo  $F_{20}$  ist, erfordert einen Aufwand, der ungefähr dem für die Anwendung des EUKLIDischen Algorithmus auf  $a$  und  $F_{20}$  entspricht, im Falle  $a = 3$  also nicht viel mehr als die Berechnung von  $F_{20} \bmod 3$ , was man im Kopf als  $(-1)^{20} + 1 = 2$  berechnen kann.

Kehren wir zurück zu Primzahltests für allgemeine Zahlen  $n$ , Wenn wir  $n - 1$  zumindest teilweise faktorisieren können, hilft gelegentlich der folgende Satz:

**Satz:** Angenommen  $n - 1 = uv$  ist das Produkt zweier teilerfremder Zahlen  $u < v$ , und wir kennen die Primzerlegung  $v = \prod q_i^{e_i}$  von  $v$ . Falls es ein  $a \in \mathbb{N}$  gibt, so daß  $a^{n-1} \equiv 1 \pmod n$  und

$$\text{ggT}(a^{(n-1)/q_i} - 1, n) = 1 \quad \text{für alle } i,$$

ist  $n$  eine Primzahl.

*Beweis:* Angenommen,  $p$  sei ein echter Primteiler von  $n$ , und  $a$  erfülle die angegebenen Bedingungen. Die Ordnung der Restklasse von  $a$  in  $(\mathbb{Z}/p)^\times$  sei  $r$ ; sie ist natürlich ein Teiler von  $p - 1$ .

Da  $a^{n-1} \equiv 1 \pmod n$  und damit erst recht modulo des Teilers  $p$ , ist  $r$  auch Teiler von  $n - 1$ ; da  $a^{(n-1)/q_i} - 1$  aber teilerfremd zu  $n$  ist, ist

$a^{(n-1)/q_i}$  nicht kongruent eins modulo  $p$ ; die Ordnung  $r$  teilt also keine der Zahlen  $(n-1)/q_i$ . Somit teilt  $r$  zwar das volle Produkt  $u \prod q_i^{e_i}$ , nicht aber das Produkt, in dem auch nur ein Exponent  $e_i$  erniedrigt wurde. Somit ist  $r$  für jedes  $i$  ein Vielfaches von  $q_i^{e_i}$  und damit auch ein Vielfaches von  $v$ . Da  $r$  ein Teiler von  $p-1$  ist, folgt insbesondere daß  $v < p$  sein muß. Nun ist aber  $n = uv$  und  $u < v$ , d.h.  $v > \sqrt{n}$ . Damit haben wir gezeigt, daß jeder echte Primteiler von  $n$  größer als  $\sqrt{n}$  sein muß. Da dies unmöglich ist, gibt es keinen echten Primteiler, d.h.  $n$  ist eine Primzahl. ■

#### §4: Der Test von Miller und Rabin

Der Test von MILLER und RABIN ist eine etwas strengere Version des Tests von FERMAT: Um zu testen, ob  $p$  eine Primzahl sein kann, schreiben wir  $p-1$  zunächst als Produkt  $2^n u$  einer Zweierpotenz und einer ungeraden Zahl; sodann berechnen wir  $a^u \bmod p$ . Falls wir das Ergebnis eins erhalten, ist erst recht  $a^{p-1} \equiv 1 \pmod p$ , und wir können nicht folgern, daß  $p$  zusammengesetzt ist.

Andernfalls quadrieren wir das Ergebnis bis zu  $n$ -mal modulo  $p$ . Falls dabei nie eine Eins erscheint, folgt nach FERMAT, daß  $p$  zusammengesetzt ist. Falls vor der ersten Eins eine von  $-1$  (bzw.  $p-1$ ) verschiedene Zahl erscheint, folgt das auch, denn im Körper  $\mathbb{F}_p$  hat die Eins nur die beiden Quadratwurzeln  $\pm 1$ . In allen anderen Fällen erfahren wir nicht mehr als bei FERMAT.

Algorithmisch funktioniert der Test also folgendermaßen:

**Schritt 0:** Wähle ein zufälliges  $a$ , schreibe  $p-1 = 2^n u$  mit einer ungeraden Zahl  $u$  und berechne  $b = a^u \bmod p$ . Falls dies gleich Eins ist, endet der Algorithmus und wir konnten nicht zeigen, daß  $p$  eine zusammengesetzte Zahl ist; sie kann prim sein.

**Schritt  $i$ ,  $1 \leq i \leq n$ :** Falls  $b \equiv -1 \pmod p$ , endet der Algorithmus und wir können nicht ausschließen, daß  $p$  prim ist. Falls  $b = 1$  ist (was frühestens im zweiten Schritt der Fall sein kann), ist  $p$  zusammengesetzt und der Algorithmus endet. Andernfalls wird  $b$  durch  $b^2 \bmod p$  ersetzt und es geht weiter mit Schritt  $i+1$ .

**Schritt  $n + 1$ :** Der Algorithmus endet mit dem Ergebnis, daß  $p$  zusammengesetzt ist.

*Beispiel:* Ist 247 eine Primzahl? Wir wählen  $a = 77$ , und da  $77^{246} \bmod 247 = 1$  ist, können wir mit FERMAT nicht ausschließen, daß 247 prim ist. Da aber  $77^{123} \bmod 247 = 77$  ist, sagt uns der Algorithmus von MILLER und RABIN im zweiten Schritt, wenn wir  $77^2 \equiv 1 \bmod 247$  betrachten, daß die Zahl zusammengesetzt sein muß.

Hätten wir allerdings mit  $a = 87$  gearbeitet, hätten wir im nullten Schritt  $87^{123} \equiv 1 \bmod 247$  berechnet und hätten  $247 = 13 \cdot 19$  nicht als zusammengesetzt erkannt.



GARY L. MILLER entwickelte diesen Test 1975 im Rahmen seiner Dissertation (in Informatik) an der Universität von Berkeley. Dabei ging es ihm nicht um einen probabilistischen Test, sondern um einen Test, der immer die richtige Antwort liefert. Er konnte zeigen, daß dies hier beim Test von hinreichend vielen geeigneten Basen der Fall ist **vorausgesetzt** die bis heute immer noch offene verallgemeinerte RIEMANN-Vermutung ist richtig. Er lehrte später zunächst einige Jahre an der University of Waterloo, inzwischen an der Carnegie Mellon University. Seine späteren Arbeiten stammen hauptsächlich aus dem Gebiet der rechnerischen Geometrie. [www.cs.cmu.edu/~glmiller](http://www.cs.cmu.edu/~glmiller)

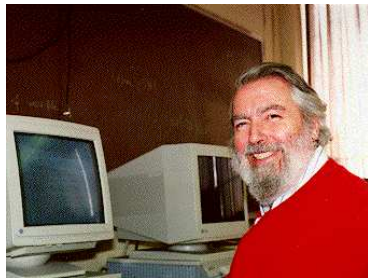


MICHAEL O. RABIN wurde 1931 in Breslau geboren. Die Familie wanderte nach Israel aus, wo er an der hebräischen Universität von Jerusalem Mathematik studierte. Nach seinem Diplom 1953 ging er nach Princeton, wo er 1957 promovierte. Seit 1958 lehrt er an der hebräischen Universität, wo er unter anderem auch Dekan der mathematischen Fakultät und Rektor war. Seit 1983 ist er zusätzlich Inhaber des THOMAS J. WATSON-Lehrstuhls für Informatik an der Harvard University (inzwischen emeritiert). Seine Forschungen, für die er u.a. 1976 den TURING-Preis erhielt, beschäftigen sich mit der Komplexität mathematischer Operationen und der Sicherheit von Informationssystemen. Seine home page in Harvard ist zu finden unter [www.seas.harvard.edu/directory/Rabin](http://www.seas.harvard.edu/directory/Rabin) .

Anscheinend wurde der Test von MILLER und RABIN bereits 1974, also vor MILLERS Veröffentlichung, von SELFRIDGE verwendet; daher sieht



man gelegentlich auch die korrektere Bezeichnung *Test von MILLER, RABIN und SELFRIDGE*.



Der amerikanische Mathematiker JOHN L. SELFRIDGE promovierte 1958 an der University of California in Los Angeles. Bis zu seiner Emeritierung lehrte er an der Northern Illinois University. Seine Arbeiten befassen sich vor allem mit der analytischen sowie der konstruktiven Zahlentheorie. Vierzehn davon schrieb er mit PAUL ERDŐS. [math.niu.edu/faculty/index.php?cmd=detail&id=91](http://math.niu.edu/faculty/index.php?cmd=detail&id=91)

## §5: Der Test von Agrawal, Kayal und Saxena

Im August 2002 fanden NEERAJ KAYAL und NITIN SAXENA, zwei Bachelor-Studenten am Indian Institute of Technology in Kanpur, zusammen mit ihrem Professor MANINDRA AGRAWAL einen Primzahltest, der ebenfalls auf dem kleinen Satz von FERMAT beruht, aber (natürlich auf Kosten eines erheblich größeren Aufwands) immer die richtige Antwort liefert; er ist inzwischen erschienen in

MANINDRA AGRAWAL, NEERAJ KAYAL, NITIN SAXENA: PRIMES is in  $P$ , *Annals of Mathematics* **160** (2004), 781-793.

Selbstverständlich war dies nicht der erste Primzahltest, der deutlich schneller als Probedivisionen zeigt, ob eine gegebene Zahl prim ist oder nicht; er ist auch bei weitem nicht der schnellste solche Test. Er hat aber gegenüber anderen solchen Tests zwei Besonderheiten:

1. Zu seinem Verständnis ist – nach einigen in der letzten Zeit gefundenen Vereinfachungen – nur elementare Zahlentheorie notwendig.
2. Es ist der bislang einzige Test, von dem man beweisen kann, daß seine Laufzeit für  $n$ -stellige Zahlen durch ein Polynom in  $n$  begrenzt werden kann.

Für uns ist vor allem der erste Punkt wichtig; der zweite ist zwar ein für Komplexitätstheoretiker sehr interessantes Ergebnis, hat aber keinerlei praktische Bedeutung: Im Buch

VICTOR SHOUP: A computational Introduction to Number Theory and Algebra, *Cambridge University Press*, <sup>2</sup>2008 (Volltext unter <http://shoup.net/ntb/>),



MANINDRA AGRAWAL erhielt 1986 seinen BTech und 1991 seinen PhD in Informatik am Indian Institute of Technology in Kanpur, wo er – abgesehen von Gastaufenthalten in Madras, Ulm, Princeton und Singapur – seither als Professor lehrt. Seine Arbeiten befassen sich hauptsächlich mit der Komplexität von Schaltungen und von Algorithmen. Für die Arbeit mit KAYAL und SAXENA erhielt er gemeinsam mit diesen unter anderem den GÖDEL-Preis 2006 für die besten Zeitschriftenveröffentlichung auf dem Gebiet der Theoretischen Informatik.

<http://www.cse.iitk.ac.in/users/manindra/>



NEERAJ KAYAL wurde 1979 geboren. Er erhielt 2002 seinen BTech und 2006 seinen PhD bei MANINDRA AGRAWAL am Indian Institute of Technology in Kanpur. Neuere Arbeiten beschäftigen unter anderem sich mit der Komplexität des Isomorphieproblems bei endlichen Ringen sowie der Lösbarkeit von bivariaten Polynomgleichungen über endlichen Körpern. Nach einem kurzen Aufenthalt an der Rutgers University arbeitet er inzwischen bei Microsoft Research.

<http://research.microsoft.com/en-us/people/neeraka/>



NITIN SAXENA wurde 1981 geboren. Er erhielt 2002 seinen Bachelor of Technology und 2006 seinen PhD bei MANINDRA AGRAWAL am Indian Institute of Technology in Kanpur. Während der Arbeit an seiner Dissertation über die Anwendung von Ringhomomorphismen auf Fragen der Komplexitätstheorie besuchte er jeweils ein Jahr lang die Universitäten Princeton und Singapur, danach arbeitete er als Postdoc in der Gruppe *Quantum Computing and Advanced Systems Research* am Centrum voor Wiskunde en Informatica in Amsterdam. 2006–2012 war er *Bonn Junior Fellow* am HAUSDORFF CENTER der Universität Bonn, seit 2013 hat er eine Professur am Indian Institute of Technology in Kanpur. Sein Interesse gilt algorithmischen Verfahren der Algebra und Zahlentheorie sowie Fragen der Komplexitätstheorie.

<http://www.cse.iitk.ac.in/users/nitin/>

dem dieser Paragraph im wesentlichen folgt, argumentiert SHOUP, daß alternative Algorithmen, so man sich auf Zahlen von weniger als  $2^{256}$  Bit beschränkt, durch eine vergleichbare Schranke abgeschätzt werden können, und natürlich sind die Zahlen, mit denen wir es üblicherweise zu tun haben, deutlich kleiner. In der Praxis sind die alternativen Algorithmen deutlich schneller.

( $2^{256}$  liegt knapp über  $10^{77}$ ; derzeitige Schätzungen für die Anzahl der Nukleonen im Universum liegen bei etwa  $10^{80}$ . Damit ist klar, daß kein Computer, der mit irgendeiner Art von heute bekannter Technologie arbeitet, je eine solche Zahl speichern kann, geschweige denn damit rechnen.)

Im folgenden wird es daher nur um eine mathematische Betrachtung des Algorithmus von AGRAWAL, KAYAL und SAXENA gehen; für einen (kurzen und elementaren) Beweis der Komplexitätsaussage sei beispielsweise auf das zitierte Buch von SHOUP verwiesen.

Die Grundidee des Algorithmus steckt im folgenden

**Satz:**  $n > 1$  sei eine natürliche Zahl und  $a \in \mathbb{N}$  sei dazu teilerfremd.  $n$  ist genau dann prim, wenn im Polynomring über  $\mathbb{Z}/n$  gilt:

$$(X + a)^n = X^n + a.$$

*Beweis:* Nach dem binomischen Lehrsatz ist

$$(X + a)^n = X^n + a^n + \sum_{i=1}^{n-1} \binom{n}{i} a^i X^{n-i}.$$

Für eine Primzahl  $n$  gilt nach dem kleinen Satz von FERMAT in  $\mathbb{Z}/n$  die Gleichung  $a^n = a$ . Außerdem ist für  $1 \leq i \leq n - 1$  der Binomialkoeffizient

$$\binom{n}{i} = \frac{n(n-1) \cdots (n-i+1)}{i!}$$

durch  $n$  teilbar, da  $n$  Faktor des Zählers, nicht aber des Nenners ist. Somit verschwinden in  $\mathbb{Z}/n$  alle diese Binomialkoeffizienten, und die Gleichung aus dem Satz ist bewiesen.

Umgekehrt sei  $n$  eine zusammengesetzte Zahl und  $p$  ein Primteiler von  $n$ . Genauer sei  $n = p^e m$  mit einer zu  $p$  teilerfremden Zahl  $m$ . Dann ist der Zähler von  $\binom{n}{p}$  genau durch  $p^e$  teilbar, denn die Faktoren  $(n-1), \dots, (n-p+1)$  sind allesamt teilerfremd zu  $p$ , und der Nenner ist genau durch  $p$  teilbar. Somit ist  $\binom{n}{p}$  zwar durch  $p^{e-1}$  teilbar, nicht aber durch  $p^e$  und damit erst recht nicht durch  $n$ . Wenn wir  $(X+a)^n$  über  $\mathbb{Z}/n$  ausmultiplizieren, kann daher der Summand  $\binom{n}{p} a^p X^{n-p}$  nicht verschwinden, und damit kann die Gleichung aus dem Satz nicht gelten. ■

In dieser Form führt der Satz allerdings noch nicht zu einem praktikablen Primzahltest: Das Ausmultiplizieren von  $(X+a)^n$  führt schließlich auf  $n+1$  Summanden, der Aufwand ist also proportional zu  $n$  und damit vergleichbar damit, daß wir für jede natürliche Zahl  $1 < m < n$  nachprüfen, ob  $n$  ohne Rest durch  $m$  teilbar ist. Die wesentliche neue Idee von AGRAWAL, KAYAL und SAXENA besteht darin zu zeigen, daß es bereits reicht, Gleichungen der im Satz genannten Art modulo eines geeigneten Polynoms  $X^r - 1$  mit einem relativ kleinen Grad  $r$  nachzuprüfen.

Konkret geht ihr Algorithmus folgendermaßen vor:

$n$  sei die zu testende natürliche Zahl, und  $\ell(n) = \lceil \log_2 n \rceil + 1$  sei die Anzahl ihrer Binärziffern.

*1. Schritt:* Stelle sicher, daß  $n$  keine Potenz einer anderen natürlichen Zahl ist.

Das läßt sich beispielsweise dadurch bewerkstelligen, daß man die Quadratwurzel, Kubikwurzel usw. von  $n$  soweit ausrechnet bis man erkennt, daß es sich um keine natürliche Zahl handelt. Der ungünstigste Fall ist offenbar der, daß  $n$  eine Zweierpotenz sein könnte; man muß also bis zur  $\lceil \log_2 n \rceil$ -ten Wurzel gehen.

*2. Schritt:* Finde die kleinste natürliche Zahl  $r > 1$  mit der Eigenschaft, daß entweder  $\text{ggT}(n, r) > 1$  ist oder aber  $\text{ggT}(n, r) = 1$  ist und  $n \bmod r$  in  $(\mathbb{Z}/r)^\times$  eine größere Ordnung als  $4\ell(n)^2$  hat.

Dies geschieht einfach dadurch, daß man die Zahlen  $r = 2, 3, \dots$  allesamt durchprobiert, bis zum ersten mal eine der beiden Bedingungen erfüllt ist. Die Bedingung über die Ordnung der Restklasse von  $n$  in

$(\mathbb{Z}/r)^\times$  prüft kann man schlimmstenfalls dadurch überprüfen, daß man nacheinander die Potenzen von  $n$  mod  $r$  berechnet, bis man entweder eine Eins gefunden hat oder aber der Exponent größer als  $4\ell(n)^2$  ist. Mit Information über die Gruppenordnung von  $(\mathbb{Z}/r)^\times$  geht es natürlich deutlich schneller, und da wir hoffen, ein kleines  $r$  zu finden, sollte  $\varphi(r)$  mit vertretbarem Aufwand berechenbar sein.

3. *Schritt*: Falls  $r = n$ , ist  $n$  prim, und der Algorithmus endet.

In der Tat: Dann haben wir für alle  $r < n$  überprüft, daß  $\text{ggT}(n, r) = 1$  ist. Wenn der Algorithmus etwas taugt, darf er natürlich höchstens für sehr kleine Werte von  $n$  mit diesem Schritt enden.

4. *Schritt*: Falls im zweiten Schritt ein  $r$  gefunden wurde, für das der  $\text{ggT}$  von  $n$  und  $r$  größer als eins ist, haben wir einen Teiler von  $n$  gefunden;  $n$  ist also zusammengesetzt, und der Algorithmus endet.

Andernfalls kennen wir nun eine zu  $n$  teilerfremde Zahl  $r$ , für die  $n$  mod  $r$  in  $(\mathbb{Z}/r)^\times$  eine größere Ordnung als  $4\ell(n)^2$  hat.

5. *Schritt*: Teste für  $j = 1, \dots, \ell \stackrel{\text{def}}{=} 2\ell(n)[\sqrt{r}] + 1$ , ob über  $\mathbb{Z}/n$

$$(X + j)^n \equiv X^n + j \pmod{(X^r - 1)}.$$

Sobald ein  $j$  gefunden wird, für das dies nicht erfüllt ist, endet der Algorithmus mit dem Ergebnis  $n$  ist zusammengesetzt.

Falls nämlich  $n$  eine Primzahl ist, stimmen  $(X + j)^n$  und  $X^n + j$  als Polynome mit Koeffizienten aus  $\mathbb{Z}/n$  nach obigem Satz überein, sind also erst recht auch gleich modulo  $(X^r - 1)$ .

6. *Schritt*: Wenn alle Tests im fünften Schritt bestanden sind, ist  $n$  eine Primzahl.

Dies zu beweisen ist die Hauptarbeit dieses Paragraphen.

Nach den Kommentaren zu den einzelnen Schritten ist klar, daß der Algorithmus für eine Primzahl  $n$  stets das richtige Ergebnis liefert; wir müssen zeigen, daß er auch zusammengesetzte Zahlen stets erkennt.

Sei also  $n$  eine zusammengesetzte Zahl. Falls  $n$  Potenz einer anderen natürlichen Zahl ist, wird dies im ersten Schritt erkannt; wir können und werden im folgenden daher annehmen, daß dies nicht der Fall ist.

Das  $r$  aus dem zweiten Schritt ist auf jeden Fall echt kleiner als  $n$ , denn als zusammengesetzte Zahl hat  $n$  insbesondere einen Teiler  $r < n$ . Der Algorithmus kann daher nicht im dritten Schritt mit der Antwort „ $n$  ist prim“ enden. Falls er im vierten Schritt endet, lieferte der zweite Schritt einen Teiler von  $n$ , und wir erhalten die richtige Antwort „ $n$  ist zusammengesetzt“.

Für den Rest des Paragraphen können wir somit annehmen, daß der zweite Schritt auf ein  $r$  führte, für das  $\text{ggT}(n, r) = 1$  ist. Wir müssen zeigen, daß einer der Tests im fünften Schritt scheitert, daß es also eine natürliche Zahl  $j$  gibt mit

$$1 \leq j \leq \ell \quad \text{und} \quad (X + j)^n \not\equiv X^n + j \pmod{(X^r - 1)} \text{ in } (\mathbb{Z}/n)[X].$$

Wir nehmen an, das sei nicht der Fall, und betrachten einen Primteiler  $p$  von  $n$ . Dieser muß größer als  $r$  sein, denn sonst hätte der Algorithmus bereits mit dem vierten Schritt spätestens bei  $r = p$  geendet.

Jede Kongruenz modulo  $n$  ist erst recht eine Kongruenz modulo  $p$ ; wir können daher davon ausgehen, daß für alle  $j$  mit  $1 \leq j \leq \ell$  gilt

$$(X + j)^n \equiv X^n + j \pmod{(X^r - 1)} \text{ in } \mathbb{F}_p[X].$$

Wenn wir zum Faktoring  $R = \mathbb{F}_p[X]/(X^r - 1)$  übergehen, ist dort also

$$(X + j)^n = X^n + j \quad \text{falls} \quad 1 \leq j \leq \ell.$$

Um diese seltsame Relation genauer zu untersuchen, betrachten wir für jede zu  $r$  teilerfremde natürliche Zahl  $k$  die Abbildung

$$\hat{\sigma}_k: \begin{cases} \mathbb{F}_p[X] \rightarrow R \\ g \mapsto g(X^k) \pmod{(X^r - 1)} \end{cases},$$

die in jedem Polynom  $g$  die Variable  $X$  überall durch  $X^k$  ersetzt.

**Lemma:**  $\hat{\sigma}_k$  ist surjektiv und sein Kern besteht genau aus den Vielfachen des Polynoms  $X^r - 1$ .

*Beweis:* Wir betrachten  $\hat{\sigma}_k$  nur für Indizes  $k$ , die zu  $r$  teilerfremd sind. Zu jedem solchen Index gibt es daher ein  $k'$ , so daß  $kk' \equiv 1 \pmod{r}$  ist,

und modulo  $X^r - 1$  ist damit  $X^{kk'} \equiv X$ . Für ein beliebiges Polynom  $g \in \mathbb{F}_p[X]$  und  $h(X) = g(X^{k'})$  ist daher in  $R$

$$\widehat{\sigma}_k(h) = h(X^k) = g(X^{kk'}) = g(X) = g,$$

die Abbildung ist also surjektiv.

Was ihren Kern betrifft, so enthält er auf jeden Fall  $X^r - 1$  und alle seine Vielfachen, denn

$$\widehat{\sigma}_k(X^r - 1) = (X^{kr} - 1) \bmod (X^r - 1) = 1^k - 1 = 0,$$

da  $X^r \equiv 1 \bmod (X^r - 1)$ .

Umgekehrt sei  $g$  irgendein Polynom aus dem Kern von  $\widehat{\sigma}_k$ . Dann ist das Polynom  $h(X) = g(X^k)$  modulo  $X^r - 1$  gleich dem Nullpolynom, ist also ein Vielfaches von  $X^r - 1$ . Konkret sei  $h = (X^r - 1)f$ . Im Faktoring  $R$  ist dann

$$g(X) = g(X^{kk'}) = h(X^{k'}) = (X^{k'r} - 1)f(X^{k'}) = 0,$$

denn wegen  $X^r = 1$  in  $R$  ist dort  $X^{k'r} - 1 = 0$ .

In  $\mathbb{F}_p[X]$  muß  $g(X)$  daher ein Vielfaches von  $X^r - 1$  sein, und genau das war die Behauptung über den Kern von  $\widehat{\sigma}_k$ . ■

Da alle Vielfachen von  $X^r - 1$  im Kern von  $\widehat{\sigma}_k$  liegen, induziert  $\widehat{\sigma}_k$  eine Abbildung  $\sigma_k$  von  $R$  nach  $R$ , die jedem Polynom  $g \bmod (X^r - 1)$  aus  $R$  das Element  $\widehat{\sigma}_k(g)$  zuordnet; nach dem gerade bewiesenen Lemma hängt dieses wirklich nur von der Restklasse  $g \bmod (X^r - 1)$  ab. Außerdem zeigt das Lemma, daß  $\sigma_k$  sowohl surjektiv als auch injektiv ist, denn der Kern von  $\widehat{\sigma}_k$  ist gleich dem Kern der Restklassenabbildung von  $\mathbb{F}_p[X]$  nach  $R$ . Damit ist  $\sigma_k$  ein bijektiver Homomorphismus von  $R$  nach  $R$ , ein sogenannter *Automorphismus* von  $R$ . Wir haben damit für jede zu  $r$  teilerfremde natürliche Zahl  $k$  einen Automorphismus  $\sigma_k: R \rightarrow R$ , der jedem Polynom in  $X$  das entsprechende Polynom in  $X^k$  zuordnet. Da wir in  $R$  rechnen, werden natürlich alle Polynome modulo  $X^r - 1$  betrachtet.

Unmittelbar aus der Definition folgt, daß die verschiedenen Automorphismen  $\sigma_k$  miteinander kommutieren; genauer ist für zwei zu  $r$  teilerfremde natürliche Zahlen  $k$  und  $k'$

$$\sigma_k \circ \sigma_{k'} = \sigma_{k'} \circ \sigma_k = \sigma_{kk'} ,$$

denn in allen drei Fällen wird im Endeffekt  $X$  durch  $X^{kk'}$  ersetzt.

Speziell für das Element  $X + j$  aus  $R$  ist  $\sigma_k(X + j) = X^k + j$ . Für  $k = n$  und  $j = 1, \dots, \ell$  ist andererseits auch

$$\sigma_n(X + j) = (X + j)^n ,$$

denn für diese  $j$  wurde ja nach unserer Annahme der Test im fünften Schritt bestanden.

Wir wollen genauer untersuchen, wann die Gleichung  $\sigma_k(f) = f^k$  erfüllt ist. Dazu definieren zwei Arten von Mengen:

$$\begin{aligned} C(f) &= \{k \in (\mathbb{Z}/r)^\times \mid \sigma_k(f) = f^k\} && \text{für alle } f \in R && \text{und} \\ D(k) &= \{f \in R \mid \sigma_k(f) = f^k\} && \text{für alle } k \in (\mathbb{Z}/r)^\times \end{aligned}$$

Beide Mengen enthalten mit zwei Elementen auch deren Produkt, denn für zwei Elemente  $k, k' \in C(f)$  ist

$$\sigma_{kk'}(f) = \sigma_k(f)\sigma_{k'}(f) = \sigma_k(\sigma_{k'}(f)) = \sigma_k(f^{k'}) = \sigma_k(f)^{k'} = f^{kk'} ,$$

und für  $f, g \in D(k)$  ist

$$\sigma_k(fg) = \sigma_k(f)\sigma_k(g) = f^k g^k = (fg)^k .$$

Der Rest des Beweises besteht darin, daß wir die „Größe“ der Menge  $D(n)$  auf zwei verschiedene Weisen abschätzen und daraus einen Widerspruch herleiten zur Annahme, daß  $n$  zusammengesetzt ist, aber trotzdem vom Algorithmus als Primzahl klassifiziert wird. Wir definieren zunächst zwei neue Zahlen:

- $s$  sei die Ordnung der Restklasse von  $p$  in  $(\mathbb{Z}/r)^\times$ . Dann ist  $r$  ein Teiler von  $p^s - 1$ , denn  $p^s \equiv 1 \pmod{r}$ .
- $t$  sei die Ordnung der von den Restklassen von  $p$  und  $n$  erzeugten Untergruppe von  $(\mathbb{Z}/r)^\times$ , d.h also die Ordnung der kleinsten Untergruppe, die beide Restklassen enthält. Da diese Untergruppe insbesondere die Restklasse von  $p$  und deren Potenzen enthält, ist  $t$  ein Vielfaches von  $s$ .



Als nächstes betrachten wir einen Körper  $K$  mit  $p^s$  Elementen. Einen solchen Körper kann man konstruieren, indem man den Vektorraum  $\mathbb{F}_p^s$  identifiziert mit dem Vektorraum aller Polynome vom Grad kleiner  $s$  mit Koeffizienten aus  $\mathbb{F}_p$  und dort eine Multiplikation einführt, die zwei Polynomen deren Produkt modulo einem festen irreduziblen Polynom vom Grad  $s$  über  $\mathbb{F}_p$  zuordnet. Man kann zeigen (siehe Algebra-Vorlesung oder entsprechendes Lehrbuch), daß es für jedes  $s$  ein solches Polynom gibt, und daß zwei verschiedene irreduzible Polynome vom Grad  $s$  zu isomorphen Körpern führen.

Aus Kapitel 1 wissen wir, daß die multiplikative Gruppe jedes endlichen Körpers zyklisch ist;  $K^\times$  ist also eine zyklische Gruppe der Ordnung  $p^s - 1$ . Diese Zahl ist, wie wir gerade gesehen haben, ein Vielfaches von  $r$ ; somit gibt es in  $K^\times$  (mindestens) ein Element  $\zeta$  der Ordnung  $r$ . Für irgendein solches Element definieren wir einen Homomorphismus

$$\widehat{\tau}: \begin{cases} \mathbb{F}_p[X] \rightarrow K \\ g \mapsto g(\zeta) \end{cases} .$$

Da  $\widehat{\tau}(X^r - 1) = \zeta^r - 1$  verschwindet, induziert  $\widehat{\tau}$  einen Ringhomomorphismus  $\tau: R \rightarrow K$ . Die angekündigten Abschätzungen der „Größe“ von  $D(n)$  beziehen sich auf die Mächtigkeit der Menge  $S = \tau(D(n))$ :

**Lemma:**  $S = \tau(D(n))$  hat höchstens  $n^{2[\sqrt{t}]}$  Elemente.

*Beweis:* Gemäß unserer Annahme ist  $n$  nicht prim, und da der Algorithmus über den ersten Schritt hinausgekommen ist, kann  $n$  auch keine Potenz einer anderen natürlichen Zahl sein, insbesondere also keine Primzahlpotenz. Daher gibt es außer dem Primteiler  $p$  noch mindestens einen weiteren Primteiler  $q \neq p$ . Wenn wir (in  $\mathbb{N}$ ) Potenzen der Form  $n^u p^v$  und  $n^{u'} p^{v'}$  mit  $u, u', v, v' \in \mathbb{N}_0$  betrachten, sind diese daher genau dann gleich, wenn  $(u, v) = (u', v')$  ist: Ist nämlich  $u \neq u'$ , so tritt  $q$  in der Primzerlegung der beiden Elemente mit verschiedenen Exponenten auf, und ist  $u = u'$ , aber  $v \neq v'$ , so gilt entsprechendes für  $p$ . Daher hat die Menge

$$I = \{n^u p^v \mid 0 \leq u, v \leq [\sqrt{t}]\}$$

mindestens  $([\sqrt{t}]+1)^2$  Elemente, und diese Zahl ist offensichtlich größer als  $t$ , die Ordnung der von  $n$  und  $p$  erzeugten Untergruppe von  $(\mathbb{Z}/r)^\times$ . Daher muß es mindestens zwei Elemente

$$k = n^u p^v \quad \text{und} \quad k' = n^{u'} p^{v'}$$

aus  $I$  geben, die modulo  $r$  dieselbe Restklasse definieren, für die also gilt:  $k \equiv k' \pmod{r}$ . Da die Exponenten  $u, u', v, v'$  höchstens gleich  $[\sqrt{t}]$  sind und  $p$  ein Teiler von  $n$  ist, können wir  $n^{2[\sqrt{t}]}$  als (sehr grobe) obere Schranke für  $k$  und  $k'$  nehmen.

Nun sei  $f \in R$  ein Element von  $D(n)$ . Nach Definition der Mengen  $C(f)$  und  $D(n)$  ist dann auch  $n$  ein Element von  $C(f)$ . Außerdem enthält  $C(f)$  stets die Eins und nach dem kleinen Satz von FERMAT auch die Primzahl  $p$ , denn Potenzieren mit  $p$  ist über  $\mathbb{F}_p$  ein Homomorphismus. Da mit zwei Elementen stets auch deren Produkt in  $C(f)$  liegt, liegen daher die Restklassen modulo  $r$  aller Elemente von  $I$  in  $C(f)$ . Insbesondere sind daher  $k$  und  $k'$  Elemente von  $C(f)$ , d.h.

$$\sigma_k(f) = f^k \quad \text{und} \quad \sigma_{k'}(f) = f^{k'}.$$

Wegen  $k \equiv k' \pmod{r}$  ist aber  $\sigma_k$  dieselbe Abbildung wie  $\sigma_{k'}$ ; daher ist  $f^k = f^{k'}$  für jedes  $f \in D(n)$ . Somit sind die Bilder  $\tau(f)$  aller  $f \in R$  Nullstellen des Polynoms  $X^k - X^{k'}$ . Dessen Grad ist das Maximum von  $k$  und  $k'$ , und da  $\tau(f)$  im Körper  $K$  liegt, gibt es höchstens so viele Nullstellen, wie der Grad angibt. Aufgrund der obigen Abschätzung für  $k$  und  $k'$  hat das Polynom daher höchstens  $n^{2[\sqrt{t}]}$  Nullstellen, und damit kann auch  $S$  nicht mehr Elemente enthalten. ■

Als untere Grenze für die Elementanzahl von  $S$  erhalten wir

**Lemma:**  $S$  enthält mindestens  $2^{\min(t,\ell)} - 1$  Elemente.

*Beweis:* Wegen der bestandenen Tests in Schritt 5 liegt  $\tau(X+j)$  in  $D(n)$  für  $j = 1, \dots, \ell$ . Da  $p > r > t \geq m$  ist, sind die Zahlen von 1 bis  $m$  auch modulo  $p$  paarweise verschieden. Die Teilmenge

$$P = \left\{ \prod_{j=1}^m (X+j)^{e_j} \mid e_j \in \{0, 1\} \text{ und } \sum_{j=1}^m e_j < m \right\}$$

von  $\mathbb{F}_p[X]$  enthält daher  $2^m - 1$  Polynome.

Aus diesen Polynomen können wir Elemente von  $R$  bzw.  $K$  machen, indem wir für die Variable  $X$  die Restklasse  $\eta = X \bmod (X^r - 1)$  bzw. das oben gewählte Element  $\zeta$  der Ordnung  $r$  einsetzen; wir erhalten Teilmengen

$$P(\eta) = \{f(\eta) \mid f \in P\} \subseteq R \quad \text{und} \quad P(\zeta) = \{f(\zeta) \mid f \in P\} \subseteq K.$$

Da sowohl  $n$  als auch  $p$  in  $D(n)$  liegen und mit zwei Elementen auch deren Produkt, liegt  $P(\eta)$  in  $D(n)$  und damit  $\tau(P(\eta)) = P(\zeta)$  in  $S$ . Das Lemma ist daher bewiesen, sobald wir gezeigt haben, daß  $P(\zeta)$  mindestens  $2^m - 1$  Elemente enthält.

Falls dies nicht der Fall wäre, müßte es in  $P$  zwei verschiedene Polynome  $g$  und  $h$  geben, für die  $g(\zeta) = h(\zeta)$  wäre. Wir müssen also zeigen, daß  $g(\zeta) = h(\zeta)$  nur dann gelten kann, wenn  $g = h$  ist.

Wie im vorigen Lemma folgt, da  $1, p$  und  $n$  alle drei sowohl in  $C(g(\eta))$  als auch in  $C(h(\eta))$  liegen, daß alle natürlichen Zahlen  $k$  der Form  $k = n^u p^v$  in diesen beiden Mengen liegen.

Da  $g(\zeta) = h(\zeta)$ , gilt für jedes solche  $k$

$$\begin{aligned} 0 &= g(\zeta)^k - h(\zeta)^k = \tau(g(\eta))^k - \tau(h(\eta))^k = \tau(g(\eta)^k) - \tau(h(\eta)^k) \\ &= \tau(g(\eta^k)) - \tau(h(\eta^k)) = g(\zeta^k) - h(\zeta^k). \end{aligned}$$

Da  $\zeta$  in  $K$  die Ordnung  $r$  hat, hängt  $\zeta^k$  nur von  $k \bmod r$  ab; die Anzahl verschiedener Restklassen der Form  $n^u p^v$  modulo  $r$  hatten wir oben mit  $t$  bezeichnet. Somit hat die Differenz  $g - h$  mindestens  $t$  Nullstellen. Andererseits sind aber  $g$  und  $h$  und damit auch ihre Differenz Polynome vom Grad höchstens  $t - 1$ , also muß  $g - h$  das Nullpolynom sein, d.h.  $g = h$ . Somit enthält  $S$  mindestens  $2^m - 1$  Elemente, wie behauptet. ■

Zum Abschluß des Beweises, daß der Test von AGRAWAL, KAYAL und SAXENA stets die richtige Antwort liefert, müssen wir nun nur noch zeigen, daß die Schranken aus den beiden letzten Lemmata, die ja unter der Voraussetzung bewiesen wurde, daß eine zusammengesetzte Zahl als prim erkannt wird, einander widersprechen, daß also die untere Schranke größer ist als die obere:

**Lemma:**  $2^{\min(t,\ell)} - 1 > n^{2\lceil\sqrt{t}\rceil}$ .

*Beweis:* Da  $\ell(n) > \log_2 n$ , genügt es zu zeigen, daß

$$2^{\min(t,\ell)} - 1 > 2^{2\ell(n)\lceil\sqrt{t}\rceil}.$$

Da beide Exponenten natürliche Zahlen sind, genügt dazu wiederum, daß  $\min(t,\ell) > 2\ell(n)\lceil\sqrt{t}\rceil$  ist, denn wenn sich die Exponenten um mindestens eins unterscheiden, ist die Differenz zwischen den Potenzen mindestens zwei. Wir müssen daher zeigen, daß sowohl  $t$  als auch  $\ell$  größer sind als  $2\ell(n)\lceil\sqrt{t}\rceil$ .

Für  $\ell = 2\ell(n)\lceil\sqrt{r}\rceil + 1$  ist das klar, da  $t$  die Ordnung einer Untergruppe von  $(\mathbb{Z}/r)^\times$  bezeichnet und damit auf jeden Fall kleiner als  $r$  ist.

Die Ungleichung  $t > 2\ell(n)\lceil\sqrt{t}\rceil$  ist sicherlich dann erfüllt, wenn sogar  $t > 2\ell(n)\sqrt{t}$  ist, und dies wiederum ist äquivalent zur Ungleichung  $t > 4\ell(n)^2$ . Nun ist aber  $t$  die Ordnung jener Untergruppe von  $(\mathbb{Z}/r)^\times$ , die von den Restklassen von  $n$  und  $p$  erzeugt wird. Da wir im zweiten Schritt des Algorithmus sichergestellt haben, daß dort allein die Ordnung der Restklasse von  $n$  schon größer ist als  $4\ell(n)^2$ , ist auch die Ungleichung für  $t$  trivial. ■

Damit ist die Korrektheit des Algorithmus vollständig bewiesen.

## Kapitel 4

### Faktorisierungsverfahren

Die MERSENNE-Zahl  $M_{67} = 2^{67} - 1$  ist keine Primzahl, denn

$$13^{M_{67}-1} \equiv 81\,868\,480\,399\,682\,966\,751 \not\equiv 1 \pmod{M_{67}}.$$

Somit ist  $M_{67}$  ein Produkt von mindestens zwei nichttrivialen Faktoren. Welche sind das?

FRANK NELSON COLE gab das Ergebnis am 31. Oktober 1903 auf einer Sitzung der American Mathematical Society bekannt: Er schrieb die Zahl

$$2^{67} - 1 = 147\,573\,952\,589\,676\,412\,927$$

auf eine der beiden Tafeln und

$$193\,707\,721 \times 761\,838\,257\,287$$

auf die andere. Dieses Produkt rechnete er wortlos aus nach der üblichen Schulmethode zur schriftlichen Multiplikation, und als er dieselbe Zahl erhielt, die auf der anderen Tafel stand, schrieb er ein Gleichheitszeichen zwischen die beiden Zahlen und setzte sich wieder. Das Ergebnis, d.h. die Faktorisierung von  $M_{67}$ , findet ein Computeralgebrasystem heute in weniger als einer Sekunde; für die damalige Zeit war es eine Sensation! COLE gab später zu, daß er drei Jahre lang jeden Sonntag Nachmittag daran gearbeitet hatte. Er versuchte  $M_{67}$  in der Form  $x^2 - y^2$  darzustellen, wobei er mit Hilfe quadratischer Reste Kongruenzbedingungen für  $x$  modulo verschiedener relativ kleiner Primzahlen aufstellte und auch verwendete, daß jeder Teiler von  $M_{67}$  kongruent eins modulo 67 und kongruent  $\pm 1$  modulo acht sein muß. Dies führte zu einer ganzen Reihe von Kongruenzen für  $x$ , die er in

$$x \equiv 1\,160\,932\,384 \pmod{1\,323\,536\,760}$$

zusammenfassen konnte. Untersuchung quadratischer Reste zeigt, daß

$$x_k = 1\,323\,536\,760k + 1\,160\,932\,384$$

frühestens für  $k = 287$  in Frage kommt, und mit  $x = x_{287}$  ist tatsächlich

$$\begin{aligned} M_{67} &= 381\,015\,982\,504^2 - 380\,822\,274\,783^2 \\ &= 193\,707\,721 \times 761\,838\,257\,287. \end{aligned}$$

Für Einzelheiten siehe

F. N. COLE: On the factoring of large numbers, *Bull. Am. Math. Soc.* **10** (1903), 134–137 oder <http://www.ams.org/bull/1903-10-03/S0002-9904-1903-01079-9/home.html>



FRANK NELSON COLE (1861–1926) wurde in Massachusetts geboren. 1882 erhielt er seinen Bachelor in Mathematik von der Harvard University; danach konnte er dank eines Stipendiums drei Jahre lang bei FELIX KLEIN in Leipzig studieren. Mit einer von KLEIN betreuten Arbeit über Gleichungen sechsten Grades wurde er 1886 in Harvard promoviert. Nach verschiedenen Positionen in Harvard und Michigan bekam er 1895 eine Professur an der Columbia University in New York, wo er bis zu seinem Tod lehrte. Seine Arbeiten befassen sich hauptsächlich mit Primzahlen und mit der Gruppentheorie.

Der Auftritt von COLE schlug selbst außerhalb der Mathematik so hohe Wellen, daß seine Faktorisierung noch fast ein Jahrhundert später vorkommt in einer New Yorker (off-Broadway) Show von RINNE GROFF mit dem Titel *The five hysterical girls theorem*. Dort bringt sich ein junger Mathematiker um, weil er in einem Beweis von der *Primzahl*  $2^{67} - 1$  ausgeht und die Tochter des Professors die obige Faktorisierung an die Tafel schreibt. Einzelheiten kann man, so man unbedingt möchte, unter <http://www.playscripts.com/play.php3?playid=551> nachlesen. (Die Show verschwand nach zwei Monaten Ende Mai 2000 in der Versenkung; sie wurde seither nur noch zweimal von Amateurgruppen aufgeführt.)

COLE konnte für seine Faktorisierung von  $M_{67}$  auf bekannte Tatsachen über die Struktur von Faktoren der MERSENNE-Zahlen zurückgreifen

und auch bei den von ihm selbst gefundenen Eigenschaften potentieller Faktoren konnte er die spezielle Struktur von  $M_{67}$  ausnutzen. Ähnlich arbeiten auch heutige Mathematiker an der Faktorisierung spezieller Zahlen, beispielsweise im Rahmen des Cunningham-Projekts zur Faktorisierung von Zahlen der Form  $b^n \pm 1$  für kleine Basen  $b$ . Für die Faktorisierung von RSA-Moduln kann man natürlich nicht mit solchen Techniken arbeiten. In diesem Kapitel soll es um Verfahren gehen, mit denen man eine zufällig gegebene Zahl ohne spezielle Struktur faktorisieren kann.

Es gibt kein „bestes“ Faktorisierungsverfahren; für Zahlen verschiedener Größenordnungen haben jeweils andere Verfahren ihre Stärken. Auch Vorwissen über die zu faktorisierende Zahl kann bei der Wahl eines geeigneten Verfahrens helfen: Bei einem RSA-Modul, der das Produkt zweier Primzahlen ähnlicher Größenordnung ist, wird man anders vorgehen als bei einer Zahl der Form  $b^n \pm 1$ . Mehr noch als bei Primzahltests gilt, daß asymptotische Komplexitätsaussagen als Auswahlkriterium nutzlos sind: Das für die Faktorisierung 150-stelliger RSA-Moduln heute optimale Verfahren, das Zahlkörpersieb, wird beim Versuch eine sechsstellige Zahl zu faktorisieren, oft nicht in der Lage sein die Faktoren zu trennen, und selbst in den Fällen, in denen es erfolgreich ist, braucht es erheblich länger als einfache Probedivisionen. Auch liefern die meisten Faktorisierungsverfahren nur *irgendeinen* Faktor; der muß nicht prim sein, und sein Kofaktor schon gar nicht. Falls also ein Faktor  $m$  einer Zahl  $N$  gefunden ist, müssen anschließend  $m$  und  $N/m$  weiter untersucht werden, und da diese Zahlen kleiner sind als  $N$ , sind dazu möglicherweise andere Verfahren besser als das zuerst angewandte.

Im folgenden sollen einige der einfachsten gebräuchlichen Verfahren vorgestellt werden.

## § 1: Die ersten Schritte

### a) Test auf Primzahl

Der schlimmste Fall für praktisch jedes Faktorisierungsverfahren tritt dann ein, wenn die zu faktorisierende Zahl eine Primzahl ist: Gerade

bei den fortgeschrittenen Verfahren gibt es oft kein anderes Abbruchkriterium als das Auffinden eines Faktors. Daher sollte (außer eventuell bei ganz kleinen Zahlen) zu Beginn einer Faktorisierung immer ein Primzahltest stehen. Da auch das Testen auf Potenzen relativ einfach ist, läßt sich eventuell auch das noch durchführen – es sei denn, daß von der Situation her (beispielsweise bei RSA-Moduln) nicht mit einer Potenz zu rechnen ist.

### b) Abdividieren kleiner Primteiler

Bei kleinen zusammengesetzten Zahlen  $n$  besteht die effizienteste Art der Faktorisierung im allgemeinen darin, einfach alle Primzahlen nacheinander durchzuprobieren, indem man sie der Reihe nach so lange abdividiert, wie es geht. Sobald der Quotient kleiner ist als das Quadrat der gerade betrachteten Primzahl, kann man sicher sein, daß auch er eine Primzahl ist und hat  $n$  vollständig faktorisiert.

Die genaue Vorgehensweise ist folgende: Wir nehmen an, daß eine Liste der Primzahlen bis zu einer gewissen Grenze zur Verfügung steht; gegebenenfalls muß diese zunächst nach ERATOSTHENES erzeugt werden.

1. *Schritt:* Setze  $M$  gleich der zu faktorisierte Zahl und  $p = 2$ .
2. *Schritt:* Solange  $M$  durch  $p$  teilbar ist, ersetze  $M$  durch  $M/p$  und notiere  $p$  als Faktor.
3. *Schritt:* Falls  $M = 1$ , sind alle Faktoren gefunden, und der Algorithmus endet. Falls  $M < p^2$ , ist  $M$  eine Primzahl und wird zur Liste der Faktoren hinzugefügt; danach endet auch in diesem Fall der Algorithmus. In allen anderen Fällen wird  $p$  auf die nächste Primzahl gesetzt und es geht zurück zum zweiten Schritt.

Als Beispiel wollen wir die Zahl 1 234 567 890 faktorisieren:

Im ersten Schritt werden  $M = 1\,234\,567\,890$  und  $p = 2$  initialisiert.

Im zweiten Schritt ist nun  $p = 2$ . Da  $M$  eine gerade Zahl ist, können wir durch  $p$  dividieren; wir notieren also die Zwei als Faktor und ersetzen  $M$  durch  $M/2 = 617\,283\,945$ . Diese Zahl ist ungerade; also geht es weiter



zum dritten Schritt, wo offensichtlich keines der Abbruchkriterien erfüllt ist. Somit wird  $p = 3$  und es geht zurück zum zweiten Schritt.

Das neue  $M$  ist durch drei teilbar, genauso auch  $M/3 = 205\,761\,315$ . Also wird nochmals durch drei dividiert, und wir erhalten den nicht mehr durch drei teilbaren Quotienten  $68\,587\,105$ . Somit werden zwei Faktoren drei notiert und es geht weiter zum dritten Schritt. Dort wird  $p = 5$  gesetzt, und es geht wieder zurück zu Schritt 2.

Das aktuelle  $M = 68\,587\,105$  ist durch fünf teilbar;  $M/5 = 13\,717\,421$ . Dies wird das neue  $M$ ; und da es nicht durch fünf teilbar ist, notieren wir nur einen Faktor fünf.

Im dritten Schritt wird wieder  $p$  erhöht und es geht zurück zum zweiten Schritt. Dort passiert nun allerdings lange Zeit nichts, denn keine der Primzahlen zwischen sieben und  $3\,593$  teilt das aktuelle  $M$ . Erst wenn  $p$  im dritten Schritt auf  $3\,607$  gesetzt wird, finden wir wieder einen Faktor. Wir notieren ihn, ersetzen  $M$  durch  $M/3\,607 = 3\,803$ , was offensichtlich nicht durch  $3\,607$  teilbar ist, und gehen weiter zum dritten Schritt.

Dort ist nun offensichtlich  $M < p^2$ , also ist  $M$  eine Primzahl, und

$$1\,234\,567\,890 = 2 \cdot 3^2 \cdot 5 \cdot 3\,607 \cdot 3\,803$$

ist vollständig faktorisiert.

Es ist klar, daß wir schon eine zwanzigstellige Zahl nur mit viel Glück auf diese Weise mit vertretbarem Aufwand vollständig faktorisieren können. Trotzdem ist Abdividieren selbst für noch viel größere Zahlen ein sinnvoller erster Schritt, denn die auf größere Faktoren spezialisierten Verfahren schaffen es im allgemeinen nicht, auch kleine Primfaktoren voneinander zu trennen.

Um Abdividieren statt zur vollständigen Faktorisierung nur zur Identifikation „kleiner“ Primfaktoren zu verwenden, ist lediglich eine kleine Modifikation des dritten Schritts notwendig: Wir legen eine Suchgrenze  $S$  fest und brechen im dritten Schritt auch dann ab, wenn  $p > S$  ist. Im letzteren Fall können wir selbstverständlich nicht behaupten, daß das verbleibende  $M$  eine Primzahl ist;  $M$  muß dann mit anderen Verfahren weiter bearbeitet werden. Bei der Wahl einer geeigneten Schranke  $S$

sollte man die Kapazität des Arbeitsspeichers und die Geschwindigkeit des verwendeten Computers berücksichtigen; ein minimaler Wert, für den der Algorithmus auf allen heutigen Computern in Sekundenbruchteilen ausgeführt werden kann, wäre etwa  $2^{16} = 65\,536$ . Bei etwas besseren Computern läßt sich auch das Abdividieren bis zu einer Million und bei derzeit aktuellen schnellen Computern auch einer Milliarde oder etwa  $2^{30}$  in weniger als einer Sekunde durchführen.

## §2: Die Verfahren von Pollard und ihre Varianten

In den Jahren um 1975 entwickelte der britische Mathematiker JOHN M. POLLARD mehrere recht einfache Algorithmen zur Faktorisierung ganzer Zahlen sowie zur Berechnung diskreter Logarithmen, die auch heute noch (teils in verbesserter Form) zu den Standardwerkzeugen der algorithmischen Zahlentheorie gehören. In diesem Paragraphen sollen die beiden bekanntesten vorgestellt werden; außerdem möchte ich zumindest kurz auf mathematisch anspruchsvollere Verallgemeinerungen eingehen. Die hier behandelten Verfahren haben im Gegensatz zu denen des nächsten Paragraphen die Eigenschaft, daß sie umso schneller zum Erfolg führen, je kleiner die gesuchten Primfaktoren sind, daß sie allerdings sehr kleine Primfaktoren oft nicht finden. Sie sind also die Verfahren der Wahl für die Weiterverarbeitung eines durch Abdividieren erhaltenen „Rests“, von dem man weiß, daß er keine allzu kleinen Primfaktoren mehr hat.

JOHN M. POLLARD ist ein britischer Mathematiker, der hauptsächlich bei British Telecom arbeitete. Er veröffentlichte zwischen 1971 und 2000 rund zwanzig mathematische Arbeiten, größtenteils auf dem Gebiet der algorithmischen Zahlentheorie. Bekannt ist er auch für seine Beiträge zur Kryptographie, für die er 1999 den RSA Award erhielt. Außer den hier vorgestellten Faktorisierungsalgorithmen entwickelte er unter anderem auch das Zahlkörpersieb, eine Variante des weiter hinten vorgestellten quadratischen Siebs, dessen Weiterentwicklungen derzeit die schnellsten Faktorisierungsalgorithmen für große Zahlen sind. Seine home page, um die er sich auch jetzt im Ruhestand noch kümmert, ist [sites.google.com/site/jmptidcott2/](http://sites.google.com/site/jmptidcott2/).

Bei den in diesem und dem nächsten Paragraphen vorgestellten Verfahren besteht das Ziel immer darin, *irgendeinen* Faktor zu finden; sobald dies erreicht ist, bricht das Verfahren ab und der gefundene Faktor sowie

sein Kofaktor werden für sich weiter untersucht – wobei natürlich immer an erster Stelle ein Primzahltest stehen sollte.

### a) Die Monte-Carlo-Methode

Monte Carlo ist ein Stadtteil von Monaco, der vor allem für seine Spielbank bekannt ist. An deren Spieltischen sollen Roulette-Schüsseln idealerweise rein zufällig für jedes Spiel von neuem eine Zahl zwischen 0 und 36 bestimmen.

Eine ähnliche Idee läßt sich auch für die Faktorisierung einer ganzen Zahl  $N$  verwenden: Ausgehend von einer Folge  $(x_i)_{i \in \mathbb{N}}$  zufällig gewählter Zahlen zwischen 1 und  $N$  (oder 0 und  $N - 1$ ) bildet man jeweils den ggT von  $x_i$  mit  $N$  in der Hoffnung, einen nichttrivialen Teiler zu finden.

Für einen Primteiler  $p$  von  $N$  können wir erwarten, daß im Mittel eine von  $p$  Zahlen  $x_i$  durch  $p$  teilbar ist. Dann ist auch  $\text{ggT}(x_i, N)$  durch  $p$  teilbar, kann aber möglicherweise größer als  $p$  sein.

Beim einfachen Abdividieren finden wir  $p$ , nachdem wir alle Primzahlen bis einschließlich  $p$  durchprobiert haben; wie wir im letzten Kapitel gesehen haben, sind dies etwa  $p / \log p$  Stück. Für jede davon brauchen wir eine Division, verglichen mit durchschnittlich  $p/2$  EUKLIDischen Algorithmen bei der obigen Methode, die nicht einmal eine Garantie dafür bietet, den Faktor zu finden. Von daher hat die neue Methode zumindest in der bislang betrachteten Form ausschließlich Nachteile und ist keine sinnvolle Alternative zum Abdividieren.

POLLARDS Idee zur Beschleunigung beruht auf dem im Anhang genauer erklärten Geburtstagsparadoxon: Die Wahrscheinlichkeit dafür, daß eine gegebene Zufallszahl durch  $p$  teilbar ist, liegt zwar nur bei  $1 : p$ , aber die Wahrscheinlichkeit, daß zwei der  $x_i$  modulo  $p$  gleich sind, steigt in der Nähe von etwa  $\sqrt{p}$  Folgengliedern ziemlich steil von nahe null zu nahe eins. Wenn wir also anstelle der größten gemeinsamen Teiler von  $N$  mit den  $x_i$  die mit den Differenzen  $x_i - x_j$  berechnen, haben wir bereits bei einer Folge der Länge um  $\sqrt{p}$  gute Chancen, einen nichttrivialen ggT zu finden.

Auch in dieser Form ist das Verfahren noch nicht praktikabel: Wenn wir ein neues  $x_i$  mit  $i \approx \sqrt{p}$  erzeugt haben, müssen wir für alle  $j < i$  den ggT von  $x_i - x_j$  berechnen, was noch einmal rund  $\sqrt{p}$  Schritte sind, so daß der Gesamtaufwand nicht proportional zu  $\sqrt{p}$  ist, sondern eher zu

$$\int_0^{\sqrt{p}} x \, dx = \frac{p}{2},$$

was keine Ersparnis ist. Dazu kommt, daß alle bereits berechneten Folgenglieder gespeichert werden müssen, der Algorithmus hat also auch einen Platzbedarf in der Größenordnung  $\sqrt{p}$ .

Dieses Problem können wir umgehen, indem wir keine echten Zufallszahlen verwenden, sondern algorithmisch eine Folge sogenannter Pseudozufallszahlen erzeugen. Typischerweise verwendet man dazu eine Rekursionsvorschrift der Form  $x_{i+1} = Q(x_i) \bmod N$  mit einem quadratischen Polynom  $Q$ . (Die bei Simulationen sehr beliebten Pseudozufallsgeneratoren nach der linearen Kongruenzmethode sind für die Monte-Carlo-Methode zur Faktorisierung nicht geeignet.) Meist nimmt man einfach Polynome der Form  $Q(x) = x^2 + c$ , wobei allerdings  $c \neq 0$  und  $c \neq -2$  sein sollte, denn eine genauere Untersuchung zeigt, daß diese Wahlen keine guten Pseudozufallszahlen liefern. Daß die anderen Wahlen von  $c$  stets gute Generatoren liefern ist zwar nicht bewiesen, aber die praktischen Erfahrungen sind positiv.

Wegen der speziellen Form der Rekursion hängt die Restklasse von  $x_{i+1}$  modulo  $p$  nur ab von  $x_i \bmod p$ ; insbesondere ist also  $x_{i+1} \equiv x_{j+1} \bmod p$ , falls  $x_i \equiv x_j \bmod p$ , und entsprechend stimmen auch für jedes  $r \geq 0$  die Zahlen  $x_{i+r}$  und  $x_{j+r}$  modulo  $p$  überein, d.h. die Folge wird modulo  $p$  periodisch mit einer Periode  $\pi$ , die  $|i - j|$  teilt.

Das Problem, Periodizität in einer Folge zu entdecken, tritt nicht nur in der Zahlentheorie auf, sondern beispielsweise auch in der Zeitreihenanalyse und anderen Anwendungen. Ein möglicher Algorithmus zu seiner Lösung, auch als Hase und Schildkröte Algorithmus bekannt, stammt von FLOYD (1967) und beruht auf folgender Beobachtung:

*Wird eine Folge  $(y_i)$  irgendwann periodisch, so gibt es Indizes  $k$  derart, daß  $y_k = y_{2k}$  ist.*

In der Tat, ist  $y_{i+\pi} = y_i$  für alle  $i \geq r$ , so können wir für  $k$  jedes Vielfache  $\ell\pi$  der Periode nehmen, das mindestens gleich  $r$  ist.



ROBERT W. FLOYD (1936–2001) beendete seine Schulausbildung bereits im Alter von 14 Jahren, um dann mit einem Stipendium an der Universität von Chicago zu studieren, wo er mit 17 einen Bachelor in *liberal arts* bekam. Danach finanzierte er sich durch Arbeit ein zweites Bachelorstudium in Physik, das er 1958 abschloß. Damit war seine akademische Ausbildung beendet; er arbeitete als Operator in einem Rechenzentrum, brachte sich selbst Programmieren bei und begann einige Jahre später mit der Publikation wissenschaftlicher Arbeiten auf dem Gebiet der Informatik. Mit 27 wurde er Assistenzprofessor in Carnegie Mellon, fünf Jahre später erhielt er einen Lehrstuhl in Stanford. Zu den vielen Entwicklungen, die er initiierte, gehört die semantische Verifikation von Programmen, Design und

Analyse von Algorithmen, Refactoring, dazu kommen Arbeiten über Graphentheorie und das FLOYD-STEINBERG dithering in der Computergraphik. 1978 erhielt er den TURING-Preis, die höchste Auszeichnung der Informatik. Stanfords Nachruf auf FLOYD ist zu finden unter [news-service.stanford.edu/news/2001/november7/floydobit-117.html](http://news-service.stanford.edu/news/2001/november7/floydobit-117.html).

Damit sieht der Grob Ablauf der Monte-Carlo-Faktorisierung einer natürlichen Zahl  $N$  folgendermaßen aus:

**Schritt 0:** Man wähle ein quadratisches Polynom  $Q$  und einen Startwert  $x_0$ . Setze  $x = y = x_0$ .

**Schritt  $i, i > 0$ :** Ersetze  $x$  durch  $Q(x) \bmod N$  und ersetze  $y$  durch  $Q(Q(y)) \bmod N$ ; berechne dann  $\text{ggT}(x - y, N)$ . Falls dieser weder eins noch  $N$  ist, wurde ein Faktor gefunden.

Man beachte, daß hier im  $i$ -ten Schritt  $x = x_i$  und  $y = x_{2i}$  ist; wir erzeugen also die Folge der  $x_i$  (Schildkröte) und die der  $x_{2i}$  (Hase) simultan, ohne Zwischenergebnisse zu speichern.

Das Teuerste an diesem Algorithmus sind die EUKLIDischen Algorithmen zur ggT-Berechnung; da wir (sofern wir kleine Primfaktoren zuvor ausgeschlossen haben) nicht wirklich erwarten, daß hier häufig ein nicht-triviales Ergebnis herauskommt, liegt es nahe, deren Anzahl möglichst zu reduzieren.

Eine Strategie dazu besteht darin, jeweils mehrere Differenzen  $x_{2i} - x_i$  modulo  $N$  aufzumultiplizieren und dann erst für das Produkt den ggT mit  $N$  zu berechnen. Die „gewisse Anzahl“ darf nicht zu groß sein, denn sonst besteht die Gefahr, daß das Produkt nicht nur durch einen, sondern gleich durch mehrere Primteiler von  $N$  teilbar ist, es sollte aber aus Effizienzgründen auch nicht zu klein sein. Wenn alle kleinen Primteiler bereits ausgeschlossen sind, zeigt die Erfahrung, daß die Zusammenfassung von etwa hundert Differenzen ein guter Kompromiß ist; wenn bereits bei den „kleinen“ Faktoren mit einer hohen Suchgrenze gearbeitet wurde, bieten sich auch höhere Werte an.

Praktisch bedeutet das, daß wir eine neue Variable  $P$  einführen mit Anfangswert eins und dann im  $i$ -ten Schritt  $P$  durch  $P \cdot (x - y) \bmod N$  ersetzen. Nur falls  $i$  durch die „gewisse Anzahl“  $m$  teilbar ist, wird anschließend der ggT von  $N$  und  $P$  berechnet; andernfalls geht es gleich weiter mit dem  $(i + 1)$ -ten Schritt.

Die Monte-Carlo-Methode wird auch als  $\rho$ -Methode bezeichnet, da die Folge der  $x_i \bmod p$  nicht von Anfang an periodisch sein muß. Sie muß aber, da es nur  $p$  Restklassen modulo  $p$  gibt, schließlich periodisch werden, d.h. sie beginnt auf dem unteren Ast des  $\rho$  und mündet irgendwann in den Kreis. Erfahrungsgemäß ist diese Methode sehr erfolgreich im Auffinden sechs- bis achtstelliger Faktoren; danach wird sie recht langsam, und kleine Faktoren kann sie oft nicht trennen.

Als Beispiel wollen wir die sechste FERMAT-Zahl

$$F_6 = 2^{64} + 1 = 18\,446\,744\,073\,709\,551\,617$$

betrachten. Mit dem quadratischen Polynom  $Q(x) = x^2 + 1$ , dem Startwert  $x_0 = 2$  und einem EUKLIDischen Algorithmus nach jeweils hundert Folgegliedern findet ein handelsüblicher PC nach 900 Iterationen in Sekundenbruchteilen den Faktor 274 177, der übrigens genau wie sein Kofaktor 67 280 421 310 721 prim ist. Damit ist  $F_6$  vollständig faktorisiert.

### **Anhang: Das Geburtstagsparadoxon**

Angenommen, in einem Raum befinden sich  $n$  Personen. Wie groß ist die

Wahrscheinlichkeit dafür, daß zwei davon am gleichen Tag Geburtstag haben?

Um diese Frage wirklich beantworten zu können, müßte man die (recht inhomogene) Verteilung der Geburtstage über das Jahr kennen; wir beschränken uns stattdessen auf ein grob vereinfachtes Modell ohne Schaltjahre mit 365 gleich wahrscheinlichen Geburtstagen. Dann ist die Wahrscheinlichkeit dafür, daß von  $n$  Personen keine zwei am gleichen Tag Geburtstage haben,

$$\prod_{k=0}^{n-1} \left(1 - \frac{k}{365}\right),$$

denn für eine Person ist das überhaupt keine Bedingung, und jede weitere Person muß die Geburtstage der schon betrachteten Personen vermeiden. (Da der Faktor mit  $k = 365$  verschwindet, wird die Wahrscheinlichkeit für  $n > 365$  zu null, wie es nach dem DIRICHLETSchen Schubfachprinzip auch sein muß.)

Nachrechnen ergibt für  $n = 23$  ungefähr den Wert 0,4927; bei 23 Personen liegt also die Wahrscheinlichkeit für zwei gleiche Geburtstage bei 50,7%. Tatsächlich dürfte sie noch deutlich höher liegen, denn bei Geburtstagen ist die Annahme einer Gleichverteilung sicherlich falsch.

Bei einer guten Folge von Zufallszahlen sollten die Restklassen modulo  $p$  in sehr guter Näherung gleichverteilt sein; die Wahrscheinlichkeit dafür, daß unter  $n$  Zufallszahlen keine zwei in der gleichen Restklasse liegen, ist somit

$$P_n = \prod_{k=0}^{n-1} \left(1 - \frac{k}{p}\right).$$

Da wir uns für einigermaßen große Werte von  $p$  interessieren (die kleinen haben wir schon abdividiert), können wir davon ausgehen, daß

$$\left(1 - \frac{1}{p}\right)^p \approx e \quad \text{und} \quad \left(1 - \frac{1}{p}\right) \approx e^{-1/p}$$

ist; für nicht zu große Werte von  $k$  ist dann auch

$$\left(1 - \frac{k}{p}\right) \approx e^{-k/p},$$

und für nicht zu große Werte von  $n$  gilt

$$P_n = \prod_{k=0}^{n-1} \left(1 - \frac{k}{p}\right) \approx \prod_{k=0}^{n-1} e^{-k/p} = e^{-\frac{1}{p} \sum_{k=0}^{n-1} k} = e^{-\frac{n(n-1)}{2p}}.$$

Für  $p = 365$  etwa ergibt dies den Näherungswert  $p_{23} \approx 0,499998$  für den korrekten Wert  $0,4927$ .

Wenn wir im Exponenten noch den Term  $n(n-1)$  durch  $n^2$  approximieren, können wir abschätzen, für welches  $n$  die Wahrscheinlichkeit  $P_n$  einen vorgegebenen Wert erreicht:

$$e^{-\frac{n^2}{2p}} = P \iff \frac{n^2}{2p} = -\ln P \iff n = \sqrt{-2p \ln P}.$$

Damit liegt  $P_n$  bei etwa 50%, falls  $n \approx \sqrt{2p \ln 2} \approx 1,177\sqrt{p}$  ist; für  $p = 365$  ergibt dies die immer noch recht gute Näherung  $22,494$ .

Für  $P = 1/1000$  ergibt sich  $n \approx 3,717\sqrt{p}$ , für  $P = 999/1000$  entsprechend  $n \approx 0,0447\sqrt{p}$ . Die Wahrscheinlichkeit dafür, daß es unter  $n$  Zufallszahlen zwei mit derselben Restklasse modulo  $p$  gibt, wechselt also bei der Größenordnung  $n \approx \sqrt{p}$  von sehr unwahrscheinlich zu sehr wahrscheinlich.

## b) Die $(p-1)$ -Methode

POLLARDS zweite Methode beruht auf dem kleinen Satz von FERMAT: Für einen Primteiler  $p$  von  $N$ , ein Vielfaches  $r$  von  $p-1$  und eine zu  $p$  teilerfremde natürliche Zahl  $a$  ist  $a^r \equiv 1 \pmod{p}$ ; der ggT von  $(a^r - 1) \pmod{N}$  und  $N$  ist also durch  $p$  teilbar.

Natürlich ist  $p-1$  nicht bekannt, wir können aber hoffen, daß  $p-1$  nur durch vergleichsweise kleine Primzahlen teilbar ist. Sei etwa  $B$  eine Schranke mit der Eigenschaft, daß  $p-1$  durch keine Primzahlpotenz größer  $B$  teilbar ist. Dann ist das Produkt  $r$  aller Primzahlpotenzen  $q^e$ , die höchstens gleich  $B$  sind, sicherlich ein Vielfaches von  $p-1$ , wenn auch ein extrem großes, das sich kaum mit realistischem Aufwand berechnen läßt. Für jedes konkrete  $a$  kann  $a^r \pmod{N}$  jedoch verhältnismäßig einfach berechnet werden: Man potenziert einfach nacheinander für jede Primzahl  $q \leq B$  modulo  $N$  mit deren größter Potenz, die



immer noch kleiner oder gleich  $B$  ist; mit dem Algorithmus zur modularen Exponentiation aus §5 des zweiten Kapitels geht das auch für sechs- bis siebenstellige Werte von  $B$  noch recht flott.

Insgesamt funktioniert POLLARDS  $(p - 1)$ -Methode zur Faktorisierung einer natürlichen Zahl  $N$  also folgendermaßen:

**Schritt 0:** Wähle eine Schranke  $B$  und eine Basis  $a$  zwischen 1 und  $N$ .

**Schritt 1:** Erstelle (z.B. nach ERATOSTHENES) eine Liste aller Primzahlen  $q \leq B$ .

**Schritt 2:** Berechne für jede dieser Primzahlen  $q$  den größten Exponenten  $e$  derart, daß auch noch  $q^e \leq B$  ist, d.h.  $e = \lfloor \log B / \log q \rfloor$ . Ersetze dann den aktuellen Wert von  $a$  durch  $a^{q^e} \bmod N$ .

**Schritt 3:** Berechne  $\text{ggT}(a - 1, N)$ . Falls ein Wert ungleich eins oder  $N$  gefunden wird, war das Verfahren erfolgreich, ansonsten nicht.

Es ist klar, daß der Erfolg dieses Verfahrens wesentlich davon abhängt, daß  $N$  einen Primteiler  $p$  hat mit der Eigenschaft, daß alle Primfaktoren von  $p - 1$  relativ klein sind. Ob dies der Fall ist, läßt sich im Voraus nicht sagen; die  $(p - 1)$ -Methode liefert daher gelegentlich ziemlich schnell sogar 20- oder 30-stellige Faktoren, während sie andererseits deutlich kleinere Faktoren oft nicht findet.

Als Beispiel betrachten wir noch einmal  $M_{67} = 2^{67} - 1$ . Wenn wir mit der Basis  $a = 17$  und der Schranke  $B = 3\,000$  arbeiten, wird  $a$  modulo  $M_{67}$  potenziert zum neuen

$$a = 111\,153\,665\,932\,902\,146\,348 \text{ mit } \text{ggT}(a - 1, M_{67}) = 193\,707\,721.$$

Damit ist (in Sekundenbruchteilen auf einem Standard-PC) eine nicht-triviale Faktorisierung gefunden, und ein Primzahltest zeigt, daß sowohl der gefundene Faktor als auch sein Komplement prim sind.

Warum die Methode Erfolg hatte, sehen wir an der Faktorisierung der um eins verminderten Faktoren:

$$193\,707\,720 = 2^3 \cdot 3^3 \cdot 5 \cdot 67 \cdot 2\,677 \quad \text{und}$$

$$761\,838\,257\,286 = 2 \cdot 3^2 \cdot 29 \cdot 67 \cdot 2\,551 \cdot 8\,539.$$

Für jede Schranke  $B \geq 2\,677$  ist also der erste Faktor ein Teiler des endgültigen  $a - 1$ , aber für  $B < 8\,539$  ist der zweite Faktor keiner.

### c) Varianten

Falls  $p - 1$  nicht nur relativ kleine Primfaktoren hat, führt die  $(p - 1)$ -Methode nicht zum Erfolg. In solchen Fällen kann man aber hoffen, daß vielleicht  $p + 1$  oder irgendeine andere Zahl in der Nähe von  $p$  nur kleine Primfaktoren hat. Wir brauchen daher Varianten der  $(p - 1)$ -Methode, bei denen es nicht auf die Primfaktoren von  $p - 1$  ankommt, sondern auf die anderer Zahlen in der Nähe von  $p$ .

Um solche Varianten zu finden, empfiehlt es sich, zunächst die  $(p - 1)$ -Methode etwas abstrakter unter gruppentheoretischen Gesichtspunkten zu betrachten.

Dort rechnen wir in der primen Restklassengruppe  $(\mathbb{Z}/N)^\times$  und damit implizit auch in  $(\mathbb{Z}/p)^\times$  für jeden Primteiler  $p$  von  $N$  – egal ob wir ihn kennen, oder nicht. In  $(\mathbb{Z}/p)^\times$  ist für jedes Element  $a$  die  $(p - 1)$ -te Potenz gleich dem Einselement; genau dasselbe gilt für jede  $r$ -te Potenz, für die der Exponent  $r$  ein Vielfaches von  $(p - 1)$  ist. Bei der  $(p - 1)$ -Methode wird ein  $r$  berechnet, das durch alle Primzahlpotenzen bis zu einer gewissen Schranke teilbar ist; falls in der Primzerlegung von  $p - 1$  keine Primzahlpotenz oberhalb der Schranke liegt, ist  $r$  ein Vielfaches von  $p - 1$ .

Allgemeiner können wir statt in  $(\mathbb{Z}/N)^\times$  und  $(\mathbb{Z}/p)^\times$  auch in einem anderen Paar von Gruppen rechnen: Wir gehen aus von einer endlichen Gruppe  $G_N$ , deren Elemente sich in irgendeiner Weise als  $r$ -tupel über  $(\mathbb{Z}/N)$  auffassen lassen; außerdem nehmen wir an, daß sich die Gruppenmultiplikation für zwei so dargestellte Elemente auf Grundrechenarten über  $\mathbb{Z}/N$  zurückführen läßt. Dann können wir die Elemente von  $G_N$  zu Tupeln über  $\mathbb{Z}/p$  reduzieren und die Menge aller so erhaltenen Tupel bildet eine Gruppe  $G_p$ . Wieder ist jede Rechnung in  $G_N$  implizit auch eine Rechnung in  $G_p$ .

Die Elementanzahl von  $G_p$  sei  $N(p)$ .

Wir wählen irgendein Element von  $G_N$  und potenzieren es mit demselben Exponenten  $r$ , mit dem wir bei der  $p - 1$ -Methode die Zahl  $a$  modulo  $N$  potenziert haben. Falls  $r$  ein Vielfaches von  $N(p)$  ist, erhalten wir ein Element  $b \in G_N$ , dessen Reduktion modulo  $p$  das Einselement

von  $G_p$  ist. Ist daher  $b_i$  die  $i$ -te Koordinate von  $b$  und  $e_i$  die von  $e$ , so muß die Differenz  $b_i - e_i$  durch  $p$  teilbar sein, und mit etwas Glück können wir  $p$  als ggT von  $n$  und  $b_i - e_i$  bestimmen.

Bleibt nur noch das Problem, geeignete Gruppen zu finden. Bei der  $(p - 1)$ -Methode ist  $G_N = (\mathbb{Z}/N)^\times$  und  $N(p) = p - 1$ .

Für die  $(p + 1)$ -Methode benutzt POLLARD die Tatsache, daß es nicht nur zu jeder Primzahl  $p$ , sondern auch zu jeder Primzahlpotenz  $p^r$  einen Körper mit entsprechender Elementanzahl. Dieser Körper  $\mathbb{F}_{p^r}$  ist natürlich verschieden vom Ring  $\mathbb{Z}/p^r$ ; er ist ein  $r$ -dimensionaler Vektorraum über  $\mathbb{F}_p$  mit geeignet definierter Multiplikation.

Speziell für  $r = 2$  hat der Körper  $\mathbb{F}_{p^2}$  eine multiplikative Gruppe  $\mathbb{F}_{p^2}^\times$  der Ordnung  $p^2 - 1 = (p + 1)(p - 1)$ . Sie hat  $\mathbb{F}_p^\times$  als Untergruppe und die Faktorgruppe  $G_p = \mathbb{F}_{p^2}^\times / \mathbb{F}_p^\times$  hat die Ordnung  $N(p) = p + 1$ . Das Rechnen in dieser Gruppe mit Repräsentanten modulo  $N$  ist etwas trickreich und benutzt die hier nicht behandelten LUCAS-Sequenzen.

Derzeit am populärsten ist eine andere Wahl von  $G_N$  und  $G_p$ : Wir nehmen für  $G_N$  eine elliptische Kurve über  $\mathbb{Z}/N$ . Dabei handelt es sich um die Menge aller Punkte  $(x, y) \in (\mathbb{Z}/N)^2$ , die einer vorgegebenen Gleichung  $y^2 = x^3 - ax - b$  genügen, wobei  $a, b$  Elemente von  $\mathbb{Z}/N$  sind, für die  $\Delta = 4a^3 - 27b^2$  teilerfremd zu  $N$  ist; dazu kommt ein weiterer Punkt  $O$ , den wir formal als  $(0, \infty)$  schreiben.  $G_p$  ist dann die entsprechende Punktmenge in  $\mathbb{F}_p^2$  zusammen mit  $O$ . Nach einem Satz von HELMUT HASSE (1898–1979) ist

$$p + 1 - 2\sqrt{p} < N(p) < p + 1 + 2\sqrt{p},$$

und wie man inzwischen weiß, kann man auch für fast jeden Wert, der diese Ungleichung erfüllt, Parameterwerte  $a$  und  $b$  finden, so daß  $N(p)$  gleich diesem Wert ist. Wenn man mit hinreichend vielen verschiedenen Kurven arbeitet, ist daher die Chance recht groß, daß der Exponent  $r$  wenigstens für eine davon ein Vielfaches von  $N(p)$  ist.

Die Multiplikation ist folgendermaßen definiert: Durch zwei Punkte  $(x_1, y_1)$  und  $(x_2, y_2)$  auf der Kurve geht genau eine Gerade; setzt man deren Gleichung  $y = mx + c$  in die Kurvengleichung ein, erhält man

ein Polynom dritten Grades in  $x$ . Dieses hat natürlich die beiden Nullstellen  $x_1, x_2$ , und daneben noch eine dritte Nullstelle  $x_3$ . Der dritte Schnittpunkt der Geraden mit der Kurve ist somit  $(x_3, mx_3 + c)$ ; als Summe der beiden Punkte definiert man aber

$$(x_1, y_1) \oplus (x_2, y_2) = (x_3, -(mx_3 + c)).$$

Man kann zeigen, daß dies die Menge der Kurvenpunkte zu einer Gruppe mit Neutralelement  $O$  macht, in der man genauso vorgehen kann wie bei der klassischen  $(p - 1)$ -Methode.

Unter den Faktorisierungsmethoden, deren Rechenzeit von der Größe des zu findenden Faktors abhängt, ist die Faktorisierung mit elliptischen Kurven die für große Zahlen derzeit beste bekannte Methode; sie fand schon Faktoren mit bis zu 67 Stellen. Produkte zweier ungefähr gleich großer Primzahlen wie beispielsweise RSA-Moduln sind für solche Methoden allerdings der schlechteste Fall; hierfür sind andere Methoden, deren Aufwand nur von der Größe der zu faktorisierenden Zahl abhängt, meist besser geeignet.

### §3: Das Verfahren von Fermat und seine Varianten

Die bisher betrachteten Verfahren funktionieren vor allem dann gut, wenn die zu faktorisierende Zahl mindestens einen relativ kleinen Primteiler hat. Das hier beschriebene Verfahren von FERMAT führt genau dann schnell ans Ziel, wenn sie sich als Produkt zweier fast gleich großer Faktoren schreiben läßt. In seiner einfachsten Form beruht es auf der dritten binomischen Formel  $x^2 - y^2 = (x + y)(x - y)$ : Ist  $N = pq$  Produkt zweier ungerader Primzahlen, so ist

$$N = (x + y)(x - y) \quad \text{mit} \quad x = \frac{p + q}{2} \quad \text{und} \quad y = \frac{p - q}{2};$$

Ausmultiplizieren führt auf die Beziehung  $N + y^2 = x^2$ .

FERMAT berechnet für  $y = 0, 1, 2, \dots$  die Zahlen  $N + y^2$ ; falls er auf ein Quadrat  $x^2$  stößt, hat er zwei Faktoren  $x \pm y$  gefunden. Da  $y$  gleich der halben Differenz der beiden Faktoren ist, kommt er umso schneller ans Ziel, je näher die beiden Faktoren beieinander liegen; dies erklärt die

Vorschrift, daß die beiden Faktoren eines RSA-Moduls zwar ungefähr gleich groß sein sollten, daß sie aber doch einen gewissen Mindestabstand einhalten müssen.

FERMAT sucht nach Zahlen  $x, y$  mit  $x^2 - y^2 = N$ . Alternativ könnten wir uns auch damit begnügen, Zahlen  $x, y$  zu finden mit  $x^2 \equiv y^2 \pmod{N}$ . Für diese gibt es ein  $k \in \mathbb{N}$ , so daß

$$kN = x^2 - y^2 = (x + y)(x - y)$$

ist, und wenn es auch keine Garantie gibt, daß mindestens eine der beiden Zahlen  $\text{ggT}(x \pm y, N)$  ein echter Faktor von  $N$  ist, sind doch bei hinreichend vielen Paaren  $(x, y)$  mit  $x^2 \equiv y^2 \pmod{N}$  die Chancen recht gut, daß mindestens eines davon einen nichttrivialen Faktor liefert.

Im letzten Jahrhundert wurden eine ganze Reihe von Verfahren entwickelt, die solche Paare  $(x, y)$  liefern; die effizientesten davon sind Varianten und Weiterentwicklungen des sogenannten *quadratischen Siebs*.

Seine Grundidee ist die folgende: Man wähle zwei Polynome  $f, g \in \mathbb{Z}[X]$  derart, daß zwar  $f(x) \equiv g(x)^2 \pmod{N}$  für alle  $x \in \mathbb{Z}$ , aber  $f(x) \neq g(x)^2$ . Der Name *quadratisches Sieb* kommt daher, daß  $f$  üblicherweise quadratisch ist wie im klassischen Beispiel

$$f(x) = \left(x + \left\lceil \sqrt{N} \right\rceil\right)^2 - N \quad \text{und} \quad g(x) = x + \left\lceil \sqrt{N} \right\rceil.$$

Da  $f(x)$  im Allgemeinen keine Quadratzahl ist, liefert die Kongruenz  $f(x) \equiv g(x)^2 \pmod{N}$  noch keinen Ansatz zur Faktorisierung.

Falls wir allerdings Werte  $x_1, x_2, \dots, x_r$  finden können, für die das Produkt der  $f(x_i)$  als Quadrat einer bekannten Zahl  $x$  dargestellt werden kann, ist

$$x^2 = \prod_{i=1}^r f(x_i) \equiv \left(\prod_{i=1}^r g(x_i)\right)^2 = y^2 \pmod{N}$$

eine Relation der gesuchten Art. Es reicht dabei natürlich, die Produkte auf der linken und auf der rechten Seite modulo  $N$  auszurechnen, denn  $\text{ggT}(x \pm y, N)$  hängt nur ab von  $x$  und  $y$  modulo  $N$ .

Um geeignete  $x_i$  zu finden, betrachten wir eine Menge  $\mathcal{B}$  von Primzahlen, die sogenannte Faktorbasis, und suchen nach Zahlen  $x_i$ , für die sich

$f(x_i)$  als Produkt von Primzahlen aus  $\mathcal{B}$  schreiben läßt, wobei eventuell auch einige Primzahlen in einer (niedrigen) Potenz auftreten können. Ist  $f(x_i) = \prod_{p \in \mathcal{B}} p^{e_{ip}}$ , so ist

$$\prod_{i=1}^r f(x_i)^{\varepsilon_i} = \prod_{p \in \mathcal{B}} p^{\sum_{i=1}^r e_{ip} \varepsilon_i}$$

genau dann ein Quadrat, wenn  $\sum_{i=1}^r e_{ip} \varepsilon_i$  für alle  $p \in \mathcal{B}$  gerade ist. Dies hängt natürlich nur ab von den  $\varepsilon_i \pmod 2$  und den  $e_{ip} \pmod 2$ ; wir können  $\varepsilon_i$  und  $e_{ip}$  daher als Elemente des Körpers mit zwei Elementen auffassen und bekommen dann über  $\mathbb{F}_2$  die Bedingungen

$$\sum_{i=1}^r e_{ip} \varepsilon_i = 0 \quad \text{für alle } p \in \mathcal{B}.$$

Betrachten wir die  $\varepsilon_i$  als Variablen, ist dies ein homogenes lineares Gleichungssystem in  $r$  Variablen mit soviel Gleichungen, wie es Primzahlen in der Faktorbasis gibt. Dieses Gleichungssystem hat nichttriviale Lösungen, falls die Anzahl der Variablen die der Gleichungen übersteigt, falls es also mehr Zahlen  $x_i$  gibt, für die  $f(x_i)$  über der Faktorbasis faktorisiert werden kann, als Primzahlen in der Faktorbasis.

Für jede nichttriviale Lösung ist

$$\prod_{i=1}^r f(x_i)^{\varepsilon_i} = \prod_{i=1}^r \left( x + \left[ \sqrt{N} \right] \right)^{2\varepsilon_i} \pmod N$$

eine Relation der Form  $x^2 \equiv y^2 \pmod N$ , die mit einer Wahrscheinlichkeit von etwa ein halb zu einer Faktorisierung von  $N$  führt. Falls wir zehn linear unabhängige Lösungen des Gleichungssystems betrachten, führt also mit einer Wahrscheinlichkeit von etwa 99,9% mindestens eine davon zu einer Faktorisierung.

Da  $\varepsilon_i$  nur die Werte 0 und 1 annimmt, stehen in obigem Produkt natürlich keine echten Potenzen: Man multipliziert einfach nur die Faktoren miteinander, für die  $\varepsilon_i = 1$  ist. Außerdem interessieren nicht die links- und rechtsstehenden Quadrate, sondern deren Quadratwurzeln;

tatsächlich also berechnet man (hier natürlich in  $\mathbb{N}_0$ )

$$x = \prod_{p \in \mathcal{B}} p^{\frac{1}{2} \sum_{i=1}^r \varepsilon_i e_{ip}} \pmod{N} \quad \text{und} \quad y = \prod_{i=1}^r \left( x + \left[ \sqrt{N} \right] \right)^{\varepsilon_i} \pmod{N}.$$

Typischerweise enthält  $\mathcal{B}$  für die Faktorisierung einer etwa hundertstelligen Zahl zwischen 100 und 120 Tausend Primzahlen, deren größte somit, wie die folgende Tabelle zeigt, im einstelligen Millionenbereich liegt.

$n$	$n$ -te Primzahl	$n$	$n$ -te Primzahl
100 000	1 299 709	600 000	8 960 453
200 000	2 750 159	700 000	10 570 841
300 000	4 256 233	800 000	12 195 257
400 000	5 800 079	900 000	13 834 103
500 000	7 368 787	1 000 000	15 485 863

$\mathcal{B}$  enthält allerdings nicht *alle* Primzahlen unterhalb einer gewissen Schranke, sondern nur die, für die  $f(x)$  eine Chance hat, durch  $p$  teilbar zu sein. Wie wir im Kapitel über quadratische Reste sehen werden, lassen sich diese im Falle eines quadratischen Polynoms leicht bestimmen.

Der Anteil der Zahlen  $x$  aus einem gegebenen Siebintervall, für die  $f(x)$  nur Primteiler aus  $\mathcal{B}$  hat, ist umso größer, je kleiner die von  $f$  angenommenen Werte in diesem Intervall sind. Dies erklärt die Wahl des obigen Beispielpolynoms  $f$ : Wir betrachten die Zahlen

$$f(x) = \left( x + \left[ \sqrt{N} \right] \right)^2 - N = x^2 + 2x \left[ \sqrt{N} \right] + \left[ \sqrt{N} \right]^2 - N.$$

Für  $x$ -Werte, die deutlich kleiner als  $\sqrt{N}$  sind (und nur mit solchen werden wir es im folgenden zu tun haben) liegen diese in der Größenordnung  $\sqrt{N}$ ; bei den meisten anderen Polynomen, die ein Quadrat minus  $N$  produzieren, wäre entweder die zu quadrierende Zahl oder deren Funktionswert deutlich größer.

Zum besseren Verständnis des Grundprinzips wollen wir versuchen, damit die Zahl 15 zu faktorisieren. Dies ist zwar eine sehr untypische Anwendung, da das quadratische Sieb üblicherweise erst für mindestens

etwa vierzigstellige Zahlen angewandt wird, aber zumindest das Prinzip sollte auch damit klarwerden.

Als Faktorbasis verwenden wir die Menge

$$\mathcal{B} = \{2, 3, 7, 11\};$$

die Primzahl fünf fehlt, da  $3 \cdot 5 = 15$  ist und daher bei einer Faktorbasis, die sowohl drei als auch fünf enthält, die Gefahr zu groß ist, daß die linke wie auch die rechte Seite der Kongruenz durch fünfzehn teilbar ist. Bei realistischen Anwendungen muß man auf solche Überlegungen keine Rücksicht nehmen, denn dann sind die Elemente der Faktorbasis höchstens siebenstellig und somit erheblich kleiner als die gesuchten Faktoren.

Wir berechnen  $f(x)$  für  $x = 1, 2, \dots$ , bis wir einige Funktionswerte haben, die über der Faktorbasis faktorisiert werden können. Die faktorisierbaren Werte sind in folgender Tabelle zusammengestellt:

$x$	$x + \left\lceil \sqrt{N} \right\rceil$	$f(x)$	Faktorisierung
1	4	1	
3	6	21	$3 \cdot 7$
5	8	49	$7^2$
6	9	66	$2 \cdot 3 \cdot 11$
10	13	154	$2 \cdot 7 \cdot 11$
54	57	3234	$2 \cdot 3 \cdot 7^2 \cdot 11$

Die erste und die dritte Zeile sind selbst schon Relationen der gesuchten Art, nämlich

$$4^2 \equiv 1 \pmod{15} \quad \text{und} \quad 8^2 \equiv 7^2 \pmod{15}.$$

Die zweite Relation ist nutzlos, denn  $8 - 7 = 1$  und  $8 + 7 = 15$ . Die erste dagegen führt zur Faktorisierung, denn

$$\text{ggT}(4 + 1, 15) = 5 \quad \text{und} \quad \text{ggT}(4 - 1, 15) = 3.$$

Da dies aber ein Zufall ist, der bei großen Werten von  $N$  so gut wie nie vorkommt, wollen wir das ignorieren und mit den Relationen zu  $x = 3, 6, 10$  und  $51$  arbeiten:

$$6^2 \equiv 3 \cdot 7 \pmod{15}$$



$$\begin{aligned} 9^2 &\equiv 2 \cdot 3 \cdot 11 \pmod{15} \\ 13^2 &\equiv 2 \cdot 7 \cdot 11 \pmod{15} \\ 57^2 &\equiv 2 \cdot 3 \cdot 7^2 \cdot 11 \pmod{15} \end{aligned}$$

Multipliziert man die ersten drei dieser Relationen miteinander, folgt

$$(6 \cdot 9 \cdot 13)^2 \equiv (2 \cdot 3 \cdot 7 \cdot 11)^2 \pmod{15}$$

oder  $702^2 \equiv 462^2 \pmod{15}$ . Da

$$\text{ggT}(702 - 462, 15) = \text{ggT}(240, 15) = 15$$

ist, bringt das leider nichts.

Wir erhalten auch dann rechts ein Quadrat, wenn wir das Produkt der ersten, dritten und vierten Relation bilden; dies führt auf

$$(6 \cdot 13 \cdot 57)^2 \equiv (2 \cdot 3 \cdot 7^2 \cdot 11)^2 \pmod{15}$$

oder  $4446^2 \equiv 3234^2 \pmod{15}$ . Hier ist

$$\text{ggT}(4446 - 3234, 15) = \text{ggT}(1212, 15) = 3,$$

womit wir die Zahl 15 faktorisiert haben – wenn auch nicht unbedingt auf die einfachstmögliche Weise.

Bei realistischen Beispielen sind die Funktionswerte  $f(x)$  deutlich größer als die Primzahlen aus der Faktorbasis; außerdem liegen die vollständig faktorisierbaren Zahlen viel dünner als hier: Bei der Faktorisierung einer hundertstelligen Zahl etwa muß man davon ausgehen, daß nur etwa jeder  $10^9$ -te Funktionswert über der Faktorbasis zerfällt.

Daher ist es wichtig, ein Verfahren zu finden, mit dem diese wenigen Funktionswerte schnell und einfach bestimmt werden können. Das ist zum Glück möglich:

Der Funktionswert  $f(x)$  ist genau dann durch  $p$  teilbar, wenn

$$f(x) \equiv 0 \pmod{p}$$

ist. Für ein Polynom  $f$  mit ganzzahligen Koeffizienten ist offensichtlich  $f(x) \equiv f(y) \pmod{p}$ , falls  $x \equiv y \pmod{p}$  ist. Daher ist für ein  $x$  mit  $f(x) \equiv 0 \pmod{p}$  auch

$$f(x + kp) \equiv 0 \pmod{p} \quad \text{für alle } k \in \mathbb{Z}.$$

Es genügt daher, im Bereich  $0 \leq x < p - 1$  nach Werten zu suchen, für die  $f(x)$  durch  $p$  teilbar ist.

Dazu kann man  $f$  auch als Polynom über dem Körper mit  $p$  Elementen betrachten und nach Nullstellen in diesem Körper suchen. Für Polynome großen Grades und große Werte von  $p$  kann dies recht aufwendig sein; hier, bei einem quadratischen Polynom, müssen wir natürlich einfach eine quadratische Gleichung lösen: In  $\mathbb{F}_p$  wie in jedem anderen Körper auch gilt

$$f(x) = \left(x + \left[\sqrt{N}\right]\right)^2 - N = 0 \iff \left(x + \left[\sqrt{N}\right]\right)^2 = N,$$

und diese Gleichung ist genau dann lösbar, wenn es ein Element  $w \in \mathbb{F}_p$  gibt mit Quadrat  $N$ , wenn also in  $\mathbb{Z}$

$$w^2 \equiv N \pmod{p}$$

ist. Für  $p > 2$  hat  $f(x) = 0$  in  $\mathbb{F}_p$  dann die beiden Nullstellen

$$x = -\left[\sqrt{N}\right] \pm w;$$

andernfalls gibt es keine Lösung.

Insbesondere kann also  $f(x)$  nur dann durch  $p$  teilbar sein, wenn  $N$  modulo  $p$  ein Quadrat ist; dies ist für etwa die Hälfte aller Primzahlen der Fall. Offensichtlich sind alle anderen Primzahlen nutzlos, und sollten daher gar nicht erst in die Faktorbasis aufgenommen werden.

Im Kapitel über quadratische Reste werden wir sehen, daß sich auch für große  $N$  und  $p$  leicht und schnell entscheiden läßt, ob  $N$  modulo  $p$  ein Quadrat ist; der Aufwand entspricht ungefähr dem eines auf  $N$  und  $p$  angewandten EUKLIDischen Algorithmus. Wir werden dort auch sehen, daß sich für solche Quadrate relativ schnell die beiden Quadratwurzeln modulo  $p$  berechnen lassen.

Das eigentliche Sieben zum Auffinden der komplett über der Faktorbasis zerlegbaren Funktionswerte  $f(x)$  geht dann folgendermaßen vor sich: Man legt ein Siebintervall  $x = 0, 1, \dots, M$  fest und speichert in einem Feld der Länge  $M + 1$  für jedes  $x$  eine Approximation von  $\log_2 |f(x)|$ .

Für jede Primzahl  $p$  aus der Faktorbasis berechnet man dann die beiden Nullstellen  $x_{1/2}$  von  $f$  modulo  $p$  im Intervall von 0 bis  $p - 1$  und

subtrahiert von jedem Feldelement mit Index der Form  $x_1 + kp$  oder  $x_2 + kp$  eine Approximation von  $\log_2 p$ .

Falls  $f(x)$  über der Faktorbasis komplett faktorisierbar ist, sollte dann am Ende der entsprechende Feldeintrag bis auf Rundungsfehler gleich null sein; um keine Fehler zu machen, untersucht man daher für alle Feldelemente, die betragsmäßig unterhalb einer gewissen Grenze liegen, durch Abdividieren, ob sie wirklich komplett faktorisieren, und man bestimmt auf diese Weise auch *wie* sie faktorisieren. Damit läßt sich dann das oben erwähnte Gleichungssystem über  $\mathbb{F}_2$  aufstellen und, falls genügend viele Relationen gefunden sind, nichttrivial so lösen, daß eine der daraus resultierenden Gleichungen  $x^2 \equiv y^2 \pmod p$  zu einer nichttrivialen Faktorisierung von  $N$  führt.

Als zwar immer noch untypisch kleines Beispiel, das besser und schneller durch Abdividieren faktorisiert werden könnte, betrachten wir die Zahl  $N = 5\,352\,499$ . Wir nehmen als Faktorbasis alle Primzahlen kleiner hundert modulo derer  $N$  ein Quadrat ist; Nachrechnen zeigt, daß

$$\mathcal{B} = \{3, 5, 11, 13, 17, 19, 23, 31, 41, 43, 53, 59, 83, 89\}$$

dann 14 Elemente enthält. Für jedes davon müssen wir die quadratische Gleichung  $f(x) \equiv 0 \pmod p$  lösen, was hier natürlich selbst durch Ausprobieren recht schnell möglich wäre. Die Lösungsmengen sind

$p =$	3	5	11	13	17	19	23
Lösungen	{1, 2}	{0, 4}	{0, 5}	{4, 11}	{6, 9}	{2, 8}	{0, 20}
$p =$	31	41	43	53	59	83	89
Lösungen	{27, 28}	{3, 4}	{26, 35}	{39, 52}	{2, 33}	{35, 70}	{23, 68}

Wenn wir damit das Intervall der natürlichen Zahlen von 1 bis 20 000 sieben, erhalten wir 18 Zahlen, die über der Faktorbasis komplett zerfallen:

$i$	$x_i$	$f(x_i)$	Faktorisierung
1	23	104397	$3 \cdot 17 \cdot 23 \cdot 89$
2	121	571857	$3 \cdot 11 \cdot 13 \cdot 31 \cdot 43$
3	533	2747217	$3 \cdot 11 \cdot 17 \cdot 59 \cdot 83$
4	635	3338205	$3 \cdot 5 \cdot 13 \cdot 17 \cdot 19 \cdot 53$
5	741	3974417	$31 \cdot 41 \cdot 53 \cdot 59$

6	895	4938765	$3 \cdot 5 \cdot 13 \cdot 19 \cdot 31 \cdot 43$
7	2013	13361777	$11 \cdot 13 \cdot 41 \cdot 43 \cdot 53$
8	2185	14879505	$3 \cdot 5 \cdot 17 \cdot 23 \cdot 43 \cdot 59$
9	2477	17591601	$3 \cdot 31 \cdot 43 \cdot 53 \cdot 83$
10	2649	19268945	$5 \cdot 19 \cdot 43 \cdot 53 \cdot 89$
11	4163	36586077	$3 \cdot 11 \cdot 19 \cdot 23 \cdot 43 \cdot 59$
12	4801	45256497	$3 \cdot 11 \cdot 13 \cdot 31 \cdot 41 \cdot 83$
13	5497	55643601	$3 \cdot 13 \cdot 17 \cdot 23 \cdot 41 \cdot 89$
14	6253	68023857	$3 \cdot 11 \cdot 19 \cdot 23 \cdot 53 \cdot 89$
15	10991	171643917	$3 \cdot 17 \cdot 23 \cdot 41 \cdot 43 \cdot 83$
16	11275	179281245	$3 \cdot 5 \cdot 11 \cdot 13 \cdot 19 \cdot 53 \cdot 83$
17	14575	279852045	$3 \cdot 5 \cdot 11 \cdot 17 \cdot 19 \cdot 59 \cdot 89$
18	18535	429286605	$3 \cdot 5 \cdot 11 \cdot 23 \cdot 31 \cdot 41 \cdot 89$

Ein Vektor  $\vec{\varepsilon} \in \mathbb{F}_2^{18}$ , für den  $\prod f(x_i)^{\varepsilon_i}$  ein Quadrat ist, löst daher über  $\mathbb{F}_2$  das lineare Gleichungssystem mit Matrix

$$\begin{pmatrix} 1 & 1 & 1 & 1 & 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 \end{pmatrix}.$$

Der GAUSS-Algorithmus führt auf die Lösungen

$$(\sigma, \mu + \rho, \mu + \rho, \mu + \rho, \sigma + \rho + \tau, \nu + \sigma, \lambda + \rho, \sigma, \\ \nu + \mu + \lambda, \nu + \sigma + \mu + \tau, \mu + \sigma + \tau, \lambda, \nu + \sigma, \mu, \nu, \rho, \sigma, \tau)$$

mit sechs freien Parametern  $\lambda, \mu, \nu, \rho, \sigma, \tau \in \mathbb{F}_2$ . Setzen wir zunächst

$\lambda = \mu = \sigma = 1, \nu = \rho = \tau = 0$ , so erhalten wir die Lösung

$$\vec{\varepsilon} = (1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 0).$$

Sie führt auf die beiden Zahlen

$$x = \prod_{p \in \mathcal{B}} p^{\frac{1}{2} \sum_{i=1}^r \varepsilon_i e_{ip}} \pmod{N} = 854\,237 \quad \text{und}$$

$$y = \prod_{i=1}^r \left( x + \left[ \sqrt{N} \right] \right)^{\varepsilon_i} \pmod{N} = 3\,827\,016.$$

Leider ist die Differenz dieser beiden Zahlen teilerfremd zu  $N$ .

Setzen wir in einem zweiten Versuch  $\nu = 1$  statt  $\nu = 0$ , so erhalten wir die weitere Lösung  $\vec{\varepsilon} = (1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0)$ , die uns die Zahlen

$$x = \prod_{p \in \mathcal{B}} p^{\frac{1}{2} \sum_{i=1}^r \varepsilon_i e_{ip}} \pmod{N} = 1\,020\,903 \quad \text{und}$$

$$y = \prod_{i=1}^r \left( x + \left[ \sqrt{N} \right] \right)^{\varepsilon_i} \pmod{N} = 4\,093\,611$$

liefert. Nun ist der ggT der Differenz mit  $N$  gleich 1 237, womit wir die Faktorisierung  $5\,352\,499 = 1237 \times 4327$  gefunden haben. In diesem Fall sind die beiden Faktoren sogar bereits Primzahlen; das wird natürlich im allgemeinen nicht der Fall sein – dieses  $N$  habe ich allerdings als Produkt zweier Primzahlen konstruiert.

Natürlich hätte uns jede der bisher behandelten Methoden dieses Ergebnis mit erheblich geringerem Aufwand und auch erheblich schneller geliefert; das quadratische Sieb entwickelt seine Stärken erst bei erheblich größeren Zahlen, für die es dann oft tagelang rechnet.

Dabei verwendet man das quadratische Sieb meist nicht in der hier vorgestellten Einfachstversion, sondern mit verschiedenen Optimierungen.

Bei realistischen Anwendungen wird der überwiegende Teil der Rechenzeit für das Sieben gebraucht. Dies läßt sich relativ einfach parallelisieren, indem man das Sieben für verschiedene Teilintervalle auf verschiedene Computer verteilt. Auf diese Weise können mehrere Tausend

Computer jeweils ein Teilintervall sieben und anschließend die gefundenen Faktorisierungen an eine Zentrale melden. Sobald genügend viele eingegangen sind, kann diese ein lineares Gleichungssystem aufstellen und dieses lösen.

Eine weitere Verbesserung, die zu kürzeren Suchintervallen und damit auch kleineren Zahlen führt, besteht darin, anstelle des einen Polynoms  $f$  mehrere Polynome zu verwenden. Auch diese können wieder auf verschiedene Computer verteilt werden. Falls einige der Polynome auch negative Werte annehmen können, muß auch das berücksichtigt werden, indem man bei der Faktorisierung der nach dem Sieben übrig gebliebenen Zahlen auch noch die  $-1$  als zusätzliche „Primzahl“ in die Faktorbasis aufnimmt.

Ab etwa 120 bis 130 Stellen wird eine Variante schneller, bei der auch mit komplizierteren als nur quadratischen Polynomen gearbeitet wird, das sogenannte Zahlkörpersieb. Es hat seinen Namen daher, daß die dahinterstehende Theorie mit algebraischen Zahlkörpern arbeitet; konkret gerechnet wird allerdings weiterhin mit ganzen Zahlen. Dieses Zahlkörpersieb ist die derzeit beste bekannte Methode zur Faktorisierung von Zahlen, die Produkte zweier Primzahlen ähnlicher Größenordnung sind; von diesem Verfahren geht also die größte Gefahr für RSA aus. Der derzeitige Rekord für dieses Verfahren ist die im Dezember 2009 gefundene Faktorisierung einer zweihundertzweiunddreißigstelligen *challenge number* der Firma RSA durch ein internationales Team; dazu wurde von August 2007 bis April 2009 auf verschiedenen Clustern von Computern gesiebt. Einzelheiten sind unter <http://eprint.iacr.org/2010/006.pdf> zu finden.

Mit heutiger Hardware und heutigen Methoden sollte eine entsprechende Faktorisierung auch für eine Zahl mit etwa Tausend Bit (dreihundert Dezimalstellen) möglich sein, allerdings hat bislang nirgends in einer öffentlich zugänglichen Quelle davon berichtet. Wahrscheinlich wäre ein entsprechendes Projekt für die beteiligten Forscher einfach zu langweilig.

# Kapitel 5

## Kettenbrüche

### §1: Der Kettenbruchalgorithmus

Der EUKLIDISCHE Algorithmus läßt sich auch verwenden, um eine reelle Zahl durch Brüche zu approximieren. Beginnen wir der Einfachheit halber mit einer rationalen Zahl  $\alpha = \frac{n}{m}$  mit  $n, m \in \mathbb{N}$ . Der erste Schritt des EUKLIDISCHEN Algorithmus dividiert  $n$  durch  $m$ :

$$n : m = q_0 \text{ Rest } r_1 \implies \alpha = \frac{n}{m} = q_0 + \frac{r_1}{m} .$$

Falls  $r_1 \neq 0$  ist, wird im zweiten Schritt  $m$  durch  $r_1$  dividiert:

$$m : r_1 = q_1 \text{ Rest } r_2 \implies \frac{m}{r_1} = q_1 + \frac{r_2}{r_1} \implies \alpha = q_0 + \frac{1}{q_1 + \frac{r_2}{r_1}} .$$

Ist auch noch  $r_2$  von Null verschieden, wird sodann  $r_1$  durch  $r_2$  dividiert:

$$r_1 : r_2 = q_2 \text{ Rest } r_3 \implies \frac{r_1}{r_2} = q_2 + \frac{r_3}{r_2} \implies \alpha = q_0 + \frac{1}{q_1 + \frac{1}{q_2 + \frac{r_3}{r_2}}} ,$$

und so weiter. Die Konstruktion muß nach endlich vielen Schritten abbrechen, denn die Folge der Reste  $r_i$  beim EUKLIDISCHEN Algorithmus ist monoton fallend und muß daher schließlich Null erreichen. Damit ist  $\alpha$  dargestellt als ein sogenannter **Kettenbruch**.

Wir können die Konstruktion auch so formulieren, daß sie nur von der Zahl  $\alpha = \frac{n}{m}$  abhängt: Der Quotient bei der Division mit Rest von  $n$

durch  $m$  ist  $q_0 = [\alpha]$ , und der durch  $m$  dividierte Rest ist  $\alpha - q_0$ . Dies führt zu folgender Formulierung des Algorithmus:

Setze zur Initialisierung  $c_0 = [\alpha]$  und schreibe

$$\alpha = c_0 + \alpha_1 \quad \text{mit} \quad 0 \leq \alpha_1 < 1.$$

Im  $i$ -ten Schritt,  $i \geq 1$ , bricht der Algorithmus ab, falls  $\alpha_i$  verschwindet; andernfalls wird  $c_i$  definiert als größte ganze Zahl kleiner oder gleich  $1/\alpha_i$  und  $\alpha_{i+1}$  so, daß gilt

$$\frac{1}{\alpha_i} = c_i + \alpha_{i+1}.$$

Offensichtlich ist dann

$$\begin{aligned} \alpha &= c_0 + \alpha_0 = c_0 + \frac{1}{c_1 + \alpha_1} = c_0 + \frac{1}{c_1 + \frac{1}{c_2 + \alpha_2}} \\ &= \dots = c_0 + \frac{1}{c_1 + \frac{1}{c_2 + \frac{1}{c_3 + \frac{1}{\ddots + \frac{1}{c_{r-1} + \frac{1}{c_r + \alpha_r}}}}} . \end{aligned}$$

Da diese Darstellung sehr viel Platz verbraucht, verwendet man dafür oft auch die kompaktere Schreibweise

$$\alpha = [c_0, c_1, \dots, c_r; \alpha_r].$$

Falls der Algorithmus mit  $\alpha_r = 0$  abbricht, steht im untersten Bruch natürlich nur  $c_r$  im Nenner, und wir schreiben den Kettenbruch kurz als  $[c_0, c_1, \dots, c_r]$ .

So, wie der Algorithmus jetzt formuliert ist, können wir ihn auch auf irrationale Zahlen  $\alpha$  anwenden. Dann kann kein  $\alpha_r$  verschwinden, denn sonst hätten wir ja eine Darstellung von  $\alpha$  als rationale Zahl. Wir können aber nach dem  $r$ -ten Schritt abbrechen und den Bruch betrachten, der entsteht, wenn wir  $\alpha_r = 0$  setzen. Diesen Bruch bezeichnen wir als die  $r$ -te **Konvergente** der Kettenbruchentwicklung von  $\alpha$ .



Als Beispiel betrachten wir  $\alpha = \sqrt{2}$ . Hier ist  $c_0 = [\sqrt{2}] = 1$  und  $\alpha_1 = \sqrt{2} - 1$ . Also ist

$$\frac{1}{\alpha_1} = \frac{1}{\sqrt{2} - 1} = \frac{\sqrt{2} + 1}{(\sqrt{2} - 1)(\sqrt{2} + 1)} = \sqrt{2} + 1,$$

d.h.  $c_1 = [1 + \sqrt{2}] = 2$  und  $\alpha_2 = 1 + \sqrt{2} - 2 = \sqrt{2} - 1 = \alpha_1$ . Damit wiederholt sich ab jetzt alles, d.h.

$$\sqrt{2} = 1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2 + \dots}}}}}$$

In Analogie zu periodischen Dezimalbrüchen schreibt man dies auch kurz in der Form

$$\sqrt{2} = [1, 2, 2, 2, \dots] = [1, \bar{2}].$$

Die ersten Partialbrüche sind

$$P_0 = 1, \quad P_1 = 1 + \frac{1}{2} = 1,5, \quad P_2 = 1 + \frac{1}{2 + \frac{1}{2}} = \frac{7}{5} = 1,4,$$

$$P_3 = 1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}} = \frac{17}{12} = 1,41\bar{6} \quad \text{und} \quad P_4 = 1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2}}}} = \frac{41}{29},$$

was ungefähr gleich 1,4137931 ist. Die Fehler  $\sqrt{2} - P_n$  sind, gerundet auf sechs Nachkommastellen, die Zahlen

$$0,414214, \quad -0,085786, \quad 0,014214, \quad -0,002453 \quad \text{und} \quad 0,000420;$$

verglichen mit den kleinen Nennern 1, 2, 5, 12 und 29 haben wir also erstaunlich gute Übereinstimmungen, und im übrigen ist auch die Kettenbruchentwicklung erheblich regelmäßiger als die Dezimalbruchdarstellung von  $\sqrt{2}$ .

Als zweites Beispiel betrachten wir  $\alpha = \pi$ ; hier erhalten wir zunächst  $c_0 = 3$  und  $\alpha_1 = \pi - 3 \approx 0,14159$ , sodann

$$c_1 = \left[ \frac{1}{\pi - 3} \right] = 7 \quad \text{und} \quad \alpha_2 \approx 0,062513285.$$

Im nächsten Schritt ist  $c_2 = \left[ \frac{1}{\alpha_2} \right] = 15$  und  $\alpha_3 \approx 0,99659976$ . Weiter geht es mit  $c_3 = 1$ ,  $c_4 = 292$ ,  $c_5 = c_6 = c_7 = 1$ ,  $c_8 = 2$  und  $c_9 = 1$ . Ein Muster ist weder erkennbar, noch ist eines bekannt.

Die Kettenbruchentwicklung von  $\pi$  ist somit

$$\pi = [3, 7, 15, 1, 292, 1, 1, 1, 2, 1, \dots].$$

Die ersten Partialbrüche und ihre Differenzen von  $\pi$  sind

$3$	$3 \frac{1}{7}$	$3 \frac{15}{106}$	$3 \frac{16}{113}$	$3 \frac{4687}{33102}$
$0,14$	$-0,0013$	$8,3 \cdot 10^{-5}$	$-2,7 \cdot 10^{-7}$	$5,8 \cdot 10^{-10}$

Auch hier haben wir wieder, verglichen mit der Größe des Nenners, exzellente Approximationseigenschaften.

## §2: Geometrische Formulierung

Wir wollen uns zunächst überlegen, daß die Konvergenten der Kettenbruchentwicklung einer irrationalen Zahl stets die bei vorgegebener Größenordnung des Nenners bestmögliche rationale Approximation dieser Zahl liefern.

Dazu betrachten wir (im wesentlichen nach dem Ansatz von HAROLD STARK in seinem Buch *An Introduction to Number Theory*, MIT Press, 1978) das Problem der rationalen Approximation von der geometrischen Seite: Zur reellen Zahl  $\alpha > 0$  haben wir die Gerade  $y = \alpha x$  durch den Nullpunkt, und offensichtlich ist  $\alpha$  genau dann rational, wenn auf dieser Geraden außer dem Nullpunkt noch ein weiterer Punkt  $(q, p)$  mit ganzzahligen Koordinaten liegt. Rationale Approximationen erhalten wir durch Punkte  $(q, p) \in \mathbb{Z} \times \mathbb{Z}$ , die in der Nähe der Geraden liegen.

(Die Reihenfolge der Koordinaten mag auf den ersten Blick verwundern; sie kommt daher, daß wir die Steigung  $\alpha$  der Geraden durch die Steigung des Ortsvektors zum Punkt  $(q, p)$  annähern wollen, und die ist  $p/q$ .)

Die folgende Konstruktion liefert Punkte  $P_n$  nahe der Geraden, die für gerade  $n$  stets unterhalb  $y = \alpha x$  liegen und für ungerade  $n$  darüber:

Wir starten mit  $P_{-2} = (1, 0)$  und  $P_{-1} = (0, 1)$ .

Zu zwei Punkten  $P = (q, p)$  und  $P' = (q', p')$ , die auf verschiedenen Seiten der Geraden liegen, gibt es stets eine nichtnegative ganze Zahl  $c \in \mathbb{N}_0$ , so daß  $P + cP'$  entweder auf der Geraden liegt oder aber auf derselben Seite wie  $P$ , während  $P + (c+1)P'$  auf der anderen Seite liegt.

Liegt nämlich beispielsweise  $P$  unterhalb der Geraden, so ist  $p/q < \alpha$ , also  $p - \alpha q < 0$ . Für den oberhalb der Geraden liegenden Punkt  $P'$  ist entsprechend  $p' - \alpha q' > 0$ . Damit ist klar, daß

$$c = \left\lfloor \frac{p - \alpha q}{p' - \alpha q'} \right\rfloor$$

das Verlangte leistet. Man überlegt sich leicht, daß diese Formel auch gilt, wenn  $P$  oberhalb und  $P'$  unterhalb der Geraden liegt.

Zähler und Nenner des obigen Bruchs lassen sich einfach geometrisch interpretieren:  $(q, \alpha q)$  hat dieselbe  $x$ -Koordinate wie  $P = (q, p)$  und liegt auf der Geraden  $y = \alpha x$ ; daher ist  $p - \alpha q$  der (gerichtete) vertikale Abstand von  $P$  zur Geraden und  $p' - \alpha q'$  entsprechend der von  $P'$ .

Ausgehend von  $P = P_{-2} = (1, 0)$  und  $P' = P_{-1} = (0, 1)$  definieren wir nun die Punkte  $P_n$  für  $n \geq 0$  mit dem wie oben definierten  $c = c_n$  aus ihren beiden Vorgängern rekursiv als

$$P_n = P_{n-2} + c_n P_{n-1}. \quad (*)$$

Dann liegt  $P_n$  auf derselben Seite der Geraden wie  $P_{n-2}$ , für gerades  $n$  also unterhalb und für ungerades oberhalb – es sei denn, irgendwann einmal liegt ein  $P_n$  auf der Geraden. In diesem Fall ist  $\alpha$  rational und wir brechen die Konstruktion ab. Für irrationales  $\alpha$  erhalten wir eine unendliche Folge von Punkten  $P_n$ .

Bezeichnen wir mit  $d_n = p_n - \alpha q_n$  den gerichteten vertikalen Abstand des Punktes  $P_n = (q_n, p_n)$  von der Geraden  $y = \alpha x$ , so ist nach obiger Formel

$$c_n = \left[ \left| \frac{d_{n-2}}{d_{n-1}} \right| \right].$$

Daher verschwindet  $c_n$  genau dann, wenn  $|d_{n-2}| < |d_{n-1}|$  ist.

Ist dagegen  $|d_{n-1}| < |d_{n-2}|$ , so ist  $c_n \geq 1$ , und da  $P_n = P_{n-2} + c_n P_{n-1}$  auf derselben Seite der Geraden liegt wie  $P_{n-2}$ , ist auch

$$\begin{aligned} d_n &= d_{n-2} + c_n d_{n-1} = d_{n-2} + \left[ \left| \frac{d_{n-2}}{d_{n-1}} \right| \right] d_{n-1} \\ &= d_{n-1} \left( \left[ \left| \frac{d_{n-2}}{d_{n-1}} \right| \right] + \frac{d_{n-2}}{d_{n-1}} \right) \end{aligned}$$

betragsmäßig kleiner als  $d_{n-1}$ . (Man beachte, daß  $d_{n-1}$  und  $d_{n-2}$  verschiedene Vorzeichen haben!) Falls daher für einen Index  $n$  der Abstand von  $P_{n-1}$  zur Geraden  $y = \alpha x$  kleiner ist als der von  $P_{n-2}$ , gilt dasselbe auch für alle folgenden Indizes, und ab dem Index  $n$  sind alle  $c_i \geq 1$ .

Die ersten beiden Abstände sind  $d_{-2} = -\alpha$  und  $d_{-1} = 1$ ; es hängt von  $\alpha$  ab, welche der beiden Zahlen den größeren Betrag hat.

Der nächste Punkt ist  $P_0 = (1, c_0)$  mit  $c_0 = [\alpha]$ , also ist  $d_0 = [\alpha] - \alpha$ , und der Betrag davon ist kleiner als  $d_{-1} = 1$ . Somit ist für alle  $n \geq 1$  der Koeffizient  $c_n$  von Null verschieden und  $|d_n| < |d_{n-1}|$ .

Wegen (\*) ist  $p_n = p_{n-2} + c_n p_{n-1}$  und  $q_n = q_{n-2} + c_n q_{n-1}$ ; daran sehen wir, daß die Folge der  $q_n$  wie auch der  $p_n$  für  $n \geq 1$  strikt monoton ansteigt, während die Folge der Differenzen

$$\left| \alpha - \frac{p_n}{q_n} \right| = \frac{|\alpha q_n - p_n|}{q_n} = \frac{|d_n|}{q_n}$$

strikt monoton fällt. Die Brüche  $p_n/q_n$  geben also immer bessere Annäherungen an  $\alpha$ .

Wir können die beiden Rekursionsformeln für  $p_n$  und für  $q_n$  zusammenfassen zur Matrixgleichung

$$\begin{pmatrix} p_n & q_n \\ p_{n-1} & q_{n-1} \end{pmatrix} = \begin{pmatrix} c_n & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} p_{n-1} & q_{n-1} \\ p_{n-2} & q_{n-2} \end{pmatrix};$$

wenden wir darauf den Multiplikationssatz für Determinanten an, erhalten wir die Formel

$$p_n q_{n-1} - q_n p_{n-1} = -(p_{n-1} q_{n-2} - q_{n-1} p_{n-2}).$$

Für  $n = 0$  ist  $p_{-1} q_{-2} - q_{-1} p_{-2} = 0 \cdot 0 - 1 \cdot 1 = -1$ ; daraus folgt induktiv die Formel  $p_n q_{n-1} - q_n p_{n-1} = (-1)^{n-1}$ . Insbesondere sind die Zahlen  $p_n$  und  $q_n$  stets teilerfremd,  $p_n/q_n$  ist also stets ein gekürzter Bruch.

Als nächstes wollen wir uns überlegen, daß die Folge dieser Brüche gegen  $\alpha$  konvergiert. Da  $P_n$  und  $P_{n+1}$  auf verschiedenen Seiten der Geraden  $y = \alpha x$  liegen, ist für  $n \geq 0$

$$\left| \alpha - \frac{p_n}{q_n} \right| \leq \left| \frac{p_{n+1}}{q_{n+1}} - \frac{p_n}{q_n} \right| = \left| \frac{p_{n+1} q_n - q_{n+1} p_n}{q_n q_{n+1}} \right| < \frac{1}{q_n^2}.$$

Da die Folge der  $q_n$  strikt monoton ansteigt, konvergiert die Folge der  $p_n/q_n$  somit gegen  $\alpha$ , und dies sogar extrem gut: Ist  $p/q$  eine rationale Approximation einer irrationalen Zahl  $\alpha$ , so kann der Fehler im allgemeinen bis zu  $1/2q$  betragen; hier ist er höchstens  $1/q^2$  und tatsächlich wohl, da wir recht grob abgeschätzt haben, meist noch kleiner. Wie wir gleich sehen werden, muß umgekehrt  $p/q$  eine Konvergente der Kettenbruchentwicklung von  $\alpha$  sein, wenn  $|\alpha - p/q| < 1/2q^2$  ist.

Zuvor müssen wir uns aber noch überlegen, daß die hier betrachteten Brüche  $p_n/q_n$  tatsächlich die Konvergenten der in §1 definierten Kettenbruchentwicklung sind und daß die hier betrachteten Zahlen  $c_i$  mit denen übereinstimmen, die der Kettenbruchalgorithmus liefert.

Dazu setzen wir

$$\alpha_n = \left| \frac{d_{n-1}}{d_{n-2}} \right| = -\frac{d_{n-1}}{d_{n-2}};$$

zumindest für  $n \geq 1$  ist dann  $\alpha_n < 1$ . Wegen  $c_n = \lceil d_{n-2}/d_{n-1} \rceil$  ist dann  $c_n = [1/\alpha_n]$ . Division der Beziehung  $d_n = d_{n-2} + c_n d_{n-1}$  durch  $d_{n-1}$  führt auf

$$\frac{d_n}{d_{n-1}} = \frac{d_{n-2}}{d_{n-1}} + c_n \quad \text{oder} \quad -\alpha_{n+1} = -\frac{1}{\alpha_n} + c_n,$$

was wir wiederum umformen können zu

$$\frac{1}{\alpha_n} = c_n + \alpha_{n+1}.$$

Da  $c_n = [1/\alpha_n]$  ist  $\alpha = c_0 + \alpha_1$ , führt dies genau auf die in §1 konstruierten Folgen der  $c_n$  und  $\alpha_n$ .

Insbesondere liefern unsere Rekursionsformeln für die Koordinaten  $p_n$  und  $q_n$  Zähler und Nenner der Konvergenten der Kettenbruchentwicklung von  $\alpha$ , was wir als Satz festhalten wollen:

**Satz:** Ist  $\alpha = [c_0, c_1, c_2, \dots]$  die Kettenbruchentwicklung einer reellen Zahl, so lassen sich die Konvergenten  $p_n/q_n$  folgendermaßen rekursiv berechnen:

$$p_0 = c_0, q_0 = 1, p_1 = c_0c_1 + 1, q_1 = c_1,$$

$$p_n = p_{n-2} + c_n p_{n-1} \quad \text{und} \quad q_n = q_{n-2} + c_n q_{n-1} \quad \text{für } n \geq 2.$$

Die so berechneten Zahlen  $p_n$  und  $q_n$  sind stets teilerfremd; genauer ist  $p_n q_{n-1} - q_n p_{n-1} = (-1)^{n-1}$  für alle  $n$ .

Zu *beweisen* gibt es hier nichts, denn wir haben alle diese Formeln bereits bewiesen für die Koordinaten  $q_n$  und  $p_n$  von  $P_n$ , und wie wir gerade gesehen haben, sind das der Nenner und der Zähler der  $n$ -ten Konvergenten. ■

Für spätere Anwendungen wollen wir noch eine Formel herleiten, wie sich  $\alpha$  aus  $\alpha_n$  sowie den Konvergenten  $p_{n-1}/q_{n-1}$  und  $p_{n-2}/q_{n-2}$  berechnet läßt: Nach der Definition beim geometrischen Zugang ist

$$\alpha_n = -\frac{d_{n-1}}{d_{n-2}} = -\frac{p_{n-1} - \alpha q_{n-1}}{p_{n-2} - \alpha q_{n-2}}.$$

Damit ist  $\alpha_n(\alpha q_{n-2} - p_{n-2}) = p_{n-1} - \alpha q_{n-1}$ , was durch Umordnung der Terme auf  $\alpha(\alpha_n q_{n-2} + \alpha q_{n-1}) = \alpha_n p_{n-2} + p_{n-1}$  führt. Also ist

$$\alpha = \frac{\alpha_n p_{n-2} + p_{n-1}}{\alpha_n q_{n-2} + q_{n-1}}.$$

### §3: Optimale Approximation

Nach den Vorbereitungen im letzten Paragraphen können wir nun beweisen, daß Kettenbrüche in der Tat bestmögliche Approximationen

sind im folgenden Sinne: Ist  $r/s$  irgendein Bruch, dessen Nenner  $s$  zwischen den Nennern  $q_{n-1}$  und  $q_n$  zweier Konvergenten der Kettenbruchentwicklung liegt, so ist  $p_{n-1}/q_{n-1}$  eine bessere Approximation als  $r/s$ :

**Lemma:**  $p_n/q_n$  seien die Konvergenten der Kettenbruchentwicklung einer reellen Zahl  $\alpha$ . Falls  $\alpha$  irrational ist oder rational mit einem Nenner echt größer  $q_n$ ,  $n \geq 2$ , so ist für jede rationale Zahl  $r/s$  mit  $s \leq q_n$  und  $r/s \notin \{p_{n-1}/q_{n-1}, p_n/q_n\}$

$$\left| \alpha - \frac{r}{s} \right| > \left| \alpha - \frac{p_{n-1}}{q_{n-1}} \right|.$$

*Beweis:* Wir betrachten die Punkte  $P_{n-1} = (q_{n-1}, p_{n-1})$ ,  $P_n = (q_n, p_n)$  und  $R = (s, r)$ . Es genügt zu zeigen, daß der vertikale Abstand von  $P_{n-1}$  zur Geraden  $y = \alpha x$  einen kleineren Betrag hat als der von  $R$ .

Wir schreiben  $R$  als ganzzahlige Linearkombination  $R = kP_{n-1} + \ell P_n$  der Punkte  $P_{n-1}$  und  $P_n$ . Das ist möglich, denn die Determinante des linearen Gleichungssystems

$$\begin{pmatrix} p_{n-1} & p_n \\ q_{n-1} & q_n \end{pmatrix} \begin{pmatrix} k \\ \ell \end{pmatrix} = \begin{pmatrix} r \\ s \end{pmatrix}$$

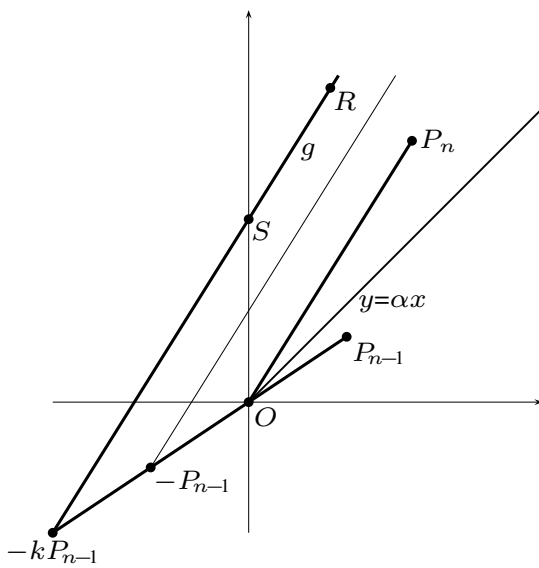
ist nach dem Satz am Ende des vorigen Paragraphen gleich

$$p_{n-1}q_n - p_nq_{n-1} = (-1)^{n-1}.$$

Die somit eindeutig bestimmte Lösung  $(k, \ell)$  des Gleichungssystems ist ganzzahlig, denn wenn wir sie nach der CRAMERSchen Regel ausdrücken, sind  $k$  und  $\ell$  Brüche mit dieser Determinante im Nenner und einer ganzzahligen Determinante im Zähler.

Für das Folgende wollen wir uns auf den Fall  $p_{n-1}/q_{n-1} < \alpha < p_n/q_n$  beschränken; der Fall  $p_{n-1}/q_{n-1} > \alpha > p_n/q_n$  geht völlig analog.

Wir betrachten die Gerade  $g$  durch  $kP_{n-1}$  mit Steigungsvektor  $\overrightarrow{OP_n}$ ; nach unserer Annahme ist ihre Steigung größer als  $\alpha$ .



Im Fall  $k < 0$  liegt der Punkt  $kP_{n-1}$  und damit die ganze Gerade  $g$  zumindest ab dem Punkt  $kP_{n-1}$  oberhalb der Geraden  $y = \alpha x$ , und wegen der größeren Steigung von  $g$  steigt der Abstand zwischen den beiden Geraden mit wachsendem  $x$ . Wir können den Abstand von  $R$  zur Geraden  $y = \alpha x$  daher nach unten abschätzen durch den Abstand des Schnittpunkts  $S$  von  $g$  mit der  $y$ -Achse. Dessen Abstand wiederum können wir nach unten abschätzen,

indem wir  $k = -1$  setzen, denn in diesem Fall ist der Abstand von  $g$  zur Geraden  $y = \alpha x$  am kleinsten. Der Punkt  $-P_{n-1}$  hat (betragsmäßig) denselben Abstand von  $y = \alpha x$  wie  $P_{n-1}$ , und da die Abszisse  $x = 0$  von  $S$  größer ist als die von  $-P_{n-1}$ , hat somit  $S$  einen größeren Abstand von  $y = \alpha x$  als  $P_{n-1}$ . Im Fall  $k < 0$  ist damit die Behauptung bewiesen.

Als nächstes betrachten wir den Fall  $k > 0$ . Dann muß  $\ell \leq 0$  sein, denn sonst wäre die  $x$ -Koordinate  $s = kq_{n-1} + \ell q_n$  von  $R$  größer als  $q_n$ . Der Punkt  $kP_{n-1}$  liegt unterhalb der Geraden  $y = \alpha x$  und die Gerade  $g$  nähert sich dieser mit steigender Abszisse immer mehr an. Da der Punkt  $R$  entweder dieselbe Abszisse wie  $kP_{n-1}$  hat oder eine kleinere, ist sein Abstand somit höchstens gleich dem von  $kP_{n-1}$ , der wiederum das  $k$ -fache des Abstands von  $P_{n-1}$  ist. Für  $k \geq 2$  erhalten wir damit die gewünschte strikte Ungleichung. Für  $k = 1$  erhalten wir auch eine, denn wegen der Voraussetzung  $R \neq P_{n-1}$  muß dann  $\ell \geq 1$  sein.

Bleibt noch der Fall  $k = 0$ . Dann ist  $R = \ell P_n$ , wobei  $\ell \neq 1$ , da  $R \neq P_n$ . Andererseits kann  $\ell$  auch nicht größer als eins sein, denn  $s \leq q_n$ . Somit kommt dieser Fall gar nicht vor. ■

Als nächstes wollen wir uns überlegen, wann gute Approximationen Konvergenten der Kettenbruchentwicklung sein *müssen*. Wir wissen be-



reits, daß für die Konvergenten gilt

$$\left| \alpha - \frac{p_n}{q_n} \right| < \frac{1}{q_n^2}.$$

Dies charakterisiert die Konvergenten allerdings noch nicht: Betrachten wir etwa die Kettenbruchentwicklung von  $\alpha = \sqrt{3}$ . Der Algorithmus liefert zunächst  $c_0 = [\sqrt{3}] = 1$  und  $\alpha_1 = \sqrt{3} - 1$ . Der Kehrwert davon ist

$$\frac{1}{\sqrt{3} - 1} = \frac{\sqrt{3} + 1}{2} \implies c_1 = 1 \quad \text{und} \quad \alpha_2 = \frac{\sqrt{3} - 1}{2}.$$

Der Kehrwert davon ist

$$\frac{2}{\sqrt{3} - 1} = \sqrt{3} + 1 \implies c_2 = 2 \quad \text{und} \quad \alpha_3 = \sqrt{3} - 1 = \alpha_1.$$

Ab hier wiederholt sich also alles periodisch, d.h.

$$\sqrt{3} = 1 + \frac{1}{1 + \frac{1}{2 + \frac{1}{1 + \frac{1}{2 + \dots}}}} = [1, \overline{1, 2}].$$

Die ersten Konvergenten der Kettenbruchentwicklung sind

$$1, \quad 2, \quad 1\frac{2}{3}, \quad 1\frac{3}{4}, \quad 1\frac{8}{11} \quad \text{und} \quad 1\frac{11}{15};$$

da die Folge der Nenner monoton steigt, gibt es also keine Konvergente mit Nenner sieben. Trotzdem ist

$$\left| \sqrt{3} - 1\frac{5}{7} \right| \approx 0,017765 < 0,2 = \frac{1}{50} < \frac{1}{49} = \frac{1}{7^2}.$$

Dafür gilt aber der folgende Satz, dessen zweite Hälfte bereits 1808 von ADRIEN-MARIE LEGENDRE (1752–1833) bewiesen wurde:

**Satz:** a) Eine irrationale Zahl  $\alpha$  erfüllt für jedes  $n \geq 2$  mindestens eine der beiden Ungleichungen

$$\left| \alpha - \frac{p_{n-1}}{q_{n-1}} \right| < \frac{1}{2q_{n-1}^2} \quad \text{oder} \quad \left| \alpha - \frac{p_n}{q_n} \right| < \frac{1}{2q_n^2}.$$

b) Erfüllen zwei ganze Zahlen  $p, q$  die Ungleichung  $\left| \alpha - \frac{p}{q} \right| < \frac{1}{2q^2}$ , so ist  $\frac{p}{q}$  eine Konvergente der Kettenbruchentwicklung von  $\alpha$ .

*Beweis:* a) Angenommen, beide Ungleichungen sind falsch. Nach Multiplikation mit  $q_{n-1}$  bzw.  $q_n$  haben wir dann die beiden Relationen

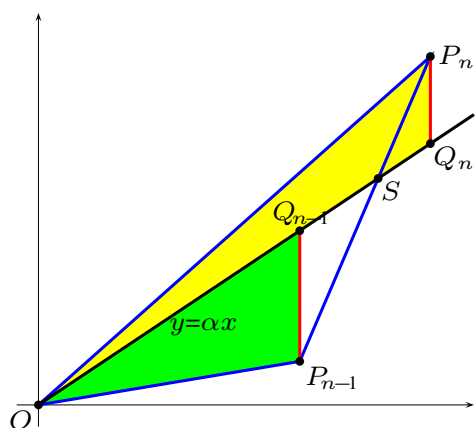
$$|q_{n-1}\alpha - p_{n-1}| \geq \frac{1}{2q_{n-1}} \quad \text{und} \quad |q_n\alpha - p_n| \geq \frac{1}{2q_n}.$$

Wir nehmen für den Beweis wieder an, daß  $p_{n-1}/q_{n-1} < \alpha < p_n/q_n$  ist; der umgekehrte Fall geht völlig analog.

Nach unserer Annahme liegt der Punkt  $P_{n-1} = (q_{n-1}, p_{n-1})$  unterhalb der Geraden  $y = \alpha x$ , und  $P_n = (q_n, p_n)$  liegt darüber.

Das Kreuzprodukt (siehe Anhang) der Vektoren  $\overrightarrow{OP_{n-1}}$  und  $\overrightarrow{OP_n}$  hat als Betrag die Fläche des davon aufgespannten Parallelogramms; das Dreieck mit Ecken  $O, P_{n-1}$  und  $P_n$  ist halb so groß. Wegen der Beziehung  $p_n q_{n-1} - q_n p_{n-1} = (-1)^{n-1}$  ist die Fläche dieses Dreiecks daher gleich  $1/2$ .

Als nächstes betrachten wir zu den Punkten  $P_i = (q_i, p_i)$  ihre Projektionen  $Q_i = (q_i, \alpha q_i)$  in  $y$ -Richtung auf die Gerade  $y = \alpha x$  und die Dreiecke  $\triangle OP_i Q_i$ . Nach Voraussetzung ist die Länge der Seite  $P_i Q_i$  für  $i = n-1$  und  $i = n$  mindestens  $1/2q_i$ . Die darauf senkrecht stehende Höhe ist  $q_i$ , also ist die Fläche jedes der beiden Dreiecks mindestens  $1/4$ .



Ist  $S$  der Schnittpunkt der Geraden  $y = \alpha x$  mit der Verbindungsstrecke von  $P_{n-1}$  und  $P_n$ , so ist das Dreieck  $\triangle OP_{n-1}P_n$  die Vereinigung der Dreiecke  $\triangle OP_{n-1}Q_{n-1}$ ,  $\triangle OP_n Q_n$  und  $\triangle P_{n-1}Q_{n-1}S$ , minus dem Dreieck  $\triangle SP_n Q_n$ . Die Dreiecke beiden  $\triangle P_{n-1}Q_{n-1}S$  und  $\triangle SP_n Q_n$  sind ähnlich, und da jede Konvergente eine bessere Approximation liefert als ihre Vorgänger,

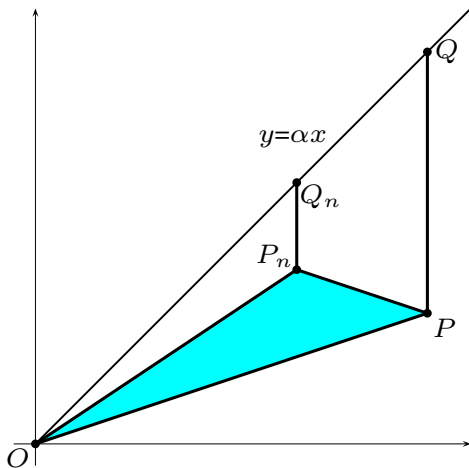
ist das zweite dieser Dreiecke das kleinere. Daher ist die Fläche des Dreiecks  $\triangle OP_{n-1}P_n$  größer als die Summe der Flächen der Dreiecke  $\triangle OP_{n-1}Q_{n-1}$  und  $\triangle OP_nQ_n$ , also größer als  $1/4 + 1/4 = 1/2$ . Dies ist ein Widerspruch zur obigen direkten Berechnung dieser Fläche.

b) Wir können natürlich voraussetzen, daß der Bruch  $p/q$  gekürzt ist, denn für jede nichtgekürzte Darstellung ist die Bedingung echt schärfer.

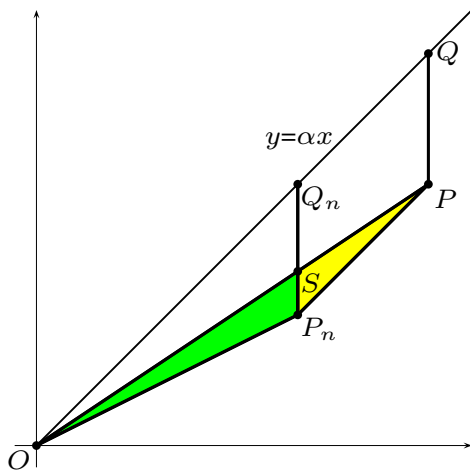
Da die Folge der Nenner  $q_n$  strikt monoton ansteigt, gibt es genau ein  $n$ , so daß  $q_n \leq q < q_{n+1}$  ist; wir müssen zeigen, daß  $p/q = p_n/q_n$  ist. Andernfalls ist  $pq_n - qp_n \neq 0$ , also – da dies eine ganze Zahl ist –  $|pq_n - qp_n| \geq 1$ . Setzen wir  $P = (q, p)$ , so ist also die Fläche des Dreiecks  $\triangle OPP_n$  mindestens gleich  $1/2$ .

Seien wieder  $Q = (q, \alpha q)$  und  $Q_n = (q_n, \alpha q_n)$  die Projektionen der betrachteten Punkte auf die Gerade  $y = \alpha x$ . Die Länge der Strecke  $\overline{PQ}$  ist  $|\alpha q - p|$ , was nach Voraussetzung kleiner als  $1/2q$  ist. Nach dem Lemma zu Beginn dieses Paragraphen ist die Strecke  $\overline{P_nQ_n}$  kürzer als  $\overline{PQ}$ , also ebenfalls kleiner als  $1/2q$  und damit erst recht kleiner als  $1/2q_n$ . Somit haben beide Dreiecke  $\triangle OPQ$  und  $\triangle OP_nQ_n$  Flächen, die kleiner sind als  $1/4$ .

Wir wollen uns überlegen, daß dann auch die Fläche des Dreiecks  $\triangle OPP_n$  kleiner als  $1/2$  sein muß, im Widerspruch zur obigen Rechnung. Die Geometrie hängt dabei stark davon ab, wie die Punkte  $P$  und  $P_n$  sowohl zueinander wie auch in Bezug auf die Gerade  $y = \alpha x$  liegen.



Betrachten wir als erstes den Fall, daß  $p_n/q_n$  zwischen  $\alpha$  und  $p/q$  liegt. Dann liegt der Punkt  $P_n$  im Innern des Dreiecks  $\triangle OPQ$ , also ist das gesamte Dreieck  $\triangle OPP_n$  im Dreieck  $\triangle OPQ$  enthalten. Da ersteres mindestens die Fläche  $1/2$  hat, letzteres aber weniger als  $1/4$ , kann dieser Fall offensichtlich nicht vorkommen.



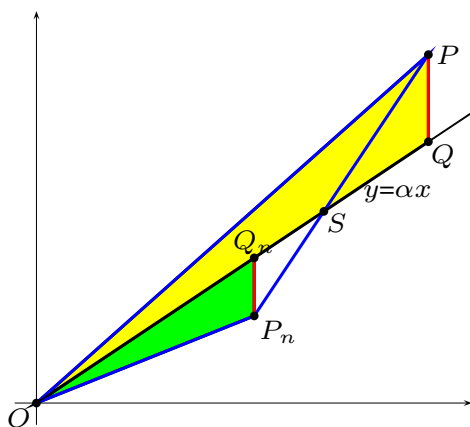
Als nächstes nehmen wir an,  $p/q$  liege zwischen  $\alpha$  und  $p_n/q_n$ . Dann schneiden sich die Strecken  $\overline{P_nQ_n}$  und  $\overline{OP}$  in einem Punkt  $S$ , und das Dreieck  $\triangle OPP_n$  ist die Vereinigung der beiden Dreiecke  $\triangle OSP_n$  und  $\triangle SPP_n$ . Zur Flächenberechnung gehen wir aus von der gemeinsamen Kante  $\overline{SP_n}$ ; die darauf senkrecht stehenden Höhen haben die Längen  $q_n$  und  $q - q_n$ . Somit ist

die verdoppelte Fläche des gesamten Dreiecks  $\triangle OPP_n$  gleich

$$|\overline{SP_n}| \cdot q_n + |\overline{SP_n}| \cdot (q - q_n) = |\overline{SP_n}| \cdot q \leq |\overline{P_nQ_n}| \cdot q \leq |\overline{PQ}| \cdot q,$$

denn da  $q$  zwischen  $q_n$  und  $q_{n+1}$  liegt, kann  $P_n$  nach obigem Lemma keinen größeren Abstand von der Geraden  $y = \alpha x$  haben als  $P$ . Rechts steht aber die verdoppelte Fläche des Dreiecks  $\triangle OPQ$ , von der wir wissen, daß sie höchstens gleich  $1/2$  ist, so daß auch dieser Fall nicht auftreten kann.

Bleibt noch der Fall, daß  $\alpha$  zwischen  $p/q$  und  $p_n/q_n$  liegt,  $P$  und  $P_n$  also auf verschiedenen Seiten der Geraden  $y = \alpha x$  liegen. Dann schneidet ihre Verbindungsstrecke  $\overline{PP_n}$  diese Gerade in einem Punkt  $S$ . Damit sind wir in einer ähnlichen Situation wie beim Beweis von a): Das Dreieck  $\triangle OPP_n$  ist gleich dem Dreieck  $\triangle OP_nQ_n$  plus dem Dreieck  $\triangle OPQ$  plus  $\triangle SP_nQ_n$  minus  $\triangle SPQ$ . Die beiden letzteren



Dreiecke sind ähnlich, und da  $\overline{PQ}$  nicht kürzer sein kann als  $\overline{P_nQ_n}$  ist das subtrahierte Dreieck mindestens genauso groß wie  $\triangle SP_nQ_n$ . Somit ist die Fläche von  $\triangle OPP_n$  höchstens gleich der Summe der Flächen von  $\triangle OPQ$  und  $\triangle OP_nQ_n$ , also kleiner als  $1/4 + 1/4 = 1/2$ . Damit haben wir auch hier einen Widerspruch, d.h.  $p/q = p_n/q_n$ . ■

## Anhang: Das Kreuzprodukt zweier Vektoren

Im  $\mathbb{R}^3$  (und nur dort) gibt es eine bilineare Verknüpfung, die zwei Vektoren einen dritten zuordnet, das (vielleicht aus der Schule bekannte) Vektorprodukt oder Kreuzprodukt. Wie schon der Name sagt, ordnet es je zwei Vektoren  $v$  und  $w$  aus  $\mathbb{R}^3$  einen *Vektor* zu, und dieser wird mit  $v \times w \in \mathbb{R}^3$  bezeichnet. Er ist festgelegt durch folgende Eigenschaften:

- $v \times w$  hat die Länge  $|v \times w| = |v| |w| |\sin \angle(v, w)|$ .  
Insbesondere ist also  $v \times w = \vec{0}$ , wenn  $v$  und  $w$  auf einer Geraden liegen, denn dann bilden sie einen Winkel von null oder 180 Grad, so daß der Sinus verschwindet.
- $v \times w$  steht senkrecht sowohl auf  $v$  als auch auf  $w$ .  
Falls  $v \times w \neq \vec{0}$  ist, spannen  $v$  und  $w$  eine Ebene auf, auf der (da wir im  $\mathbb{R}^3$  sind) genau ein eindimensionaler Unterraum senkrecht steht. Darin gibt es allerdings für jede vorgegebene positive Länge zwei Vektoren, die sich durch ihr Vorzeichen unterscheiden. Um  $v \times w$  eindeutig festzulegen, brauchen wir daher noch eine weitere Bedingung:
- Die drei Vektoren  $v$ ,  $w$  und  $v \times w$  bilden ein Rechtssystem, d.h. wenn man die Finger der *rechten* Hand so ausrichtet, daß der Daumen in Richtung von  $v$  zeigt und der Zeigefinger in Richtung von  $w$ , so zeigt und der Mittelfinger in Richtung von  $v \times w$ .

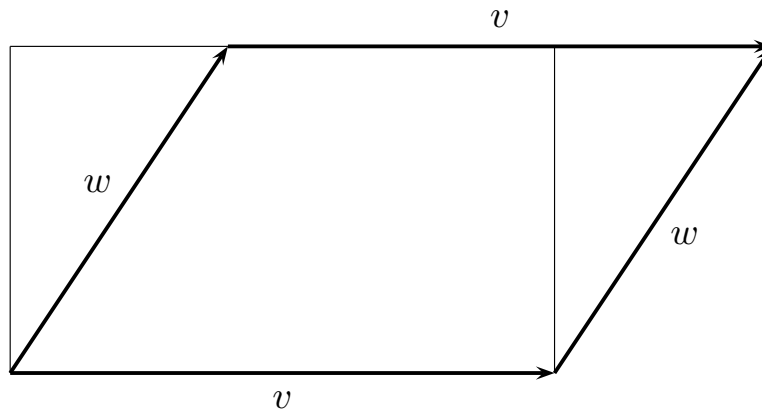
Alternativ kann man ein Rechtssystem auch so definieren, daß sich, ein von  $v$  nach  $w$  gedrehter Korkenzieher in Richtung  $v \times w$  in den Kork bohrt. Ähnlich geht es auch mit Schrauben; da es allerdings neben den (üblichen) Rechtsschrauben auch die (seltenen) Linksschrauben gibt, ist diese Definition eventuell zirkulär: Alles hängt davon ab, wie man Rechtsschrauben definiert.

Aus jeder dieser Regeln folgt sofort die *Antikommutativität* des Vektorprodukts:

$$v \times w = -w \times v.$$

Weitere Rechenregeln lassen sich leicht geometrisch ableiten: Da der Sinus eines Winkels gleich Gegenkathete durch Hypotenuse ist, ist in der von  $v$  und  $w$  aufgespannten Ebenen  $|w| |\sin \angle(v, w)|$  gleich der Länge

des auf die senkrecht auf  $v$  stehenden Geraden projizierten Vektors  $w$ , das heißt also gleich der Höhe des in der Abbildung eingezeichneten Rechtecks. Die Länge des Vektors  $v \times w$  ist daher gleich dem Flächeninhalt dieses Rechtecks und damit – wie eine Scherung zeigt – gleich der Fläche des von  $v$  und  $w$  aufgespannten Parallelogramms.



Daraus folgt nun sofort das Distributivgesetz

$$v \times (w + u) = v \times w + v \times u$$

für den zweiten Faktor, und wegen der Antikommutativität folgt daraus wiederum das für den ersten:

$$(u + v) \times w = u \times w + v \times w .$$

Um das Vektorprodukt in Koordinaten ausrechnen zu können, müssen wir zunächst die Produkte der Koordinateneinheitsvektoren  $e_i$  kennen. Da sie allesamt die Länge eins haben und paarweise aufeinander senkrecht stehen, ist klar, daß das Produkt zweier verschiedener dieser Vektoren bis aufs Vorzeichen gleich dem dritten ist; das Vorzeichen hängt ab von der Orientierung des Koordinatensystems. Das Produkt eines Vektors  $e_i$  mit sich selbst ist natürlich, wie jedes Produkt eines Vektors mit sich selbst, gleich dem Nullvektor, denn der eingeschlossene Winkel ist null Grad.

Für die folgende Rechnung wollen wir annehmen, daß  $e_1, e_2$  und  $e_3$  in dieser Reihenfolge ein Rechtssystem bilden; das ist beispielsweise dann der Fall, wenn  $e_1$  nach rechts,  $e_2$  nach vorne und  $e_3$  nach oben zeigt.

Dann folgt sofort, daß

$$e_1 \times e_2 = e_3$$

ist, und nach einigen Fingerübungen auch findet man auch die Formeln

$$e_2 \times e_3 = e_1 \quad \text{und} \quad e_1 \times e_3 = -e_2.$$

Die Produkte mit vertauschten Faktoren sind natürlich gerade das negative davon, und  $e_i \times e_i = 0$ . Für

$$v = \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} \quad \text{und} \quad w = \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix}$$

ist also

$$v \times w = (v_1 e_1 + v_2 e_2 + v_3 e_3) \times (w_1 e_1 + w_2 e_2 + w_3 e_3),$$

nach den obigen Rechenregeln gleich

$$\begin{aligned} & \sum_{i=1}^3 \sum_{j=1}^3 v_i w_j e_i \times e_j \\ &= (v_2 w_3 - v_3 w_2) e_1 + (v_3 w_1 - v_1 w_3) e_2 + (v_1 w_2 - v_2 w_1) e_3, \end{aligned}$$

d.h.

$$\begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} \times \begin{pmatrix} w_1 \\ w_2 \\ w_3 \end{pmatrix} = \begin{pmatrix} v_2 w_3 - v_3 w_2 \\ v_3 w_1 - v_1 w_3 \\ v_1 w_2 - v_2 w_1 \end{pmatrix}.$$

Dies läßt sich dadurch merken, daß man im Schema

$$\begin{array}{ccccc} e_1 & & e_2 & & e_3 & & e_1 & & e_2 \\ & \searrow & & \times & & \times & & \swarrow & \\ v_1 & & v_2 & & v_3 & & v_1 & & v_2 \\ & \swarrow & & \times & & \times & & \searrow & \\ w_1 & & w_2 & & w_3 & & w_1 & & w_2 \end{array}$$

von  $e_i$  ausgeht und als dessen Koeffizient das Zweierprodukt entlang der schrägen Linie nach rechts unten *positiv* und das entlang der schrägen Linie nach links unten *negativ* nimmt; man wendet also die SARRUSSche Regel an auf die „Determinante“

$$\begin{vmatrix} e_1 & e_2 & e_3 \\ v_1 & v_2 & v_3 \\ w_1 & w_2 & w_3 \end{vmatrix}.$$

## §4: Kettenbrüche und Kalender

Schon in den ältesten bekannten Kulturen richtete sich die Zeitrechnung nach astronomischen Gesetzmäßigkeiten: dem Umlauf der Erde um die Sonne, dem Umlauf des Mondes um die Erde sowie der Drehung der Erde um sich selbst.

Der Tag als Zeiteinheit ist ein so selbstverständlicher Teil unseres Lebensrhythmus, daß er als Zeiteinheit nie zur Debatte stand. Ebenso selbstverständlich war, daß als Tag nicht die Dauer einer vollen Drehung der Erde um ihre Achse genommen wurde, sondern der etwa vier Minuten längere Zeitraum, bis sie der Sonne wieder dieselbe Stelle zuwendet.

Als nächstgrößere Einheit führten wahrscheinlich die Babylonier die Woche ein; ob sie sich dabei vom ungefähren Abstand zwischen zwei Mondphasen leiten ließen, ist unbekannt; vielleicht wurde die Dauer von sieben Tagen auch einfach deshalb gewählt, weil die Sieben als heilige Zahl galt.

Definitiv vom Mond abgeleitet ist der Monat. Der Mond dreht sich bekanntlich um die Erde; die Zeit für einen vollständigen Umlauf beträgt ungefähr 27,3 Tage. Dieser sogenannte *siderische* Monat spielt allerdings für die Kalenderrechnung keine Rolle; für die Zeitbestimmung wurden seit Alters her die gut und einfach zu beobachtenden Mondphasen verwendet. Da der Mond nicht selbst leuchtet, sondern das Sonnenlicht reflektiert, hängen diese ab vom Winkelabstand zwischen Sonne und Mond; für den Kalender relevant ist daher der sogenannte *synodische* Monat von 29,53 Tagen, nach dem sich dieser Winkelabstand wiederholt. (Der tatsächliche Abstand zwischen zwei Neumonden ist wegen der komplizierten Mondbewegung keine Konstante; erst im Mittel kommt man auf den synodischen Monat.)

Da 29,53 keine ganze Zahl ist, lassen sich Monate nicht einfach als eine feste Anzahl von Tagen definieren sondern müssen in einem lunaren Kalender mal 29, mal 30 Tage haben.

Einer der einfachsten und zugleich einer der jüngsten dieser Kalender ist der islamische: Sobald mindestens zwei vertrauenswürdige Männer den



neuen Mond gesehen haben, beginnt ein neuer Monat, bei bewölktem Himmel unabhängig davon dreißig Tage nach dem letzten Monatsanfang. Alle zwölf Monate beginnt ein neues Jahr.

Es ist klar, daß bei einer solchen Festlegung die Länge der Monate sowohl innerhalb eines Jahres als auch von Jahr zu Jahr schwankt, außerdem sind die Jahre nicht synchron zum Umlauf der Erde um die Sonne. Da der Kalender auf der arabischen Halbinsel entstand, wo Jahreszeiten keine Rolle spielen und auch die Landwirtschaft das ganze Jahr über konstante Bedingungen vorfindet, ist letzteres dort kein Nachteil.

Für Regionen mit ausgeprägten Jahreszeiten oder jährlich wiederkehrenden Ereignissen ist die Synchronisation des Kalenders mit der Sonne wichtiger als die mit dem Mond. Das wohl älteste Beispiel eines reinen Sonnenjahrs bietet der ägyptische Kalender. Da die für die Landwirtschaft fundamentalen jährlichen Nilüberschwemmungen ungefähr mit der ersten Sichtung des Sterns Sirius übereinstimmten, wurde dieses Ereignis als Beginn des neuen Jahrs genommen. Dies ist das sogenannte *siderische* Jahr mit einer Länge von 365,256 Tagen. Obwohl die Ägypter wußten, daß diese Länge ungefähr  $365\frac{1}{4}$  Tage beträgt, legten Sie doch fest, daß jedes Jahr genau 365 Tage haben sollte, verteilt auf zwölf Monate zu jeweils dreißig Tagen sowie fünf Zusatztage.

Ein Jahr, das wirklich synchron zu den Jahreszeiten ist, sollte allerdings nicht anhand des Fixsternhimmels definiert werden, sondern anhand jahreszeitlicher Phänomene wie beispielsweise der Tag- und Nachtgleiche oder dem Durchgang der Sonne durch den Frühlingspunkt. Das ist (im Mittel) das sogenannte *tropische* Jahr mit einer Länge von 365,2422 Tagen. Der Unterschied zum siderischen Jahr ist zwar gering, aber – wie wir gleich sehen werden – trotzdem relevant.

Viel bedeutender als dieser Unterschied war aber zunächst einmal die Tatsache, daß ein Jahr mit exakt 365 Tagen natürlich im Laufe der Jahrhunderte zu einem Verlust der Synchronisation des Kalenders mit den Jahreszeiten führt. Aus diesem Grund beauftragte GAIUS JULIUS CAESAR (100–44) den alexandrinischen Astronomen SOSIGENES mit einer Kalenderreform, die die damit verbundene Verschiebung des Jahresanfangs (die sich in nur 120 Jahren auf einen Monat summiert) kompensieren sollte.

Das Ergebnis, der *Julianischer Kalender*, wurde offiziell eingeführt zum 1. Januar des Jahres 709 *ab urbe condita*, d.h. nach Gründung der Stadt Rom. In unserer heutigen Zeitrechnung handelt es sich dabei um das Jahr 45 v.Chr. Die Monatsnamen und -längen des Julianischen Kalenders sind die heute noch gebräuchlichen. Er unterscheidet sich vom klassischen ägyptischen Kalender durch die zusätzliche Regel, daß in jedem vierten Jahr ein Schalttag eingeführt wird, der 29. Februar.

Dabei blieb es bis ins sechzehnte Jahrhundert. Bis dahin hatte das astronomisch etwas zu lange Julianische Jahr zu einer Verschiebung beispielsweise der Frühjahrs-sonnenwende um rund elf Tage geführt. Um dies zu korrigieren, setzte Papst GREGOR XIII (UGO BONCOMPAGNI, 1502–1585, Papst ab 1572) eine Kommission ein, auf Grund von deren Empfehlungen in den Jahren, die durch hundert aber nicht durch vierhundert teilbar sind, auf den Schalttag verzichtet wird. Dieser *Gregorianische Kalender* trat in den katholischen Ländern am Freitag, dem 15. Oktober 1582 in Kraft; um die bis dahin akkumulierten Fehler des Julianischen Kalenders zu kompensieren, folgte dieser Tag auf den noch Julianischen Donnerstag, den 4. Oktober 1582. In nichtkatholischen Ländern galt der Julianische Kalender noch länger, wurde aber schließlich (mit den verschiedensten Übergangsregeln) überall durch den Gregorianischen ersetzt – zuletzt 1927 in der Türkei.

Um die Neuerung des Gregorianischen Kalenders zu verstehen, betrachten wir die Kettenbruchentwicklung von  $365,2422$ ; sie ist

$$[365, 4, 7, 1, 3, 4, 1, 1, 1, 2]$$

und hat die Konvergenten

$$365, \quad 365\frac{1}{4}, \quad 365\frac{7}{29}, \quad 365\frac{8}{33}, \quad 365\frac{31}{128}, \quad 365\frac{132}{545}, \quad \dots$$

Der Julianische Kalender verwendet die einfach zu realisierende Konvergente  $365\frac{1}{4}$ . Unter den folgenden Konvergenten finden wir keine mit einem Nenner, der sich gut für eine einfache Kalenderregel eignen würde. Am ehesten kommt vielleicht noch der Nenner 33 in Frage, da er ungefähr ein Drittel von hundert ist. Arbeitet man mit der Approximation  $365\frac{8}{33}$ , so sollten unter 33 Jahren acht Schaltjahre sein, also

24 pro 99 Jahre. Nimmt man stattdessen 24 Schaltjahre pro Jahrhundert, was der Regel entspricht, daß durch hundert teilbaren Jahre *keine* Schaltjahre sind, so sorgt der Unterschied zwischen 99 und 100 dafür, daß nach jeweils 400 Jahren eine Vierjahresperiode fehlt. Auch diese hat Anspruch auf ein Schaltjahr, daher die Gregorianische Regel, daß durch hundert teilbare Jahre *keine* Schaltjahre sind, es sei denn, die Jahreszahl sei sogar durch vierhundert teilbar. Innerhalb einer jeden Periode von 400 Jahren gibt es also  $100 - 3 = 97$  Schaltjahre; das Gregorianische Jahr hat somit eine Länge von  $365 \frac{97}{400}$  Tagen, eine praktikable Zahl in der Nähe der Konvergente  $365 \frac{8}{33}$ .

Zum Einstieg in die Kalenderrechnung beginnen wir mit dem einfachsten Problem, den Wochentagen, und überlegen wir uns, auf welchen Wochentag der  $T$ -te Tag des  $M$ -ten Monats im Jahr  $J$  fällt.

Alle Wochen haben exakt sieben Tage und die Jahre haben nach einer recht klaren Regel 365 oder 366 Tage; es ist daher relativ einfach, den Wochentag für den  $i$ -ten Tag des Jahres zu berechnen, sofern man ihn für irgendeinen anderen Tag dieses Jahres kennt: Er hängt innerhalb eines Jahres schließlich nur ab von  $i \bmod 7$ .

Um auch die Abhängigkeit vom Jahr noch zu berücksichtigen, ist es am einfachsten, die Tage nicht vom 1. Januar des jeweils betrachteten Jahres aus zu zählen, sondern ab irgendeinem festen Datum. Historisch sinnvoll wäre hier beispielsweise das Datum der Einführung des Gregorianischen Kalenders; da hier aber Jahr, Monat und Tag „krumme“ Zahlen sind, würde dies zu unnötig komplizierten mathematischen Formeln mit zu vielen willkürlich erscheinenden Konstanten führen.

Für die Mathematik ist es unerheblich, ob zum fiktiven Anfangspunkt bereits der Gregorianische Kalender in Gebrauch war oder nicht; wir können daher beispielsweise ausgehen von einem 1. Januar eines fiktiven Jahres Null. (Die Zählung der Jahre ab Christi Geburt wurde im sechsten Jahrhundert initiiert von DIONYSIUS EXIGUUS, der allerdings ein falsches Geburtsjahr 1 berechnete, auf das wir uns heute noch beziehen. Jahre davor interessierten ihn nicht; der erste der auch Jahre zuvor in Bezug auf Christi Geburt datierte, war wohl der angelsächsische Theologe und Historiker BEDA VENERABILIS (673–735), der – da er keine Null

kannte – das Jahr vor dem Jahr eins *nach* Christus als eins *vor* Christus bezeichnete.)

Wenn wir dem fiktiven 1. Januar 0 die Nummer eins geben und der Einfachheit halber davon ausgehen, daß das Jahr Null *kein* Schaltjahr war, können wir die Nummer des 31. Dezembers des Jahres  $J - 1$  folgendermaßen berechnen: Bis dahin sind  $J - 1$  Jahre verflossen, von denen jedes mindestens 365 Tage hatte; damit kommen schon einmal  $365(J - 1)$  Tage zusammen. Was noch fehlt sind die Schalttage: Im Julianischen Kalender hätte es davon  $[(J - 1)/4]$  gegeben, im Gregorianischen sind aber die durch 100, nicht aber durch 400 teilbaren Jahre keine Schaltjahre, also müssen wir  $[(J - 1)/100] - [(J - 1)/400]$  subtrahieren. Der  $i$ -te Tag des Jahres  $J$  hat somit die Nummer

$$365(J - 1) + [(J - 1)/4] - [(J - 1)/100] + [(J - 1)/400] + i .$$

Um den zugehörigen Wochentag zu finden, müssen wir nun nur noch den Wochentag des Tags Nummer eins bestimmen. Dazu können wir von irgendeinem bekannten Datum ausgehen: Dienstag, der 8. Mai 2018 ist der  $(31 + 28 + 31 + 30 + 8) = 128$ -te Tag des Jahres 2018, hat also die Nummer

$$365 \cdot 2017 + 504 - 20 + 5 + 128 = 736\,8222 .$$

Diese Zahl ist kongruent zwei modulo sieben; somit war Tag eins ein Montag. Geben wir den Wochentagen, wie es die DIN- und ISO-Normen vorsehen, von Montag ausgehend die Nummern eins bis sieben, so fällt der Tag mit Nummer  $i$  also auf den Wochentag mit Nummer  $i \bmod 7$ , wobei die Null dem normgemäß mit 7 bezeichneten Sonntag entspricht.

Tatsächlich hätten wir die obige Rechnung etwas vereinfachen können: Da  $365 \equiv 1 \pmod{7}$  ist, reicht es, wenn wir

$$(J - 1) + [(J - 1)/4] - [(J - 1)/100] + [(J - 1)/400] + i \pmod{7}$$

berechnen, im Beispiel also

$$2017 + 504 - 20 + 5 + 128 = 2634 \equiv 2 \pmod{7} .$$

Um den Wochentag zu einem vorgegebene Datum bestimmen zu können, müssen wir immer noch berechnen, der wievielte Tag des Jahres

der  $T$ -te Tag des  $M$ -ten Monats ist. Eine sehr einfache Methode besteht darin, daß wir zählen, wie viele Tage vor dem Ersten des jeweiligen Monats bereits vergangen sind. Dabei müssen wir natürlich zwischen Schaltjahren und gewöhnlichen Jahren unterscheiden: Für letztere seien dies  $t_M$  Tage, für erste  $s_M$ . Dann haben wir folgende Tabelle:

$M =$	1	2	3	4	5	6	7	8	9	10	11	12
$t_M =$	0	31	59	90	120	151	181	212	243	273	304	334
$s_M =$	0	31	60	91	121	152	182	213	244	274	305	335

Dann fällt der  $T$ -te Tag des  $M$ -ten Monats des Jahrs  $J$  auf den Wochentag mit der Nummer

$$(J - 1) + [(J - 1)/4] - [(J - 1)/100] + [(J - 1)/400] + t_M + T$$

modulo sieben, falls  $J$  kein Schaltjahr ist; andernfalls muß  $t_M$  durch  $s_M$  ersetzt werden. Es genügt natürlich, die Zahlen  $t_M$  oder  $s_M$  modulo 7 einzusetzen, also

$M =$		1	2	3	4	5	6	7	8	9	10	11	12
$t_M \bmod 7 =$		0	3	3	6	1	4	6	2	5	0	3	5
$s_M \bmod 7 =$		1	4	4	0	2	5	0	3	6	1	4	6

Da wohl niemand eine dieser beiden Tabellen auswendig lernen möchte, stellt sich die Frage, ob es vielleicht auch eine geschlossene Formel gibt. Dazu ignorieren wir zunächst einmal die historisch überkommenen Monatslängen und tun so, als könnte ein Mathematiker am grünen Tisch festlegen, wie er 365 Tage auf zwölf Monate verteilt.

Für ihn wäre die am wenigsten irreguläre Verteilung der Monatslängen wohl die, bei der Tag  $i$  eines Jahres mit  $N$  Tagen genau dann im  $k$ -ten Monat liegt, wenn gilt

$$\frac{(k-1)N}{12} < i \leq \frac{kN}{12} \quad \text{oder} \quad k-1 < \frac{12i}{N} \leq k,$$

d.h.  $k$  ist die kleinste ganze Zahl größer oder gleich  $12i/N$ . Für ein Jahr mit  $N = 365$  Tagen würde dies auf die Monatslängen

$$30, 30, 31, 30, 31, 30, 30, 31, 30, 31, 30, 31$$

führen, in einem Schaltjahr mit  $N = 366$  hätten alle ungeraden Monate dreißig und alle geraden Monate 31 Tage. Ein Februar mit 28 oder 29 Tagen kann bei einer derartigen Strategie natürlich nie vorkommen.

Trotzdem läßt sich auch unser chaotisches Monatssystem fast auf eine solche Formel bringen: Nehmen wir an, wir hätten ein Jahr mit 367 Tagen und zwölf Monaten. Dann liefert uns die obige Vorgehensweise Monate der Längen

$$30, 31, 30, 31, 30, 31, 31, 30, 31, 30, 31, 31,$$

wir haben also wie im wirklichen Kalender an zwei Stellen aufeinanderfolgende Monate mit 31 Tagen. Im Kalender sind dies Juli/August und Dezember/Januar, hier sind es die Monate 6 und 7 sowie 11 und 12. Wenn wir zyklisch um eine Position verschieben, so daß die hintere 31 an der ersten Stelle steht, stimmen die beiden Positionen überein, und abgesehen vom Februar, der hier dreißig Tage hat, haben wir genau die Folge der Monatslängen.

(Kurioserweise gab es 1712 in Schweden sogar ein Jahr mit 367 Tagen und einem 30. Februar: 1699 wurde beschlossen, langsam zum Gregorianischen Kalender überzugehen und dazu als erstes den Schalttag 1700 zu streichen. Danach wurde der Beschluß aufgegeben, und um wieder synchron zum Julianischen Kalender zu werden, gab es 1712 einen 30. Februar als zweiten Schalttag. 1753 wurde der Gregorianische Kalender dann endgültig und abrupt eingeführt.)

Die Anzahl der Tage vor dem Ersten des  $M$ -ten Monats wäre in unserem hypothetischen Kalender einfach gleich  $[367M/12]$ ; durch die zyklische Verschiebung wird diese Formel freilich zerstört. Sie kann aber gerettet werden durch eine Verschiebung im Zähler: Wie explizites Nachrechnen zeigt, sind in einem Jahr mit 367 Tagen, in dem der Februar dreißig Tage hat, vor dem Ersten des  $M$ -ten Monats gleich  $[(367M - 362)/12]$  Tage vergangen. Somit ist für unseren realen Kalender

$$t_M = \begin{cases} \left[ \frac{367M-362}{12} \right] & \text{für } M \leq 2 \\ \left[ \frac{367M-362}{12} \right] - 2 & \text{für } M \geq 3 \end{cases} \quad \text{und}$$

$$s_M = \begin{cases} \left[ \frac{367M-362}{12} \right] & \text{für } M \leq 2 \\ \left[ \frac{367M-362}{12} \right] - 1 & \text{für } M \geq 3 \end{cases} .$$

Der Wochentag des  $T$ -ten Tags im  $M$ -ten Monat des Jahrs  $J$  ist somit

$$(J - 1) + \left[ \frac{J - 1}{4} \right] - \left[ \frac{J - 1}{100} \right] + \left[ \frac{J - 1}{400} \right] + \left[ \frac{367M - 362}{12} \right] + \delta_M + T$$

modulo sieben mit

$$\delta_M = \begin{cases} 0 & \text{falls } M \leq 2 \\ -1 & \text{falls } M \geq 3 \text{ und } J \text{ Schaltjahr} \\ -2 & \text{falls } M \geq 3 \text{ und } J \text{ kein Schaltjahr} \end{cases} .$$

Diese Formel gilt selbstverständlich nur für Daten nach dem Gregorianischen Kalender; bei älteren Daten muß man zunächst wissen, auf welchem Kalender und welchem Jahresanfang sie beruhen.

Als beispielsweise der amerikanische Naturwissenschaftler, Philosoph und Politiker BENJAMIN FRANKLIN geboren wurde, zeigten die Kalender in seiner Heimatstadt Boston den 6. Januar 1705. Massachusetts war zu der Zeit noch britische Kolonie, und da Großbritannien den Gregorianischen Kalender erst 1752 einführt, ist das ein Julianisches Datum. Gregorianisch ist sein Geburtstag elf Tage später, d.h. am 17. Januar. Allerdings handelt es sich dabei nicht um den 17. Januar 1705, sondern um den des Jahres 1706: In Großbritannien begann das neue Jahr damals nämlich nicht am ersten Januar, sondern am 25. März. Auf den 31. Dezember 1705 folgte also der 1. Januar 1705 und auf den 24. März 1705 der 25. März 1706. Solche Besonderheiten bei der Interpretation alter Datumsangaben gibt es viele; hier ist also Vorsicht geboten.

Mindestens genauso wichtig wie eine verbesserte Schaltjahrregel war für Papst Gregor das Datum des Osterfests; auch darum sollte sich seine Kommission kümmern. Damit kam nun plötzlich auch der Mond in den Kalender, denn 325 beschloß das Konzil von Nicäa (bei Konstantinopel), daß Ostern stets am ersten Sonntag nach dem ersten Vollmond am oder nach der Frühlings-Tag-und-Nacht-Gleiche zu feiern sei. (Man beachte, daß das erste *nach* im Sinne eines  $>$ , das zweite im Sinne eines  $\geq$  definiert ist. Der Grund für das  $>$  lag darin, daß so Ostern nur sehr selten gleichzeitig mit dem jüdischen Pascha-Fest begangen wird.)

Der erste Vollmond am oder nach der Frühlings-Tag-und-Nacht-Gleiche kann nicht einfach nach einer ähnlichen Regel wie der Monatsanfang im islamischen Kalender bestimmt werden, also etwa dann, wenn ihn

mindestens zwei vertrauenswürdige Kardinäle gesehen haben, denn der österliche Festkreis beginnt bereits siebenzig Tage vor Ostern. Das Datum mußte daher im Voraus berechnet werden. Der Mathematiker und Informatiker DONALD E. KNUTH sagt in Abschnitt 1.3.2, Aufgabe 14 seiner *Art of Computer Programming: There are many indications that the sole important application of arithmetic in Europe during the Middle Ages was the calculation of the Easter date, and so such algorithms are historically significant.* So ordnete etwa KARL DER GROSSE (747/8–814) bei seiner Neuordnung des Bildungssystems mehrfach an, daß in jeder Diözese mindestens ein Geistlicher in der Lage sein müsse, das Osterdatum zuverlässig zu berechnen. Schauen wir uns also an, wie Papst Gregor das Osterdatum berechnen ließ.

Wir brauchen Informationen über die Wochentage, über die Mondphasen und über die Tag-Nacht-Gleiche im Frühling. Letztere ist, da der Gregorianische Kalender das tropische Jahr recht genau approximiert, noch recht lange konstant am 21. März jedes Jahres; bei der bei der Berechnung des Osterdatums geht man daher stets von diesem Tag aus. Wie man Wochentage bestimmt, haben wir uns gerade überlegt; bleibt also noch das Problem mit den Mondphasen.

Die mittlere Zeitspanne zwischen zwei Neumonden, die synodische Umlaufzeit des Mondes, beträgt etwa 29,5306 Tage; ein tropisches Jahr mit seinen 365,2422 Tagen besteht also aus

$$365,2422 : 29,5306 \approx 12,3679$$

solchen Zykeln. Die Kettenbruchentwicklung dieses Quotienten ist

$$[12, 2, 1, 2, 1, 1, 4, 1, 81]$$

mit Konvergenten

$$12, \quad 12\frac{1}{3}, \quad 12\frac{3}{8}, \quad 12\frac{4}{11}, \quad 12\frac{7}{19}, \quad 12\frac{32}{87}, \quad 12\frac{39}{106}, \quad \dots$$

Für einen Kalender, der sowohl mit der Sonne als auch mit dem Mond synchronisiert ist, könnte man also Jahre mit 12 und 13 Monaten kombinieren, wobei in erster Näherung jedes dritte Jahr 13 Monate hätte. Tatsächlich war man im fünften vorchristlichen Jahrhundert bereits erheblich weiter: Der um 440 v.Chr. lebende Athener Mathematiker und



Astronomen METON verwendete die Konvergente mit Nenner 19. Ein Metonischer Zyklus besteht demnach aus 19 Jahren, darunter zwölf *Gemeinjahren* aus zwölf Monaten und sieben *Schaltjahren* aus 13 Monaten. Die Monate hatten teils 29, teils 30 Tage. Der darauf basierende Kalender wurde in Griechenland bis 46 v.Chr. verwendet. Die Synchronisation zwischen Sonnen- und Mondzyklen ist fast perfekt:

$$19 \cdot 365,2422 = 6939,6018 \quad \text{und} \quad 235 \cdot 29,5306 = 6939,691,$$

der Fehler pro Zyklus liegt also bei nur etwa zwei Stunden. Seit Einführung des Gregorianischen Kalenders sind etwas über 22 Metonische Zykeln vergangen; der akkumulierte Fehler liegt also noch unter zwei Tagen.

Der Gregorianische Kalender geht deshalb bei der Bestimmung des Osterdatums nicht von astronomischen Beobachtungen aus, sondern von Metonischen Zykeln, allerdings mit einer Korrektur für den akkumulierten Fehler. Ebenfalls unberücksichtigt bleiben die Irregularitäten der realen Mondbewegung; gerechnet wird mit einer Approximation der *mittleren* Mondbewegung. Auf den ersten Blick seltsam erscheinen mag auch die Tatsache, daß bei der Fehlerkorrektur mit der *Julianischen* Jahreslänge von  $365\frac{1}{4}$  Tagen gerechnet wird; der Grund lag wohl vor allem darin, daß Papst Gregor bisherige Praktiken nicht mehr als unbedingt notwendig ändern wollte.

Die wesentliche Größe, mit der die Mondphasen in unseren an der Sonne orientierten Kalender gebracht werden, ist der sogenannte *Epakt*. Mit diesem Wort bezeichneten die Griechen die Anzahl der Tage, die an Neujahr seit dem letzten Neumond des alten Jahres vergangen waren. Gemäß dem Metonischen Zyklus sollte diese Zahl sich alle 19 Jahre wiederholen; in der Kalenderrechnung wird daher die um eins vermehrte Restklasse modulo 19 der Jahreszahl als die „Goldene Zahl“ bezeichnet. (Die Addition der Eins kommt natürlich daher, daß zur Zeit ihrer Einführung die Null in der europäischen Mathematik noch nicht vorkam.)

Wenn jedes Jahr genau 365 Tage hätte, könnten wir einfach mit den  $12 \times 29,5 = 354$  Tagen eines Mondjahrs vergleichen und wüßten dann,

daß sich die Mondphase an einem festen Datum jedes Jahr um elf Tage verschiebt. Als *Mondphase* bezeichnen wir dabei die Anzahl von Tagen, die seit dem letzten Neumond vergangen sind.

Eine der vielen Vereinfachungen in der Berechnung des Osterdatums liegt darin, daß man innerhalb des aktuellen Metonischen Zyklus im wesentlichen von dieser Formel ausgeht, die Schalttage also ignoriert.

Da die Schalttage Ende Februar eingeführt werden, wir uns aber für den Vollmond am oder nach dem 21. März interessieren, sollten wir nicht mit dem klassischen Epakt, der Mondphase des 1. Januar, rechnen: Sonst gäbe es schließlich algorithmisch unangenehme Fallunterscheidungen für die Schaltjahre. Aus Effizienzgründen bietet sich an, stattdessen mit einem *verschobenen* Epakt zu rechnen, d.h. mit der Mondphase eines geeigneten Datums, das näher bei Ostern liegt.

Die Länge eines lunaren Zyklus liegt bei ungefähr 29,5 Tagen; zum einfacheren Rechnen sollten wir das zumindest für den einen Zyklus, in den Ostern fällt, auf den ganzzahligen Wert 30 runden. Der erste Vollmond nach dem 21. März ist dann der letzte Vollmond vor dem 19. April, und sein Abstand zum 19. April ist, wenn wir den Vollmond als Tag mit Mondphase 14 betrachten, gleich dem Abstand des letzten Neumonds vor dem 5. April zum 5. April, also die Mondphase des 5. Aprils. Somit bietet sich an, als verschobenen Epakt die Mondphase des 5. Aprils zu nehmen. Gemäß unserer vereinfachenden Annahmen sollte auch sie sich alle 19 Jahre wiederholen und sollte sich zumindest innerhalb eines Metonischen Zyklus von Jahr zu Jahr modulo 30 um elf verschieben.

Damit brauchen wir nur noch für ein Jahr des Metonischen Zyklus den tatsächlichen Wert der Mondphase des fünften Aprils kennen, um den verschobenen Epakt allgemein berechnen zu können. Die vor der Gregorianischen Reform gebräuchliche Formel berechnet ihn für das Jahr  $J$  als

$$E = (14 + 11 \cdot (J \bmod 19)) \bmod 30 .$$

Der erste Vollmond am oder nach dem 21. März lag somit  $E$  Tage vor dem 19. April, und Ostern war der (echt) darauf folgende Sonntag. So

wird Ostern noch heute in fast allen orthodoxen Kirchen berechnet; die einzige Ausnahme ist die finnische.

Der Gregorianische Kalender modifiziert diese Formel durch drei zusätzliche Terme: Zunächst berücksichtigt er, daß der Metonische Zyklus nicht wirklich exakt ist, insbesondere dann nicht, wenn man mit dem Julianischen Jahr arbeitet:

$$19 \cdot 365,25 = 6939,750 \quad \text{und} \quad 235 \cdot 29,5306 = 6939,691 ;$$

hier beträgt die Differenz also 0,059 Tage pro Zyklus und

$$0,059 \times \frac{100}{19} \approx 0,31$$

Tage pro Jahrhundert. Die Gregorianische Osterformel approximiert dies durch  $8/25 = 0,32$ , addiert allerdings im  $h$ -ten Jahrhundert nicht  $[8h/25]$ , sondern  $[(5 + 8h)/25]$ . Diese Modifikation soll in erster Linie dafür sorgen, daß Ostern möglichst selten mit dem jüdischen Paschafest zusammenfällt. Für das Jahrhundert wird dabei die gleiche Konvention benutzt wie für die Feier des Jahrtausendanfangs am 1. Januar 2000: Das  $h$ -te Jahrhundert beginnt mit dem Jahr  $100(h - 1)$ , d.h.  $h = [J/100] + 1$ .

Da der Gregorianische Kalender bei der Korrektur der Metonischen Zyklen mit Julianischen Jahren arbeitet, am Ende aber ein Gregorianisches Datum braucht, müssen als nächstes die unterschiedlichen Anzahlen von Schalttagen berücksichtigt werden, d.h. die „ausfallenden“ Schalttage des Gregorianischen Kalenders müssen subtrahiert werden. Das sind drei Stück pro 400 Jahre, also wird  $[3h/4]$  subtrahiert. Dies ergäbe die neue Formel

$$E = \left( 14 + 11 \cdot (J \bmod 19) + \left[ \frac{5 + 8h}{25} \right] - \left[ \frac{3h}{4} \right] \right) \bmod 30 .$$

Tatsächlich gibt es noch eine weitere Modifikation, die dafür sorgen soll, daß die 19 Epakte eines Metonischen Zyklus alle verschieden sind und  $E = 0$  nicht auftritt: Falls  $E = 0$  ist oder falls  $E = 1$  ist und  $J \bmod 19 > 10$ , wird  $E$  um eins erhöht. Der (berechnete) Vollmond ist dann  $E$  Tage vor dem 19. April, und Ostern wird weiterhin am darauf folgenden Sonntag gefeiert.

Das Jahr  $J = 2018$  liegt im  $h = 21$ . Jahrhundert und  $2018 \equiv 4 \pmod{19}$ . Somit ist

$$E = \left( 14 + 11 \cdot 4 + \left[ \frac{173}{25} \right] - \left[ \frac{63}{4} \right] \right) \pmod{30} = 19$$

Der rechnerische Vollmond ist daher am 31. März, und an diesem Tag um 14 Uhr 37 wurde er auch tatsächlich erreicht. Der darauf folgende Sonntag ist der 1. April 2018, also wurde Ostern an diesem Tag gefeiert.

## §5: Eine kryptographische Anwendung

Beim RSA-Verfahren wählt man den öffentlichen Exponenten  $e$  oft ziemlich klein, z.B.  $e = 2^{16} + 1$ . Dies hat den Vorteil, daß zumindest die Verschlüsselung ziemlich schnell geht und man nur zur Entschlüsselung mit einem Exponenten in der Größenordnung des Moduls arbeiten muß.

Für jemanden, der RSA hauptsächlich für elektronische Unterschriften verwendet, würde sich anbieten, stattdessen den privaten Exponenten  $d$  relativ klein zu wählen. Dann könnte er schnell viele Dokumente unterschreiben, und falls jeder Empfänger nur eines davon bekommt, fällt dessen höherer Aufwand bei der Überprüfung nicht so sehr ins Gewicht.

Natürlich kann man nicht  $d = 3$  oder  $d = 2^{16} + 1$  wählen: Der private Exponent muß schließlich geheim sein und es darf nicht möglich sein, ihn durch Probieren zu erraten.

Andererseits geht man heute bei symmetrischen Kryptoverfahren davon aus, daß ein Verfahren sicher ist, falls ein Gegner mindestens  $2^{128}$  Möglichkeiten durchprobieren muß, so daß gängige Verfahren wie AES mit einer Schlüssellänge von 128 Bit auskommen. Verglichen damit erscheinen 2048 Bit für einen privaten Entschlüsselungsexponenten recht hoch.

Trotzdem läßt sich hier nicht wesentlich sparen, denn ein Gegner kann kurze private Exponenten nicht nur durch Ausprobieren bestimmen, sondern auch wesentlich schneller nach dem Kettenbruchalgorithmus.

Wir gehen aus von einem öffentlichen RSA-Schlüssel  $(N, e)$  sowie dem zugehörigen privaten Exponenten  $d$ . Dann gibt es bekanntlich eine

natürliche Zahl  $k$ , so daß  $ed - k\varphi(N) = 1$  ist. Dies können wir umschreiben als

$$\frac{e}{\varphi(N)} - \frac{k}{d} = \frac{1}{d\varphi(N)}.$$

Falls  $d$  sehr viel kleiner ist als  $\varphi(N)$  haben wir hier einen Bruch mit dem großen Nenner  $\varphi(N)$  sehr gut angenähert durch einen Bruch mit dem sehr viel kleineren Nenner  $d$ . Für hinreichend kleines  $d$  ist das nur möglich, wenn  $k/d$  eine Konvergente der Kettenbruchentwicklung von  $e/\varphi(N)$  ist.

Das mag zunächst harmlos erscheinen, denn die Sicherheit von RSA beruht ja gerade darauf, daß niemand außer dem Inhaber des privaten Schlüssels  $d$  die Faktorisierung  $N = pq$  und damit den Wert von

$$\varphi(N) = (p - 1)(q - 1) = N - (p + q) + 1$$

kennt. Dafür kennt aber jeder den Wert von  $N$ , und wie die obige Gleichung zeigt, liegt der recht nahe bei  $\varphi(N)$ : Die Primzahlen  $p$  und  $q$  sind schließlich nur von der Größenordnung  $\sqrt{N}$ . Damit sollte  $k/d$  auch eine gute Approximation für  $e/N$  liefern.

In der Tat zeigte Kryptologe MICHAEL JAMES WIENER 1990 ein Resultat, wonach insbesondere der folgende Satz gilt:

**Satz:** Ist  $N = pq$  Produkt zweier Primzahlen  $p$  und  $q$  mit  $p < q < 2q$ , und ist  $d < \frac{1}{3} \sqrt[4]{N}$  der private Exponent zum öffentlichen Exponenten  $e < \varphi(N)$ , so ist  $d$  Nenner einer Konvergenten der Kettenbruchentwicklung von  $e/N$ .

*Beweis:* Wegen  $ed \equiv 1 \pmod{\varphi(N)}$  gibt es ein  $k \in \mathbb{N}$ ; so daß  $ed - k\varphi(N) = 1$  ist; wegen  $e < \varphi(N)$  ist dabei  $k < d$ . Nach dem Satz von LEGENDRE aus §3 reicht es, wenn wir zeigen können, daß

$$\left| \frac{e}{N} - \frac{k}{d} \right| < \frac{1}{2d^2}$$

ist, denn dann ist  $k/d$  eine Konvergente der Kettenbruchentwicklung

von  $e/N$ .

$$\begin{aligned}
 \left| \frac{e}{N} - \frac{k}{d} \right| &= \left| \frac{ed - kN}{dN} \right| \\
 &= \left| \frac{(ed - k\varphi(N)) + k\varphi(N) - kN}{dN} \right| \\
 &= \left| \frac{1 + k(\varphi(N) - N)}{dN} \right| \\
 &= \left| \frac{1 + k(1 - p - q)}{dN} \right| = \frac{k(p + q) - (k + 1)}{dN} \\
 &< \frac{k(p + q)}{dN}.
 \end{aligned}$$

Natürlich ist  $p < \sqrt{N}$ , und wegen der Voraussetzung  $q < 2p$  folgt  $p + q < 3\sqrt{N}$ . Außerdem ist  $k < d < \frac{1}{3} \sqrt[4]{N}$ , also

$$\left| \frac{e}{N} - \frac{k}{d} \right| < \frac{3k\sqrt{N}}{dN} = \frac{3k}{d\sqrt{N}} < \frac{\sqrt[4]{N}}{d\sqrt{N}} = \frac{1}{d\sqrt[4]{N}}.$$

Dies ist genau dann kleiner als  $1/2d^2$ , wenn  $d < \frac{1}{2} \sqrt[4]{N}$  ist. Nach Voraussetzung ist aber  $d$  sogar kleiner als  $\frac{1}{3} \sqrt[4]{N}$ , womit der Satz bewiesen wäre. ■

Um  $d$  zu berechnen, müssen wir daher nur so lange Konvergenten  $p_n/q_n$  bestimmen, bis für einen der Nenner  $q_n$  die Exponentiation mit  $q_n$  modulo  $N$  invers ist zu der mit  $e$ . Falsche Kandidaten sollten dabei praktisch immer bereits beim ersten Versuch erkannt werden.

Tatsächlich gibt es Algorithmen, mit denen man sogar private Exponenten  $d < N^{0,289}$  rekonstruieren kann, und manche Fachleute meinen, daß man vielleicht sogar in vielen Fällen mit  $d < \sqrt{N}$  mit geeigneten Algorithmen eine realistische Erfolgchance haben könnte; bei diesen Attacken arbeitet man allerdings nicht mit Kettenbrüchen, sondern mit anderen Verfahren zur diophantischen Approximation.

Private Exponenten müssen somit immer groß sein. Falls man von einem vorgegebenen öffentlichen Exponenten ausgeht, ist das für reali-

stische  $N$  mit an Sicherheit grenzender Wahrscheinlichkeit erfüllt; Vorsicht ist nur geboten, wenn man mit dem privaten Exponenten startet. Daher verlangen auch die Vorschriften der Bundesnetzagentur, daß man immer vom öffentlichen Exponenten  $e$  ausgehen muß, und erst daraus einen privaten Exponenten berechnet.

## §6: Die Kettenbruchentwicklung der Eulerschen Zahl

Aufgabe 1b) des neunten Übungsblatts läßt eine erstaunliche Regelmäßigkeit in der Kettenbruchentwicklung von  $e$  vermuten:

$$e = [2, 1, 2, 1, 1, 4, 1, 1, 6, 1, 1, 8, 1 \dots].$$

Diese Entwicklung ist bereits im 18. Kapitel der 1748 erschienenen *Introductio in analysin infinitorum* von EULER enthalten; HERMITE bewies sie 1873 im Rahmen seiner Arbeit über die Transzendenz von  $e$  mit anderen Methoden die zusammenhängen mit der Approximation der Exponentialfunktion durch rationale Funktionen. Sein Schüler PADÉ entwickelte später eine systematische Theorie solcher Approximationen, die PADÉ-Approximanten, die in der Numerik eine große Rolle spielen für die näherungsweise Berechnung von Standardfunktionen. Durch Kombination solcher Ideen kamen verschiedene Mathematiker zu immer einfacheren Beweisen; der hier wiedergegebene Beweis von HENRY COHN erschien 2006 im *American Mathematical Monthly*; direkt dahinter folgt eine Arbeit von THOMAS J. OSLER, der den Beweis so verallgemeinert, daß er auch die Kettenbruchentwicklungen der Wurzeln aus  $e$  liefert.

Wir können die obige Kettenbruchentwicklung noch etwas regelmäßiger schreiben, indem wir beachten, daß für alle  $x \in \mathbb{R}$  gilt

$$1 + \frac{1}{0 + \frac{1}{1+x}} = 1 + (1+x) = 2+x;$$

der obige Kettenbruch kann also auch geschrieben werden als

$$[1, 0, 1, 1, 2, 1, 1, 4, 1, 1, 6, 1, 1, 8, 1, \dots].$$

Hier läßt sich der  $n$ -te Koeffizient  $c_i$  völlig regelmäßig durch  $n$  ausdrücken: Für  $n = 3k + 1$  mit  $k \in \mathbb{N}_0$  ist er  $2k$ , ansonsten eins.

Für die Kettenbruchentwicklung der  $M$ -ten Wurzel aus  $e$  müssen wir daran nur wenig ändern: Hier wollen wir sehen, daß

$$c_{3k} = c_{3k+2} = 1 \quad \text{und} \quad c_{3k+1} = (2k+1)M - 1$$

ist für alle  $k \in \mathbb{N}_0$ , d.h.

$$\sqrt[M]{e} = [1, M-1, 1, 1, 3M-1, 1, 1, 5M-1, 1, 1, 7M-1, 1, \dots].$$

Wir gehen aus von diesem Kettenbruch und wollen zeigen, daß er gegen  $\sqrt[M]{e}$  konvergiert. Nach dem Satz am Ende von §2 lassen sich der Zähler  $p_n$  und der Nenner  $q_n$  der  $n$ -ten Konvergente eines Kettenbruchs  $[c_0, c_1, c_2, \dots]$  rekursiv berechnen nach den Formeln

$$p_0 = c_0, q_0 = 1, p_1 = c_0c_1 + 1, q_1 = c_1,$$

$$p_n = p_{n-2} + c_n p_{n-1} \quad \text{und} \quad q_n = q_{n-2} + c_n q_{n-1} \quad \text{für } n \geq 2.$$

Speziell für die hier betrachtete Kettenbruchentwicklung haben wir also die Anfangsterme  $p_0 = q_0 = p_1 = 1$  und  $q_1 = M - 1$ , was insbesondere bedeutet, daß die erste Konvergente im Falle  $M = 1$ , bei der Kettenbruchentwicklung von  $e$  also, nicht definiert ist. Weiter geht es nach obiger Rekursionsformel; da die  $c_n$  von  $n \bmod 3$  abhängen, bekommen wir je nach Restklasse von  $n$  drei verschiedene Formeln:

$$\begin{aligned} p_{3k} &= p_{3k-2} + p_{3k-1} & q_{3k} &= q_{3k-2} + q_{3k-1} \\ p_{3k+1} &= p_{3k-1} + ((2k+1) - 1)M p_{3k} & q_{3k+1} &= q_{3k-1} + ((2k+1)M - 1)q_{3k} \\ p_{3k+2} &= p_{3k} + p_{3k+1} & q_{3k+2} &= q_{3k} + q_{3k+1} \end{aligned}$$

Wir müssen zeigen, daß die Folge der Quotienten  $p_n/q_n$  gegen  $\sqrt[N]{e}$  konvergiert.

Der Trick dazu hängt mit PADÉ-Approximanten zusammen; ich möchte darauf nicht eingehen, sondern ohne Begründung einfach die drei Integrale

$$A_k = \int_0^1 \frac{x^k(x-1)^k}{k!M^{k+1}} e^{x/M} dx, \quad B_k = \int_0^1 \frac{x^{k+1}(x-1)^k}{k!M^{k+1}} e^{x/M} dx$$

$$\text{und} \quad C_k = \int_0^1 \frac{x^k(x-1)^{k+1}}{k!M^{k+1}} e^{x/M} dx$$

betrachten.



**Satz:** Für alle  $k \in \mathbb{N}_0$  gilt:

$$\begin{aligned} p_{3k} - q_{3k} \sqrt[M]{e} &= -A_k \\ p_{3k+1} - q_{3k+1} \sqrt[M]{e} &= B_k \\ p_{3k+2} - q_{3k+2} \sqrt[M]{e} &= C_k \end{aligned}$$

Da in allen drei Integranden der Zähler kleiner als eins und die Exponentialfunktion höchstens  $\sqrt[M]{e}$  ist, während der Nenner für  $k \rightarrow \infty$  gegen  $\infty$  geht, ist

$$\lim_{k \rightarrow \infty} A_k = \lim_{k \rightarrow \infty} B_k = \lim_{k \rightarrow \infty} C_k = 0;$$

daher folgt aus diesem Satz sofort

**Korollar:**  $\lim_{n \rightarrow \infty} \frac{p_n}{q_n} = \sqrt[M]{e}$ , d.h.

$$\sqrt[M]{e} = [1, M-1, 1, 1, 3M-1, 1, 1, 5M-1, 1, 1, 7M-1, 1, \dots].$$

Insbesondere ist  $e = [1, 0, 1, 1, 2, 1, 1, 4, 1, \dots] = [2, 1, 2, 1, 1, 4, 1, \dots]$ . ■

Der obige Satz wird durch Induktion bewiesen. Für  $k = 0$  ist

$$A_0 = \int_0^1 \frac{1}{M} e^{x/M} dx = e^{x/M} \Big|_0^1 = \sqrt[M]{e} - 1$$

$$B_0 = \int_0^1 \frac{x}{M} e^{x/M} dx = (x - M) e^{x/M} \Big|_0^1 = (1 - M) \sqrt[M]{e} + M$$

$$C_0 = \int_0^1 \frac{x-1}{M} e^{x/M} dx = (x - 1 - M) e^{x/M} \Big|_0^1 = -M \sqrt[M]{e} + M + 1$$

Nach den eingangs angegebenen Rekursionsformeln für die  $p_n$  und  $q_n$  ist

$$p_0 = q_0 = 1, \quad p_1 = M, \quad q_1 = M - 1, \quad p_2 = 1 + M \quad \text{und} \quad q_2 = M.$$

Die drei Formeln aus dem Satz werden also für  $k = 0$  zu den Gleichungen

$$\begin{aligned} 1 - \sqrt[M]{e} &= 1 - \sqrt[M]{e} \\ M - (M-1)\sqrt[M]{e} &= (1-M)\sqrt[M]{e} + M \\ 1 + M - M\sqrt[M]{e} &= -M\sqrt[M]{e} + M + 1, \end{aligned}$$

die offensichtlich alle drei richtig sind.

Für den Induktionsschritt brauchen wir Beziehungen zwischen den Integralen  $A_k$ ,  $B_k$  und  $C_k$ . Hier gilt für alle  $k \in \mathbb{N}$

$$\begin{aligned} a) \quad A_k &= -B_{k-1} - C_{k-1} \\ b) \quad B_k &= -((2k+1)M-1)A_k + C_{k-1} \\ c) \quad C_k &= B_k - A_k \end{aligned}$$

Zum *Beweis* von *a)* wenden wir die LEIBNIZ-Regel zur Ableitung eines Produkts an auf das Produkt der drei Faktoren  $x^k$ ,  $(x-1)^k$  und  $e^{x/M}$ . Für ein solches Dreierprodukt ist  $(uvw)' = u'vw + uv'w + uvw'$ , also ist

$$\begin{aligned} &\frac{d}{dx} x^k (x-1)^k e^{x/M} \\ &= kx^{k-1}(x-1)^k e^{x/M} + kx^k(x-1)^{k-1} e^{x/M} + \frac{x^k(x-1)^k}{M} e^{x/M}. \end{aligned}$$

Division durch  $k!M^k$  macht daraus

$$\begin{aligned} &\frac{d}{dx} \frac{x^k(x-1)^k}{k!M^k} e^{x/M} \\ &= \frac{x^{k-1}(x-1)^k e^{x/M}}{(k-1)!M^k} + \frac{x^k(x-1)^{k-1} e^{x/M}}{(k-1)!M^k} + \frac{x^k(x-1)^k}{k!M^{k+1}} e^{x/M}. \end{aligned}$$

Integrieren wir beide Seiten von 0 bis 1, so erhalten wir auf der linken Seite den Wert null, da die Stammfunktion des Integranden an beiden Intervallenden verschwindet. Rechts erhalten wir die Summe der Integrale  $C_{k-1}$ ,  $B_{k-1}$  und  $A_k$ ; somit ist  $A_k + B_{k-1} + C_{k-1} = 0$ , was *a)* beweist.

Der Beweis von *b)* geht ähnlich: Wir berechnen zunächst die Ableitung von  $x^k(x-1)^{k+1} e^{x/M}$  und dividieren wieder durch  $k!M^k$ ; wir erhalten

$$\frac{d}{dx} \frac{x^k(x-1)^{k+1}}{k!M^k} e^{x/M}$$

$$\begin{aligned}
&= \frac{x^{k-1}(x-1)^{k+1}}{(k-1)!M^k} e^{x/M} + \frac{(k+1)x^k(x-1)^k}{k!M^k} e^{x/M} + \frac{x^k(x-1)^{k+1}}{k!M^{k+1}} e^{x/M} \\
&= \frac{kMx^{k-1}(x-1)^{k+1} + M(k+1)x^k(x-1)^k + x^k(x-1)^{k+1}}{k!M^{k+1}} e^{x/M} \\
&= \frac{x^{k-1}(x-1)^k (kM(x-1) + (k+1)Mx + x(x-1))}{k!M^{k+1}} e^{x/M} \\
&= \frac{x^{k-1}(x-1)^k \left( (2k+1)M - 1 \right) x - kM + x^2}{k!M^{k+1}} e^{x/M} \\
&= ((2k+1)M - 1) \frac{x^k(x-1)^k}{k!M^{k+1}} e^{x/M} - \frac{x^{k-1}(x-1)^k}{(k-1)!M^k} e^{x/M} \\
&\quad + \frac{x^{k+1}(x-1)^k}{k!M^{k+1}} e^{x/M}.
\end{aligned}$$

Wenn wir die linke Seite dieser Gleichung von 0 bis 1 integrieren, erhalten wir wieder den Wert null, bei der rechten erhalten wir

$$((2k+1)M - 1)A_{k-1} - B_k + C_{k-1}.$$

Auflösen nach  $B_k$  zeigt die Behauptung *b*)

Zum Beweis von *c*) schließlich gehen wir aus von der Gleichung  $(x-1)^2 = x(x-1) - (x-1)$  und multiplizieren diese mit

$$\frac{x^k(x-1)^{k-1}}{k!M^{k+1}} e^{x/M}.$$

Integration von 0 bis 1 führt auf die Gleichung  $C_k = B_k - A_k$ .

Damit sind alle drei Relationen bewiesen, und wir können mit dem Induktionsschritt zum Beweis unseres zentralen Satzes beginnen. Sei also  $k \geq 1$ ; wir nehmen an, daß die drei Gleichungen für  $k-1$  gelten.

Als erstes wollen wir zeigen, daß

$$p_{3k} - q_{3k} \sqrt[M]{e} = -A_k$$

ist. Nach den Rekursionsformeln für Zähler und Nenner der Konvergen-ten ist

$$p_{3k} = p_{3k-2} + p_{3k-1} \quad \text{und} \quad q_{3k} = q_{3k-2} + q_{3k-1};$$

also ist

$$\begin{aligned} p_{3k} - q_{3k} \sqrt[M]{e} &= (p_{3k-2} - q_{3k-2} \sqrt[M]{e}) + (p_{3k-1} - q_{3k-1} \sqrt[M]{e}) \\ &= B_{k-1} + C_{k-1} = -A_k \end{aligned}$$

nach Induktionsannahme und der Beziehung  $A_k = -B_{k-1} - C_{k-1}$ .

Genauso können wir auch bei den anderen beiden Gleichungen vorgehen:

$$\begin{aligned} p_{3k+1} - q_{3k+1} \sqrt[M]{e} &= (p_{3k-1} - q_{3k-1} \sqrt[M]{e}) \\ &\quad + ((2k+1)M - 1)(p_{3k} - q_{3k} \sqrt[M]{e}) \\ &= C_{k-1} - ((2k+1)M - 1)A_k = B_k \end{aligned}$$

und

$$\begin{aligned} p_{3k+2} - q_{3k+2} \sqrt[M]{e} &= (p_{3k} - q_{3k} \sqrt[M]{e}) + (p_{3k+1} - q_{3k+1} \sqrt[M]{e}) \\ &= -A_k + B_k = C_k \end{aligned}$$

Somit gelten alle drei Beziehungen auch für  $k$ , also für alle  $k \in \mathbb{N}_0$ . Dies beweist den Satz sowie die Kettenbruchentwicklungen für  $e$  und seine Wurzeln. ■

Wer sich genauer dafür interessiert, wie man auf die hier einfach hingeschriebenen Integrale  $A_k$ ,  $B_k$  und  $C_k$  kommt, sollte die verwendeten Originalarbeiten (und eventuell auch die dort zitierte Literatur) konsultieren:

HENRY COHN: A Short Proof of the Simple Continued Fraction Expansion of  $e$ , *American Mathematical Monthly* **113** (2006), 56–62

und

THOMAS J. OSLER: A Proof of the Continued Fraction Expansion of  $e^{1/M}$ , *American Mathematical Monthly* **113** (2006), 62–66

Das *American Mathematical Monthly*, eine Mitgliederzeitschrift der *Mathematical Association of America*, ist im Internet frei verfügbar.

## Kapitel 6

# Gaußsche Zahlen und Quaternionen

Die algebraische Zahlentheorie untersucht neben den „klassischen“ ganzen Zahlen auch ganze Zahlen in Erweiterungskörpern von  $\mathbb{Q}$ ; teilweise führt dies auch zu einem besseren Verständnis von Sätzen der elementaren Zahlentheorie. Hier wollen wir uns auf ein einziges Beispiel beschränken, die sogenannten GAUSSschen Zahlen, und dazu noch eine nichtkommutative Verallgemeinerung betrachten, die Quaternionen.

### § 1: Der Ring der Gaußschen Zahlen

GAUSSsche Zahlen sind komplexe Zahlen  $x + iy$  mit  $x, y \in \mathbb{Z}$ ; wir bezeichnen die Menge aller dieser Zahlen mit  $\mathbb{Z}[i] = \mathbb{Z} \oplus i\mathbb{Z}$ . Da Summen und Produkte GAUSSscher Zahlen wieder GAUSSsche Zahlen sind, überlegt man sich leicht, daß  $\mathbb{Z}[i]$  ein Ring ist.

**Definition:**  $z = x + iy$  sei eine GAUSSsche Zahl.

a)  $\bar{z} = x - iy$  heißt die zu  $z$  konjugiert komplexe Zahl.

b)  $N(z) = z\bar{z} = x^2 + y^2$  heißt die Norm von  $z$ .

c)  $z \in \mathbb{Z}[i]$  heißt *Einheit*, wenn es eine Zahl  $w \in \mathbb{Z}[i]$  gibt, so daß  $zw = 1$  ist.

d)  $z \in \mathbb{Z}[i]$  heißt *irreduzibel*, falls gilt:  $z$  ist keine Einheit, und ist  $z = wq$  Produkt zweier GAUSSscher Zahlen, so muß  $w$  oder  $q$  eine Einheit sein.

**Lemma:** a) Für  $z, w \in \mathbb{Z}[i]$  ist  $N(zw) = N(z)N(w)$ .

b)  $z \in \mathbb{Z}[i]$  ist genau dann eine Einheit, wenn  $N(z) = 1$  ist.

c) Die Einheiten sind genau die vier Zahlen  $1, -1, i$  und  $-i$ .

d) Ist  $N(z)$  eine Primzahl, so ist  $z$  irreduzibel.

*Beweis:* a)  $N(zw) = zw \cdot \overline{zw} = zw\overline{zw} = z\overline{z}w\overline{w} = N(z)N(w)$ .

b) Ist  $z$  eine Einheit, so gibt es ein  $w \in \mathbb{Z}[i]$  mit  $zw = 1$ , und nach a) ist  $N(z)N(w) = N(zw) = N(1) = 1$ . Da die Norm einer GAUSSschen Zahl in  $\mathbb{N}_0$  liegt, folgt  $N(z) = N(w) = 1$ . Ist umgekehrt  $N(z) = 1$ , so ist nach Definition von  $N(z) = z\overline{z}$  das Produkt  $z\overline{z} = 1$ , d.h.  $z$  ist eine Einheit.

c)  $z = x+iy$  sei nach b) genau dann eine Einheit, wenn  $N(z) = x^2+y^2 = 1$  ist. Da  $x$  und  $y$  ganze Zahlen sind, muß eine der beiden verschwinden und die andere  $\pm 1$  sein, was genau auf die vier angegebenen Fälle führt.

d) Ist  $z = wq$ , so ist nach a) auch  $N(z) = N(w)N(q)$ . Da alle Normen nichtnegative ganze Zahlen sind und  $N(z)$  eine Primzahl, muß  $N(w)$  oder  $N(q)$  gleich eins sein, d.h.  $w$  oder  $q$  ist eine Einheit. ■

Die Umkehrung von d) gilt nicht: Beispielsweise ist  $N(3) = 3 \cdot 3 = 9$ , aber trotzdem ist die Drei in  $\mathbb{Z}[i]$  irreduzibel: Ist nämlich  $3 = zw$ , so ist  $N(3) = N(z)N(w)$ ; falls weder  $z$  noch  $w$  eine Einheit ist, müssen also  $N(z) = N(w) = 3$  sein. Es gibt aber keine GAUSSsche Zahl der Norm drei, denn die Gleichung  $x^2 + y^2 = 3$  hat keine ganzzahligen Lösungen.

## §2: Euklidische Ringe

In Kapitel I bewiesen wir die eindeutige Primzerlegung in  $\mathbb{Z}$  mit Hilfe des EUKLIDischen Algorithmus. Um zu sehen, ob wir ähnliches auch für die GAUSSschen Zahlen beweisen können, liegt es daher nahe, Ringen zu untersuchen, in denen es einen EUKLIDischen Algorithmus gibt. Solche Ringe heißen EUKLIDische Ringe.

Wie wir gesehen haben, ist die Division mit Rest das wichtigste Werkzeug beim EUKLIDischen Algorithmus, und wie sich in diesem Abschnitt herausstellen wird, brauchen wir kein weiteres. Wir definieren daher

**Definition:** Ein EUKLIDischer Ring ist ein Integritätsbereich  $R$  zusammen mit einer Abbildung  $\nu: R \setminus \{0\} \rightarrow \mathbb{N}_0$ , so daß gilt: Ist  $x|y$ , so ist

$\nu(x) \leq \nu(y)$ , und zu je zwei Elementen  $x, y \in R$  gibt es Elemente  $q, r \in R$  mit

$$x = qy + r \quad \text{und} \quad r = 0 \quad \text{oder} \quad \nu(r) < \nu(y).$$

Wir schreiben auch  $x : y = q$  Rest  $r$  und bezeichnen  $r$  als Divisionsrest bei der Division von  $x$  durch  $y$ .

Das Standardbeispiel ist natürlich der Ring  $\mathbb{Z}$  der ganzen Zahlen mit  $\nu(x) = |x|$ . Ein anderes Beispiel ist der Polynomring  $k[X]$  über einem Körper  $k$ : Hier können wir  $\nu(f)$  für ein Polynom  $f \neq 0$  als den Grad von  $f$  definieren; dann erfüllt auch die Polynomdivision mit Rest die Forderung an einen EUKLIDischen Ring.

Man beachte, daß weder der Quotient noch der Divisionsrest eindeutig bestimmt sein muß: Beispielsweise ist schon in  $\mathbb{Z}$  einerseits  $15 : 4 = 3$  Rest  $3$ , andererseits aber auch  $4$  Rest  $-1$ , wobei letzteres im EUKLIDischen Algorithmus möglicherweise sogar schneller ans Ziel führt.

**Lemma:** In einem EUKLIDischen Ring  $R$  gibt es zu je zwei Elementen  $x, y \in R$  einen ggT. Dieser kann nach dem EUKLIDischen Algorithmus berechnet werden und läßt sich als Linearkombination von  $x$  und  $y$  mit Koeffizienten aus  $R$  darstellen.

*Beweis:* In jedem Integritätsbereich folgt aus der Gleichung  $x = qy + r$  mit  $x, y, q, r \in R$ , daß die gemeinsamen Teiler von  $x$  und  $y$  gleich denen von  $y$  und  $r$  sind. Speziell in einem EUKLIDischen Ring können wir dabei  $r$  als Divisionsrest wählen und, wie beim klassischen EUKLIDischen Algorithmus, danach  $y$  durch  $r$  dividieren usw., wobei wir eine Folge von Divisionsresten  $r_i$  erhalten mit der Eigenschaft, daß in jedem Schritt die gemeinsamen Teiler von  $x$  und  $y$  gleich denen von  $r_{i-1}$  und  $r_i$  sind. Außerdem ist stets entweder  $r_i = 0$  oder  $\nu(r_i) < \nu(r_{i-1})$ , so daß die Folge nach endlich vielen Schritten mit einem  $r_n = 0$  abbrechen muß. Auch hier sind die gemeinsamen Teiler von  $r_{n-1}$  und  $r_n = 0$  genau die gemeinsamen Teiler von  $x$  und  $y$ . Da jede Zahl Teiler der Null ist, sind die gemeinsamen Teiler von  $r_{n-1}$  und Null aber genau die Teiler von  $r_{n-1}$ , und unter diesen gibt es natürlich einen größten, nämlich  $r_{n-1}$  selbst. Somit haben auch  $x$  und  $y$  einen größten gemeinsamen Teiler,

nämlich den nach dem EUKLIDischen Algorithmus berechneten letzten von Null verschiedenen Divisionsrest  $r_{n-1}$ .

Auch die lineare Kombinierbarkeit folgt wie im klassischen Fall: Bei jeder Division mit Rest ist der Divisionsrest als Linearkombination von Dividend und Divisor darstellbar; beim EUKLIDischen Algorithmus beginnen wir mit Dividend  $x$  und Divisor  $y$ , die natürlich beide als Linearkombinationen von  $x$  und  $y$  darstellbar sind, und induktiv folgt, daß auch alle folgenden Dividenden und Divisoren sind als Reste einer vorangegangenen Division Linearkombinationen von  $x$  und  $y$  sind, also ist es auch ihr Divisionsrest. Insbesondere ist auch der ggT als letzter nichtverschwindender Divisionsrest Linearkombination von  $x$  und  $y$ , und die Koeffizienten können wie in Kapitel I mit dem erweiterten EUKLIDischen Algorithmus berechnet werden. ■

**Lemma:** *a)* In einem EUKLIDischen Ring  $R$  ist jedes Element  $x \neq 0$  mit  $\nu(x) = 0$  eine Einheit.

*b)* Ist  $x = yz \neq 0$ , wobei  $y, z$  keine Einheiten sind, so ist  $\nu(y) < \nu(x)$  und  $\nu(z) < \nu(x)$ .

*Beweis:* *a)* Wir dividieren eins durch  $x$  mit Rest:  $1 : x = q$  Rest  $r$ . Dann ist entweder  $r = 0$  oder aber  $\nu(r) < \nu(x) = 0$ . Letzteres ist nicht möglich, also ist  $qx = 1$  und  $x$  eine Einheit.

*b)* Da  $y$  und  $z$  Teiler von  $x$  sind, sind  $\nu(y), \nu(z) \leq \nu(x)$ . Um zu zeigen, daß  $\nu(y)$  echt kleiner als  $\nu(x)$  ist, dividieren wir  $y$  mit Rest durch  $x$ ; das Ergebnis sei  $q$  Rest  $r$ , d.h.  $y = qx + r$  mit  $r = 0$  oder  $\nu(r) < \nu(x)$ . Wäre  $r = 0$ , wäre  $y$  ein Vielfaches von  $x$ , es gäbe also ein  $u \in R$  mit  $y = ux = u(yz) = (uz)y$ . Damit wäre  $uz = 1$ , also  $z$  eine Einheit, im Widerspruch zur Annahme. Somit ist  $\nu(r) < \nu(x)$ .

Als Teiler von  $x$  ist  $y$  auch Teiler von  $r = y - qx = y(1 - qx)$ , also muß  $\nu(y) \leq \nu(r) < \nu(x)$  sein. Genauso folgt, daß auch  $\nu(z) < \nu(x)$  ist. ■

**Satz:** Jeder EUKLIDische Ring ist faktoriell.

*Beweis:* Wir müssen zeigen, daß jedes Element  $x \neq 0$  aus  $R$  bis auf Reihenfolge und Assoziiertheit eindeutig als Produkt aus einer Einheit und



geeigneten Potenzen irreduzibler Elemente geschrieben werden kann. Wir beginnen damit, daß sich  $x$  überhaupt in dieser Weise darstellen läßt.

Dazu benutzen wir die Betragsfunktion  $\nu: R \setminus \{0\} \rightarrow \mathbb{N}_0$  des EUKLIDischen Rings  $R$  und beweisen induktiv, daß für  $n \in \mathbb{N}_0$  alle  $x \neq 0$  mit  $\nu(x) \leq n$  in der gewünschten Weise darstellbar sind.

Ist  $\nu(x) = 0$ , so ist  $x$  nach obigem Lemma eine Einheit. Diese kann als sich selbst mal dem leeren Produkt von Potenzen irreduzibler Elemente geschrieben werden.

Für  $n > 1$  unterscheiden wir zwei Fälle: Ist  $x$  irreduzibel, so ist  $x = x$  eine Darstellung der gewünschten Form, und wir sind fertig.

Andernfalls läßt sich  $x = yz$  als Produkt zweier Elemente schreiben, die beide keine Einheiten sind. Somit sind nach obigem Lemma  $\nu(y) < \nu(x)$  und  $\nu(z) < \nu(x)$ , beide lassen sich also nach Induktionsvoraussetzung als Produkte von Einheiten und Potenzen irreduzibler Elemente schreiben. Damit läßt sich auch  $x = yz$  so darstellen.

Als nächstes müssen wir uns überlegen, daß diese Darstellung bis auf Reihenfolge und Einheiten eindeutig ist. Das wesentliche Hilfsmittel hierzu ist die folgende Zwischenbehauptung:

*Falls ein irreduzibles Element  $p$  ein Produkt  $xy$  teilt, teilt es mindestens einen der beiden Faktoren.*

Zum *Beweis* betrachten wir den ggT von  $x$  und  $p$ . Dieser ist insbesondere ein Teiler von  $p$ , also bis auf Assoziiertheit entweder  $p$  oder 1. Im ersten Fall ist  $p$  Teiler von  $x$  und wir sind fertig; andernfalls können wir

$$1 = \alpha p + \beta x$$

als Linearkombination von  $p$  und  $x$  schreiben. Multiplikation mit  $y$  macht daraus  $y = \alpha p x + \beta x y$ , und hier sind beide Summanden auf der rechten Seite durch  $p$  teilbar: Bei  $\alpha p x$  ist das klar, und bei  $\beta x y$  folgt es daraus, daß nach Voraussetzung  $p$  ein Teiler von  $x y$  ist. Also ist  $p$  Teiler von  $y$ , und die Zwischenbehauptung ist bewiesen.

Induktiv folgt sofort:

Falls ein irreduzibles Element  $p$  ein Produkt  $\prod_{i=1}^r x_i$  teilt, teilt es mindestens einen der Faktoren  $x_i$ .

Um den Beweis des Satzes zu beenden, zeigen wir induktiv, daß für jedes  $n \in \mathbb{N}_0$  alle Elemente mit  $\nu(x) \leq n$  eine bis auf Reihenfolge und Einheiten eindeutige Primfaktorzerlegung haben.

Für  $n = 0$  ist  $x$  eine Einheit; hier ist die Zerlegung  $x = x$  eindeutig.

Seien nun

$$x = u \prod_{i=1}^r p_i^{e_i} = v \prod_{j=1}^s q_j^{f_j}$$

zwei Zerlegungen eines Elements  $x \in R$ , wobei wir annehmen können, daß alle  $e_i, f_j \geq 1$  sind. Dann ist  $p_1$  trivialerweise Teiler des ersten Produkts, also auch des zweiten. Wegen der Zwischenbehauptung teilt  $p_1$  also mindestens eines der Elemente  $q_j$ , d.h.  $p_1 = wq_j$  ist bis auf eine Einheit  $w$  gleich  $q_j$ . Da  $p_1$  keine Einheit ist, ist  $\nu(x/p_1) < \nu(x)$ ; nach Induktionsannahme hat also  $x/p_1 = x/(wq_j)$  eine bis auf Reihenfolge und Einheiten eindeutige Zerlegung in irreduzible Elemente. Damit hat auch  $x$  diese Eigenschaft. ■

*Bemerkung:* Die Umkehrung dieses Satzes gilt nicht: Beispielsweise sind nach einem Satz von GAUSS auch  $\mathbb{Z}[X]$  sowie Polynomringe in mehr als einer Veränderlichen über  $\mathbb{Z}$  oder einem Körper faktoriell, aber keiner dieser Ringe ist EUKLIDISCH, da sich weder der ggT eins von zwei und  $X$  in  $\mathbb{Z}[X]$  noch der ggT eins von  $X$  und  $Y$  in  $k[X, Y]$  als Linearkombination der Ausgangselemente schreiben läßt.

Wir interessieren uns in diesem Kapitel vor allem für die GAUSSSchen Zahlen und wollen uns überlegen, daß auch diese einen EUKLIDISCHEN Ring bilden.

Dazu brauchen wir zunächst eine Abbildung  $\nu$  nach  $\mathbb{N}_0$ . Für  $\mathbb{Z}$  konnten wir einfach den Betrag nehmen; für die GAUSSSchen Zahlen können wir unser Glück versuchen mit der Norm.

Falls  $\mathbb{Z}[i]$  zusammen mit der Norm ein EUKLIDISCHER Ring ist, muß es zu je zwei Elementen  $r, s \in \mathbb{Z}[i]$  mit  $s \neq 0$  ein Element  $q \in \mathbb{Z}[i]$

geben, so daß  $N(r - sq) < N(s)$  ist. Division durch  $s$  macht daraus die Ungleichung

$$N\left(\frac{r}{s} - q\right) < N(1) = 1.$$

Da sich jedes Element von  $\mathbb{Q}[i]$  als so ein Quotient  $r/s$  mit  $r, s \in \mathbb{Z}[i]$  darstellen läßt, muß es also zu jedem  $z \in \mathbb{Q}[i]$  ein  $w \in \mathbb{Z}[i]$  geben, so daß  $N(z-w) < 1$  ist. Dazu schreiben wir  $z = x+iy$  mit  $x, y \in \mathbb{Q}$  und wählen ganze Zahlen  $u, v$  derart, daß  $|x - u|$  und  $|y - v|$  kleiner oder gleich  $\frac{1}{2}$  sind. Für  $w = u + iv$  ist dann  $N(z - w) = (x - u)^2 + (y - v)^2 \leq \frac{1}{2} < 1$ . Damit ist gezeigt, daß  $\mathbb{Z}[i]$  ein EUKLIDischer und damit auch faktorieller Ring ist.

Betrachten wir als Beispiel die Division von  $23 + 9i$  durch  $2 - 3i$ . In  $\mathbb{Q}[i]$  ist

$$\frac{23 + 9i}{2 - 3i} = \frac{(23 + 9i)(2 + 3i)}{13} = \frac{19}{13} + \frac{87}{13}i.$$

Da  $19 : 13 = 1$  Rest 6 und  $87 : 13 = 6$  Rest 9 ist, liegt das Element  $1 + 7i$  aus  $\mathbb{Z}[i]$  am nächsten bei dieser Zahl. Die Norm von

$$\frac{19}{13} + \frac{87}{13}i - (1 + 7i) = \frac{6}{13} - \frac{4}{13}i$$

ist  $(6^2 + 4^2)/13^2 = 52/169$  und damit deutlich kleiner als eins. Somit ist

$$(23 + 9i) : (2 - 3i) = (1 + 7i) \text{ Rest } -2i$$

ein mögliches Ergebnis der Division mit Rest. Ein anderes wäre

$$(23 + 9i) : (2 - 3i) = (1 + 6i) \text{ Rest } 3,$$

denn auch die Norm von 3 ist kleiner als die von  $2 + 3i$ . (Der Rest wurde jeweils als Dividend minus Divisor mal Quotient berechnet.)

### §3: Quaternionen

Nachdem durch die komplexen Zahlen  $\mathbb{R}^2$  mit der Struktur eines Körpers versehen war, versuchten viele Mathematiker ähnliches auch für  $\mathbb{R}^3$  zu erreichen. Natürlich kann weder  $\mathbb{R}^3$  noch sonst ein  $\mathbb{R}^n$  mit  $n > 2$  zu einem Körper gemacht werden, denn ein solcher Körper wäre eine

algebraische Erweiterung von  $\mathbb{R}$ ; da aber der algebraische Abschluß von  $\mathbb{R}$  gleich  $\mathbb{C}$  ist, muß dann  $n = 1$  oder  $n = 2$  sein.

Die damaligen Mathematiker waren jedoch bescheidener: Ihnen genügte es, einfach irgendeine Art von Multiplikation zu finden, die nicht unbedingt allen Körperaxiomen genügte – von Körpern sprach damals ohnehin noch niemand.

Erst 1940 konnte HEINZ HOPF (1894–1971) (auf dem Umweg über Vektorfelder auf Sphären) zeigen, daß das nicht möglich ist: Selbst eine bilineare Abbildung  $\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  kann nur dann existieren, wenn  $n$  eine Zweierpotenz ist, und 1958 zeigten dann unabhängig voneinander und mit verschiedenen Methoden JOHN MILNOR und MICHEL KERVAIRE, daß auch noch  $n \leq 8$  sein muß, so daß nur die vier Möglichkeiten  $n = 1, 2, 4$  und  $8$  in Frage kommen. Genau in diesen Fällen waren auch bereits entsprechende Produkte bekannt: Für  $n = 1$  und  $2$  haben wir natürlich die reelle bzw. komplexe Multiplikation. Den Fall  $n = 4$  löste HAMILTON 1843: Er fand eine Multiplikation auf  $\mathbb{R}^4$ , die zwar nicht kommutativ ist, ansonsten aber alle Körperaxiome erfüllt. Man spricht in so einem Fall von einem *Schiefkörper* oder, in der neueren Literatur, einer *Divisionsalgebra*. HAMILTON bezeichnete seine vierdimensionalen Zahlen als *Quaternionen*. Kurz danach konstruierte ARTHUR CAYLEY (1821–1895) ein nicht-assoziatives Produkt auf  $\mathbb{R}^8$ ; die so erhaltenen „Zahlen“ nannte er *Oktaven*.



WILLIAM ROWEN HAMILTON (1805–1865) wurde in Dublin geboren; bereits mit fünf Jahren sprach er Latein, Griechisch und Hebräisch. Mit dreizehn begann er, mathematische Literatur zu lesen, mit 21 wurde er, noch als Student, Professor der Astronomie am Trinity College in Dublin. Er verlor allerdings schon bald sein Interesse an der Astronomie und beschäftigte sich stattdessen mit mathematischen und physikalischen Problemen. Am bekanntesten ist er für seine Entdeckung der Quaternionen, vorher publizierte er aber auch bedeutende Arbeiten über Optik, Dynamik und Algebra.

HAMILTON wählte eine Basis von  $\mathbb{H} = \mathbb{R}^4$ , die aus der Eins sowie drei „imaginären Einheiten“  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  besteht, d.h.  $\mathbf{i}^2 = \mathbf{j}^2 = \mathbf{k}^2 = -1$ . Außerdem postulierte er, daß  $\mathbf{ij} = -\mathbf{ji} = \mathbf{k}$  sein sollte; daraus lassen sich dann über

das Assoziativgesetz auch die anderen Produkte imaginärer Einheiten berechnen.

Damit ist, wenn man die Gültigkeit des Distributivgesetzes postuliert, eine Multiplikation auf  $\mathbb{R}^4$  definiert; der Beweis, daß hierbei alle Körperaxiome außer der Kommutativität der Multiplikation erfüllt sind, enthält wie üblich nur einen etwas schwierigeren Punkt, die Existenz von Inversen; der Rest ist mühsame Abhakerei.

Zum Glück fand CAYLEY 1858 einen einfacheren Weg: Die vier komplexen  $2 \times 2$ -Matrizen

$$E = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad I = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \quad J = \begin{pmatrix} 0 & i \\ i & 0 \end{pmatrix} \quad \text{und} \quad K = \begin{pmatrix} i & 0 \\ 0 & -i \end{pmatrix}$$

erfüllen dieselben Relationen

$$I^2 = J^2 = K^2 = -E \quad \text{und} \quad IJ = -JI = K;$$

wir können also die Quaternion  $a + bi + cj + dk$  identifizieren mit der Matrix

$$aE + bI + cJ + dK = \begin{pmatrix} a + di & b + ci \\ -b + ci & a - di \end{pmatrix} \in \mathbb{C}^{2 \times 2}.$$

Da für Matrizen das Assoziativgesetz wie auch das Distributivgesetz gelten, ist klar, daß das Produkt zweier solcher Matrizen wieder von derselben Form ist und daß auch die Quaternionenmultiplikation Assoziativ- und Distributivgesetz erfüllt.

Die Quaternionen entsprechen somit genau den komplexen  $2 \times 2$ -Matrizen der Form

$$\begin{pmatrix} \alpha & \beta \\ -\bar{\beta} & \bar{\alpha} \end{pmatrix} \quad \text{mit} \quad \alpha = a + di, \beta = b + ci.$$

Die Determinante dieser Matrix ist  $\alpha\bar{\alpha} + \beta\bar{\beta} = a^2 + b^2 + c^2 + d^2$ .

Definieren wir in Analogie zum Fall der quadratischen Zahlkörper wieder das konjugierte Element zu  $\gamma = a + bi + cj + dk$  als die Quaternion  $\bar{\gamma} = a - bi - cj - dk$ , so entspricht  $\bar{\gamma}$  der Matrix

$$\begin{pmatrix} \bar{\alpha} & -\beta \\ \bar{\beta} & \alpha \end{pmatrix} \quad \text{und} \quad \begin{pmatrix} \alpha & \beta \\ -\bar{\beta} & \bar{\alpha} \end{pmatrix} \begin{pmatrix} \bar{\alpha} & -\beta \\ \bar{\beta} & \alpha \end{pmatrix} = (\alpha\bar{\alpha} + \beta\bar{\beta})E.$$

Damit folgt insbesondere, daß  $\gamma\bar{\gamma}$  eine reelle Zahl ist, die genau dann verschwindet, wenn  $\gamma = 0$  ist. Wir bezeichnen diese Zahl wieder als die *Norm*  $N(\gamma)$  der Quaternion  $\gamma$ , und wieder ist  $\bar{\gamma}/N(\gamma)$  das multiplikative Inverse zu  $\gamma$  – sowohl für die Links- als auch die Rechtsmultiplikation.

$N(\gamma)$  ist gleichzeitig die Determinante der  $\gamma$  zugeordneten Matrix; aus dem Multiplikationssatz für Determinanten folgt daher sofort die Formel

$$N(\gamma\delta) = N(\gamma)N(\delta).$$

# Kapitel 7

## Quadratische Formen

Eine quadratische Form ist ein Ausdruck der Form

$$F(x, y) = Ax^2 + Bxy + Cy^2 \quad \text{mit } A, B, C \in \mathbb{Z};$$

die Zahlentheorie interessiert sich vor allem dafür, welche Werte  $F(x, y)$  für  $x, y \in \mathbb{Z}$  annehmen kann.

### § 1: Summen zweier Quadrate

Der einfachste Fall ist die Form  $F(x, y) = x^2 + y^2$ , die offensichtlich keine negativen Werte annehmen kann. Sie hängt eng zusammen mit dem Ring  $\mathbb{Z}[i] = \mathbb{Z} + i\mathbb{Z}$  der GAUSSSchen Zahlen, d.h. der komplexen Zahlen mit ganzzahligem Real- und Imaginärteil, denn

$$x^2 + y^2 = (x + iy)(x - iy).$$

Eine Zahl  $n \in \mathbb{N}_0$  ist also genau dann als Summe zweier Quadrate darstellbar, wenn es eine GAUSSSche ganze Zahl  $x + iy \in \mathbb{Z}[i]$  gibt, so daß  $n$  das Produkt von  $x + iy$  mit seiner konjugiert komplexen Zahl  $\overline{x + iy} = x - iy$  ist.

Das Quadrat einer geraden Zahl ist durch vier teilbar, das einer ungeraden Zahl  $2k + 1$  ist  $4k^2 + 4k + 1 \equiv 1 \pmod{4}$ ; somit ist jede Summe zweier Quadrate kongruent null, eins oder zwei modulo vier. Eine Zahl kongruent drei modulo vier kann also nicht als Summe zweier Quadratzahlen auftreten.

Auf der Suche nach positiven Ergebnissen können wir uns auf Primzahlen beschränken, denn wie FIBONACCI bereits im dreizehnten Jahrhundert zeigte, gilt:

**Lemma:** Sind zwei Zahlen  $n, m \in \mathbb{N}$  darstellbar als Summen zweier Quadrate, so gilt dasselbe für ihr Produkt  $nm$ .

*Beweis:* Wenn  $n = a^2 + b^2$  und  $m = c^2 + d^2$  als Summen zweier Quadrate darstellbar sind, gilt für  $\alpha = a + ib$  und  $\beta = c + id \in \mathbb{Z}[i]$ , so daß  $n = \alpha\bar{\alpha}$  und  $m = \beta\bar{\beta}$  ist. Dann ist

$$nm = (\alpha\bar{\alpha})(\beta\bar{\beta}) = (\alpha\beta)(\overline{\alpha\beta}).$$

Wegen  $\alpha\beta = (ac - bd) + i(ad + bc)$  ist also  $nm = (ac - bd)^2 + (ad + bc)^2$ . ■

FIBONACCI bewies dieses Lemma natürlich nicht auf dem Umweg über GAUSSsche Zahlen; er fand die obige Formel wahrscheinlich durch geschicktes Probieren,

$2 = 1^2 + 1^2$  ist als Summe zweier Quadrate darstellbar; wir müssen daher nur die ungeraden Primzahlen untersuchen. Hier wissen wir bereits, daß Zahlen kongruent drei modulo vier keine Summen zweier Quadrate sein können.

**Satz:** Eine ungerade Primzahl  $p$  ist genau dann darstellbar als Summe zweier Quadrate, wenn  $p \equiv 1 \pmod{4}$ . Diese Darstellung ist (abgesehen von den Vorzeichen) eindeutig bis auf die Reihenfolge der Summanden.

*Beweis:* Aus Kapitel I, §8 wissen wir, daß die multiplikative Gruppe des Körper  $\mathbb{F}_p$  von einem einzigen Element  $g$  erzeugt wird. Für  $p = 4k + 1$  ist dann  $g^{4k} = 1$ , also  $g^{2k} = -1$ . Somit ist  $-1 = p - 1$  in  $\mathbb{F}_p$  ein Quadrat.

In  $\mathbb{Z}$  gibt es daher Zahlen  $x$ , für die  $x^2 \equiv -1 \pmod{p}$  ist oder, anders ausgedrückt,  $x^2 + 1 = \ell p$  für ein  $\ell \in \mathbb{N}$ . Da jede Restklasse modulo  $p$  einen Vertreter mit Betrag kleiner  $p/2$  enthält, können wir dabei annehmen, daß  $|x| < p/2$  ist; dann ist mit einer geeigneten natürlichen Zahl  $\ell$

$$x^2 + 1 = \ell p < \frac{p^2}{4} + 1 < \frac{p^2}{2} \implies \ell < p.$$

Es gibt also ein  $\ell < p$ , so daß  $\ell p$  Summe zweier Quadrate ist. Das kleinste solche  $\ell$  sei  $m$ ; wir müssen zeigen, daß  $m = 1$  ist.

Zunächst ist klar, daß  $m$  eine ungerade Zahl sein muß, denn aus der Formel  $x^2 + y^2 = mp$  mit geradem  $m$  folgt, daß  $x$  und  $y$  entweder beide



gerade oder beide ungerade sind;  $x \pm y$  sind also gerade und

$$\left(\frac{x+y}{2}\right)^2 + \left(\frac{x-y}{2}\right)^2 = \frac{x^2+y^2}{2} = \frac{m}{2}p,$$

im Widerspruch zur Minimalität von  $m$ .

Falls die Behauptung falsch wäre, müßte somit  $m \geq 3$  sein. Wir definieren dann zwei neue Zahlen  $u, v \in \mathbb{Z}$  durch die Bedingungen

$$|u| < \frac{m}{2}, \quad |v| < \frac{m}{2}, \quad u \equiv y \pmod{m} \quad \text{und} \quad v \equiv x \pmod{m}.$$

Offensichtlich können nicht beide dieser Zahlen verschwinden, denn sonst wären  $x$  und  $y$  beide durch  $m$  teilbar, also wäre  $x^2 + y^2 = mp$  durch  $m^2$  teilbar und  $p$  durch  $m$ . Das kann aber nicht sein, denn  $p$  ist prim und  $1 < m < p$ . Weiter ist

$$u^2 + v^2 \equiv x^2 + y^2 = mp \equiv 0 \pmod{m},$$

also gibt es eine natürliche Zahl  $r$ , so daß  $u^2 + v^2 = rm$  ist. Da  $u^2 + v^2$  kleiner ist als  $\frac{1}{2}m^2$ , ist  $r < \frac{m}{2}$ .

Nach der zu Beginn des Paragraphen zitierten Formel von FIBONACCI, d.h. also durch explizite Berechnung von  $(u+iv)(x+iy)$  und Berechnung der Norm davon, erhalten wir die Formel.

$$(rm)(mp) = (u^2 + v^2)(x^2 + y^2) = (xu - yv)^2 + (xv + yu)^2.$$

Dabei ist nach Definition von  $u$  und  $v$

$$xu - yv \equiv xy - yx = 0 \pmod{m} \quad \text{und} \quad xv + yu \equiv x^2 + y^2 \equiv 0 \pmod{m},$$

beide Zahlen sind also durch  $m$  teilbar. Somit gibt es natürliche Zahlen  $X, Y$  mit

$$(rm)(mp) = m^2rp = (mX)^2 + (mY)^2 \quad \text{oder} \quad rp = X^2 + Y^2.$$

Da  $r < \frac{m}{2}$ , widerspricht dies der Minimalität von  $m$ .

Damit haben wir gezeigt, daß  $m = 1$  sein muß, d.h.  $p$  läßt sich als Summe zweier Quadrate darstellen. Wir müssen uns noch überlegen, daß diese Darstellung bis auf die Reihenfolge der Faktoren eindeutig ist.

Angenommen, es gibt zwei Darstellungen  $p = x^2 + y^2 = u^2 + v^2$ . In  $\mathbb{Z}[i]$  ist dann

$$p = (x + iy)(x - iy) = (u + iv)(u - iv).$$

Alle Faktoren haben Norm  $p$  und sind somit irreduzibel, und aus §5 des vorigen Kapitels wissen wir, daß  $\mathbb{Z}[i]$  ein EUKLIDischer, insbesondere also faktorieller Ring ist. Daher unterscheiden sich die beiden Zerlegungen nur durch Einheiten von  $\mathbb{Z}[i]$ . Auch diese kennen wir aus Kapitel 6: Nach dem Lemma aus §6 sind es genau die Elemente  $\pm 1$  und  $\pm i$ . Somit ist entweder  $x^2 = u^2$  und  $y^2 = v^2$  oder umgekehrt, womit die Eindeutigkeit bis auf Reihenfolge und Vorzeichen bewiesen wäre. ■

Als erste Anwendung davon können wir die Primzahlen im Ring  $\mathbb{Z}[i]$  der GAUSSschen Zahlen bestimmen:

**Korollar:** Eine Primzahl  $p \in \mathbb{N}$  ist genau dann irreduzibel in  $\mathbb{Z}[i]$ , wenn  $p \equiv 3 \pmod{4}$ . Andernfalls zerfällt sie in das Produkt zweier konjugiert komplexer irreduzibler Elemente  $r \pm is$  mit  $r^2 + s^2 = p$ .

*Beweis:*  $p = 2 = (1 + i)(1 - i)$  zerfällt offensichtlich, und dies ist bereits die Primzerlegung, denn  $N(1 \pm i) = 2$  hat keine echten Teiler.

Falls eine ungerade Primzahl  $p$  einen echten Teiler  $r + is$  hat, ist sie auch durch  $r - is$  teilbar. Da die Norm von  $p$  gleich  $p^2$  ist und  $r \pm is$  keine Einheiten, muß  $N(r \pm is) = p$  sein. Damit folgt zunächst, daß  $r \pm is$  prim sind, denn ein echter Teiler müßte als Norm einen echten Teiler von  $p$  haben. Außerdem folgt, daß sich  $(r + is)(r - is) = r^2 + s^2$  höchstens durch eine Einheit von  $p$  unterscheidet. Da beides positive Zahlen sind, muß diese gleich eins sein, d.h. die Primzerlegung von  $p$  in  $\mathbb{Z}[i]$  ist

$$p = (r + is)(r - is) = r^2 + s^2.$$

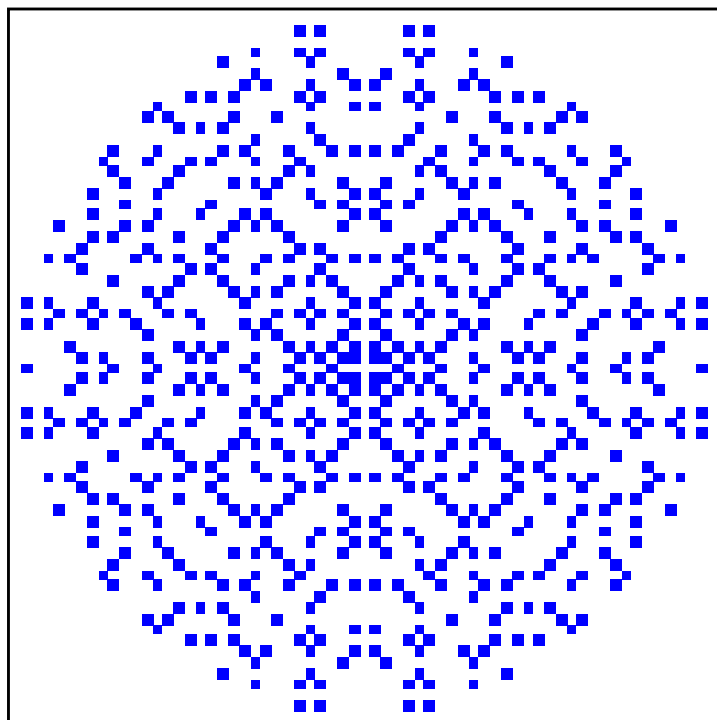
Nach dem Satz ist daher  $p \equiv 1 \pmod{4}$ .

Ist umgekehrt  $p \equiv 1 \pmod{4}$ , so gibt es nach dem Satz zwei ganze Zahlen  $r, s$ , so daß  $p = r^2 + s^2$  ist, d.h.  $p = (r + is)(r - is)$  zerfällt in  $\mathbb{Z}[i]$ , und das Argument aus dem vorigen Abschnitt zeigt, daß dies die Primzerlegung ist.

Somit zerfallen genau die Primzahlen  $p \equiv 1 \pmod{4}$  und die Zwei, d.h. genau die  $p \equiv 3 \pmod{4}$  bleiben prim. ■

In der Abbildung sind die GAUSSschen Primzahlen  $a + ib$  der Norm höchstens 1000 durch Quadrate um den Punkt  $(a, b) \in \mathbb{R}^2$  dargestellt.

Mancher Leser wird hier ein gelegentlich von Designern verwendetes Muster erkennen.



Kehren wir zurück zur Ausgangsfrage: Wann kann eine vorgegebene natürliche Zahl als Summe zweier Quadrate dargestellt werden?

**Satz:** Eine natürliche Zahl  $n$  läßt sich genau dann als Summe zweier Quadrate schreiben, wenn jeder Primteiler  $p \equiv 3 \pmod{4}$  in der Primzerlegung von  $n$  mit einer geraden Potenz auftritt.

*Beweis:* Zunächst ist die Bedingung hinreichend, denn da mit  $n$  auch jedes Produkt  $c^2 n$  als Summe zweier Quadrate darstellbar ist, können wir die Primteiler  $p \equiv 3 \pmod{4}$  ignorieren. Nach dem gerade bewiesenen Satz wissen wir, daß jede Primzahl  $p \equiv 1 \pmod{4}$  Summe zweier Quadrate ist, und natürlich gilt dies auch für  $2 = 1^2 + 1^2$ . Damit ist nach dem obigen Lemma auch jedes Produkt solcher Primzahlen als Summe zweier Quadrate darstellbar.

Umgekehrt sei  $n = x^2 + y^2$  und  $d = \text{ggT}(x, y)$ . Mit  $x = du$ ,  $y = dv$  und  $n = d^2 m$  ist dann  $m = u^2 + v^2$ , und  $m$  enthält genau dann einen Primteiler  $p \equiv 3 \pmod{4}$  in ungerader Potenz, wenn dies für  $n$  der Fall ist.

Ein solcher Primteiler  $p$  teilt auch  $u^2 + v^2 = (u + iv)(u - iv)$  im Ring  $\mathbb{Z}[i]$  der GAUSSSchen Zahlen. Falls  $p$  auch dort eine Primzahl ist, muß  $p$  mindestens einen der beiden Faktoren teilen; komplexe Konjugation zeigt, daß es dann auch den anderen teilt. Damit teilt es auch deren Summe  $2u$  und Differenz  $2iv$ ; da  $p$  ungerade ist und  $i$  eine Einheit, teilt  $p$  also die zueinander teilerfremden Zahlen  $u$  und  $v$ , ein Widerspruch.

Somit ist  $p$  in  $\mathbb{Z}[i]$  keine Primzahl; nach obigem Korollar muß daher  $p = 2$  oder  $p \equiv 1 \pmod{4}$  sein. Damit ist jeder Primteiler  $p \equiv 3 \pmod{4}$  von  $n$  zugleich ein Teiler von  $d$  und tritt in  $n$  daher mit einer geraden Potenz auf. ■

Für zusammengesetzte Zahlen ist die Darstellung als Summe zweier Quadrate im allgemeinen nicht mehr eindeutig. Über die Primzerlegung in  $\mathbb{Z}[i]$  läßt sich die Anzahl verschiedener Darstellungen leicht erkennen: Natürlich entsprechen auch für eine beliebige natürliche Zahl  $n$  die Darstellungen als Summe zweier Quadrate den Darstellungen von  $n$  als Norm eines Elements von  $\mathbb{Z}[i]$ , wobei assoziierte Elemente bis auf Reihenfolge auf dieselbe Zerlegung führen.

Aus der Primzerlegung von  $n$  in  $\mathbb{Z}$  können wir leicht auf die Primzerlegung in  $\mathbb{Z}[i]$  schließen: Primzahlen kongruent drei modulo vier bleiben nach obigem Korollar auch in  $\mathbb{Z}[i]$  irreduzibel, die kongruent eins modulo vier sind Produkte zweier konjugierter Elemente  $x \pm iy$ . Die beiden Faktoren sind nicht assoziiert, denn sonst wäre  $|x| = |y|$  und  $p = x^2 + y^2$  wäre gerade. Die Zwei schließlich ist Produkt der beiden irreduziblen Elemente  $1 \pm i$ , und die sind assoziiert zueinander, denn  $(1 - i) \cdot i = 1 + i$ .

Wir sortieren daher in der Primzerlegung von  $n$  nach den Kongruenzklassen modulo vier der Primfaktoren:

$$n = 2^e \prod_{j=1}^r p_j^{f_j} \prod_{k=1}^s q_k^{2g_k} \quad \text{mit} \quad p_j \equiv 1 \pmod{4}, \quad q_k \equiv 3 \pmod{4}.$$

Für jedes  $p_j$  wählen wir ein  $\pi_j \in \mathbb{Z}[i]$  derart, daß  $\pi_j \cdot \bar{\pi}_j = p_j$  ist; dann

ist  $n$  in  $\mathbb{Z}[i]$  assoziiert zu

$$(1+i)^{2e} \prod_{j=1}^r \pi_j^{f_j} \prod_{j=1}^r \bar{\pi}_j^{f_j} \prod_{k=1}^s q_k^{2g_k}.$$

Ein Element  $\alpha \in \mathbb{Z}[i]$ , für das  $N(\alpha) = n$  sein soll, hat daher bis auf eine Einheit die Form

$$\alpha = (1+i)^e \prod_{j=1}^r \pi_j^{h_j} \prod_{j=1}^r \bar{\pi}_j^{f_j - h_j} \prod_{k=1}^s q_k^{g_k},$$

mit  $0 \leq h_j \leq f_j$ . Die Anzahl verschiedener Möglichkeiten ist somit gleich dem Produkt der  $(f_j + 1)$ , wobei hier allerdings die Darstellungen  $n = x^2 + y^2$  und  $n = y^2 + x^2$  für  $x \neq y$  als verschieden gezählt werden.

Die im Vergleich zur Größe von  $n$  meisten verschiedenen Darstellungen gibt es offenbar dann, wenn  $n$  ein Produkt verschiedener Primzahlen ist, die allesamt kongruent eins modulo vier sind. In diesem Fall ist die Anzahl der Darstellungen gleich zwei hoch Anzahl der Faktoren.

## §2: Anwendung auf die Berechnung von $\pi$

Aus der Analysis I ist bekannt, daß gilt

$$\arctan x = x - \frac{x^3}{3} + \frac{x^5}{5} - \frac{x^7}{7} + \frac{x^9}{9} - \frac{x^{11}}{11} + \frac{x^{13}}{13} - \frac{x^{15}}{15} + \dots;$$

falls es jemand nicht mehr weiß: Die Ableitung des Arkustangens ist  $1/(1+x^2)$ , und nach der Summenformel für die geometrische Reihe ist

$$\frac{1}{1+x^2} = 1 - x^2 + x^4 - x^6 + x^8 - x^{10} + x^{12} - x^{14} + \dots.$$

Durch gliedweise Integration folgt wegen  $\arctan 0 = 0$  die obige Formel. Eine bekannte Anwendung davon ist der Spezialfall  $x = 1$ :

$$\frac{\pi}{4} = \arctan 1 = 1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \frac{1}{11} + \frac{1}{13} - \frac{1}{15} + \dots.$$

Zur praktischen Berechnung von  $\pi$  ist diese Formel allerdings völlig unbrauchbar und der Alptraum eines jeden Numerikers: Zunächst einmal sind alternierende Summen grundsätzlich problematisch, allerdings ist

das hier vergleichsweise harmlos: Wenn wir jeden negativen Summanden von seinem Vorgänger subtrahieren, bekommen wir eine Reihe

$$\frac{\pi}{4} = \frac{2}{1 \cdot 3} + \frac{2}{5 \cdot 7} + \frac{2}{9 \cdot 11} + \frac{2}{13 \cdot 15} + \dots$$

mit lauter positiven Gliedern. Die Summanden sind jedoch immer noch monoton fallend, so daß die Rundungsfehler der ersten Additionen bei hinreichend langer Summation größer sind als die hinteren Summanden. Man muß also, wenn man eine endliche Teilsumme berechnen will, von hinten nach vorne summieren und damit bereits vor Beginn der Rechnung die Anzahl der Terme festlegen. Bei jeder Erhöhung der Anzahl der Summanden muß die gesamte Rechnung von vorne beginnen.

Dazu kommt, daß die Reihe extrem langsam konvergiert: Dividieren wir obige Gleichung durch zwei und berechnen für

$$\frac{\pi}{8} = \sum_{n=0}^{\infty} \frac{1}{(4n+1)(4n+3)}$$

die Teilsummen

$$S_N = \sum_{n=0}^N \frac{1}{(4n+1)(4n+3)},$$

so erhalten wir für die ersten Zehnerpotenzen  $N$  die folgenden Fehler:

$N$	10	100	1 000	10 000
$\pi - 8S_N$	$4,5 \cdot 10^{-2}$	$5,0 \cdot 10^{-3}$	$5,0 \cdot 10^{-4}$	$5,0 \cdot 10^{-5}$
$N$	100 000	1 000 000	10 000 000	100 000 000
$\pi - 8S_N$	$5,0 \cdot 10^{-6}$	$5,0 \cdot 10^{-7}$	$5,0 \cdot 10^{-8}$	$5,0 \cdot 10^{-9}$

Für eine zusätzliche Dezimalstelle muß also der Rechenaufwand ziemlich genau verzehnfacht werden. Angesichts der Tatsache, daß heute mehrere Billionen Ziffern von  $\pi$  bekannt sind, ist klar, daß es bessere Wege zur Berechnung von  $\pi$  geben muß.

Einer davon benutzt Zahlen mit einer großen Anzahl verschiedener Darstellungen als Summen zweier Quadrate. Die Reihe für den Arkustangens konvergiert sicherlich umso besser, je kleiner der Wert von  $x$  ist. Wenn wir also den Winkel  $\frac{\pi}{4}$  aufteilen können in mehrere kleine Winkel,

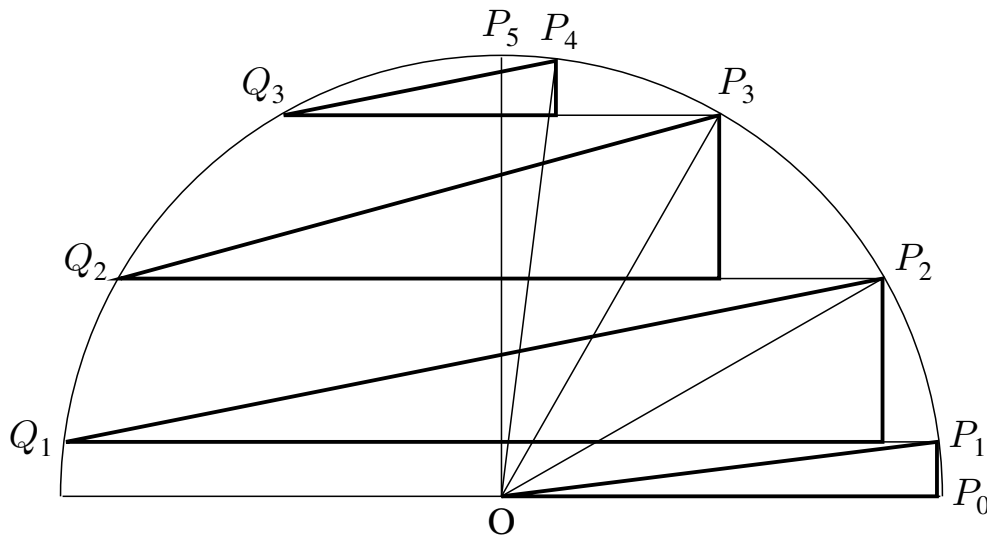
deren Tangens wir kennen, sollten bessere Ergebnisse zu erwarten sein. Genau das können wir mit solchen Zahlen erreichen.

Angenommen, wir haben für eine Zahl  $n$  die  $r$  verschiedenen Darstellungen

$$n = x_1^2 + y_1^2 = \cdots = x_r^2 + y_r^2$$

als Summen von Quadraten, wobei  $y_1 < \cdots < y_r$  sei. Dann ist  $x_i = y_{r-i}$ , denn wir können ja in jeder Darstellung die Reihenfolge der Faktoren vertauschen. Wir wollen außerdem voraussetzen, daß  $n$  nicht das Doppelte eines Quadrats ist, so daß stets  $x_i \neq y_i$  und somit  $r$  eine gerade Zahl ist.

Die Punkte  $P_i = (x_i, y_i)$  und  $Q_i = (-x_i, y_i)$  für  $i = 1, \dots, r$  liegen auf der Kreislinie  $x^2 + y^2 = n$  um den Nullpunkt  $O$ , genauso die drei Punkte  $P_0 = (\sqrt{n}, 0)$ ,  $Q_0 = (-\sqrt{n}, 0)$  und  $P_{r+1} = (0, \sqrt{n})$ .



Da die  $y$ -Koordinaten  $y_i$  der  $P_i$  der Größe nach geordnet sind, ist

$$\frac{\pi}{2} = \sum_{i=0}^r \angle OP_i P_{i+1} = 2 \sum_{i=0}^{r/2-1} \angle OP_i P_{i+1} + \angle OP_{r/2} P_{r/2+1}.$$

Leider ist keines der Dreiecke  $\triangle OP_i P_{i+1}$  rechtwinklig, so daß uns die ganzzahligen Koordinaten der (meisten)  $P_i$  bei der Berechnung der Winkel  $\angle OP_i P_{i+1}$  nichts nützen.

Nun lehrt uns aber ein Satz der Elementargeometrie, der (im Anhang zu diesem Paragraphen bewiesene) Satz vom Innenwinkel, daß der Winkel  $\angle OP_i P_{i+1}$  doppelt so groß ist wie der Winkels  $\angle Q_i P_i P_{i+1}$ . Letzterer gehört zu einem rechtwinkligen Dreieck, denn natürlich ändert sich nichts am Winkel, wenn wir den Punkt  $P_i$  ersetzen durch die senkrechte Projektion  $P'_i = (x_{i+1}, y_i)$  von  $P_{i+1}$  auf die Gerade  $Q_i P_i$ . Somit ist

$$\frac{\pi}{2} = 2\angle OP'_0 P_1 + 4 \sum_{i=1}^{r/2-1} \angle Q_i P'_i P_{i+1} + 2\angle Q_{r/2} P'_{r/2} P_{r/2+1}.$$

Division durch zwei macht daraus

$$\frac{\pi}{4} = \angle OP'_0 P_1 + 2 \sum_{i=1}^{r/2-1} \angle Q_i P'_i P_{i+1} + \angle Q_{r/2} P'_{r/2} P_{r/2+1}.$$

In dieser Darstellung sind die drei Punkte, die den Winkel bestimmen, in allen Fällen die Eckpunkte eines rechtwinkligen Dreiecks, sie haben alle samt ganzzahlige Koordinaten, und zumindest die Katheten der Dreiecke haben ganzzahlige Längen. Somit können wir alle auftretenden Winkel ausdrücken durch Arkustangenswerte rationaler Zahlen.

Als Beispiel betrachten wir das kleinste Produkt dreier verschiedener Primzahlen kongruent eins modulo vier, also  $n = 5 \cdot 13 \cdot 17 = 1105$ . Aus den Darstellungen

$$5 = 1^2 + 2^2, \quad 13 = 2^2 + 3^2 \quad \text{und} \quad 17 = 1^2 + 4^2$$

verschafft man sich leicht die vier Darstellungen

$$1105 = 4^2 + 33^2 = 9^2 + 32^2 = 12^2 + 31^2 = 23^2 + 24^2,$$

zu denen natürlich auch noch vier mit vertauschten Faktoren kommen. Wir haben also

$$P_1 = (33, 4), \quad P_2 = (32, 9), \quad P_3 = (31, 12), \quad P_4 = (24, 23), \\ P_8 = (4, 33), \quad P_7 = (9, 32), \quad P_6 = (12, 31), \quad P_5 = (23, 24);$$

dazu kommen noch die beiden Randpunkte  $P_0 = (\sqrt{1105}, 0)$  sowie  $P_9 = (0, \sqrt{1105})$ .



Die  $Q_i$  für  $1 \leq i \leq 8$  unterscheiden sich von den  $P_i$  nur durch das Vorzeichen der Abszisse. Damit können wir die Tangenten aller Winkel bei  $O$  berechnen:

$$\tan \angle OP_0P_1 = \tan \angle OP_8P_9 = \frac{y_1}{x_1} = \frac{4}{33}$$

$$\tan \angle OP_1P_2 = \tan \angle OP_7P_8 = \tan 2\angle Q_1P_1P_2 = \frac{y_2 - y_1}{x_1 + x_2} = \frac{5}{65} = \frac{1}{13}$$

$$\tan \angle OP_2P_3 = \tan \angle OP_6P_7 = \tan 2\angle Q_2P_2P_3 = \frac{y_3 - y_2}{x_2 + x_3} = \frac{3}{63} = \frac{1}{21}$$

$$\tan \angle OP_3P_4 = \tan \angle OP_5P_6 = \tan 2\angle Q_3P_3P_4 = \frac{y_4 - y_3}{x_3 + x_4} = \frac{11}{55} = \frac{1}{5}$$

$$\tan \angle OP_4P_5 = \tan 2\angle Q_4P_4P_5 = \frac{y_5 - y_4}{x_4 + x_5} = \frac{1}{47}$$

Die Summe aller dieser Winkel ist

$$\frac{\pi}{4} = \arctan \frac{4}{33} + 2 \arctan \frac{1}{13} + 2 \arctan \frac{1}{21} + 2 \arctan \frac{1}{5} + \arctan \frac{1}{47}.$$

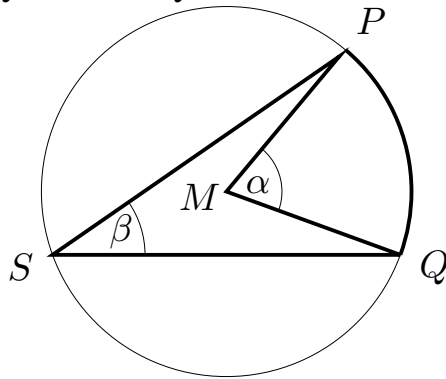
Approximieren wir dies, indem wir jede der TAYLOR-Reihen durch das TAYLOR-Polynom vom Grad  $n$  ersetzen, erhalten wir die folgenden betragsmäßigen Abweichungen  $\Delta_n$  zwischen  $\pi$  und dem Vierfachen dieser Summe:

$n$	1	3	5	7	9
$\Delta_n$	$2,5 \cdot 10^{-2}$	$5,2 \cdot 10^{-4}$	$1,4 \cdot 10^{-5}$	$4,4 \cdot 10^{-7}$	$1,4 \cdot 10^{-8}$
$n$	11	13	15	17	19
$\Delta_n$	$5 \cdot 10^{-10}$	$4,9 \cdot 10^{-10}$	$6 \cdot 10^{-13}$	$2,1 \cdot 10^{-14}$	$7,7 \cdot 10^{-16}$

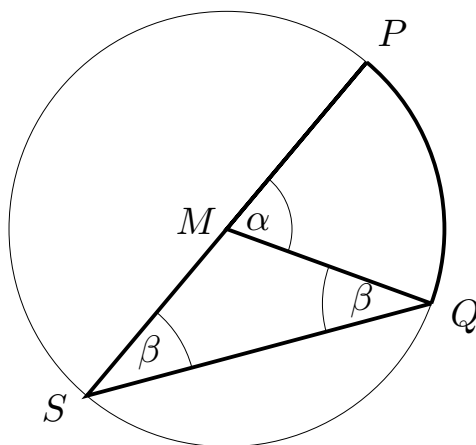
Die Verbesserung gegenüber der Berechnung via  $\frac{\pi}{4} = \arctan 1$  ist dramatisch: Die dort betrachtete Teilsumme  $S_N$  entspricht der Auswertung des TAYLOR-Polynoms vom Grad  $n = 4N + 3$ , und selbst wenn wir  $N$  auf hundert Millionen setzen, haben wir noch einen Fehler von  $5 \cdot 10^{-7}$ . Mit dem neuen Ansatz kommen wir bereits mit TAYLOR-Polynomen vom Grad neun auf einen Fehler, der gerade mal ein Zehntel davon beträgt. An Stelle von hundert Millionen Summanden mußten wir dazu nur fünf TAYLOR-Polynome mit jeweils fünf Summanden auswerten.

### Anhang: Der Satz vom Innenwinkel

**Satz:**  $P, Q, S$  seien Punkte auf einer Kreislinie mit Mittelpunkt  $M$ . Dann ist  $\angle MPQ = 2\angle SPQ$ .

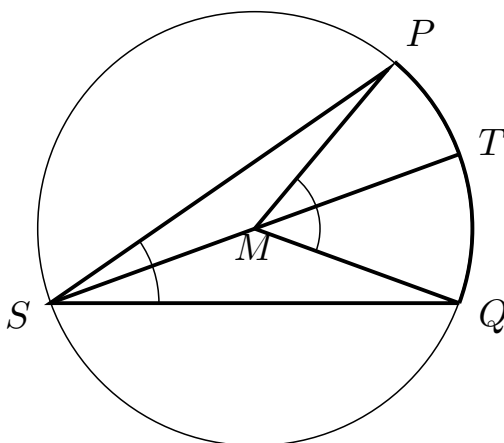


*Beweis:* Am einfachsten ist der Fall, daß  $M$  auf der Verbindungsstrecke



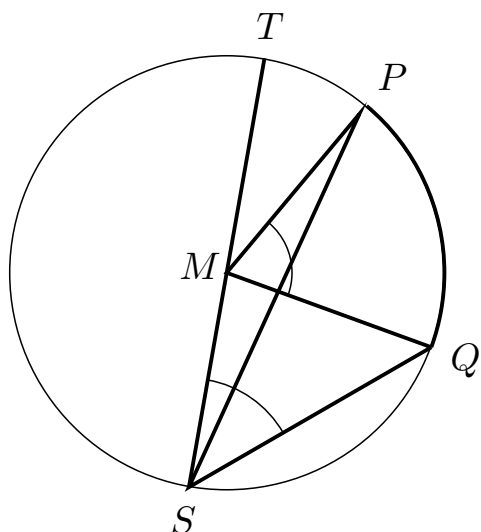
von  $S$  mit einem der beiden Punkte  $P$  und  $Q$  liegt; wir nehmen an, er liege auf  $\overline{SP}$ . (Der andere Fall ist spiegelsymmetrisch dazu und geht genauso.) Dann ist das Dreieck  $\triangle MSQ$  gleichschenkelig, d.h. wir haben bei  $S$  und bei  $Q$  denselben Winkel  $\beta$ . Der verbleibende Dreieckswinkel bei  $M$  ist somit  $\pi - 2\beta$ . Andererseits ist dies aber der Komplementärwinkel zu  $\alpha = \angle MPQ$ , also ist  $\alpha = 2\beta$ , wie behauptet.

Der allgemeine Fall kann auf diesen Spezialfall zurückgeführt werden: Liegen  $P$  und  $Q$  auf verschiedenen



Seiten des Durchmessers durch  $S$ , dessen anderer Endpunkt  $T$  sei, so erfüllen auch die Punkte  $S, P, T, M$  sowie die Punkte  $S, Q, T, M$  die Voraussetzung des Satzes, und in beiden Fällen sind wir in der Situation des bereits bewiesenen Spezialfalls. Addition der Ergebnisse für diese beiden Fälle liefert das Ergebnis für die Punkte  $S, P, Q, M$ .

Bleibt noch der Fall, daß  $P$  und  $A$  auf derselben Seite des Durchmessers  $\overline{ST}$  liegen. Auch in diesem Fall erfüllen



wieder sowohl die Punkte  $S, P, T, M$  als auch die Punkte  $S, Q, T, M$  die Voraussetzungen des Satzes, und beides Mal sind wir in der Situation des eingangs bewiesenen Spezialfalls. Dieses Mal führt die Subtraktion dieser beiden Ergebnisse zum gewünschten Resultat für die Ausgangssituation mit den Punkten  $S, P, Q, M$ .

Damit ist der Satz vollständig bewiesen. ■

### §3: Der Satz von Lagrange

Es ist nicht möglich, eine beliebige natürliche Zahl als Summe von höchstens drei Quadratzahlen zu schreiben; das kleinste Gegenbeispiel ist die Sieben. Wie EULER vermutete und LAGRANGE bewies, kommt man aber immer mit höchstens vier Quadratzahlen aus.

Einer der vielen Beweise dieses Satzes ist recht ähnlich zu dem des Zweiquadratesatzes aus §1; statt mit dem Ring  $\mathbb{Z}[i]$  der GAUSSschen Zahlen arbeiten wir aber mit dem Ring

$$\mathbb{Z} \oplus \mathbb{Z}i \oplus \mathbb{Z}j \oplus \mathbb{Z}k$$

der ganzen Quaternionen. Auch hier haben wir eine Normabbildung, und eine ganze Zahl  $n$  ist offensichtlich genau dann als Summe von vier Quadraten darstellbar, wenn sie Norm einer ganzen Quaternion ist. Wegen der Multiplikativität der Norm reicht es also wieder, wenn wir Primzahlen  $p$  betrachten.

Zur Vorbereitung zeigen wir zunächst

**Lemma:** Zu jeder Primzahl  $p$  gibt es ganze Zahlen  $x, y, z \in \mathbb{Z}$  und eine natürliche Zahl  $m < p$ , so daß gilt:  $mp = x^2 + y^2 + z^2$

*Beweis:* Für  $p = 2$  ist  $2 = 1^2 + 1^2 + 0^2$ ; sei also  $p$  ungerade.

Von den Zahlen  $a^2$  mit  $0 \leq a \leq \frac{1}{2}(p-1)$  sind keine zwei kongruent modulo  $p$ , denn  $a^2 - b^2 = (a+b)(a-b)$ , und falls  $0 \leq a, b < \frac{1}{2}(p-1)$  sind beide Faktoren kleiner als  $p$ . Damit gibt es auch in den Mengen

$$\mathcal{M}_1 = \{-a^2 \mid 0 \leq a \leq \frac{1}{2}(p-1)\}$$

und

$$\mathcal{M}_2 = \{1 + a^2 \mid 0 \leq a \leq \frac{1}{2}(p-1)\}$$

keine zwei Elemente, die modulo  $p$  kongruent sind. Da die beiden Mengen disjunkt sind und jede davon  $\frac{1}{2}(p+1)$  Elemente hat, enthält ihre Vereinigung  $p+1$  Elemente; hier muß es also mindestens zwei Elemente geben, die modulo  $p$  kongruent sind. Es gibt also Zahlen  $x, y \in \mathbb{Z}$  mit  $-x^2 \equiv 1 + y^2 \pmod{p}$ , d.h.  $x^2 + y^2 + 1^2 = mp$  ist durch  $p$  teilbar. Da  $x, y \leq \frac{1}{2}(p-1)$ , ist dabei  $m < p$  und das Lemma ist bewiesen. ■

**Lemma:** Jede Primzahl  $p$  läßt sich als Summe von höchstens vier Quadraten schreiben.

*Beweis:* Für  $p = 2$  wissen wir das; sei also  $p$  wieder ungerade. Nach dem vorigen Lemma gibt es eine natürliche Zahl  $m < p$  derart, daß  $mp$  als Summe von sogar höchstens drei Quadraten darstellbar ist;  $\ell$  sei die kleinste natürliche Zahl, für die  $\ell p$  als Summe von höchstens vier Quadraten darstellbar ist. Natürlich ist dann auch  $\ell < p$ .

Wäre  $\ell$  eine gerade Zahl, so wäre auch die Summe der vier Quadrate gerade, und dazu gibt es drei Möglichkeiten: Entweder alle Summanden sind gerade, oder alle sind ungerade, oder zwei davon sind gerade, der Rest ungerade. Im letzteren Fall wollen wir die vier Zahlen  $w, x, y, z$  so anordnen, daß  $w$  und  $x$  gerade sind,  $y$  und  $z$  dagegen ungerade. Dann sind in allen drei Fällen  $w \pm x$  und  $y \pm z$  gerade, und

$$\left(\frac{w+x}{2}\right)^2 + \left(\frac{w-x}{2}\right)^2 + \left(\frac{y+z}{2}\right)^2 + \left(\frac{y-z}{2}\right)^2 = \frac{w^2 + x^2 + y^2 + z^2}{2} = \frac{\ell}{2}p,$$

im Widerspruch zur Minimalität von  $\ell$ . Also ist  $\ell$  ungerade, und falls das Lemma falsch wäre, müßte  $\ell \geq 3$  sein.

Wir betrachten die modulo  $\ell$  zu  $w, x, y, z$  kongruenten ganzen Zahlen  $W, X, Y, Z$  vom Betrag kleiner  $\ell/2$ . Wie schon beim Zwei-Quadrate-Satz können diese nicht allesamt verschwinden, denn sonst wären  $w, x, y, z$  durch  $\ell$  teilbar, also ihre Quadratsumme  $\ell p$  durch  $\ell^2$ , was wegen  $\ell < p$  für eine Primzahl  $p$  nicht möglich ist.

Somit ist  $0 < W^2 + X^2 + Y^2 + Z^2 < 4 \cdot \left(\frac{\ell}{2}\right)^2 = \ell^2$ . Andererseits ist aber

$$W^2 + X^2 + Y^2 + Z^2 \equiv w^2 + x^2 + y^2 + z^2 \equiv 0 \pmod{\ell};$$

also ist

$$W^2 + X^2 + Y^2 + Z^2 = \ell m \quad \text{mit} \quad 1 \leq m < \ell.$$

Damit haben die Quaternionen

$$q = w + \mathbf{i}x + \mathbf{j}y + \mathbf{k}z \quad \text{und} \quad Q = W + \mathbf{i}X + \mathbf{j}Y + \mathbf{k}Z$$

die Normen  $N(q) = \ell p$  und  $N(Q) = \ell m$ , ihr Produkt hat also die Norm  $\ell^2 mp$ . Zumindest von der Norm her spricht also nichts dagegen, daß dieses Produkt durch  $\ell$  teilbar sein könnte.

Tatsächlich ist  $q\bar{Q}$  durch  $\ell$  teilbar, und das sieht man am schnellsten durch brutales Nachrechnen: In

$$\begin{aligned} q\bar{Q} = & (wW + xX + yY + zZ) + (-wX + xW - yZ + zY)\mathbf{i} \\ & + (-wY + yW - zX + xZ)\mathbf{j} + (-wZ + zW - xY + yX)\mathbf{k} \end{aligned}$$

sind alle vier Klammern durch  $\ell$  teilbar, denn modulo  $\ell$  sind alle Großbuchstaben gleich den entsprechenden Kleinbuchstaben, so daß die Koeffizienten von  $\mathbf{i}, \mathbf{j}, \mathbf{k}$  trivialerweise modulo  $\ell$  verschwinden, und für den Realteil haben wir

$$wW + xX + yY + zZ \equiv w^2 + x^2 + y^2 + z^2 = \ell p \equiv 0 \pmod{\ell}.$$

Somit ist

$$\frac{q\bar{Q}}{\ell} = A + B\mathbf{i} + C\mathbf{j} + D\mathbf{k}$$

eine Quaternion mit ganzzahligen Koeffizienten, und

$$A^2 + B^2 + C^2 + D^2 = N\left(\frac{q\bar{Q}}{\ell}\right) = \frac{N(q\bar{Q})}{\ell^2} = \frac{N(q)N(Q)}{\ell^2} = mp.$$

Dies widerspricht aber der Minimalität von  $\ell$ .

Somit muß  $\ell = 1$  sein, und der Satz ist bewiesen. ■

**Satz (LAGRANGE):** Jede natürliche Zahl läßt sich als Summe von höchstens vier Quadraten schreiben.

*Beweis:* Wie wir in Kapitel 6, §7 gesehen haben, läßt sich eine Zahl  $n$  genau dann als Summe von höchstens vier Quadraten schreiben, wenn sie Norm einer ganzen Quaternion ist. Da wir gerade gesehen haben, daß sich jede Primzahl als Summe von höchstens vier Quadraten schreiben läßt (und die Eins natürlich auch), folgt die Behauptung aus der Multiplikativität der Norm. ■

# Kapitel 8

## Quadratische Reste

### § 1: Das Legendre-Symbol

**Definition:** Für eine Primzahl  $p$  und eine nicht durch  $p$  teilbare natürliche Zahl  $a$  ist das LEGENDRE-Symbol

$$\left(\frac{a}{p}\right) = \begin{cases} +1 & \text{falls es ein } x \in \mathbb{N} \text{ gibt mit } x^2 \equiv a \pmod{p} \\ -1 & \text{sonst} \end{cases}$$

Im ersten Fall bezeichnen wir  $a$  als *quadratischen Rest* modulo  $p$ , andernfalls als quadratischen Nichtrest. Für eine durch  $p$  teilbare Zahl  $a$  setzen wir  $\left(\frac{a}{p}\right) = 0$ .

Sind  $a, b$  zwei modulo  $p$  kongruente natürliche Zahlen, so ist offensichtlich  $\left(\frac{a}{p}\right) = \left(\frac{b}{p}\right)$ . Wir haben daher auch für  $a \in \mathbb{F}_p^\times$  ein wohldefiniertes LEGENDRE-Symbol  $\left(\frac{a}{p}\right)$ , das durch die Vorschrift  $\left(\frac{0}{p}\right) = 0$  auf ganz  $\mathbb{F}_p$  fortgesetzt wird.



ADRIEN-MARIE LEGENDRE (1752–1833) wurde in Toulouse oder Paris geboren; jedenfalls ging er in Paris zur Schule und studierte Mathematik und Physik am dortigen Collège Mazarin. Ab 1775 lehrte er an der Ecole Militaire und gewann einen Preis der Berliner Akademie für eine Arbeit über die Bahn von Kanonenkugeln. Andere Arbeiten befaßten sich mit der Anziehung von Ellipsoiden und der Himmelsmechanik. Ab etwa 1785 publizierte er auch Arbeiten über Zahlentheorie, in denen er z.B. das quadratische Reziprozitätsgesetz bewies sowie die Irrationalität von  $\pi$  und  $\pi^2$ .

**Lemma:** Das LEGENDRE-Symbol definiert einen Gruppenhomomorphismus

$$\left(\frac{\cdot}{p}\right) : \begin{cases} \mathbb{F}_p^\times \rightarrow \{+1, -1\} \\ a \mapsto \left(\frac{a}{p}\right) \end{cases} .$$

Für  $p = 2$  ist dies der triviale Homomorphismus, für ungerade  $p$  ist er surjektiv. Insbesondere gibt es dann jeweils  $\frac{p-1}{2}$  quadratische Reste und Nichtreste.

*Beweis:* Für  $p = 2$  ist  $\mathbb{F}_2^\times = \{1\}$ , und  $1 = 1^2$  ist ein quadratischer Rest.

Sei nun  $p$  ungerade. Der Homomorphismus

$$\begin{cases} \mathbb{F}_p^\times \rightarrow \mathbb{F}_p^\times \\ x \mapsto x^2 \end{cases}$$

hat den Kern  $\{+1, -1\}$ , also besteht das Bild aus  $\frac{p-1}{2}$  Elementen, den quadratischen Resten.

Trivialerweise ist das Produkt zweier quadratischer Reste wieder ein quadratischer Rest. Ist  $a = x^2$  ein quadratischer Rest und  $b$  ein Nichtrest, so ist auch  $ab$  ein quadratischer Nichtrest, denn wäre  $ab = y^2$ , wäre  $b = (yx^{-1})^2$  ein quadratischer Rest. Da Multiplikation mit  $b$  injektiv ist, folgt, daß sich jeder quadratische Nichtrest in der Form  $bc$  darstellen läßt, wobei  $c$  ein quadratischer Rest ist. Damit folgt, daß das Produkt zweier quadratischer Nichtreste ein quadratischer Rest ist, denn  $bc \cdot bd = b^2 cd$ , wobei  $c$  und  $d$  Quadrate in  $\mathbb{F}_p^\times$  sind. ■

**Lemma (EULER):** Ist  $p$  eine ungerade Primzahl und kein Teiler von  $a$ , so ist  $\left(\frac{a}{p}\right) \equiv a^{\frac{p-1}{2}} \pmod{p}$ .

*Beweis:*  $g$  sei ein erzeugendes Element von  $\mathbb{F}_p^\times$ . Dann ist offensichtlich jede Potenz  $g^r$  mit geradem  $r$  ein quadratischer Rest, und da es genau  $\frac{p-1}{2}$  verschiedene solcher Potenzen gibt, sind das auch *alle* quadratischen Reste. Somit ist  $g^r$  genau dann ein quadratischer Rest, wenn  $r$  gerade ist.



Da  $g$  ein erzeugendes Element ist, kann  $g^{(p-1)/2}$  nicht gleich eins sein; da nach dem kleinen Satz von FERMAT aber sein Quadrat  $g^{p-1} = 1$  ist, folgt  $g^{(p-1)/2} = -1$ . Für  $a = g^r$  ist somit

$$a^{\frac{p-1}{2}} = (g^r)^{\frac{p-1}{2}} = \left(g^{\frac{p-1}{2}}\right)^r = (-1)^r$$

genau dann gleich eins, wenn  $a$  ein quadratischer Rest ist, und  $-1$  sonst. ■

**Korollar:** Für ungerades  $p$  ist

$$\left(\frac{-1}{p}\right) = (-1)^{\frac{p-1}{2}} = \begin{cases} +1 & \text{falls } p \equiv 1 \pmod{4} \\ -1 & \text{falls } p \equiv 3 \pmod{4} \end{cases} .$$

## §2: Das quadratische Reziprozitätsgesetz

**Quadratisches Reziprozitätsgesetz:** Für zwei verschiedene ungerade Primzahlen  $p, q$  ist

$$\left(\frac{p}{q}\right) \left(\frac{q}{p}\right) = (-1)^{\frac{(p-1)(q-1)}{4}} \quad \text{und} \quad \left(\frac{2}{p}\right) = (-1)^{\frac{p^2-1}{8}} .$$

Zum *Beweis* betrachten wir ein zum Nullpunkt symmetrisches Vertreter-system von  $\mathbb{F}_p^\times$  in  $\mathbb{Z}$ , nämlich

$$R = \{-h, \dots, -1, 1, \dots, h\} \quad \text{mit} \quad h = \frac{p-1}{2} .$$

Weiter sei  $S = \{q, 2q, \dots, hq\}$ . Da  $p$  und  $q$  teilerfremd sind, haben zwei verschiedene Elemente von  $S$  verschiedene Restklassen modulo  $p$ .

*1. Schritt (GAUSS):*  $q$  sei eine beliebige Primzahl und  $p \neq q$  eine ungerade Primzahl. Dann ist  $\left(\frac{q}{p}\right) = (-1)^m$ , wobei  $m$  die Anzahl jener Elemente von  $S$  bezeichnet, die modulo  $p$  kongruent sind zu einem negativen Element von  $R$ .

*Beweis:*  $a_1, \dots, a_m$  seien die negativen Elemente von  $R$ , die zu Elementen aus  $S$  kongruent sind,  $b_1, \dots, b_n$  die positiven. Dann ist

$$a_1 \cdots a_m \cdot b_1 \cdots b_n \equiv \prod_{i=1}^h (iq) = h!q^h \pmod{p} .$$

Natürlich sind  $a_i$  und  $a_j$  für  $i \neq j$  zwei verschiedene Zahlen, genauso auch  $b_i$  und  $b_j$ . Außerdem kann auch nie  $|a_i| = |b_j|$  sein, denn sonst wäre einerseits  $a_i + b_j = 0$ , andererseits gäbe es aber Zahlen  $1 \leq k, \ell \leq h$ , so daß  $a_i \equiv kq$  und  $b_j \equiv \ell q \pmod{p}$ . Also wäre  $(k + \ell)q$  durch  $p$  teilbar, was nicht möglich ist, denn  $k + \ell \leq 2h = p - 1$ . Damit sind die Beträge der  $a_i$  und der  $b_j$  genau die Zahlen von 1 bis  $h$ , d.h.

$$a_1 \cdots a_m b_1 \cdots b_n = (-1)^m h!.$$

Vergleich mit der obigen Kongruenz zeigt, daß dann  $q^h \equiv (-1)^m \pmod{p}$  ist, also nach dem vorigen Lemma  $\left(\frac{q}{p}\right) = (-1)^m$ . ■

2. Schritt (GAUSS): Für zwei ungerade Primzahlen  $p \neq q$  ist

$$\left(\frac{q}{p}\right) = (-1)^M \quad \text{mit} \quad M = \sum_{i=1}^h \left[\frac{iq}{p}\right], \quad \text{und} \quad \left(\frac{2}{p}\right) = (-1)^{\frac{p^2-1}{8}}.$$

Im *Beweis* sei zunächst auch noch der Fall  $q = 2$  zugelassen. Für  $i \leq h$  sei  $r_i = iq - p \cdot \left[\frac{iq}{p}\right]$ ; dann ist  $0 \leq r_i < p$  und  $iq = p \cdot \left[\frac{iq}{p}\right] + r_i$ . Falls  $iq$  modulo  $p$  kongruent ist zu einem negativen Element  $a_j \in R$ , ist also  $r_i = p + a_j$ ; falls  $r_i \equiv b_j > 0$  ist dagegen  $r_i = b_j$ . Somit ist

$$\sum_{i=1}^h iq = p \sum_{i=1}^h \left[\frac{iq}{p}\right] + \sum_{i=1}^m (a_i + p) + \sum_{i=1}^n b_i = pM + mp + \sum_{i=1}^m a_i + \sum_{i=1}^n b_i.$$

Andererseits ist

$$\sum_{i=1}^h iq = \frac{h(h+1)}{2} \cdot q = \frac{1}{2} \cdot \frac{p-1}{2} \cdot \frac{p+1}{2} \cdot q = \frac{p^2-1}{8} \cdot q.$$

Außerdem wissen wir aus dem ersten Schritt, daß

$$\{-a_1, \dots, -a_m, b_1, \dots, b_n\} = \{1, \dots, h\}$$

ist, d.h.

$$-\sum_{i=1}^m a_i + \sum_{i=1}^n b_i = \sum_{i=1}^h i = \frac{h(h+1)}{2} = \frac{p^2-1}{8}$$

und damit ist  $\sum_{i=1}^n b_i = \frac{p^2-1}{8} + \sum_{i=1}^m a_i$ . Setzen wir das alles in die obige Formel ein, erhalten wir die Beziehung

$$\frac{p^2-1}{8} \cdot q = (M+m)p + \frac{p^2-1}{8} + 2 \sum_{i=1}^m a_i$$

oder

$$\frac{p^2-1}{8} \cdot (q-1) = (M+m)p + 2 \sum_{i=1}^m a_i.$$

Im Falle einer ungeraden Primzahl  $q$  steht rechts eine gerade Zahl; damit muß auch  $M+m$  gerade sein, d.h.  $(-1)^M = (-1)^m$ , und die Behauptung folgt aus dem ersten Schritt.

Für  $q=2$  ist  $M=0$ , da  $\left[\frac{2i}{p}\right]$  für alle  $i \leq h$  verschwindet. Modulo zwei wird die obige Beziehung daher zu

$$\frac{p^2-1}{8} \equiv mp \equiv m \pmod{2},$$

so daß die Behauptung auch hier aus dem ersten Schritt folgt. ■

3. Schritt (EISENSTEIN):  $p$  und  $q$  seien ungerade Primzahlen,

$$h = \frac{p-1}{2}, \quad k = \frac{q-1}{2}, \quad M = \sum_{i=1}^h \left[\frac{iq}{p}\right] \quad \text{und} \quad N = \sum_{i=1}^k \left[\frac{ip}{q}\right].$$

Dann ist  $M+N = hk$ .

*Beweis:* Im Innern des Rechtecks mit Ecken  $(0,0)$ ,  $(\frac{p}{2}, 0)$ ,  $(0, \frac{q}{2})$  und  $(\frac{p}{2}, \frac{q}{2})$  liegen  $hk$  Gitterpunkte, nämlich die Punkte  $(i, j)$  mit  $1 \leq i \leq h$  und  $1 \leq j \leq k$ .

Die Diagonale des Rechtecks liegt auf der Geraden  $y = \frac{q}{p}x$  und enthält keine Gitterpunkte. Unterhalb der Diagonalen liegen  $\left[\frac{iq}{p}\right]$  Punkte mit Abszisse  $i$ , insgesamt also  $M$  Punkte. Darüber liegen  $\left[\frac{ip}{q}\right]$  Punkte mit Ordinate  $i$ , insgesamt also  $N$  Punkte. Somit ist  $hk = M+N$ . ■

Zum *Beweis* des quadratischen Reziprozitätsgesetzes müssen wir nun nur noch alles kombinieren: Nach dem zweiten und dem dritten Schritt ist

$$\left(\frac{q}{p}\right) \left(\frac{p}{q}\right) = (-1)^M \cdot (-1)^N = (-1)^{M+N} = (-1)^{hk} = (-1)^{\frac{p-1}{2} \frac{q-1}{2}}.$$



CARL FRIEDRICH GAUSS (1777–1855) leistete wesentliche Beiträge zur Zahlentheorie, zur nichteuklidischen Geometrie, zur Funktionentheorie, zur Differentialgeometrie und Kartographie, zur Fehlerrechnung und Statistik, zur Astronomie und Geophysik *usw.* Als Direktor der Göttinger Sternwarte baute er zusammen mit dem Physiker Weber den ersten Telegraphen. Er leitete die erste Vermessung und Kartierung des Königreichs Hannover und zeitweise auch den Witwenfond der Universität Göttingen; seine hierbei gewonnene Erfahrung benutzte er für erfolgreiche Spekulationen mit Aktien. Seine 1801 veröffentlichten *Disquisitiones arithmeticae* sind auch noch heute fundamental für die Zahlentheorie.



FERDINAND GOTTHOLD MAX EISENSTEIN (1823–1852), genannt Gotthold, wurde in Berlin geboren. Als einziges seiner sechs Geschwister starb er nicht bereits während der Kindheit an Meningitis. Im Alter von 17 Jahren, noch als Schüler, besuchte er Mathematikvorlesungen der Universität, unter anderem bei DIRICHLET. Ab 1842 las er die *Disquisitiones arithmeticae* von GAUSS, den er 1844 in Göttingen besuchte. Trotz zahlreicher wichtiger Arbeiten erhielt er nie eine gut bezahlte Position und überlebte vor allem dank der Unterstützung durch ALEXANDER VON HUMBOLDT. 1847 habilitierte er sich in Berlin und hatte dort unter anderem RIEMANN als Studenten. Er starb 29-jährig an Tuberkulose.

**Bemerkung:** Die rechten Seiten der Gleichungen im quadratischen Reziprozitätsgesetz lassen sich auch durch Kongruenzbedingungen ausdrücken:  $(p-1)/2$  ist genau dann gerade, wenn  $p \equiv 1 \pmod{4}$ , entsprechend für  $q$ . Somit ist

$$\left(\frac{p}{q}\right) \left(\frac{q}{p}\right) = \begin{cases} +1 & \text{falls } p \equiv 1 \pmod{4} \text{ oder } q \equiv 1 \pmod{4} \\ -1 & \text{falls } p \equiv q \equiv 3 \pmod{4} \end{cases}.$$

Ist  $p = 8r + k$ , so ist  $p^2 = 64r^2 + 16r + k^2 \equiv k^2 \pmod{16}$ , also ist  $p^2 - 1 \equiv k^2 - 1 \pmod{16}$ . Für  $k = \pm 1$  ist dies null, für  $k = \pm 3$  acht. Somit ist

$$\left(\frac{2}{p}\right) = (-1)^{\frac{p^2-1}{8}} = \begin{cases} +1 & \text{falls } p \equiv \pm 1 \pmod{8} \\ -1 & \text{falls } p \equiv \pm 3 \pmod{8} \end{cases}.$$

Das quadratische Reziprozitätsgesetz läßt sich gelegentlich dazu verwenden, um ein LEGENDRE-Symbol einfach zu berechnen. Wenn wir beispielsweise entscheiden wollen, ob sieben ein quadratischer Rest modulo 17 ist, sagt es uns (da  $17 \equiv 1 \pmod{4}$ ), daß  $\left(\frac{7}{17}\right) = \left(\frac{17}{7}\right)$  ist. Letzteres ist gleich  $\left(\frac{3}{7}\right)$ , da  $17 \equiv 3 \pmod{7}$ . Hier haben wir zwei Primzahlen, die beide kongruent drei modulo vier sind, also ist  $\left(\frac{3}{7}\right) = -\left(\frac{7}{3}\right) = -\left(\frac{1}{3}\right) = -1$ , denn die Eins ist natürlich modulo jeder Primzahl ein quadratischer Rest. Also ist sieben modulo 17 ein quadratischer Nichtrest.

Genauso können wir auch leicht feststellen, ob 13 quadratischer Rest modulo 1 000 003 ist: Da  $13 \equiv 1 \pmod{4}$ , ist  $\left(\frac{13}{1\,000\,003}\right) = \left(\frac{1\,000\,003}{13}\right)$ . Da  $1\,000\,003 \equiv 4 \pmod{13}$ , ist dies gleich  $\left(\frac{4}{13}\right)$ , und das ist natürlich eins, da  $4 = 2^2$  modulo jeder Primzahl ein Quadrat ist. Somit ist auch 13 ein Quadrat modulo 1 000 003.

Das Problem bei dieser Vorgehensweise besteht darin, daß wir normalerweise nicht soviel Glück haben wie hier und als Reduktionen stets Primzahlen erhalten. Wir sollten daher ein quadratisches Reziprozitätsgesetz haben, das auch funktioniert, wenn die beteiligten Zahlen nicht prim sind.

### §3: Das Jacobi-Symbol

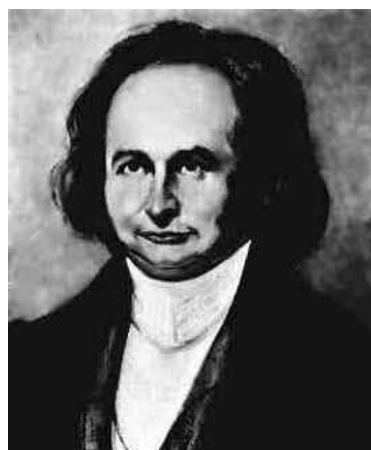
Wie wir in §1 gesehen haben, definiert das LEGENDRE-Symbol in Bezug auf seinen „Zähler“ einen Homomorphismus; wir können versuchen, es zu erweitern, indem wir dasselbe auch für den „Nenner“ postulieren:

**Definition:** Ist  $n = \prod_{i=1}^r p_i^{e_i}$  eine ungerade Zahl und  $m$  eine zu  $n$  teiler-

fremde Zahl, so ist das JACOBI-Symbol definiert als

$$\left(\frac{m}{n}\right) = \prod_{i=1}^r \left(\frac{m}{p_i}\right)^{e_i}.$$

Falls  $m$  und  $n$  nicht teilerfremd sind, setzen wir  $\left(\frac{m}{n}\right) = 0$ .



CARL GUSTAV JACOB JACOBI (1804–1851) wurde in Potsdam als Sohn eines jüdischen Bankiers geboren und erhielt den Vornamen Jacques Simon. Im Alter von zwölf Jahren bestand er sein Abitur, mußte aber noch vier Jahre in der Abschlußklasse des Gymnasiums bleiben, da die Berliner Universität nur Studenten mit mindestens 16 Jahren aufnahm. 1824 beendete er seine Studien mit dem Staatsexamen für Mathematik, Griechisch und Latein und wurde Lehrer. Außerdem promovierte er 1825 und begann mit seiner Habilitation. Etwa gleichzeitig konvertierte er zum Christentum, so daß er ab 1825 an der Universität Berlin und ab 1826 in

Königsberg lehren konnte. 1832 wurde er dort Professor. Zehn Jahre später mußte er aus gesundheitlichen Gründen das rauhe Klima Königsbergs verlassen und lebte zunächst in Italien, danach für den Rest seines Lebens in Berlin. Er ist vor allem berühmt durch seine Arbeiten zur Zahlentheorie und über elliptische Integrale.

Für eine Primzahl  $n$  und ein nicht dadurch teilbares  $m$  stimmt das JACOBI-Symbol natürlich mit dem LEGENDRE-Symbol überein, und man kann sich fragen, ob man hier wirklich einen neuen Namen braucht. Dieser ist gerechtfertigt, weil es einen ganz wesentlichen Unterschied zwischen den beiden Symbolen gibt: Beispielsweise ist

$$\left(\frac{2}{15}\right) = \left(\frac{2}{3}\right) \left(\frac{2}{5}\right) = (-1)^{\frac{3^2-1}{8}} \cdot (-1)^{\frac{5^2-1}{8}} = (-1) \cdot (-1) = 1,$$

aber zwei ist offensichtlich kein quadratischer Rest modulo 15: Sonst müßte es schließlich erst recht quadratischer Rest modulo drei und modulo fünf sein, aber die entsprechenden LEGENDRE-Symbole sind  $-1$ . In der Tat gibt es modulo 15 nur vier quadratische Reste: 1, 4, 6 und 10.

Das JACOBI-Symbol gibt daher keine Auskunft darüber, ob eine Zahl quadratischer Rest ist oder nicht; lediglich wenn es gleich  $-1$  ist, können wir sicher sein, daß wir es mit einem quadratischen Nichtrest zu tun haben, denn dann muß ja auch schon für mindestens einen Primteiler

des „Nenners“ das LEGENDRE-Symbol gleich  $-1$  sein, während ein quadratischer Rest modulo einer Zahl  $n$  erst recht quadratischer Rest modulo eines jeden Teilers von  $n$  sein muß.

Die Nützlichkeit des JACOBI-Symbols kommt in erster Linie daher, daß auch dafür das quadratische Reziprozitätsgesetz gilt und es somit zur Berechnung von LEGENDRE-Symbolen verwendet werden kann:

**Satz:** Für zwei ungerade Zahlen  $m, n$  mit  $\text{ggT}(m, n) = 1$  ist

$$\left(\frac{n}{m}\right) \left(\frac{m}{n}\right) = (-1)^{\frac{n-1}{2} \frac{m-1}{2}} \quad \text{und} \quad \left(\frac{2}{m}\right) = (-1)^{\frac{m^2-1}{8}}.$$

*Beweis:* Sei  $n = \prod_{i=1}^r p_i^{e_i}$  und  $m = \prod_{j=1}^s q_j^{f_j}$ . Nach Definition des JACOBI-Symbols und weil das LEGENDRE-Symbol bei festgehaltenem „Nenner“ einen Homomorphismus definiert, ist dann

$$\left(\frac{n}{m}\right) = \prod_{j=1}^s \left(\frac{n}{q_j}\right)^{f_j} = \prod_{j=1}^s \prod_{i=1}^r \left(\frac{p_i}{q_j}\right)^{e_i f_j} \quad \text{und} \quad \left(\frac{m}{n}\right) = \prod_{j=1}^s \prod_{i=1}^r \left(\frac{q_j}{p_i}\right)^{e_i f_j}.$$

Nach dem quadratischen Reziprozitätsgesetz aus §2 ist daher

$$\begin{aligned} \left(\frac{n}{m}\right) \left(\frac{m}{n}\right) &= \prod_{i=1}^r \prod_{j=1}^s \left((-1)^{\frac{p_i-1}{2} \frac{q_j-1}{2}}\right)^{e_i f_j} = (-1)^{\sum_{i=1}^r \sum_{j=1}^s \frac{p_i-1}{2} \frac{q_j-1}{2} e_i f_j} \\ &= (-1)^{\left(\sum_{i=1}^r \frac{p_i-1}{2} e_i\right) \left(\sum_{j=1}^s \frac{q_j-1}{2} f_j\right)} = \left((-1)^{\sum_{i=1}^r \frac{p_i-1}{2} e_i}\right)^{\sum_{j=1}^s \frac{q_j-1}{2} f_j}. \end{aligned}$$

Dies ist genau dann gleich  $+1$ , wenn mindestens einer der beiden Exponenten gerade ist; andernfalls ist es gleich  $-1$ .

Im Produkt

$$(-1)^{\sum_{i=1}^r \frac{p_i-1}{2} e_i} = \prod_{i=1}^r (-1)^{\frac{p_i-1}{2} e_i}$$

können wir alle Faktoren weglassen, für die  $e_i$  gerade ist oder aber  $p_i \equiv 1 \pmod{4}$ . Das Produkt ist also gleich  $(-1)^N$  mit

$N =$  Anzahl der Indizes  $i$  mit  $p_i \equiv 3 \pmod{4}$  und  $e_i$  ungerade.

Die Faktoren  $p_i^{e_i}$  sind genau dann kongruent eins modulo vier, wenn  $p_i \equiv 1 \pmod{4}$  oder  $e_i$  gerade ist, denn  $3^2 \equiv 1 \pmod{4}$ . Andernfalls ist  $p_i^{e_i} \equiv 3 \equiv -1 \pmod{4}$ . Somit ist auch  $n \equiv (-1)^N \pmod{4}$ , also

$$(-1)^{\sum_{i=1}^r \frac{p_i-1}{2} e_i} = (-1)^N = (-1)^{\frac{n-1}{2}}.$$

Ist dies gleich +1, so ist die rechte Seite der Gleichung für  $\left(\frac{n}{m}\right) \left(\frac{m}{n}\right)$  ebenfalls +1, andernfalls zeigt das gleiche Argument für  $m$ , daß sie gleich  $(-1)^{(m-1)/2}$  ist. In jedem Fall erhalten wir daher die gewünschte Formel

$$\left(\frac{m}{n}\right) \left(\frac{n}{m}\right) = (-1)^{\frac{n-1}{2} \frac{m-1}{2}}.$$

Genauso folgt auch, daß  $\left(\frac{2}{m}\right) = (-1)^{(m^2-1)/8}$  ist, denn dies ist +1 für  $m \equiv \pm 1 \pmod{8}$  und -1 für  $m \equiv \pm 3 \pmod{8}$ . Das Produkt zweier Primzahlen kongruent  $\pm 1$  modulo acht ist wieder kongruent  $\pm 1$ , genauso das zweier Primzahlen kongruent  $\pm 3$  modulo acht. Damit führt dieselbe Argumentation wie oben zum Ziel. ■

Als Anwendung können wir uns überlegen, modulo welcher Primzahlen eine vorgegebene Zahl  $a$  quadratischer Rest ist. Modulo seiner Primteiler verschwindet  $a$  und ist somit ein Quadrat. Sei also  $p$  kein Teiler von  $a$ .

Für  $a = 2$  haben wir gesehen, daß  $\left(\frac{2}{p}\right)$  nur von der Kongruenzklasse  $p \pmod{8}$  abhängt; wegen der Multiplikativität des JACOBI-Symbols reicht es also, wenn wir ungerade  $a$  betrachten. Nach dem gerade bewiesenen Gesetz ist dann

$$\left(\frac{a}{p}\right) = (-1)^{\frac{a-1}{2} \frac{p-1}{2}} \left(\frac{p}{a}\right).$$

Für festes  $a$  ist  $(a-1)/2$  ein konstanter Wert,  $(p-1)/2$  hängt nur ab von  $p \pmod{4}$ , und  $\left(\frac{p}{a}\right)$  hängt ab von  $p \pmod{a}$ . Insgesamt hängt es also nur ab von  $p \pmod{4a}$ , ob  $a$  ein quadratischer Rest oder Nichtrest modulo  $p$  ist.

Betrachten wir als Beispiel den Fall  $a = 3$ . Hier ist  $(a-1)/2 = 1$ , also

$$(-1)^{\frac{(a-1)}{2} \frac{p-1}{2}} = (-1)^{\frac{p-1}{2}} = \begin{cases} +1 & \text{falls } p \equiv 1 \pmod{4} \\ -1 & \text{falls } p \equiv 3 \pmod{4} \end{cases},$$



und

$$\left(\frac{p}{3}\right) = \begin{cases} +1 & \text{falls } p \equiv 1 \pmod{3} \\ -1 & \text{falls } p \equiv 2 \pmod{3} \end{cases}.$$

Somit ist für eine Primzahl  $p > 3$

$$\left(\frac{3}{p}\right) = \begin{cases} +1 & \text{falls } p \pmod{12} \in \{1, 11\} \\ -1 & \text{falls } p \pmod{12} \in \{5, 7\} \end{cases}.$$

Für  $a = 5$  ist  $(a - 1)/2 = 2$  gerade, also

$$\left(\frac{5}{p}\right) = \left(\frac{p}{5}\right) = \begin{cases} +1 & \text{falls } p \pmod{5} \in \{1, 4\} \\ -1 & \text{falls } p \pmod{5} \in \{2, 3\} \end{cases},$$

## §4: Anwendungen quadratischer Reste

Zum Abschluß dieses Kapitels sollen kurz noch einige Anwendungen quadratischer Reste vorgestellt werden:

### a) Münzwurf per Telephon

A und B können sich nicht einigen, wer von ihnen eine dringend notwendige aber unangenehme Arbeit übernehmen soll. Also werfen sie eine Münze. Vorher entscheidet sich etwa A für „Wappen“, B für „Zahl“, dann wirft A die Münze in die Luft. Wenn sie mit Wappen nach oben auf den Boden fällt, hat er gewonnen, andernfalls B.

Stellen wir uns nun aber vor, A und B stehen nicht nebeneinander, sondern befinden sich an verschiedenen Orten und diskutieren per Telephon, wer was machen soll. Auch hier könnte A wieder eine Münze werfen, allerdings sieht jetzt nur A, wie sie zu Boden fällt; wenn er gewinnt, muß B sehr viel Vertrauen in ihn haben, um das zu glauben.

Mit Hilfe von quadratischen Resten läßt sich der Münzwurf so simulieren, daß *beide* den Ausgang überprüfen können und jeder mit der gleichen Wahrscheinlichkeit gewinnt.

Dazu wählt sich A zwei Primzahlen  $p, q \equiv 3 \pmod{4}$ , die so groß sind, daß B das Produkt  $N = pq$  nicht mit einem Aufwand von nur wenigen Minuten faktorisieren kann. ( $p$  und  $q$  können also deutlich kleiner sein als

bei RSA, wo man mit Gegnern rechnen muß, die monatelang rechnen.) Dieses  $N$  schickt er an B.

B wählt sich nun eine zufällige Zahl  $x$  zwischen eins und  $N$  und schickt deren Quadrat  $y = x^2 \bmod N$  an A.

A kennt die Faktorisierung von  $N$  und kann die Gleichungen

$$z^2 \equiv y \pmod{p} \quad \text{und} \quad z^2 \equiv y \pmod{q}$$

lösen: Wie wir von Aufgabe vier des dritten Übungsblatts wissen, sind  $\pm y^{(p+1)/4}$  und  $\pm y^{(q+1)/4}$  die Lösungen. Nach dem chinesischen Restesatz kann er sich somit vier Zahlen zwischen null und  $N - 1$  konstruieren, die allesamt die Kongruenz  $u^2 \equiv y \pmod{N}$  erfüllen. Er entscheidet sich zufällig für eine dieser vier Möglichkeiten (dies entspricht dem Münzwurf) und schickt das entsprechende  $u$  an B.

B kennt nun zwei Zahlen  $x$  und  $u$ , die beide das Quadrat  $y$  haben. Möglicherweise ist  $u = x$ ; in diesem Fall hat er keine neue Information bekommen, und er hat verloren. Das gleiche gilt im Fall  $u \equiv -x \pmod{N}$ , d.h.  $u = N - x$ .

Ist aber  $u \neq \pm x$ , was mit 50%-iger Wahrscheinlichkeit eintritt, hat B gewonnen und muß das nun gegenüber A beweisen. Da  $u^2 \equiv y \pmod{N}$  ist erst recht  $u^2 \equiv x \pmod{p}$  und  $u^2 \equiv x \pmod{q}$ . Da quadratische Gleichungen in einem Körper höchstens zwei Lösungen haben, ist daher  $u \equiv \pm x \pmod{p}$  und  $u \equiv \pm x \pmod{q}$ . Falls in beiden Gleichungen das gleiche Vorzeichen steht, ist  $u \equiv \pm x$ ; andernfalls ist  $x - u$  durch genau eine der beiden Primzahlen teilbar, und B kann diese als ggT von  $N$  und  $x - u$  ausrechnen. Damit hat er  $N$  faktorisiert und schickt als Beweis die Faktoren an A.

Wenn B sich nicht an die Regeln hält und ein  $y$  an A schickt, das kein Quadrat modulo  $N$  ist, merkt A dies bei der Berechnung der modularen Quadratwurzeln; falls A ein  $u$  schickt, dessen Quadrat von  $y$  verschieden ist, kann B dies leicht feststellen, denn wenn er verloren hat, muß  $u = x$  oder  $u = N - x$  sein. (Er kann natürlich auch  $u^2 \bmod N$  berechnen.)

## b) Akustik von Konzerthallen

Alte Konzerthallen waren zwangsläufig sehr hoch: Andernfalls wäre die Luft während eines längeren Konzerts bei voll besetztem Saal zu schnell verbraucht gewesen. Mit den Fortschritten der Lüftungstechnik verschwand diese Notwendigkeit; dafür sorgten steigende Bau- und Heizungskosten für immer niedrigere Säle. Auf die Luftqualität hatte das keinen nennenswerten Einfluß; die Akustik der Hallen allerdings wurde deutlich schlechter.

Der Grund dafür ist intuitiv recht klar und wurde auch durch Messungen und Hörerbefragungen in einer Reihe von Konzertsälen experimentell bestätigt: Die Hörer bevorzugen Schall, der von den Seitenwänden kommt und daher mit verschiedener Stärke bei den beiden Ohren eintrifft gegenüber Schall von oben, der beide Ohren mit gleicher Stärke erreicht und somit keinen räumlichen Eindruck hinterläßt.

Eine mögliche Abhilfe bestünde darin, die Decken aus absorbierendem Material zu bauen. Dem steht entgegen, daß in einem großen Konzertsaal aller Schall, der von der Bühne kommt, den Hörer auch wirklich erreichen sollte: Ansonsten müßte der Schall aus Lautsprechern kommen und man könnte sich das Konzert genauso gut daheim per Radio oder CD anhören.

Der Schall muß daher von der Decke reflektiert werden, darf die Ohren der Zuhörer aber nicht von oben erreichen. Er sollte daher beispielsweise möglichst diffus zu den Seitenwänden hin gestreut werden, so daß der größte Teil der Energie die Zuhörer über die Seitenwände erreicht.

Der Einfachheit halber wollen wir uns auf eindimensionale Wellen beschränken und damit auch nur diffuse Reflektion in einer Richtung betrachten, der Querrichtung des Konzertsaals.

Eine Welle hat eine räumliche wie auch zeitliche Periodizität. Zeitlich periodische Funktionen sind beispielsweise Sinus und Kosinus; wie die FOURIER-Analyse lehrt, läßt sich jede stückweise stetige zeitlich periodische Funktion (bis auf sogenannte Nullfunktionen) aus Sinus- und Kosinusfunktionen zusammensetzen, so daß es reicht, solche Funktionen zu betrachten.

Da der Umgang mit den Additionstheoremen für trigonometrische Funktionen recht umständlich ist, schreibt man Wellen allerdings meist komplex in der Form  $f(t) = Ae^{i\omega t}$  mit der Maßgabe, daß nur der Realteil dieser Funktion physikalische Realität beschreibt. Aufgrund der EULERSchen Formel  $e^{i\varphi} = \cos \varphi + i \sin \varphi$  lassen sich so, falls man für  $A$  beliebige komplexe Konstanten zuläßt, alle Funktionen der Art  $a \cos \omega t + b \sin \omega t$  als Realteile erhalten, und da beispielsweise

$$\cos(\alpha + \beta) = \Re e^{i(\alpha+\beta)} = \Re(e^{i\alpha} e^{i\beta}) = \cos \alpha \cos \beta - \sin \alpha \sin \beta$$

ist, lassen sich auf diese Weise auch die Additionstheoreme auf einfache Multiplikationen von Exponentialfunktionen zurückführen.

Auch die räumliche Periodizität läßt sich mit trigonometrischen oder – besser – Exponentialfunktionen ausdrücken; hier schreiben wir entsprechend  $g(x) = Be^{ikx}$ .

Um einen räumlich und zeitlich periodischen Vorgang zu beschreiben, kombinieren wir die beiden Ansätze und betrachten beispielsweise die Funktion

$$\psi(x, t) = Ae^{i(\omega t - kx)} = Ae^{ik(\frac{\omega}{k}t - x)}.$$

Wie man der zweiten Form ansieht, hängt  $\psi(x, t)$  nur ab von  $x - \frac{\omega}{k}t$ , was wir auch so interpretieren können, daß

$$v = \frac{\omega}{k} = \frac{\lambda}{T} = \frac{\lambda\omega}{2\pi}$$

die Ausbreitungsgeschwindigkeit der Welle ist; denn eine Änderung der Zeit um  $\Delta t$  hat denselben Effekt wie eine Änderung des Orts um  $v \cdot \Delta t$ .

Da Sinus und Kosinus die Periode  $2\pi$  haben, müssen wir für eine Schwingung der Frequenz  $\nu$  den Parameter  $\omega$  gleich  $2\pi\nu$  wählen, denn dann fallen  $1/\nu$  Perioden in das Intervall  $0 \leq t \leq 1$ . Aus diesem Grund wird  $\omega = 2\pi\nu$  als die *Kreisfrequenz* der Schwingung bezeichnet. In der räumlichen Dimension nimmt die Wellenlänge  $\lambda$  die Rolle der zeitlichen Periode ein; dementsprechend muß hier  $k = 2\pi/\lambda$  gesetzt werden. Diese Konstante wird als *Wellenzahl* bezeichnet.

Schallwellen breiten sich bei  $20^\circ \text{ C}$  in Luft mit einer Geschwindigkeit von etwa  $v = 343 \text{ m/s}$  aus; der hörbare Frequenzbereich beginnt bei

$\nu = 16$  Hz und kann bis zu etwa  $\nu = 20$  kHz gehen. Die Wellenlängen, mit denen wir es zu tun haben, variieren also zwischen etwa  $\lambda = 21,5$  m und  $\lambda = 1,75$  cm. Der Kammerton  $a'$  mit 440 Hz hat eine Wellenlänge von knapp 78 cm.

Bei einer Reflektion können wir nach HUYGENS annehmen, daß von jedem Punkt der reflektierenden Fläche eine neue Welle ausgeht; ihre Amplitude ist gleich der Amplitude der dort eintreffenden Welle mal einem Reflektionsfaktor  $\rho(x)$ , der im Idealfall gleich eins ist.



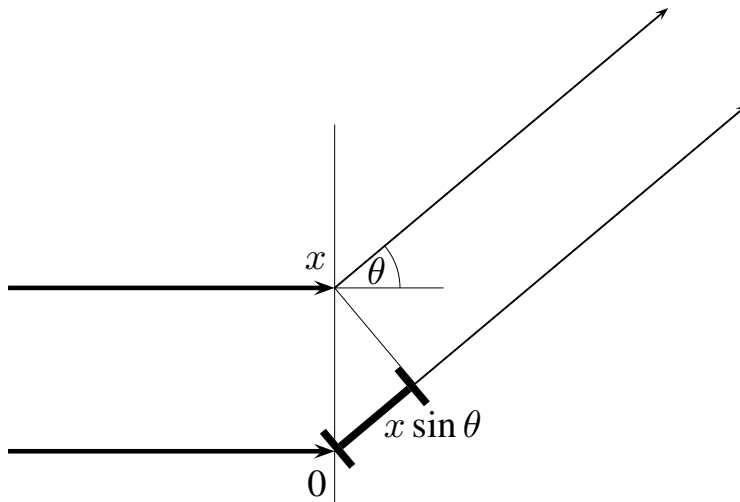
CHRISTIAAN HUYGENS (1629–1695) kam aus einer niederländischen Diplomatenfamilie. Dadurch und später auch durch seine Arbeit hatte er Kontakte zu führenden europäischen Wissenschaftlern wie DESCARTES und PASCAL. Nach seinem Studium der Mathematik und Jura arbeitete er teilweise auch selbst als Diplomat, interessierte sich aber bald vor allem für Astronomie und den Bau der dazu notwendigen Instrumente. Er entwickelte eine neue Methode zum Schleifen von Linsen und erhielt ein Patent für die erste Pendeluhr. Trotz des französisch-niederländischen Kriegs arbeitete er einen großen Teil seines Lebens an der *Académie Royale des Sciences* in Paris, wo beispielsweise LEIBNIZ viel Mathematik bei ihm lernte. HUYGENS war ein scharfer

Kritiker sowohl von NEWTONS Theorie des Lichts als auch seiner Gravitationstheorie, die er für absurd und nutzlos hielt. Gegen Ende seines Lebens beschäftigte er sich mit der Möglichkeit außerirdischen Lebens.

Da es uns nur um den mittleren Schalldruck, nicht aber um seine Variation geht, können wir den  $\omega t$ -Term ignorieren und einfach mit der Funktion  $Ae^{-ikx}$  arbeiten. Wir interessieren uns, wieviel Schall unter welchem Winkel reflektiert wird.

Die Schallwellen die von zwei verschiedenen Punkten unter einem Winkel  $\theta$  ausgehen haben, wie die Zeichnung zeigt, einen Laufwegunterschied von  $x \sin \theta$ , wobei  $x$  den Abstand der beiden Punkte bezeichnet.

Der Laufwegunterschied von  $x \sin \theta$  entspricht einem Phasenfaktor  $e^{-ikx \sin \theta}$ . Wählen wir also die Phase im Nullpunkt als Referenz (die wir in den zu ignorierenden Phasenfaktor der einfallenden Welle hineinziehen können), ist die Summe aller unter dem Winkel  $\theta$  abge-



henden Strahlen gleich

$$\int_{-\infty}^{\infty} \rho(x) e^{-ikx \sin \theta} dx ;$$

das ist die sogenannte FOURIER-Transformierte von  $\rho(x)$ , ausgewertet im Punkt  $u = k \sin \theta$ . Wenn wir den Schall möglichst gleichmäßig verteilen wollen, müssen wir die Funktion  $\rho$  daher so wählen, daß ihre FOURIER-Transformierte möglichst konstant ist.

Eine Möglichkeit dazu sind das, was Physiker als *Reflektions-Phasengitter* bezeichnen: Die Decke besteht aus einem Material mit konstantem, möglichst großem Reflektionsgrad, aber die Höhe der Decke variiert stufenförmig mit dem Querschnitt. Wenn die Höhe der einer festen Stelle um den Betrag  $h$  über der Nulllinie liegt, muß der dort reflektierte Schall gegenüber dem an der Nulllinie reflektierten den zusätzlichen Weg  $2h$  zurücklegen; dies kann man formal so ausdrücken, daß man in der Reflektionsfunktion  $r(x)$  den zusätzlichen Faktor  $e^{2i\omega h}$  einfügt.

Bei den sogenannten SCHROEDER-Reflektoren werden die Abstände zur Nulllinie so gewählt, daß die Längen  $2\omega h$  gleich den quadratischen Resten modulo einer ungeraden Primzahl sind, die Decke ist also treppenförmig aufgebaut, wobei die  $n$ -te Stufe eine Höhe proportional zu  $n^2 \bmod p$  hat. Das obige FOURIER-Integral läßt sich dann approximieren durch die diskrete FOURIER-Transformierte

$$\hat{r}(m) = \frac{1}{\sqrt{p}} \sum_{n=0}^{p-1} e^{2\pi i n^2/p} e^{-2\pi i n m t} = \frac{1}{\sqrt{p}} \sum_{n=0}^{p-1} e^{2\pi i n(n-m)/p} .$$

Ihr Betragsquadrat ist

$$\begin{aligned}
 |\widehat{r}(m)|^2 &= \frac{1}{\sqrt{p}} \sum_{n=0}^{p-1} e^{2\pi i n(n-m)/p} \cdot \frac{1}{\sqrt{p}} \sum_{n=0}^{p-1} e^{-2\pi i n(n-m)/p} \\
 &= \frac{1}{p} \sum_{n=0}^{p-1} \sum_{k=0}^{p-1} e^{2\pi i n(n-m)/p} e^{-2\pi i k(k-m)/p} \\
 &= \frac{1}{p} \sum_{n=0}^{p-1} \sum_{k=0}^{p-1} e^{2\pi i (n^2 - k^2 - (n-k)m)/p}
 \end{aligned}$$

Die Summanden hängen nur ab von den Restklassen modulo  $p$  der Indizes  $k$  und  $n$ , und für festes  $n$  durchläuft mit  $k$  auch  $n - k$  alle diese Restklassen. Daher können wir dies weiter ausrechnen als

$$\begin{aligned}
 |\widehat{r}(m)|^2 &= \frac{1}{p} \sum_{n=0}^{p-1} \sum_{k=0}^{p-1} e^{2\pi i (n^2 - (n-k)^2 - km)/p} \\
 &= \frac{1}{p} \sum_{k=0}^{p-1} e^{-2\pi i km/p} \sum_{n=0}^{p-1} e^{2\pi i ((n^2 - (n-k)^2)/p)} \\
 &= \frac{1}{p} \sum_{k=0}^{p-1} e^{-2\pi i km/p} \sum_{n=0}^{p-1} e^{2\pi i (2kn - k^2)/p} .
 \end{aligned}$$

Die zweite Summe können wir schreiben als

$$e^{-2\pi i k^2} \sum_{n=0}^{p-1} e^{4\pi i kn/p} .$$

Für  $k = 0$  ist sowohl der Vorfaktor wie auch jeder der Summanden gleich eins, wir erhalten also insgesamt  $p$ . Für  $k \neq 0$  und  $k < p$  ist die Summe aber gleich null, denn

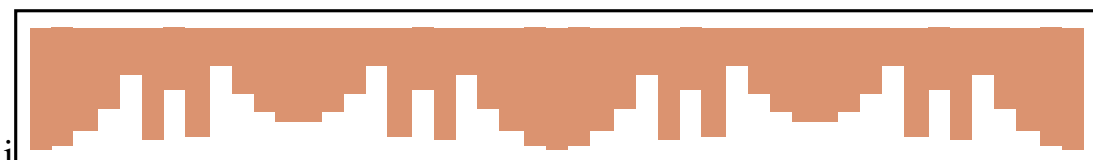
$$e^{4\pi i k/p} \sum_{n=0}^{p-1} e^{4\pi i kn/p} = \sum_{n=0}^{p-1} e^{4\pi i (k+1)n/p} = \sum_{n=1}^p e^{4\pi i kn/p} = \sum_{n=0}^{p-1} e^{4\pi i kn/p} ,$$

die Summe ändert also ihren Wert nicht, wenn man sie mit der von eins verschiedenen Zahl  $e^{4\pi i k/p}$  multipliziert, und damit muß sie verschwin-

den. Somit ist

$$|\widehat{r}(m)|^2 = \frac{1}{p} e^0 \cdot p = 1$$

für alle  $m$ , wir haben also die gewünschte Diffusionseigenschaft.



Die obige Abbildung zeigt den Querschnitt über ein solches Phasengitter, hier für  $p = 23$ . Entsprechende SCHROEDER-Reflektoren zu den verschiedensten Primzahlen gibt es in vielen Konzertsälen und Opernhäusern, oft allerdings verborgen hinter schalldurchlässigem Material.



MANFRED ROBERT SCHROEDER (1926–2009) wurde in Ahlen geboren. Er studierte Physik an der Universität Göttingen, wo er 1952 promovierte. Danach arbeitete er bei den AT & T Bell Laboratories in Murray Hill, New Jersey auf dem Gebiet der Akustik; diese Arbeit führte unter anderem zu 45 Patenten. 1969 wechselte er als Professor für Akustik an die Universität Göttingen, wo er bis zu seiner Emeritierung lehrte. Er schrieb mehrere Bücher, unter anderem

*Number theory in Science and Communication* und *Fractals, Chaos, Power Laws*. Der Inhalt dieses Abschnitts ist kurz im ersten dieser Bücher dargestellt sowie ausführlich in M.R. SCHROEDER: Binaural dissimilarity and optimum ceilings for concert halls: More lateral sound diffusion, *J. Acoust. Soc. Am.* **65** (4), 1979

[www.physik3.gwdg.de/~mrs](http://www.physik3.gwdg.de/~mrs)



# Kapitel 9

## Die Fermat-Vermutung für Zahlen und für Polynome

### § 1: Zahlen und Funktionen

Zum Beweis der eindeutigen Primzerlegung in der Hauptordnung eines Zahlkörpers mußten wir nur nachweisen, daß diese ein EUKLIDischer Ring ist, denn wie wir aus Kapitel 6, §5 wissen, ist jeder EUKLIDische Ring faktoriell. Da der Polynomring in einer Veränderlichen über einem Körper EUKLIDisch ist, gilt also auch dort das Gesetz von der eindeutigen Primzerlegung, wobei die irreduziblen Polynome die Rolle der Primzahlen einnehmen.

Besonders einfach ist die Situation, wenn der Grundkörper  $k$  algebraisch abgeschlossen ist: Dann hat jedes nichtkonstante Polynom  $f \in k[x]$  eine Nullstelle  $a$  und ist damit durch  $x - a$  teilbar. In diesem Fall sind also alle irreduziblen Polynome linear.

Die Zerlegung in irreduzible Elemente ist bekanntlich nur eindeutig bis auf Einheiten, und die Einheiten eines Polynomrings sind nach §1 aus Kapitel 6 gerade die des Koeffizientenrings, hier also die von Null verschiedenen Elemente von  $k$ . Durch Multiplikation mit einem solchen Element kann man den höchsten Koeffizienten eines jeden Polynoms zu eins machen; die irreduziblen Polynome über einem algebraisch abgeschlossenen Körper sind also bis auf Assoziiertheit genau die Polynome der Form  $x - a$  mit  $a \in k$  und sie entsprechen eindeutig den Elementen von  $k$ .

Den Primzahlen von  $\mathbb{Z}$  entsprechen daher im Polynomring über einem algebraisch abgeschlossenen Körper die Punkte der affinen Geraden über  $k$ , wir können also geometrisch argumentieren.

Natürlich gibt es – zum Teil beträchtliche – Unterschiede zwischen  $\mathbb{Z}$  und dem Polynomring über einem Körper, aber gerade das macht die Analogie so interessant: Da es für jeden der beiden Ringe ein eigenes Instrumentarium gibt, kann man versuchen die damit bewiesenen Resultate auf den jeweils anderen Fall zu übertragen, was idealerweise zu neuen Sätzen und sonst zumindest zu interessanten Vermutungen führt.

Als Beispiel für Parallelen und Unterschiede zwischen den beiden Situationen wollen wir die FERMAT-Vermutung betrachten. FERMAT schrieb bekanntlich um 1637 an den Rand seiner Arithmetik des DIOPHANTOS von Alexandrien, daß die Gleichung

$$x^n + y^n = z^n$$

für  $n \geq 3$  keine Lösung in ganzen Zahlen habe – außer natürlich den trivialen Lösungen, bei denen eine der beiden Variablen verschwindet. (Die französische Übersetzung der Arithmetik, die er dabei benutzte, stammt übrigens von BACHET DE MÉZIRIAC, denn wir bereits als Entdecker des erweiterten EUKLIDISCHEN Algorithmus kennen. Bekannt wurde FERMATs Randbemerkung erst, als dessen Sohn CLÉMENT-SAMUEL DE FERMAT 1670 die Arithmetik mit den Randbemerkungen seines fünf Jahre zuvor gestorbenen Vaters veröffentlichte.)

Die direkte Verallgemeinerung auf Polynomringe ist sicherlich falsch: Die Gleichung  $f^n + g^n = h^n$  ist zumindest für *konstante* Polynome über einem algebraisch abgeschlossenen Körper immer lösbar: Für beliebig vorgegebene Konstanten  $f, g \in k$  muß man einfach  $h = \sqrt[n]{f^n + g^n}$  setzen. Das sind allerdings, wenn wir uns wirklich für Polynome interessieren, uninteressante Lösungen, vergleichbar den Lösungen  $x^n + 0^n = x^n$  der klassischen FERMAT-Gleichung.

Auch wenn wir verlangen, daß die Grade aller beteiligter Polynome positiv sein sollen, gibt es triviale Lösungen: Ist  $f$  irgendein beliebiges Polynom und sind  $a, b, c \in k$  so, daß gilt  $a^n + b^n = c^n$ , ist natürlich auch  $(af)^n + (bf)^n = (cf)^n$ . Was wir höchstens erwarten können ist also das folgende Analogon zur klassischen FERMAT-Vermutung:

*Es ist nicht möglich, einen Kubus in zwei Kuben oder ein Biquadrat in zwei Biquadraten und ganz allgemein irgendeine der unendlich vielen Potenzen jenseits des Quadrats in zwei ebensolche zu teilen. Ich habe einen wunderbaren Beweis hierfür gefunden, aber der Rand ist zu schmal, um ihn zu fassen.*

Für  $n \geq 3$  gibt es keine paarweise teilerfremden Polynome  $f, g, h$  mit positivem Grad, so daß  $f^n + g^n = h^n$  ist.

Für Körper positiver Charakteristik ist selbst das noch falsch: Über einen Körper der Charakteristik  $p$  ist schließlich  $f^p + g^p = (f + g)^p$  für beliebige Polynome  $f$  und  $g$ , und dasselbe gilt auch wenn man den Exponenten  $p$  durch eine seiner Potenzen ersetzt. Wir können also höchstens für Körper der Charakteristik null erwarten, daß diese Vermutung für alle Exponenten  $n \geq 3$  richtig ist, und genau das werden wir im übernächsten Paragraphen (zumindest für den Körper der komplexen Zahlen) auch beweisen. Zunächst aber wollen wir schauen, was im bei FERMAT ausgeschlossenen Fall des Exponenten zwei passiert.

## §2: Pythagoräische Tripel

Betrachten wir zunächst den Fall der Polynome, wobei wir uns der Einfachheit halber gleich auf Polynome mit komplexen Koeffizienten beschränken wollen. Ist  $f^2 + g^2 = h^2$ , so ist

$$f^2 = h^2 - g^2 = (h + g)(h - g).$$

Wenn wir  $f$  und  $g$  als teilerfremd voraussetzen, sind auch  $g$  und  $h$  teilerfremd und somit auch  $h + g$  und  $h - g$ , denn jeder gemeinsame Teiler dieser beiden Polynome wäre auch ein Teiler ihrer Summe  $2h$  sowie ihrer Differenz  $2g$ .

Wenn wir die Zerlegung von  $f^2$  in irreduzible Faktoren vergleichen mit der von  $h + g$  und  $h - g$  folgt somit, daß jeder irreduzible Faktor von  $f$  entweder in  $h + g$  oder in  $h - g$  in gerader Potenz auftreten muß. Da jede komplexe Zahl ein Quadrat ist, können wir auch eine eventuell auftretende Einheit als Quadrat schreiben; somit gibt es zwei Polynome  $u, v \in \mathbb{C}[x]$  derart, daß

$$u^2 = h + g, \quad v^2 = h - g \quad \text{und} \quad uv = f$$

ist, d.h.

$$f = uv, \quad g = \frac{u^2 - v^2}{2} \quad \text{und} \quad h = \frac{u^2 + v^2}{2}.$$

Starten wir umgekehrt mit zwei beliebigen teilerfremden Polynomen  $u, v \in \mathbb{C}[x]$ , erhalten mit Hilfe dieser Formeln Lösungen der Gleichung  $f^2 + g^2 = h^2$ . Damit kennen wir alle teilerfremden Lösungen, und die restlichen erhalten wir, indem wir alle drei Polynome mit einem gemeinsamen Faktor multiplizieren.

Nehmen wir als ein einfaches Beispiel  $u = 2x$  und  $v = 2$ , ist also

$$f = 4x, \quad g = x^2 - 1 \quad \text{und} \quad h = x^2 + 1;$$

in der Tat ist

$$(2x)^2 + (x^2 - 1)^2 = (x^2 + 1)^2.$$

Hier erhalten wir also mit geringem Aufwand eine vollständige Übersicht über alle Lösungen.

Versuchen wir das gleiche auch für den klassischen Fall! Wegen des Satzes von PYTHAGORAS bezeichnet man ein Tripel  $(x, y, z)$  ganzer Zahlen mit  $x^2 + y^2 = z^2$  als pythagoräisches Tripel; für jedes solche Tripel gibt es ein rechtwinkliges Dreieck mit Seitenlängen  $|x|$ ,  $|y|$  und  $|z|$ .

Wir bezeichnen das Tripel  $(x, y, z)$  als *primitiv*, wenn sie sich nicht als Vielfaches eines anderen schreiben läßt, wenn also die Zahlen  $x, y, z$  keinen gemeinsamen Teiler haben. Sobald wir alle primitiven Lösungen kennen, können wir daraus die gesamte Lösungsmenge konstruieren, denn jede nichtprimitive Lösung ist Vielfaches einer primitiven.

Wie im Fall der Polynome gehen wir aus von einem primitiven Tripel  $(x, y, z)$  und wenden die dritte binomische Formel an:

$$x^2 = z^2 - y^2 = (z + y)(z - y).$$

Hier können wir leider nicht mehr ohne weiteres folgern, daß  $z + y$  und  $z - y$  teilerfremd und damit Quadrate sind: Sind  $y$  und  $z$  beide ungerade, so sind ihre Summe und ihre Differenz beide gerade, also durch zwei teilbar. Wir müssen uns also zunächst über die Paritäten von  $y$  und  $z$  klarwerden.

Für eine primitive Lösung  $(x, y, z)$  müssen bereits  $x$  und  $y$  teilerfremd sein, denn ist  $d$  ein gemeinsamer Teiler von  $x$  und  $y$ , so sind  $x^2$  und  $y^2$  beide durch  $d^2$  teilbar, also auch ihre Summe  $z^2$ . Wegen der eindeutigen

Zerlegbarkeit einer natürlichen Zahl in Primfaktoren ist dann auch  $z$  durch  $d$  teilbar, d.h.  $d$  ist ein gemeinsamer Teiler von  $x$ ,  $y$  und  $z$ .

Insbesondere können daher  $x$  und  $y$  nicht beide gerade sein; mindestens eine der beiden Zahlen muß ungerade sein. Andererseits können aber auch nicht beide Zahlen ungerade sein: Wäre nämlich  $x = 2u + 1$  und  $y = 2v + 1$ , so wäre

$$z^2 = (2u + 1)^2 + (2v + 1)^2 = 4u^2 + 4u + 1 + 4v^2 + 4v + 1 \equiv 2 \pmod{4},$$

was unmöglich ist, da modulo vier nur null und eins Quadrate sind.

Somit muß in einem primitiven pythagoräischen Tripel  $(x, y, z)$  eine der beiden Zahlen  $x, y$  gerade sein und die andere ungerade. Da mit  $(x, y, z)$  auch  $(y, x, z)$  ein primitives pythagoräisches Tripel ist, genügt es, wenn wir diejenigen Tripel betrachten, in denen  $x$  gerade ist und  $y$  ungerade. Offensichtlich ist dann auch  $z$  ungerade.

Für so ein Tripel steht in der Gleichung

$$x^2 = z^2 - y^2 = (z + y)(z - y)$$

rechts das Produkt zweier gerader Zahlen. Im Gegensatz zur Situation bei den Polynomen haben wir hier also keine teilerfremden Faktoren.

Dividieren wir aber durch zwei, so können wir wie oben argumentieren, daß  $\frac{1}{2}(z+y)$  und  $\frac{1}{2}(z-y)$  teilerfremd sind, denn jeder gemeinsame Teiler wäre Teiler ihrer Summe  $z$  und ihrer Differenz  $y$ .

Jetzt können wir wieder die eindeutige Primzerlegung anwenden: Da

$$x^2 = 2^2 \cdot \left(\frac{z+y}{2}\right) \cdot \left(\frac{z-y}{2}\right),$$

wobei die beiden Klammern teilerfremd sind, gibt es ganze Zahlen  $u, v$ , so daß

$$u^2 = \left(\frac{z+y}{2}\right), \quad v^2 = \left(\frac{z-y}{2}\right) \quad \text{und} \quad 2uv = x$$

ist, also

$$x = 2uv, \quad y = u^2 - v^2 \quad \text{und} \quad z = u^2 + v^2.$$

Umgekehrt definieren diese Formeln für beliebige ganze Zahlen  $u$  und  $v$  ein primitives pythagoräisches Tripel, und bis auf die Reihenfolge von

$x$  und  $y$  erhalten wir so auch jedes dieser Tripel, für  $u = 2$  und  $v = 1$  etwa das seit Jahrtausenden bekannte Tripel  $(4, 3, 5)$ . Da sich in alten Sakralbauten und -anlagen auch deutlich kompliziertere pythagoräische Tripel nachweisen lassen, steht zu vermuten, daß die obige Konstruktion möglicherweise bereits in einigen steinzeitlichen Kulturen zumindest teilweise bekannt waren; entsprechende Thesen vertritt beispielsweise bekannte algebraische Geometer B.L. VAN DER WAERDEN (1903–1996), der nach seiner Emeritierung auch mehrere Bücher über die Geschichte der Mathematik veröffentlichte. Mit pythagoräischen Tripel beschäftigt er sich ausführlich in

B.L. VAN DER WAERDEN: *Geometry and algebra in ancient civilizations*, Springer, 1983

### §3: Der Satz von Mason

In diesem Abschnitt wollen wir, wie bereits angekündigt, sehen, daß es für  $n \geq 3$  keine zueinander teilerfremden Polynome positiven Grades  $f, g, h \in \mathbb{C}[x]$  gibt mit  $f^n + g^n = h^n$ .

Der *Beweis* beruht darauf, daß die Polynome  $f^n$  und  $g^n$  dieselben Nullstellen wie  $f$  und  $g$  haben, aber mit  $n$ -facher Vielfachheit. Ist  $f^n + g^n = h^n$ , so hat auch die Summe dieser beiden Potenzen im Vergleich zum Grad relativ wenige Nullstellen, diese aber mit mindestens  $n$ -facher Vielfachheit. Nach einem 1983 von R.C. MASON bewiesenen Satz können in einer solchen Situation aber  $f^n, g^n$  und  $h^n$  nicht zu wenige verschiedene Nullstellen haben:

**Satz:** Bezeichnet  $n_0(f)$  die Anzahl verschiedener (komplexer) Nullstellen eines Polynoms  $f$ , so gilt für drei nichtkonstante, teilerfremde Polynome  $f, g, h$  mit  $f + g = h$

$$n_0(fgh) \geq \max(\deg f, \deg g, \deg h) + 1 .$$

Bevor wir diesen Satz beweisen, wollen wir uns zunächst überlegen, daß daraus wirklich die FERMAT-Vermutung für Polynome folgt:

Für drei nichtkonstante teilerfremde Polynome  $f, g, h$  mit  $f^n + g^n = h^n$  ist nach dem Satz von MASON

$$\begin{aligned} n_0(f^n g^n h^n) &\geq \max(\deg f^n, \deg g^n, \deg h^n) + 1 \\ &= n \max(\deg f, \deg g, \deg h) + 1. \end{aligned}$$

Andererseits ist aber

$$\begin{aligned} n_0(f^n g^n h^n) &= n_0(fgh) \\ &\leq \deg f + \deg g + \deg h \\ &\leq 3 \max(\deg u(x), \deg v(x), \deg w(x)), \end{aligned}$$

denn die Anzahl *verschiedener* Nullstellen einer Potenz eines Polynoms ist gleich der Anzahl verschiedener Nullstellen des Polynoms selbst, und die Nullstellenanzahl eines Polynom kann nicht größer sein als der Grad.

Damit haben wir insgesamt die Ungleichung

$$\begin{aligned} 3 \max(\deg f, \deg g, \deg h) \\ &\geq n_0(f^n g^n h^n) \\ &\geq n \max(\deg f, \deg g, \deg h) + 1, \end{aligned}$$

die bei nichtkonstanten Polynomen nur für  $n \leq 2$  gelten kann. Somit gibt es für  $n \geq 3$  keine nichtkonstanten teilerfremden Polynome, für die  $f^n + g^n = h^n$  ist.

Zu einem vollständigen Beweis der FERMAT-Vermutung für Polynome fehlt nun nur noch der Beweis des Satzes von MASON. Die Idee dazu ist folgende: Ist  $f + g = h$ , so betrachten wir den Quotienten  $g/f$  im rationalen Funktionenkörper  $\mathbb{C}(x)$ . Da  $f$  und  $g$  teilerfremd sind, ist das ein gekürzter Bruch. Falls wir diesen auch in der Form  $g/f = G/F$  schreiben können mit Polynomen  $F, G$  vom Grad höchstens  $n_0(fgh) - 1$  schreiben können, so haben auch  $f$  und  $g$  höchstens den Grad  $n_0(fgh) - 1$ . Wegen  $f + g = h$  gilt dasselbe auch für  $h$ , und damit wäre der Satz bewiesen.

Um  $g/f$  als Quotienten zweier neuer Polynome auszudrücken, schreiben wir zunächst

$$\frac{g}{f} = \frac{S}{R} \quad \text{mit} \quad R = \frac{f}{h} \quad \text{und} \quad S = \frac{g}{h}.$$

Dabei ist  $R + S = 1$ , die Summe  $R' + S'$  der Ableitungen verschwindet also. Aus der Gleichung

$$R' + S' = \frac{R'}{R}R + \frac{S'}{S}S$$

folgt die neue Darstellung

$$\frac{g}{f} = \frac{S}{R} = -\frac{R'/R}{S'/S}.$$

Rechts stehen die logarithmischen Ableitungen von  $R$  und  $S$  im Zähler und Nenner, und damit lassen sich gut die Nullstellen von  $f$ ,  $g$  und  $h$  ins Spiel bringen: Nach der LEIBNIZ-Regel ist bekanntlich

$$(uv)' = u'v + uv', \quad \text{also} \quad \frac{(uv)'}{uv} = \frac{u'}{u} + \frac{v'}{v},$$

die logarithmische Ableitung eines Produkts ist also einfach die Summe der logarithmischen Ableitungen der Faktoren. Daraus folgt sofort, daß die logarithmische Ableitung eines Quotienten gleich der Differenz aus logarithmischer Ableitung des Zählers und logarithmischer Ableitung des Nenners ist. Schreiben wir

$$f = f_0 \prod_{i=1}^r (x-a_i)^{n_i}, \quad g = g_0 \prod_{j=1}^s (x-b_j)^{m_j} \quad \text{und} \quad h = h_0 \prod_{k=1}^t (x-c_k)^{p_k},$$

so ist also

$$\begin{aligned} \frac{R'}{R} &= \frac{f'}{f} - \frac{h'}{h} = \sum_{i=1}^r \frac{n_i}{x-a_i} - \sum_{k=1}^t \frac{p_k}{x-c_k}, \\ \frac{S'}{S} &= \frac{g'}{g} - \frac{h'}{h} = \sum_{j=1}^s \frac{m_j}{x-b_j} - \sum_{k=1}^t \frac{p_k}{x-c_k} \\ \text{und} \quad \frac{g}{f} &= \frac{R'/R}{S'/S} = -\frac{\sum_{i=1}^r \frac{n_i}{x-a_i} - \sum_{k=1}^t \frac{p_k}{x-c_k}}{\sum_{j=1}^s \frac{m_j}{x-b_j} - \sum_{k=1}^t \frac{p_k}{x-c_k}}. \end{aligned}$$



Erweitern wir Zähler und Nenner mit dem Hauptnenner aller Summanden erweitern, d.h. mit dem Polynom vom Grad  $r + s + t = n_0(fgh)$

$$H = \prod_{i=1}^r (x - a_i) \cdot \prod_{j=1}^s (x - b_j) \cdot \prod_{k=1}^t (x - c_k),$$

so erhalten wir im Zähler wie auch im Nenner Summen von Polynomen vom Grad  $n_0(fgh) - 1$ , als Polynome vom Grad höchstens  $n_0(fgh) - 1$ , wie gewünscht. Damit ist der Satz von MASON bewiesen.

#### §4: Die abc-Vermutung

Der Erfolg des Satzes von MASON beim Beweis der FERMAT-Vermutung für Polynome legt es nahe, etwas ähnliches auch im klassischen Fall zu versuchen.

Da natürliche Zahlen weder Grade noch Nullstellen haben, müssen wir dazu den Satz von MASON zunächst einmal so umformulieren, daß wir eine Aussage bekommen, die ein sinnvolles Analogon für natürliche Zahlen hat.

Dazu ordnen wir einem Polynom  $f$  anstelle der Anzahl  $n_0(f)$  seiner (verschiedenen) Nullstellen ein Polynom  $N_0(f)$  dazu, das genau diese Nullstellen mit jeweils der Vielfachheit eins haben soll: Für

$$f = f_0 \prod_{i=1}^r (x - a_i)^{n_i} \quad \text{sei} \quad N_0(f) \stackrel{\text{def}}{=} \prod_{i=1}^r (x - a_i),$$

so daß der Grad von  $N_0(f)$  gerade die im vorigen Paragraphen definierte Zahl  $n_0(f)$  ist.

Der Vorteil des Polynoms  $N_0(f)$  besteht darin, daß wir eine analoge Definition leicht auch für natürliche Zahlen hinschreiben können: Für

$$n = \prod_{i=1}^r p_i^{e_i} \quad \text{setzen wir} \quad N_0(n) \stackrel{\text{def}}{=} \prod_{i=1}^r p_i.$$

Mit Hilfe der Polynome  $N_0(f)$  läßt sich der Satz von MASON folgendermaßen umformulieren:

*Gilt für drei teilerfremde Polynome  $f, g$  und  $h$  die Gleichung  $f + g = h$ , so hat jedes der drei Polynome einen kleineren Grad als das Polynom  $N_0(fgh)$ .*

In dieser Formulierung kommt immer noch der Grad vor, für den wir bei natürlichen Zahlen keine Verwendung haben. Betrachten wir aber den Grad lediglich als eine Methode, einem Polynom eine Zahl aus  $\mathbb{N}_0$  zuzuordnen, so können wir, wenn wir bereits natürliche Zahlen haben, einfach ganz auf ihn verzichten; falls wir ganze Zahlen betrachten, liegt es nahe, ihn durch den Betrag zu ersetzen.

Gemäß dieser Philosophie können wir nun probeweise die folgende Aussage formulieren:

**A1:** *Ist  $c = a+b$  für drei zueinander teilerfremde natürliche Zahlen  $a, b, c$ , so ist jede der drei Zahlen kleiner als  $N_0(abc)$ .*

Damit haben wir eine sinnvolle Aussage über natürliche Zahlen gefunden, die – falls sie korrekt ist – sofort die FERMAT-Vermutung impliziert: Gibt es nämlich drei natürliche Zahlen  $x, y, z$  mit der Eigenschaft, daß  $x^n + y^n = z^n$  für ein  $n \geq 3$ , so gibt es auch drei zueinander teilerfremde Zahlen  $x, y, z$  mit dieser Eigenschaft: Wir müssen einfach die drei Zahlen durch ihren größten gemeinsamen Teiler kürzen. Alsdann muß, falls obige Aussage richtig ist, jede der drei Potenzen  $x^n, y^n, z^n$  kleiner sein als  $N_0(x^n y^n z^n)$ . Nun ist aber

$$N_0(x^n y^n z^n) = N_0(xyz) \leq xyz,$$

d.h. jede der drei Zahlen  $x^n, y^n, z^n$  wäre kleiner als  $xyz$ . Damit wäre

$$(xyz)^n = x^n y^n z^n < (xyz)^3,$$

was für  $n \geq 3$  offensichtlich nicht möglich ist.

Angesichts der Komplexität des WILESSchen Beweises fällt es schwer, an einen so einfachen Beweis zu glauben, und in der Tat ist die obige Aussage in dieser Form falsch:

Betrachten wir etwa die Gleichung  $8 + 1 = 9$ . Offensichtlich sind die drei Summanden teilerfremd zueinander, aber sowohl 8 als auch 9 sind größer als  $N_0(8 \cdot 1 \cdot 9) = 2 \cdot 3 = 6$ . Ganz so einfach geht es also nicht.

Da der Grad eines Polynoms nicht durch konstante Faktoren beeinflußt wird, könnte man versuchen, als „richtiges“ Analogon zum Satz von MASON eine abgeschwächte Aussage zu nehmen, die nur eine Abschätzung bis auf einen konstanten Faktor enthält, etwa

**A2:** *Ist  $c = a+b$  für drei zueinander teilerfremde natürliche Zahlen  $a, b, c$ , so gibt es eine Konstante  $K$  derart, daß jede der drei Zahlen kleiner ist als  $K \cdot N_0(abc)$ .*

Diese Aussage ist trivialerweise richtig: Wir müssen nur eine Konstante  $K$  wählen, die größer ist als das Maximum von  $a, b$  und  $c$ . Leider ist sie auch völlig nutzlos, denn solange die Konstante von  $a, b$  und  $c$  abhängen darf, haben wir keine Chance, damit die FERMAT-Vermutung zu beweisen.

Wir müssen die Aussage also noch einmal umformulieren:

**A3:** *Es gibt eine Konstante  $K$ , so daß gilt: Ist  $c = a+b$  für drei zueinander teilerfremde natürliche Zahlen  $a, b, c$ , so ist jede der drei Zahlen kleiner als  $K \cdot N_0(abc)$ .*

Wie wir gleich sehen werden, würde hieraus die FERMAT-Vermutung zumindest für alle hinreichend großen Exponenten  $n$  folgen, allerdings ist die Aussage, so wie sie dasteht, leider immer noch falsch:

Betrachten wir die Gleichung

$$a_n + b_n = c_n \quad \text{mit} \quad a_n = 3^{2^n} - 1, \quad b_n = 1 \quad \text{und} \quad c_n = 3^{2^n}. \quad (*)$$

Wäre sie richtig, müßte für jedes  $n$  gelten:

$$3^{2^n} \leq K N_0((3^{2^n} - 1)3^{2^n}) = K \cdot 3 \cdot N_0(3^{2^n} - 1).$$

Um  $N_0(3^{2^n} - 1)$  abzuschätzen, beachten wir, daß gilt

$$3^{2^n} = (3^{2^{n-1}})^2 \quad \text{und} \quad 3^{2^n} - 1 = (3^{2^{n-1}} + 1)(3^{2^{n-1}} - 1)$$

nach der dritten binomischen Formel. Wenden wir dies mehrfach an,

erhalten wir

$$\begin{aligned}
 3^{2^n} - 1 &= (3^{2^{n-1}} + 1)(3^{2^{n-1}} - 1) \\
 &= (3^{2^{n-1}} + 1)(3^{2^{n-2}} + 1)(3^{2^{n-2}} - 1) \\
 &= (3^{2^{n-1}} + 1)(3^{2^{n-2}} + 1)(3^{2^{n-3}} + 1)(3^{2^{n-3}} - 1) \\
 &= \dots \\
 &= (3^{2^{n-1}} + 1)(3^{2^{n-2}} + 1) \dots (3^2 + 1)(3^1 + 1)(3^1 - 1).
 \end{aligned}$$

In der letzten Zeile steht ein Produkt aus  $n + 1$  geraden Zahlen; somit ist  $3^{2^n} - 1$  durch  $2^{n+1}$  teilbar. Das Produkt  $N_0(3^{2^n} - 1)$  aller *verschiedener* Primteiler von  $3^{2^n} - 1$  erfüllt daher die Ungleichung

$$N_0(3^{2^n} - 1) \leq 2 \cdot \frac{3^{2^n} - 1}{2^{n+1}} = \frac{3^{2^n} - 1}{2^n},$$

denn das Produkt aller ungerader Primteiler kann höchstens gleich  $(3^{2^n} - 1)/2^n$  sein.

Falls **A3** korrekt wäre, müßte nach Gleichung (\*) also gelten

$$3^{2^n} \leq \frac{3K}{2^n}(3^{2^n} - 1) \quad \text{für alle } n.$$

Das kann aber unmöglich der Fall sein, denn für hinreichend große  $n$  ist der Faktor  $\frac{3K}{2^n}$  kleiner als eins, so daß  $3^{2^n}$  echt kleiner als sich selbst sein müßte.

Auf der Suche nach einem Analogon für den Satz von MASON müssen wir daher noch weiter abschwächen. *Eine* Möglichkeit dazu ist die 1986 aufgestellte

**abc-Vermutung** von MASSER und OESTERLÉ: Zu jedem  $\varepsilon > 0$  gibt es eine Konstante  $K(\varepsilon)$ , so daß für alle teilerfremden natürlichen Zahlen  $a, b, c$  mit  $a + b = c$  gilt: Jede der drei Zahlen  $a, b, c$  ist kleiner oder gleich  $K(\varepsilon) \cdot N_0(abc)^{1+\varepsilon}$ .

Diese Vermutung ist, im Gegensatz zur FERMAT-Vermutung, bis heute offen.

Wir wollen uns überlegen, daß sie zumindest für große Exponenten  $n$  die FERMAT-Vermutung impliziert.

Dazu betrachten wir eine Lösung  $x^n + y^n = z^n$  mit o.B.d.A. teilerfremden natürlichen Zahlen  $x, y, z$  und wählen uns irgendein  $\varepsilon > 0$ . Nach der *abc*-Vermutung gibt es dazu eine Konstante  $K(\varepsilon)$ , so daß  $x^n, y^n$  und  $z^n$  allesamt höchstens gleich

$$K(\varepsilon)N_0(x^n y^n z^n)^{1+\varepsilon} = K(\varepsilon)N_0(xyz)^{1+\varepsilon} \leq K(\varepsilon)(xyz)^{1+\varepsilon}$$

sind. Für ihr Produkt gilt daher

$$x^n y^n z^n \leq K(\varepsilon)^3 (xyz)^{3(1+\varepsilon)} \quad \text{oder} \quad (xyz)^{n-3-3\varepsilon} \leq K(\varepsilon)^3.$$

$K(\varepsilon)^3$  ist eine feste Zahl; es gibt daher einen Exponenten  $m$  derart, daß  $2^m > K(\varepsilon)^3$  ist. Da das Produkt  $xyz$  auf jeden Fall nicht kleiner als zwei sein kann, ist dann für  $n - 3 - 3\varepsilon > m$  oder  $n > m + 3 + 3\varepsilon$  insbesondere

$$(xyz)^{n-3-3\varepsilon} > K(\varepsilon)^3.$$

Für Exponenten  $n > m + 3 + 3\varepsilon$  kann daher die FERMAT-Gleichung keine Lösung in natürlichen Zahlen haben.

Ob und gegebenenfalls welche konkreten Schranken für  $n$  man damit erreichen kann, hängt natürlich davon ab, wie  $K(\varepsilon)$  von  $\varepsilon$  abhängt. Dazu gibt es im Augenblick nicht einmal Vermutungen.

Für weitere Informationen zu §3 und §4 sei auf einen Vortrag verwiesen, den SERGE LANG in Zürich vor einem „allgemeinen“ Publikum hielt und dem ich hier im wesentlichen gefolgt bin:

SERGE LANG: Die *abc*-Vermutung, *Elemente der Mathematik* **48** (1993), 89-99

Der Artikel ist (wie die gesamte Zeitschrift *Elemente der Mathematik*) unter <http://www.bibliothek.uni-regensburg.de/ezeit/?2135837> frei zugänglich.

## §5: Die Frey-Kurve

Da die FERMAT-Vermutung seit 1994 bewiesen ist, die *abc*-Vermutung aber immer noch offen, mußte der Beweis der FERMAT-Vermutung

natürlich andere Wege gehen. Die meisten dieser Wege führen in Gebiete, die weit jenseits dessen liegen, was selbst ein guter auf Zahlentheorie spezialisierter Diplom-Mathematiker im Laufe seines Studiums lernen kann, aber zumindest die Grundidee der *abc*-Vermutung, daß man nämlich Summenbeziehungen zwischen großen Zahlen nicht ohne ein gewisses Minimum an verschiedenen Primfaktoren realisieren kann, spielt in modifizierter Weise in der Tat eine große Rolle.

Der Anstoß kam 1984 von GERHARD FREY, damals Professor an der Universität Saarbrücken, wo er auf dem Gebiet der Arithmetik elliptischer Kurven arbeitete. Heute leitet er die Arbeitsgruppe Zahlentheorie am Institut für experimentelle Mathematik der (inzwischen mit Duisburg vereinigten) Universität Essen und beschäftigt sich mit der Anwendung elliptischer (und anderer) Kurven in der Kryptologie.

Elliptische Kurven sind ebene Kurven, die durch eine Gleichung der Form  $y^2 = f_3(x)$  beschrieben werden mit einem Polynom  $f_3(x)$  vom Grad drei mit drei verschiedenen Nullstellen. Da das Quadrat einer reellen Zahl nicht negativ sein kann, gibt es im Reellen nur Punkte mit  $x$ -Koordinaten, für die  $f_3(x) \geq 0$  ist. Im Falle  $f_3(x) > 0$  erfüllt mit  $y$  auch  $-y$  die obige Gleichung, die Kurve ist also symmetrisch zur  $x$ -Achse.

Falls  $f_3(x)$  nur zwei verschiedene Nullstellen hat, muß eine der Nullstellen doppelt sein, und bei diesem  $x$ -Wert überkreuzt sich die Kurve; wir reden dann von einer Knotenkurve.

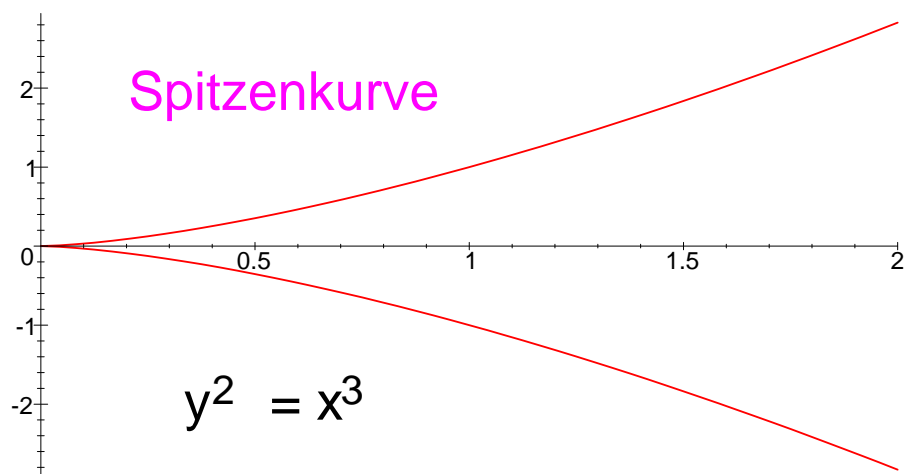
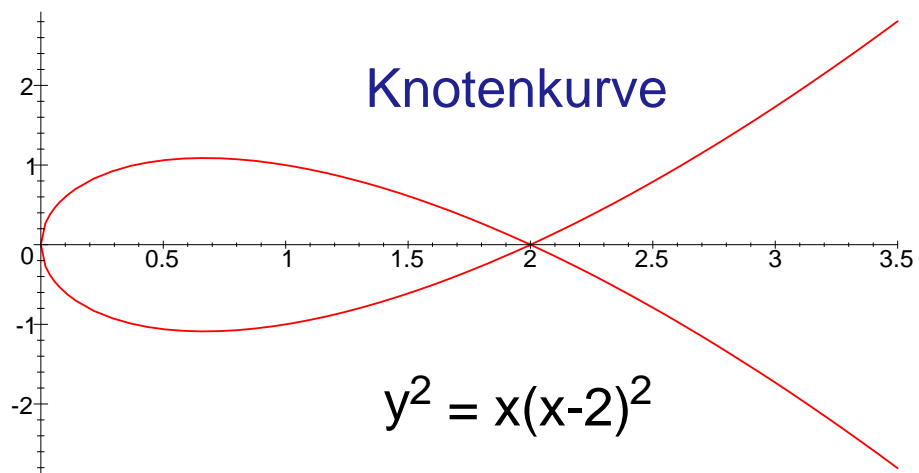
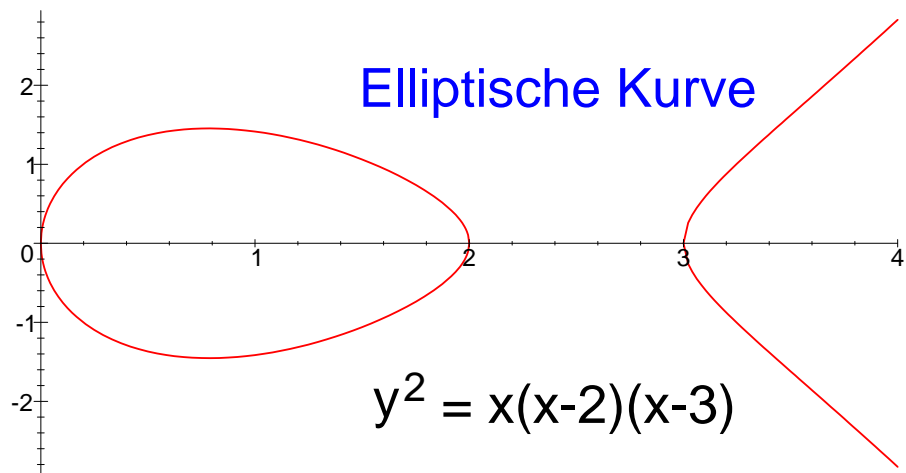
Hat schließlich  $f_3(x)$  nur eine, dafür aber dreifache Nullstelle, entsteht eine Spitzenkurve.

FREY betrachtete eine (hypothetische) Lösung

$$x^n + y^n = z^n$$

der FERMAT-Gleichung mit teilerfremden natürlichen Zahlen  $x, y, z$  und  $n \geq 5$ . (Den Fall  $n = 4$  hat, wie wir gesehen haben, bereits FERMAT selbst gelöst, den Fall  $n = 3$  nicht viel später EULER.) Wenn es eine solche Lösung gibt, dann gibt es auch eine Lösung für einen Primzahlexponenten  $\ell$ , denn ist  $\ell$  ein Primteiler von  $n$  und  $n = \ell m$ , so ist auch

$$a^\ell + b^\ell = c^\ell \quad \text{mit} \quad a = x^m, \quad b = y^m \quad \text{und} \quad c = z^m$$



eine Lösung, und auch  $a, b, c$  sind teilerfremd. Auch für  $\ell$  genügt es, den Fall  $\ell \geq 5$  zu betrachten, denn wenn wir für  $\ell$  den größten Primteiler von  $n$  nehmen, bedeutet  $\ell = 2$ , daß  $n$  eine Zweierpotenz sein muß, was für  $n = 2$  kein Widerspruch zur FERMAT-Vermutung ist und für  $n = 4$  und damit auch jede höhere Zweierpotenz nach FERMATs Beweis ausgeschlossen ist. Für den Fall  $\ell = 3$  kann wieder auf EULER verwiesen werden.

Zur obigen Lösung definiert FREY die elliptische Kurve

$$y^2 = x(x - a^\ell)(x + b^\ell)$$

zu, die er aber nicht nur über den reellen oder komplexen Zahlen betrachtet, sondern auch über den ganzen Zahlen modulo einer Primzahl  $p$ :

Ist allgemein

$$y^2 = (x - x_1)(x - x_2)(x - x_3)$$

eine Kurvengleichung mit ganzen Zahlen  $x_1, x_2, x_3$ , so können wir für  $x$  und  $y$  auch ganze Zahlen einsetzen und diese Gleichung modulo  $p$  betrachten. Wir sprechen wieder von einer elliptischen Kurve, einer Knotenkurve oder einer Spitzenkurve je nachdem wie viele der Nullstellen  $x_1, x_2$  und  $x_3$  modulo  $p$  noch verschieden sind.

Die obige Gleichung definiert genau dann eine elliptische Kurve, wenn alle drei Nullstellen verschieden sind, wenn also die sogenannte Diskriminante

$$\Delta = (x_1 - x_2)(x_1 - x_3)(x_2 - x_3)$$

von Null verschieden ist. Modulo  $p$  definiert sie eine elliptische Kurve, wenn  $\Delta$  auch modulo  $p$  noch von Null verschieden ist, wenn also  $p$  kein Teiler von  $\Delta$  ist.

Speziell für die FREY-Kurve  $y^2 = x(x - a^\ell)(x + b^\ell)$  ist die Diskriminante

$$\Delta = (0 - a^\ell)(0 - b^\ell)(a^\ell - (-b)^\ell) = a^\ell b^\ell (a^\ell + b^\ell) = a^\ell b^\ell c^\ell = (abc)^\ell$$

stets von Null verschieden; modulo  $p$  verschwindet sie genau dann, wenn  $p$  ein Teiler von  $\Delta$  ist, d.h., wenn  $p$  eine der drei Zahlen  $a, b, c$  teilt.

Da die Diskriminante als  $\ell$ -te Potenz von  $abc$  verglichen mit  $a, b, c$  ziemlich groß ist, heißt das, daß es im Verhältnis zur Größe der Diskriminante



erstaunlich wenige Primzahlen gibt, modulo derer wir *keine* elliptische Kurve erhalten; wir sind also wieder einer ähnlichen Situation wie bei der *abc*-Vermutung. Die FREYSche Kurve sieht damit so aus, als sei sie fast zu schön, um wirklich zu existieren.

Einen Anhaltspunkt zum Beweis dieser Nichtexistenz liefert eine Vermutung, die auf um 1955 durchgeführte Rechnungen und Spekulationen des japanischen Mathematikers TANIYAMA zurückgeht und heute je nach Autor mit irgendeiner Kombination der drei Namen TANIYAMA, SHIMURA und WEIL bezeichnet wird. Danach sollte es zu einer elliptischen Kurve  $E$  mit ganzzahligen Koeffizienten eine surjektive Abbildung  $X_0(N) \rightarrow E$  von einer sogenannten Modulkurve  $X_0(N)$  auf  $E$  geben, wobei  $N$  im wesentlichen das Produkt aller Primzahlen  $p$  ist, modulo derer  $E$  keine elliptische Kurve mehr ist. Wie FREYS Rechnungen zeigen, hat seine Kurve vor diesem Hintergrund sehr seltsame Eigenschaften.

Als er damals hier in Mannheim über seine Resultate vortrug, meinte er noch, er glaube nicht, daß die FERMAT-Vermutung so bewiesen werde; er veröffentlichte sein Ergebnis auch nicht in einer der großen internationalen Fachzeitschriften, sondern als Band 1, Heft 1 einer gerade neu gestarteten Schriftenreihe der Universität Saarbrücken, in einfachster Aufmachung xerographiert mit einem nur schwarz-weiß gestalteten Karton als Umschlag:

GERHARD FREY: Links between stable elliptic curves and certain diophantine equations, *Annales Universitatis Saraviensis, Series Mathematicae*, **1** (1), 1986

1987 verschärfte der französische Mathematiker JEAN-PIERRE SERRE die TANIYAMA-Vermutung, und aus dieser stärkeren Vermutung folgt in der Tat, daß die FREY-Kurve nicht existieren kann. Leider ist die SERRESche Vermutung bis heute noch nicht bewiesen.

SERRE erhielt übrigens 2002 den ersten der vom norwegischen Parlament gestifteten ABEL-Preise, die seither zur Erinnerung an den norwegischen Mathematiker NIELS HENRIK ABEL (1802–1829) jedes Jahr in gleicher Weise und gleicher Ausstattung wie die Nobel-Preise für hervorragende Leistungen auf dem Gebiet der Mathematik vergeben werden.

SERRE stellte jedoch noch zusätzlich seine sogenannte  $\varepsilon$ -Vermutung auf,

und auch aus der TANIYAMA-Vermutung zusammen mit der  $\varepsilon$ -Vermutung folgt die Nichtexistenz der FREY-Kurve und damit die FERMAT-Vermutung. Diese  $\varepsilon$ -Vermutung bewies KENNETH RIBET von der Universität Berkeley 1990. Die Grundidee seines Beweises läßt sich interpretieren als eine Art zweidimensionale Version eines Beweises von ERNST EDUARD KUMMER (1810–1893), der die FERMAT-Vermutung 1846 für sogenannte reguläre Primzahlen als Exponenten bewies. (Eine Primzahl  $p$  heißt regulär, wenn die Hauptordnung des Körpers  $\mathbb{Q}[\zeta_p]$  der  $p$ -ten Einheitswurzeln faktoriell ist.) Der Beweis von RIBET ist allerdings erheblich aufwendiger.

Damit war also die FERMAT-Vermutung zurückgeführt auf die TANIYAMA-Vermutung. Diese Vermutung schließlich (die für die weitere mathematische Forschung erheblich wichtiger ist als die FERMAT-Vermutung) bewies WILES 1994.