

# MATHEMATICS MAGAZINE



- Trigonometric Series and Theories of Integration
- Mathematical Certificates
- An Advanced Calculus Approach to Finding the Fermat Point

## EDITORIAL POLICY

The aim of *Mathematics Magazine* is to provide lively and appealing mathematical exposition. This is not a research journal and, in general, the terse style appropriate for such a journal (lemma-theorem-proof-corollary) is not appropriate for an article for the *Magazine*. Articles should include examples, applications, historical background, and illustrations, where appropriate. They should be attractive and accessible to undergraduates and would, ideally, be helpful in supplementing undergraduate courses or in stimulating student investigations. Articles on pedagogy alone, unaccompanied by interesting mathematics, are not suitable. Neither are articles consisting mainly of computer programs unless these are essential to the presentation of some good mathematics. Manuscripts on history are especially welcome, as are those showing relationships between various branches of mathematics and between mathematics and other disciplines.

The full statement of editorial policy appears in this *Magazine*, Vol. 64, pp. 71–72, and is available from the Editor. Manuscripts to be submitted should not be concurrently submitted to, accepted for publication by, nor published by another journal or publisher.

Send new manuscripts to: Martha Siegel, Editor, *Mathematics Magazine*, Towson State University, Towson, MD 21204. Manuscripts should be typewritten and double spaced and prepared in a style consistent with the format of *Mathematics Magazine*. Authors should submit the original and two copies and keep one copy. In addition, authors should supply the full five-symbol Mathematics Subject Classification number, as described in *Mathematical Reviews*, 1980 and later. Illustrations should be carefully prepared on separate sheets in black ink, the original without lettering and two copies with lettering added. Do not use staples.

## AUTHORS

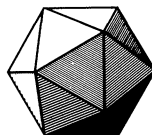
**Alan D. Gluchoff** received his Ph.D. at the University of Wisconsin, Madison, in 1981, and has been teaching at Villanova University since then. His research interests lie mostly in the area of functions of one complex variable. The present article had its origins in a question that nagged at him as he instructed advanced calculus and variables, namely, "Why do we really need the Riemann/Lebesgue Integral?" The discovery that these integrals were tied closely by their creators to trigonometric series was a revelation whose appeal he hopes is shared by others who read the article.

**John Higgins** is currently Professor of Computer Science at Brigham Young University, where he was previously Professor of Mathematics. He received his Ph.D. in Mathematics in 1966 from the University of California and has taught mathematics and computer science at universities in the United States and Japan. Current research interests are in cryptography and computational graph theory. The article on prime certificates was a result of research in cryptography.

**Douglas Campbell** received his Ph.D. from the University of North Carolina in 1971. After joining the Computer Science Department at Brigham Young University in 1982, his interest turned to provably efficient ways to obtain knowledge. He combines this with an interest in algorithms for improving written English. He is author of the program *Thelma Thistleblossom*, an integrated grammar, spelling, style, punctuation, and cliché checker.

Vol. 67 No. 1 February 1994

---



# MATHEMATICS MAGAZINE

## EDITOR

Martha J. Siegel  
*Towson State University*

## ASSOCIATE EDITORS

Donna Beers  
*Simmons College*

Douglas M. Campbell  
*Brigham Young University*

Paul J. Campbell  
*Beloit College*

Underwood Dudley  
*DePauw University*

Susanna Epp  
*DePaul University*

George Gilbert  
*Texas Christian University*

Judith V. Grabiner  
*Pitzer College*

David James  
*Howard University*

Dan Kalman  
*American University*

Loren C. Larson  
*St. Olaf College*

Thomas L. Moore  
*Grinnell College*

Bruce Reznick  
*University of Illinois*

Kenneth A. Ross  
*University of Oregon*

Doris Schattschneider  
*Moravian College*

Harry Waldman  
*MAA, Washington, DC*

## EDITORIAL ASSISTANT

Dianne R. McCann

The *MATHEMATICS MAGAZINE* (ISSN 0025-570X) is published by the Mathematical Association of America at 1529 Eighteenth Street, N.W., Washington, D.C. 20036 and Montpelier, VT, bimonthly except July/August.

The annual subscription price for the *MATHEMATICS MAGAZINE* to an individual member of the Association is \$16 included as part of the annual dues. (Annual dues for regular members, exclusive of annual subscription prices for MAA journals, are \$64. Student and unemployed members receive a 66% dues discount; emeritus members receive a 50% discount; and new members receive a 40% dues discount for the first two years of membership.) The nonmember/library subscription price is \$68 per year.

Subscription correspondence and notice of change of address should be sent to the Membership/Subscriptions Department, Mathematical Association of America, 1529 Eighteenth Street, N.W., Washington, D.C. 20036. Microfilmed issues may be obtained from University Microfilms International, Serials Bid Coordinator, 300 North Zeeb Road, Ann Arbor, MI 48106.

Advertising correspondence should be addressed to Ms. Elaine Pedreira, Advertising Manager, The Mathematical Association of America, 1529 Eighteenth Street, N.W., Washington, D.C. 20036.

Copyright © by the Mathematical Association of America (Incorporated), 1994, including rights to this journal issue as a whole and, except where otherwise noted, rights to each individual contribution. Reprint permission should be requested from Marcia P. Sward, Executive Director, Mathematical Association of America, 1529 Eighteenth Street, N.W., Washington, D.C. 20036. General permission is granted to Institutional Members of the MAA for noncommercial reproduction in limited quantities of individual articles (in whole or in part) provided a complete reference is made to the source.

Second class postage paid at Washington, D.C. and additional mailing offices.

Postmaster: Send address changes to Mathematics Magazine Membership/Subscriptions Department, Mathematical Association of America, 1529 Eighteenth Street, N.W., Washington, D.C. 20036-1385.

PRINTED IN THE UNITED STATES OF AMERICA

---

# ARTICLES

---

## Trigonometric Series and Theories of Integration

ALAN D. GLUCHOFF

Villanova University  
Villanova, PA 19085

### Introduction

This essay was written in response to my experience in teaching the definite integral of a function of one variable. In covering the topic first in freshman calculus, then in advanced calculus, and finally in a course on real variables, these questions would occur to me repeatedly: What is the need for more general theories of integration? To what mathematical use can these theories be put? Is there a way that we can justify the time spent on developing the theories of Riemann and Lebesgue integration, a context in which their development seems natural and answers previously unanswered questions?

The usual justification for the development of these theories is the desire for generality, the tendency to see how far concepts can be pushed in an effort to get at their essence. If this motivating principle is acceptable to both students and instructor then the above questions are perhaps less urgent. But my classroom experiences have suggested to me that something more would be helpful, namely an idea of the necessity for the new concepts within a mathematical framework.

As an example, consider the transition from the Riemann to the Lebesgue integral. The usual view of the inadequacy of the Riemann integral is that the theory is incomplete because one can have a sequence of functions, each of which is Riemann integrable, but whose pointwise limit function is not. From an advanced viewpoint, say, that of function space theory, that is indeed a drawback, but should it be considered as one by a student of advanced calculus? After all, limit functions do not always share the properties of their defining sequences—continuity is an example of this. A related problem is that of reversing the order of integration and summation, and of convergence theorems in general: Why all the concern over the weakest conditions that assure that  $\lim_{n \rightarrow \infty} \int_a^b f_n = \int_a^b \lim_{n \rightarrow \infty} f_n$ ? Of what use are these theorems? (Some answers will be given later in this paper.)

In doing some background preparation for a seminar in Fourier series, I learned about their history, and I believe that the connection with trigonometric series can be useful in putting these integration issues into a context. The three main integration theories that we teach can roughly be attributed to Cauchy, Riemann, and Lebesgue, and each of these men was very interested in trigonometric series. Cauchy devoted a paper to the subject [7], Riemann introduced his integral in a classic on the series, [19], and Lebesgue devoted several papers and a short book to them [15], [16]. (Books [3], [12], and [14] are excellent general references for this material.)

This paper is an attempt to show how we may pedagogically relate integration theories and trigonometric series. Let us begin with a simple question seemingly unrelated to integration: Is there a closed-form analytic representation for an arbitrary real-valued function defined on  $[-\pi, \pi]$ ? That is, can every function defined on this

interval be expressed as an algebraic combination of well-known elementary functions like polynomials, trigonometric functions, or exponentials, or limits of such combinations? The Taylor series comes to mind immediately, but, in the charming words of one writer [21] "... the power series obtained by combining terms of the form  $c_n x^n$ " define the most civilized members of mathematical society—the so-called analytic functions—which are most orderly in their behavior, being continuous throughout their domains, possessing derivatives of all orders." Thus, for instance, discontinuous functions cannot be expressed by power series. We then turn to trigonometric series, that is, series of the form  $\sum_{k=0}^{\infty} A_k \cos kx + B_k \sin kx$ ; see FIGURE 1. It is reasonable to suspect that by combining functions of this sort, where the higher frequency sines and cosines are very choppy, we may obtain functions that are highly irregular. What can we say about the class of functions that are expressible by trigonometric series?

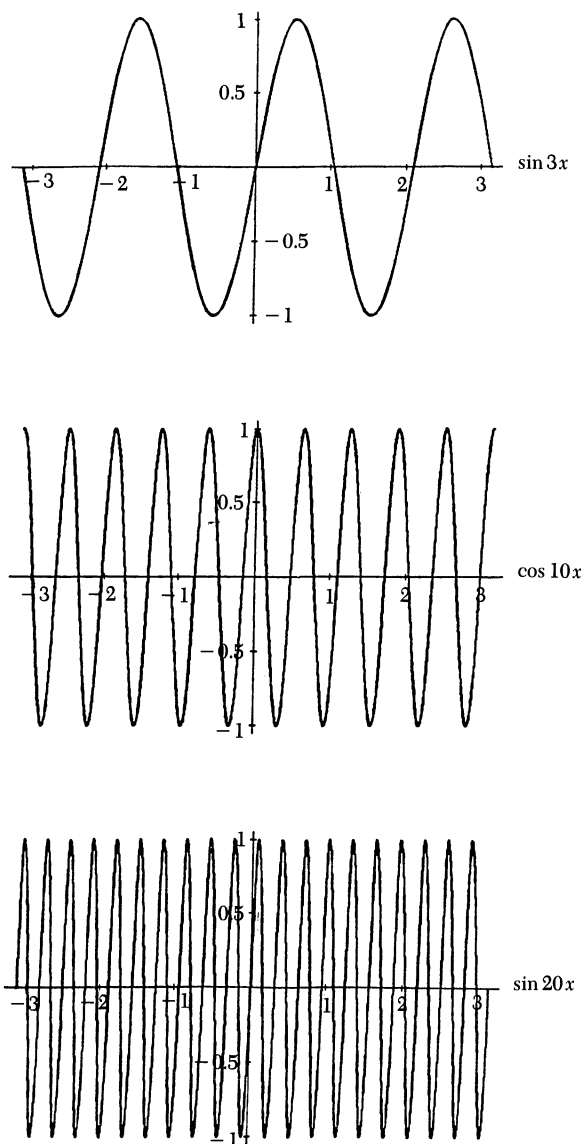


FIGURE 1

A judicious use of specific series and their computer graphics plots would confirm that many kinds of behavior occur with these series: discontinuous, very smooth, and fractal-like graphs can be produced; see FIGURE 2 (a) through (c). The last example in the figure is one of an interesting set of functions, investigated by Weierstrass in 1872, given by  $f(x) = \sum_{n=0}^{\infty} b^n \cos(a^n x)$  for  $0 < b < 1$ ,  $a$  a positive integer. These functions are continuous for all  $x$  but differentiable *nowhere* if  $ab > 1$ . The continuity follows easily from the Weierstrass M-test, and nondifferentiability, though not elementary to prove, is at least plausible when one looks at the partial sums in FIGURE 2(c): The frequencies of the cosines increase in such a pattern that their superposition forces the monotone portions of the graphs to become increasingly vertical, while greater oscillation is introduced with each additional term. (The mathematician Paul du Bois-Reymond, who published Weierstrass' proof of nondifferentiability in 1874, considered these functions "equally too strange for immediate perception as well as for critical understanding".) The "density" in the plane of the graphs of similar series led Besicovich and Ursell in 1937 to investigate the "dimension" of these graphs using a generalization of ordinary Euclidean dimension to fractional numbers. Some graphs turn out to have a fractional dimension between one and two! See [18].

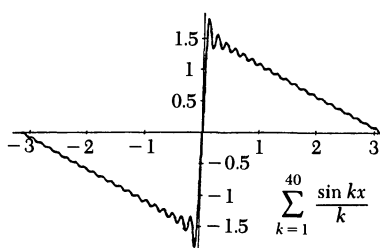
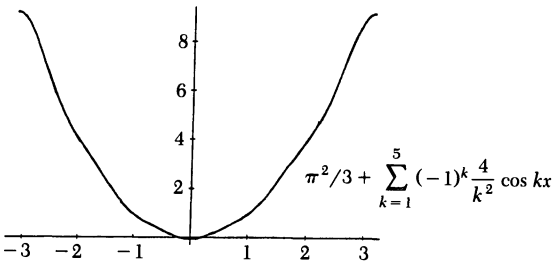
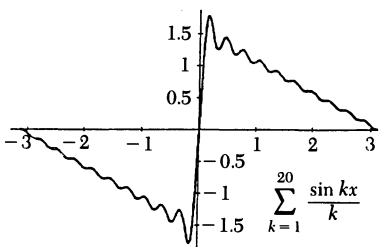
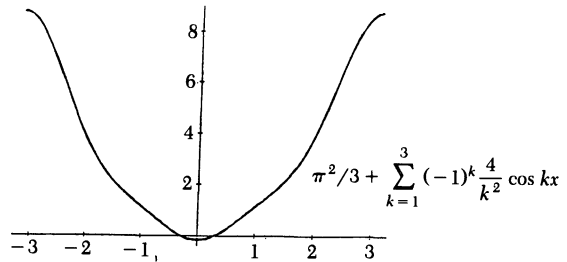
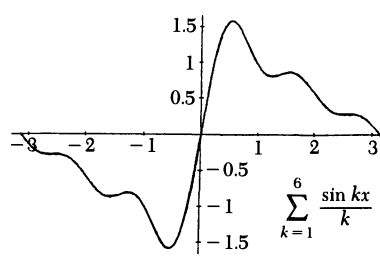


FIGURE 2a

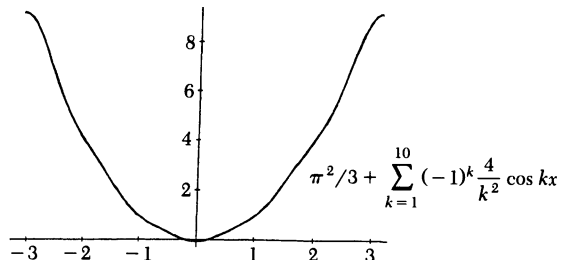


FIGURE 2b

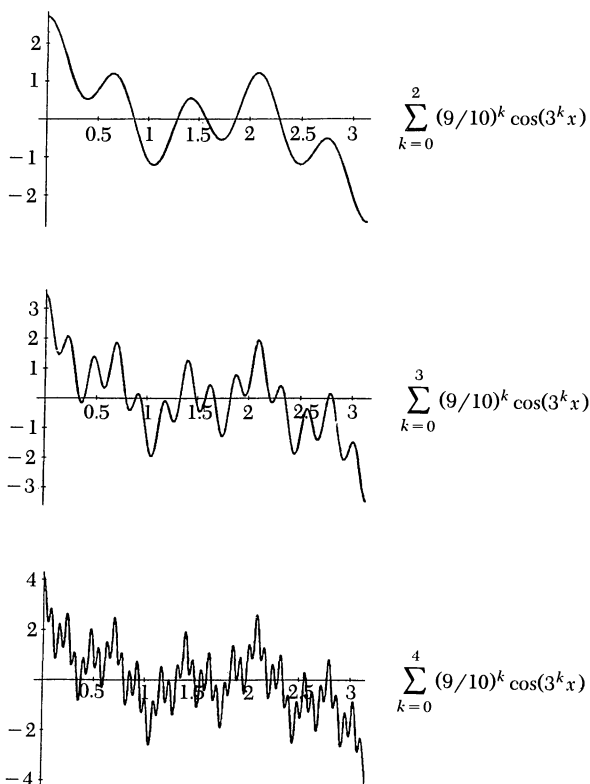


FIGURE 2c

Since such a wide variety of behavior seems possible with these series, the need for a rigorous mathematical investigation into the representation of an arbitrary function is now felt, and we can address it by beginning with a function and trying to find a trigonometric representation of it. This approach begins to involve us in theories of integration as will be seen in the following section, where we start with a naïve conception of the integral and review the classical argument of Fourier and Euler to arrive at the Fourier series. It is questions about the meaning of the coefficients in this series, the generality of representation theorems, and the uniqueness of Fourier series representation that motivate the refinement of the integral concept presented in succeeding sections. My hope is that seeing the inadequacies of each theory in this context will make the motivation for the next theory seem compelling.

This presentation will be patterned after historical development, and we will often stop to introduce other ways in which trigonometric series influenced the development of concepts in analysis, but this is not a recapitulation of history. The job of disentangling the precise notions of integral, continuity, and convergence that earlier mathematicians had and distinguishing them from their modern counterparts is a formidable, though fascinating task, but is one better left to the above mentioned references, in which students may find a more complete story.

In what follows, all functions  $f$  will be defined on  $[-\pi, \pi]$  and are assumed *bounded*, for simplicity. In each section we assume familiarity with the definition of the appropriate integral and the main theorems concerning that theory, but the reader should imagine himself at different levels of sophistication in each section: in Section



I, in freshman calculus having done integration theory; in Section II, in advanced calculus having covered theorems on continuous functions and uniform convergence; in Section III, in advanced calculus having covered the Riemann integral (knowledge of bounded variation functions would also be useful here); and in Section IV, having completed a study of measure theory and the Lebesgue integral on the line.

## I. Representation of Functions by Trigonometric Series— Some Basic Questions

We begin by assuming only a basic knowledge of calculus, say at the college freshman level, and deal with elementary functions. What is our procedure for representing an elementary function by trigonometric series?

We follow the classical argument given by Fourier in 1822 [13] using the orthogonality of the family of sines and cosines. Let us consider, for example,  $f(x) = x^2$  on the interval  $[-\pi, \pi]$ , assume that it can be represented by a trigonometric series and attempt to solve for the coefficients. If  $x^2 = a_0/2 + \sum_{k=1}^{\infty} (a_k \cos kx + b_k \sin kx)$  then we can determine any  $a_n$  as follows: Multiply both sides of the equation by  $\cos nx$ , and integrate both sides from  $-\pi$  to  $\pi$ :

$$\int_{-\pi}^{\pi} x^2 \cos nx \, dx = \int_{-\pi}^{\pi} \left[ (a_0/2) \cos nx + \sum_{k=1}^{\infty} (a_k \cos kx \cos nx + b_k \sin kx \cos nx) \right] dx.$$

Now assuming that, on the right, the integral of the sum is the sum of the integrals, we have

$$\begin{aligned} \int_{-\pi}^{\pi} x^2 \cos nx \, dx &= a_0/2 \int_{-\pi}^{\pi} \cos nx \, dx \\ &+ \sum_{k=1}^{\infty} \left[ a_k \int_{-\pi}^{\pi} \cos kx \cos nx \, dx + b_k \int_{-\pi}^{\pi} \sin kx \cos nx \, dx \right]. \end{aligned}$$

By using some trigonometric identities one can prove that

$$\int_{-\pi}^{\pi} \cos kx \cos nx \, dx = \begin{cases} 0 & k \neq n \\ \pi & k = n \end{cases}$$

and

$$\int_{-\pi}^{\pi} \sin kx \cos nx \, dx = 0 \text{ for all } n \text{ and } k.$$

(See [20], pp. 11–12 for details of this.) Thus we have

$$a_n = \frac{1}{\pi} \int_{-\pi}^{\pi} x^2 \cos nx \, dx = \frac{(-1)^n 4}{n^2}$$

after two integrations by parts.

Similarly,  $b_n$  can be found by multiplying both sides by  $\sin nx$ ; the result is  $b_n = 0$  for all  $n$ . Putting all this together gives the tentative conclusion that

$$x^2 = \pi^2/3 + \sum_{k=1}^{\infty} (-1)^k \frac{4}{k^2} \cos kx.$$

How do we check the validity of the argument? A nonrigorous approach would be to plot partial sums  $\pi^2/3 + \sum_{k=1}^n (-1)^k (4/k^2) \cos kx$  and verify that they approach  $x^2$  on  $[-\pi, \pi]$ . See FIGURE 2 (b). More examples, computations, and plots will validate this procedure for other elementary functions.

Does this idea work for any  $f$  defined on  $[-\pi, \pi]$ , not necessarily elementary? Again, we write  $f(x) = a_0/2 + \sum_{k=1}^{\infty} (a_k \cos kx + b_k \sin kx)$ , and repeat the above procedure to find that

$$a_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos nx \, dx \quad \text{and} \quad b_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin nx \, dx. \quad (1)$$

With these numbers as coefficients, we then need to check somehow that the series converges to  $f(x)$ .

We must now ask if this argument is valid. It depends on these points:

- We first assumed what was to be proved, namely that a function defined on  $[-\pi, \pi]$  has a representation of the required form; there is no a priori reason to believe this is true.
- We assumed that the integrals in (1) *have meaning*; that they can either be evaluated by the Fundamental Theorem of Calculus, or have independent meaning as definite integrals.
- We assumed that we can integrate an infinite series by integrating each term and adding up the results. This is surely true for a finite series, but is it true for an infinite series?
- We then assumed that the series converges to  $f(x)$ ; this has not been *proved* at all, even in the elementary examples.

Thus the meaning of the symbol  $\int_a^b f(x) \, dx$  is related in an unexpected way to the representation of a function  $f$  by a trigonometric series: One candidate for such a representation has (1) as its coefficients. So, if this candidate is going to have any chance at being a representation, we must find an explicit meaning for those integrals.

Fourier himself had little problem accepting the above argument; his attitude seemed to be that if one can solve for the above coefficients in any manner, then the resulting representation must be valid ([14], p. 7). But the points (a)–(d) above constitute the modern criticism of the argument and must be answered before we accept the conclusion. In honor of Fourier's work with these series, the coefficients (1) are called the Fourier coefficients of  $f$ , and the series formed using these coefficients is called the Fourier series associated with  $f$ . As mentioned above, it is not a priori clear that the Fourier series for  $f$  actually converges to  $f$ ; nor is it clear that it is the only such trigonometric series that might do so.

To return to point (b), how do we assign a meaning to the symbols in (1)? In our example  $f(x) = x^2$ , we evaluated the integrals by the Fundamental Theorem of Calculus; this procedure, if it were to generalize, would require showing that every function  $f$  has an antiderivative. But is this obvious? An alternative is to go back to the definition of  $\int_a^b f(x) \, dx$  as a limit of sums, and make sure that *it* is meaningful: Is it true that every function  $f$  defined on  $[-\pi, \pi]$  has a well-defined number  $\lim_{\Delta x \rightarrow 0} \sum_{i=1}^n f(x_i)(x_i - x_{i-1})$  associated with it?

In (c) we see an impetus to exploring the question of term-by-term integration of an infinite series. This topic occurs in another related way. Suppose we ask whether there is more than one way to represent a given function by trigonometric series, i.e., suppose

$$f(x) = \sum_{k=0}^{\infty} A_k \cos kx + B_k \sin kx = \sum_{k=0}^{\infty} C_k \cos kx + D_k \sin kx \quad (2)$$

for all  $x$  in  $[-\pi, \pi]$ . Must  $A_k = C_k$  and  $B_k = D_k$  for all  $k$ ? This question is of interest not just for “theoretical” reasons but for a “practical” one—if one is to use these expansions in any computations it is important to know whether there is more than one representation for a given function. One approach to answering this question is again to multiply both sides by  $\cos nx$  and integrate over  $[-\pi, \pi]$ . *Assuming we may integrate term-by-term*, we get  $A_n = C_n$ ;  $B_n = D_n$  to follow similarly. Thus it would appear that only one such representation is possible, but again the argument assumes that the integrals all have meaning, and that term-by-term integration is possible, points (b) and (c).

This question of uniqueness of representation has an interesting historical sideline. In the early 1870s Georg Cantor concerned himself with this matter, and in doing so was led to considerations about sets of points on the real line in the following way. Cantor proved that if (2) holds for all  $x$ , then indeed  $A_n = C_n$  and  $B_n = D_n$  for all  $n$ . His proof avoided the above questionable integrations. He was then able to obtain the same conclusion if (2) holds for all but a finite number of points, the so-called exceptional set. This led to the natural question of whether uniqueness follows if (2) holds for all but an *infinite* exceptional set; in 1872 he was able to show that this is true provided that the exceptional sets are distributed in certain ways, see [3], pp. 275–280. Point-set theory was at this time so undeveloped that Cantor had to spend much of his energies proposing a theory of real numbers and point sets, including notions of Cauchy sequences, accumulation points, and a Bolzano-Weierstrass-type theorem, in order to prove his theorems. Subsequent to this work Cantor’s interests in fact remained with transfinite set theory and the theory of the continuum of real numbers. See [8] for a nice account of this.

## II. The Cauchy Integral and Continuous Functions

Our first theory of integration uses the definition of the integral given by Cauchy [6] in 1823, *where the function to be integrated is assumed continuous on  $[a, b]$* . “As early as 1814 he [Cauchy] came to recognize the importance of continuity in the modern sense and of the sum conception of the definite integral. After 1818 when Cauchy learned of Fourier’s work—particularly Fourier’s interpretation of the coefficients  $a_n$  and  $b_n$  (see (1)) as integrals—he was undoubtedly all the more convinced of the importance of these notions” [14], p. 9. See also [11], chapter 6, for this section.

Cauchy’s integral is defined by the usual limit of sums, which, by virtue of uniform continuity of the integrand, is guaranteed to exist. (It is interesting to note that in his proof of the existence of his integral for continuous functions, Cauchy did not appreciate the fact that *uniform* continuity, not merely point continuity, of the function was the property needed to guarantee convergence of the sums to a limit. He apparently did not realize the distinction between the concepts at this time.) This is the integral that students see in a freshman calculus course, although rigorous proofs of its convergence are perhaps not given until advanced calculus. It has some advantages over the notion in the preceding section. It can be shown to exist for the class of continuous functions, and using it one can prove that any  $g$  continuous on  $[a, b]$  has an antiderivative  $G$ : Just define  $G(x) = \int_a^x g(t) dt$ . Thus there are two choices for the values in (1), either the limiting sums or the value  $F(b) - F(a)$ , for  $F$  an antiderivative of, say,  $f(x) \cos nx$ . The Fundamental Theorem can be rigorously proved using the Cauchy integrals; these facts were all demonstrated in [6].

Theorems involving this integral allow a proof showing, that, under certain mild restrictions on  $f$ , the Fourier series for  $f$  converges to  $f$  at all points of  $[-\pi, \pi]$ . We

will outline such a proof now; it follows the argument first used by Bonnet in 1849 [2], who reworked Dirichlet's earlier treatment of 1829 [9] by using integral mean-value theorems. Dirichlet's paper represents the first rigorous treatment of conditions assuring the convergence of the Fourier series of  $f$  to  $f$ . The proof given here should be accessible to advanced calculus students and would ideally be introduced prior to the Riemann integral, to give a sense of analysis using only the "theory of the continuous". This proof, although perhaps difficult at first, is worth studying for its use of integration facts and techniques like a) and c) below. Such facts unfortunately often are presented as isolated examples in analysis books, unrelated to a mathematical goal.

We assume that  $f(\pi) = f(-\pi)$ , thus  $f$  can be extended periodically on the real line by  $f(x + 2\pi) = f(x)$ ; the extended function is continuous for all  $x$ . We also assume that  $f$  has only finitely many maxima and minima in  $[-\pi, \pi]$ , and will prove that  $f$  is representable by its Fourier series. Now let  $x \in [-\pi, \pi]$  be fixed, define

$$s_n(x) = a_0/2 + \sum_{k=1}^n (a_k \cos kx + b_k \sin kx) \quad (3)$$

where  $a_k$  and  $b_k$  are as in (1). Then

$$\begin{aligned} s_n(x) &= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(y) dy \\ &\quad + \frac{1}{\pi} \sum_{k=1}^n \left[ \left( \int_{-\pi}^{\pi} f(y) \cos ky dy \right) \cos kx + \left( \int_{-\pi}^{\pi} f(y) \sin ky dy \right) \sin kx \right] \\ &= \frac{1}{\pi} \int_{-\pi}^{\pi} f(y) \left[ \frac{1}{2} + \sum_{k=1}^n \cos k(y-x) \right] dy. \end{aligned}$$

An elementary argument with trigonometric identities ([20], p. 71) shows that we may write this as

$$s_n(x) = \frac{1}{\pi} \int_{-\pi}^{\pi} f(y) \frac{\sin[(n+1/2)(y-x)]}{2 \sin[(1/2)(y-x)]} dy.$$

A change of variable  $t = (1/2)(y-x)$  and an adjustment of the limits of integration gives

$$s_n(x) = \frac{1}{\pi} \int_{-\pi/2}^{\pi/2} f(x+2t) \frac{\sin(2n+1)t}{\sin t} dt.$$

Our argument will show

$$\lim_{n \rightarrow \infty} \frac{1}{\pi} \int_0^{\pi/2} f(x+2t) \frac{\sin(2n+1)t}{\sin t} dt = \frac{f(x)}{2}; \quad (4)$$

a similar argument will show that the integral over  $[-\pi/2, 0]$  approaches the same limit as well, giving the desired result. As a final preliminary step, it can be shown that, since  $\sin t$  "behaves like  $t$ " over  $[0, \pi/2]$ , we may replace (4) with

$$\lim_{n \rightarrow \infty} \frac{1}{\pi} \int_0^{\pi/2} f(x+2t) \frac{\sin(2n+1)t}{t} dt = \frac{f(x)}{2}. \quad (5)$$

What is our next step? Intuitively, we proceed as follows: Since  $f$  is continuous at  $x$ , then near 0,  $f(x+2t)$  "looks like"  $f(x)$ . We may then get the idea to break the

integral in (5) into two pieces: one over  $[0, \delta]$ , where  $f(x + 2t) \approx f(x)$ , and one over  $[\delta, \pi/2]$ :

$$\begin{aligned} \frac{1}{\pi} \int_0^{\pi/2} f(x + 2t) \frac{\sin(2n + 1)t}{t} dt &\approx \frac{1}{\pi} f(x) \int_0^{\delta} \frac{\sin(2n + 1)t}{t} dt \\ &+ \frac{1}{\pi} \int_{\delta}^{\pi/2} f(x + 2t) \frac{\sin(2n + 1)t}{t} dt. \end{aligned}$$

As  $n \rightarrow \infty$ , we will show that the second integral on the right vanishes and the first approaches  $\pi/2$ . Let us list some facts needed for the demonstration.

A change of variables in the first integral on the right above shows that knowledge of  $\int_0^{\infty} (\sin t/t) dt$  is needed:

$$(a) \quad \lim_{M \rightarrow \infty} \int_0^M \frac{\sin x}{x} dx = \int_0^{\infty} \frac{\sin x}{x} dx = \pi/2.$$

This can be proved in several ways: See [5], p. 135 or p. 202. Furthermore,

(b) if  $0 < p < q$ , we have  $|\int_p^q (\sin x/x) dx| < \pi$ . This follows from an analysis of the curve  $(\sin x)/x$  and (a) ([5], p. 222). Looking at the second integral on the right we see that it would be nice to be able to remove  $f(x + 2t)$  and focus on an integral expression related to  $\int_0^{\infty} (\sin x/x) dx$  again.

Our method for doing this is using (c) a mean-value theorem for integrals whose earliest form is due to Bonnet [2]. Let  $\phi$  and  $\psi$  be continuous on  $[a, b]$  with  $\phi$  monotonic and  $\psi$  changing signs no more than a finite number of times on  $[a, b]$ . Then there is a  $\xi$ ,  $a \leq \xi \leq b$  with  $\int_a^b \phi(x)\psi(x) dx = \phi(a)\int_a^{\xi} \psi(x) dx + \phi(b)\int_{\xi}^b \psi(x) dx$ . This is not the most general form of this theorem; but it has the advantage that its proof ([5], p. 105–108) involves only the intermediate-value property of continuous functions and other elementary ideas, thus staying within the confines of our Cauchy integral theory for continuous functions.

Here is the remainder of the proof:  $[0, \pi/2]$  may be partitioned by  $\mathcal{P} = \{a_k\}$  with  $f(x + 2t)$  monotone in  $t$  in each  $[a_{k-1}, a_k]$ . Let  $\varepsilon > 0$  be given, then there is a  $\delta > 0$  with  $0 < \delta < a_1$ , and  $|f(x + 2t) - f(x)| < \varepsilon/4$  if  $|t| < \delta$ . We now write

$$\begin{aligned} &\left| \frac{1}{\pi} \int_0^{\pi/2} f(x + 2t) \frac{\sin(2n + 1)t}{t} dt - \pi \frac{f(x)}{2} \right| \\ &\leq \left| \frac{1}{\pi} \int_0^{\delta} [f(x + 2t) - f(x)] \frac{\sin(2n + 1)t}{t} dt \right| \\ &\quad + \left| \frac{1}{\pi} f(x) \int_0^{\delta} \frac{\sin(2n + 1)t}{t} dt - f(x) \frac{\pi}{2} \right| + \left| \frac{1}{\pi} \int_{\delta}^{a_1} f(x + 2t) \frac{\sin(2n + 1)t}{t} dt \right| \\ &\quad + \left| \frac{1}{\pi} \sum_{k=2}^n \int_{a_{k-1}}^{a_k} f(x + 2t) \frac{\sin(2n + 1)t}{t} dt \right| = \text{I} + \text{II} + \text{III} + \text{IV}. \end{aligned}$$

A change of variables in II shows that, as  $n \rightarrow \infty$ ,  $\text{II} \rightarrow 0$  by (a). Applying (c) to III gives

$$\begin{aligned} \int_{\delta}^{a_1} f(x + 2t) \frac{\sin(2n + 1)t}{t} dt &= f(x + 2\delta) \int_{\delta}^{\xi(n)} \frac{\sin(2n + 1)t}{t} dt \\ &\quad + f(x + 2a_1) \int_{\xi(n)}^{a_1} \frac{\sin(2n + 1)t}{t} dt, \end{aligned}$$

where  $\delta < \xi(n) < a_1$ . A change of variable in each integral on the right and an appeal

to the convergence of  $\int_0^\infty (\sin t/t) dt$  shows that, as  $n \rightarrow \infty$ ,  $\text{III} \rightarrow 0$ . The same argument applied to each summand in IV shows that it, too, goes to 0 with increasing  $n$ . Finally, using (c) on I gives

$$I = \left| \frac{1}{\pi} [f(x + 2\delta) - f(x)] \int_{\xi(n)}^\delta \frac{\sin(2n+1)t}{t} dt \right|,$$

and so by (b) and a change of variable, this gives  $|I| < \varepsilon/4$  for all  $n$ . Thus for  $n$  sufficiently large,

$$\left| \frac{1}{\pi} \int_0^{\pi/2} f(x+2t) \frac{\sin(2n+1)t}{t} dt - \pi \frac{f(x)}{2} \right| < \varepsilon/4 + \varepsilon/4 + \varepsilon/4 + \varepsilon/4,$$

and we are done.

Thus, from a theory of integration for continuous functions we get a theorem that proves the representability of certain functions by trigonometric series, namely the Fourier series. Our earlier intuition about the general approach of Fourier series is borne out and the validity of our example  $f(x) = x^2$  is confirmed.

It is now natural to ask if the proof can be extended to more general, and especially discontinuous,  $f$ . The requirement that  $f$  have only finitely many maxima and minima is an obvious target for removal. Note that we needed this assumption so that we could break the interval of integration up into subintervals on which  $f$  is monotone and apply the mean-value theorem on each subinterval. It is not at all obvious how to modify the proof for a continuous function with infinitely many maxima and minima in an interval—for instance,  $f(x) = x \sin(1/x)$ . A different problem is posed by the removal of the continuity restriction on  $f$ , since *removing the continuity restriction on  $f$  requires a new definition of the integral*. This is easy if there is only a finite number of discontinuities: If  $f$  is discontinuous only at  $a = x_0 < x_1 < x_2 < x_3 < \cdots < x_n = b$ , then one naturally defines  $\int_a^b f(x) dx$  as  $\sum_{i=1}^n \int_{x_{i-1}}^{x_i} f(x) dx$ . A modification of the proof is possible in this case (assuming, in addition, the extrema restriction on each subinterval) giving convergence of  $s_n(x)$  to  $\frac{1}{2}[\lim_{x \rightarrow x_i^+} f(x) + \lim_{x \rightarrow x_i^-} f(x)]$  at each  $x_i$ . This was the case considered by Dirichlet in [9]. The next logical step is to ask whether the proof can be extended to a function with infinitely many discontinuities in  $[-\pi, \pi]$ . The “worst possible scenario” is that these discontinuities are densely distributed. The function

$$\chi_Q(x) = \begin{cases} 0 & x \text{ irrational} \\ 1 & x \text{ rational} \end{cases}$$

shows that this may happen. In fact  $\chi_Q$  was first introduced by Dirichlet as an instance of a function not satisfying the above criteria for convergence of Fourier series. “The function so defined has finite and determinate values for every value of  $x$ , and yet one does not know how to substitute it in the series, *seeing that the various integrals [(1)] that enter into this series will lose all meaning in this case*” [9, p. 132]. Again, a more general definition of the integral is necessary.

In his doctoral thesis in 1864 Rudolf Lipschitz attempted such extensions of the integral to more general classes of functions in an effort to generalize Dirichlet’s results (Dirichlet himself attempted such an extension earlier), but in doing so was led into errors involving point-set topology on the line. It remained for Riemann to propose a new definition that would cover these cases.

Finally, let us note that it is within this theory that we begin to get some answers to the other questions raised earlier, that of the uniqueness of representation by trigonometric series. If (2) holds where both series converge *uniformly* on  $[-\pi, \pi]$ ,

then both sums are continuous, their integrals are defined, and term-by-term integration is possible. The argument that follows (2) then shows  $A_k = C_k$ ,  $B_k = D_k$  for all  $k$ , so there is at most one way of representing a function by a *uniformly* convergent trigonometric series. The same argument shows that, if a continuous function  $f$  is represented by a uniformly convergent trigonometric series, then the coefficients of  $f$  must be the Fourier coefficients. [This is the “proof” from section I.] These arguments fail, however, if the series do not converge uniformly; the series cannot even be assumed to be integrable in this case. In section IV we will see how this situation can be remedied.

The concept of uniform convergence itself has an interesting history that involves trigonometric series in a different way. Cauchy included in his *Cours d'analyse* a theorem stating that if a series of functions continuous at a point  $x$  converges to a limit function, then the limit function is also continuous at  $x$ . This was shown to be false by Abel in 1826, and his counterexample was the trigonometric series  $\sum_{n=1}^{\infty} (-1)^{n+1} (\sin(nx)/n)$ , which has a discontinuity similar to that in FIGURE 1(a). Abel also showed that certain series of continuous functions converge to continuous limits. But in these examples uniform convergence was present, though Abel did not realize this. Weierstrass in 1841 was the first to understand the importance of this concept, and he demonstrated the validity of term-by-term integration of uniformly convergent series, the device used in the preceding paragraph. The existence of non-uniformly convergent Fourier series leads also to the possibility of a given function being representable by different trigonometric series, since the preceding uniqueness argument is not valid. This was observed by Heine in 1870. He in turn induced Cantor to turn his attention to the series, with results noted earlier.

### III. Riemann's Theory of Integration

It is in the context of representation of functions by trigonometric series that Riemann introduced his definition of what is now known as the Riemann integral. In “Über die Darstellbarkeit einer Function durch eine trigonometrische Reihe” [19] in 1854 (published 1867) he states: “The uncertainty which still prevails in several fundamental points of the theory of definite integrals obliges us to set down a few things about the idea of the definite integral and the extent of its validity.” Presumably he had in mind the difficulties of Dirichlet in extending Cauchy's integral. He then proceeds to define the integral and develop a necessary and sufficient condition that an integrable function be representable by a trigonometric series.

As is usually pointed out in advanced calculus courses, the main difference between Riemann's definition and Cauchy's is that the function to be integrated is not necessarily assumed to be continuous; the definition simultaneously defines a class for which integration is possible—those for which the limiting sums  $\sum_i f(x_i)(x_i - x_{i-1})$  converge in the required way—and gives the value of the integral as the limit of these sums. If  $f$  is continuous, the integrals agree in value.

One immediate result is that there are many well-known functions that are integrable in the new sense but not in the old, for example,  $f(x) = \sum_{x_n < x} 1/2^n$  where  $x_n$  is an enumeration of the rationals in  $[a, b]$  is monotone, hence integrable. Note this function has discontinuities precisely at the rationals, a dense set. Riemann himself gave a similar example of an integrable function with densely distributed discontinuities [19, p. 242]. These examples indicate that even a “pathological” function now stands a chance at being representable by a trigonometric series since the integrals in (1) are well defined for any Riemann integrable  $f$ .

And in fact many such functions can now be proved to be so representable. A more general version of the mean-value theorem for integrals mentioned in the preceding section is available: The conclusion stated there is valid if  $\phi$  is monotone only and  $\psi$  is Riemann integrable. The appropriate generalization of the class of functions treated in section II is the class of functions of bounded variation. These functions are Riemann integrable since they are differences of monotone functions [5, p. 80]. It can be shown that if  $f$  is such a function, the Fourier series of  $f$  converges to  $\frac{1}{2}[\lim_{t \rightarrow x^+} f(t) + \lim_{t \rightarrow x^-} f(t)]$  at each  $x$ . The proof method is a straightforward adaptation of that in II, with all the integrals understood in the Riemann sense. The proof uses the decomposition of a bounded variation function into a difference of monotones, to which the mean-value theorem is applied. In particular, this device begins to address the earlier problem of dealing with continuous functions having infinitely many maxima and minima, since each a function will be representable by its Fourier series provided that it is of bounded variation. Of course, the concept of bounded variation could have been introduced for continuous functions only and this difficulty overcome earlier, but it is interesting that historically this concept was introduced by Jordan in 1881, well after Riemann's theory was developed. Perhaps the focus on the integrability of arbitrary monotone functions shifted attention to this more general class at this time. Interestingly, the class first appeared in an attempt to extend Dirichlet's proof of convergence to a wider class of functions.

The criteria for Riemann integrability can be stated in many equivalent ways. Let us single out one:  $f$  can be shown to be Riemann integrable if, and only if, there are two sequences of step functions  $g_n, h_n$  with  $g_n \leq f \leq h_n$  on  $[a, b]$  and  $\int_a^b (h_n - g_n) \rightarrow 0$  as  $n \rightarrow \infty$  [17, p. 185]. Thus in some sense  $f$  is Riemann integrable if, and only if, it can be well approximated by step functions. The idea of step-function approximation of integrable functions is a useful tool in proving theorems about them, including theorems about Fourier series. One well-known result, the Riemann-Lebesgue Lemma, states that if  $f$  is Riemann integrable, then  $\{a_k\}$  and  $\{b_k\}$  in (1) must tend to 0 as  $k \rightarrow \infty$ . This follows easily if  $f$  is a step function, since  $f(x) = \sum_{i=1}^n c_i \chi_{(x_{i-1}, x_i)}$  where

$$\chi_{(x_{i-1}, x_i)} = \begin{cases} 1 & x \in (x_{i-1}, x_i) \\ 0 & \text{elsewhere} \end{cases}$$

implies, for example, that

$$b_k = \sum_{i=1}^n c_i \int_{x_{i-1}}^{x_i} \sin kx \, dx = \sum_{i=1}^n c_i \left[ \frac{\cos kx_{i-1} - \cos kx_i}{k} \right],$$

which  $\rightarrow 0$  as  $k \rightarrow \infty$ . The case for general  $f$  then follows by the approximation by step functions. The use of step function approximation may be considered a product of the style of thinking involved in Riemann integration; although this kind of argument is possible with the more elementary Cauchy integral, it perhaps is more natural in the Riemann setting.

Other theorems describing conditions for convergence of Fourier series at a point, as opposed to all points in an interval, can be proved using the Riemann-Lebesgue Lemma. For example, we know that  $s_n(x) \rightarrow f(x)$  if, and only if,

$$\lim_{n \rightarrow \infty} \frac{1}{\pi} \int_{-\pi}^{\pi} [f(x+u) - f(x)] \frac{\sin(n+1/2)u}{u} \, du = 0.$$

(See [20], pp. 75–76 for details, which are similar to the work in section I.) Letting  $\delta > 0$  be small, we may split the integral into portions on  $(-\delta, \delta)$ ,  $[\delta, \pi]$ , and  $[-\pi, -\delta]$ , and use the Riemann-Lebesgue Lemma to show that, as  $n \rightarrow \infty$ , the latter



two pieces vanish for any  $f$ . If  $f$  is assumed to be, say, differentiable at  $x$ , then the remaining integral over  $(-\delta, \delta)$  also approaches 0 as  $n \rightarrow \infty$  by the Lemma. Differentiability can be replaced by weaker conditions such as Lipschitz continuity, yet another concept introduced specifically to deal with criteria for convergence of Fourier series. (It appeared in Lipschitz's previously mentioned 1864 work.)

With all these advances, it is perhaps difficult to see why the Riemann integral is not sufficient to deal with the problems involved with trigonometric series representation. But there are problems remaining. There are functions that still are not integrable even in this new sense;  $\chi_Q$  is one. We might wish to know about representation for these functions. Also, the new theory offers no new theorems to help us deal easily with the uniqueness questions at the end of section I. This is a special case of the issue of integrability of limits of Riemann integrable functions and the general questions of reversal of limiting operations raised in the introduction. We give two more examples of how this issue may arise.

For  $f$  bounded and integrable, and  $s_n$  as in (3) straightforward manipulations of the integral expressions involved lead to the identity

$$\frac{1}{\pi} \int_{-\pi}^{\pi} f^2 = a_0^2/2 + \sum_{k=1}^n (a_k^2 + b_k^2) + \frac{1}{\pi} \int_{-\pi}^{\pi} (f - s_n)^2.$$

What happens as  $n \rightarrow \infty$ ? Since the integral on the right is positive, we have

$$a_0^2/2 + \sum_{k=1}^n a_k^2 + b_k^2 \leq \frac{1}{\pi} \int_{-\pi}^{\pi} f^2,$$

so the series  $a_0^2/2 + \sum_{k=1}^{\infty} a_k^2 + b_k^2$  is convergent. Now if  $f$  is sufficiently nice throughout  $[-\pi, \pi]$  we might have  $s_n(x) \rightarrow f(x)$  for all  $x$ , and if we were allowed to conclude from this that  $\int_{-\pi}^{\pi} (f - s_n)^2 \rightarrow 0$ , we would have a nice identity

$$\frac{1}{\pi} \int_{-\pi}^{\pi} f^2 = a_0^2/2 + \sum_{k=1}^{\infty} a_k^2 + b_k^2. \quad (6)$$

Can we say that  $\int_{-\pi}^{\pi} (f - s_n)^2 \rightarrow 0$ ? Equation (6) is in fact true; it is called Parseval's identity. For example, choosing  $f(x) = x^2$  as in section I gives  $\sum_{n=1}^{\infty} 1/n^4 = \pi^4/90$ !

Another example: Suppose  $f$  is not Riemann integrable on  $[-\pi, \pi]$ . We may attempt a trigonometric series representation of  $f$  as follows. Motivated by the idea that integrable functions are approximable by step functions, we may try to obtain a sequence of step functions  $f_n$  converging to  $f$ , and define  $a_k = \lim_{n \rightarrow \infty} \int_{-\pi}^{\pi} f_n(x) \cos kx \, dx$ ,  $b_k = \lim_{n \rightarrow \infty} \int_{-\pi}^{\pi} f_n(x) \sin kx \, dx$  for the coefficients of a trigonometric series. This would be a good candidate, but how do we know the limits exist, even assuming the  $f_n$  can be found?

## IV. The Lebesgue Integral

In the theory of the Lebesgue integral we once again encounter an extension of the concept of integrability that was related by its creator to problems of representation by trigonometric series. As Lebesgue states in his first extended treatment of this topic in 1902 [15]: "In concerning myself with trigonometric series, I had as a special goal that of showing the use of the notion of the integral which I introduced in my thesis in the study of discontinuous functions of a real variable."

He then outlines his definition of the integral, which is heavily dependent on the idea of the measure of a set  $E \subset [a, b]$ . This concept played no role in Riemann's

integral. The exterior measure  $m_e(E)$  is the greatest lower bound of the numbers  $\sum_{i=1}^{\infty} |b_i - a_i|$ , where  $\bigcup_{i=1}^{\infty} (a_i, b_i) \supset E$ . If  $E \subset (a, b)$ , then the interior measure  $m_i(E)$  is defined as  $m_i(E) = |b - a| - m_e((a, b) - E)$ ;  $E$  is said to be measurable if  $m_e(E) = m_i(E)$ , and its measure is defined as  $m(E) = m_e(E) = m_i(E)$ . The Lebesgue integral is then defined for certain kinds of bounded functions  $f$  on  $[a, b]$ :  $f$  is called measurable if  $f^{-1}[c, d]$  is a measurable set for every interval  $[c, d] \subset [a, b]$ . Now if  $m \leq f(x) \leq M$  for all  $x \in [a, b]$  and  $m = y_0 < y_1 < y_2 < \cdots < y_n = M$  is any subdivision of the range of  $f$ , we define  $E_i = f^{-1}[y_{i-1}, y_i]$ ,  $i = 1, \dots, n$ ; then  $E_i$  are measurable sets. Lebesgue shows that the sums  $\sum_{i=1}^n y_i m(E_i)$  tend toward a definite limit as  $|y_i - y_{i-1}| \rightarrow 0$  uniformly; this limiting value is the Lebesgue integral of  $f$ . In some sense the Lebesgue integral replaces the "Riemann sums"  $\sum_{i=1}^n f(x_i)(x_i - x_{i-1})$  with a more general kind of sum, one involving measures of sets  $E_i$  into which  $[a, b]$  is subdivided and on which  $f$  maintains values in a certain subdivision of its range. His integral is an extension of Riemann's in the sense that any function integrable in Riemann's sense is integrable in Lebesgue's, with integrals agreeing. But, for example,  $\chi_Q$  is Lebesgue integrable but not Riemann integrable.

The theory developed by Lebesgue and his successors is very extensive, but one theorem in particular is important for us and was often used by Lebesgue in his early papers. If  $f_n$  is a sequence of measurable functions on  $[a, b]$ , and if there is an  $M \geq 0$  such that  $|f_n(x)| \leq M$  for all  $x$  and all  $n$ , and if  $f_n(x) \rightarrow f(x)$  for all  $x$  (convergence almost everywhere suffices), then  $f$  is measurable and  $\int_a^b f_n \rightarrow \int_a^b f$ . This is usually called the Bounded Convergence Theorem (BCT).

One immediate consequence of the broader definition of the integral is that the theorems that were stated earlier about convergence of the Fourier series of  $f$  to  $f$  can now be stated for a more general class of functions, since those proofs go through with the integrals interpreted in the Lebesgue sense. "Now that we know the form for the coefficients [the Fourier coefficients expressed as Lebesgue integrals], we can hope to study the convergence of the [Fourier] development by ordinary methods. And, in fact, it was sufficient for me to modify very little the reasonings previously employed in order to find a wide set of cases of convergence," Lebesgue noted in [15]. But this is relatively minor compared with the use of his new theorems.

One can see this immediately in the question of uniqueness. If

$$f(x) = A_0/2 + \sum_{k=1}^{\infty} A_k \cos kx + B_k \sin kx, \quad (7)$$

then  $f$  is automatically measurable, being the pointwise limit of a sequence of measurable functions, hence integrable. (Recall that we are assuming throughout that  $f$  is bounded.) Furthermore if the partial sums are uniformly bounded, i.e.,  $|A_0/2 + \sum_{k=1}^n A_k \cos kx + B_k \sin kx| \leq M$  for all  $x$  and  $n$ , then by BCT the argument outlined in section I is valid if the integrals are interpreted in the Lebesgue sense. Thus  $A_k$  and  $B_k$  must be the Fourier coefficients of  $f$ . Uniform convergence of the partial sums is thus replaceable by the weaker hypothesis of uniform boundedness and no additional assumption on  $f$  is required.

But even the hypothesis of uniform boundedness may be omitted; this was the first result given by Lebesgue in [15]. We give an outline of his proof, which depends on a theorem proved earlier by Riemann (given as (a) below).

Let us assume for simplicity that  $f$  is continuous, (7) holds, and the goal is to show that  $A_k$  and  $B_k$  are the Fourier coefficients of  $f$ . Here are two needed results.

(a) Define

$$F(x) = (A_0/2)x^2 - \sum_{k=1}^{\infty} \frac{1}{k^2} (A_k \cos kx + B_k \sin kx);$$

$F$  is the series obtained from the original one by “formally” integrating it twice.

Define

$$\frac{\Delta^2 F(x)}{t^2} = \frac{F(x+t) + F(x-t) - 2F(x)}{t^2},$$

then  $\lim_{t \rightarrow 0} (\Delta^2 F/t^2)(x) = f(x)$  for all  $x$ . (The limiting quotient defined here is a generalization of double differentiation usually called the second Schwartz derivative.) Thus Riemann showed that when the series for  $f$  is formally integrated twice and “Schwartz differentiated” twice, the result is  $f$  again, and

$$(b) \quad \inf_{s \in [x-t, x+t]} f(s) \leq \frac{F(x+t) + F(x-t) - 2F(x)}{t^2} \leq \sup_{s \in [x-t, x+t]} f(s), \quad \text{for any fixed } x \text{ and } t.$$

This was proved by Lebesgue in section 1 of [15].

Lebesgue’s proof now proceeds as follows: Assume (7), then since  $f$  is bounded,  $\Delta^2 F(x)/t^2$  is uniformly bounded in  $x$  and  $t$ , by (b), we may integrate twice to get

$$\int_0^\phi \int_0^\theta f(x) \, dx \, d\theta = \int_0^\phi \int_0^\theta \lim_{t \rightarrow 0} \frac{\Delta^2 F}{t^2}(x) \, dx \, d\theta = \lim_{t \rightarrow 0} \int_0^\phi \int_0^\theta \frac{\Delta^2 F}{t^2}(x) \, dx \, d\theta \quad \text{by BCT.}$$

An evaluation of the iterated integral gives

$$\int_0^\phi \int_0^\theta f(x) \, dx \, d\theta = F(\phi) - F(0) - \phi \lim_{t \rightarrow 0} \frac{\Delta^2 F_1(0)}{t^2},$$

where  $F_1$  is an antiderivative for  $F$ . Thus there are constants  $A$  and  $B$  such that  $F(\phi) = \int_0^\phi \int_0^\theta f(t) \, dt \, d\theta + A\phi + B$ . Now we apply the procedure outlined in section I to determine the coefficients of the trigonometric series for  $F(\phi) - (A_0/4)\phi^2$  in terms of integrals of  $f$ . The procedure is valid since the series for  $F(\phi)$  is uniformly convergent. Lebesgue obtains three triple integral formulae relating  $A_n$ ,  $B_n$  and  $f$ ; for instance,

$$-A_n/n^2 = \frac{1}{\pi} \int_{-\pi}^\pi \int_0^\phi \int_0^\theta f(t) \cos n\phi \, dt \, d\theta \, d\phi - \frac{2A_0}{n^2}.$$

Reversing the orders of integration and solving gives

$$A_n = \frac{1}{\pi} \int_{-\pi}^\pi f(t) \cos nt \, dt + 2 \left[ A_0 - \frac{1}{2\pi} \int_0^{2\pi} f(t) \, dt \right].$$

Now using the convergence of  $\sum_{n=1}^\infty A_n$  to get  $A_n \rightarrow 0$ , and the Riemann-Lebesgue Lemma, we let  $n \rightarrow \infty$ . So the bracketed quantity is 0, yielding  $A_n = (1/\pi) \int_{-\pi}^\pi f(t) \cos nt \, dt$ . The companion formula for the other coefficient follows from a similar argument. Thus the focus on Fourier series as the “best candidate” for representing a function by a trigonometric series is validated: Any trigonometric representation of a continuous function must be its Fourier series.

This proof is actually valid for measurable  $f$ . We required that  $f$  be continuous to simplify the assumptions needed for reversing the orders of integration. Lebesgue showed in [15] how this restriction may be removed.

As another example of the use of BCT, we return to Parseval’s identity, (6). In [16] is a short proof of this identity for  $f$  measurable, which Lebesgue obtained by

working with  $\sigma_n(x) = [s_0(x) + s_1(x) + \cdots + s_{n-1}(x)]/n$ , the Cesàro means of  $s_n$ . These are “smoothed-out” versions of (3) and have the useful property that  $|\sigma_n(x)| \leq \sup_{-\pi \leq t \leq \pi} |f(t)|$  for all  $n$  and  $x$ , i.e., they are uniformly bounded. Fejér showed in [10] that  $\sigma_n \rightarrow f$  at each point of continuity of  $f$ . A deeper result of Lebesgue’s theory actually shows  $\sigma_n \rightarrow f$  almost everywhere.

Thus if we assume  $f$  is measurable, BCT gives

$$\frac{1}{\pi} \int_{-\pi}^{\pi} f^2 = \lim_{n \rightarrow \infty} \frac{1}{\pi} \int_{-\pi}^{\pi} \sigma_n^2 = \lim_{n \rightarrow \infty} \left[ (1/2) a_0^2 + \sum_{k=1}^n (a_k^2 + b_k^2) (1 - k/n)^2 \right],$$

where the last equality is another straightforward calculation. But

$$\lim_{n \rightarrow \infty} \left[ (1/2) a_0^2 + \sum_{k=1}^n (a_k^2 + b_k^2) (1 - k/n)^2 \right] \leq \lim_{n \rightarrow \infty} \left[ (1/2) a_0^2 + \sum_{k=1}^n a_k^2 + b_k^2 \right] \leq \int_{-\pi}^{\pi} f^2,$$

and the Parseval identity follows.

As a final example of the use of the BCT, we give a quick proof of what is known as the Cantor-Lebesgue Lemma. We have seen that if  $f$  is integrable in Riemann’s sense, then the Fourier coefficients tend to zero. The same is true if  $f$  is Lebesgue integrable; this is also proved in [15]. A similar problem is: If  $\sum_{n=0}^{\infty} c_n \cos nx + d_n \sin nx$  is *any* convergent trigonometric series, do we have  $c_n, d_n \rightarrow 0$ ? The answer is yes, as was first proved by Cantor [4]. Here is a short proof: Since the given series converges, we must have  $c_n \cos nx + d_n \sin nx \rightarrow 0$  as  $n \rightarrow \infty$  for all  $x$  in  $[-\pi, \pi]$ . Now if the lemma is false, there is an  $m$  such that  $c_{n_j}^2 + d_{n_j}^2 \geq m$  for a subsequence  $n_j$ . Defining  $f_{n_j}(x) = (c_{n_j} \cos n_j x + d_{n_j} \sin n_j x)^2 / (c_{n_j}^2 + d_{n_j}^2)$ , we have that  $f_{n_j}$  is uniformly bounded (the quantity in parenthesis in the numerator can be written as a single sine or cosine with amplitude equal to the square root of the denominator) and  $f_{n_j} \rightarrow 0$  for all  $x$  in  $[-\pi, \pi]$ . Thus  $\int_{-\pi}^{\pi} f_{n_j}(x) dx \rightarrow 0$  by BCT; but an easy calculation shows that these integrals all equal  $\frac{1}{2}$ , giving a contradiction.

This lemma, far from being of isolated interest, is the first step in Cantor’s famous proof of the uniqueness of trigonometric series expansion; using it, one can prove that any two trigonometric series that are equal for all values of  $x$  in  $[-\pi, \pi]$  must have the same coefficients. This answers affirmatively the question at the end of section I. (See [1] for details on this proof.)

In fairness it perhaps should be said that similar convergence theorems to the BCT can be proved for sequences of continuous or Riemann-integrable functions without the full development of the Lebesgue measure theory ([17], p. 209, for example), but these theorems require an additional hypothesis about the continuity or Riemann-integrability of the limit function, as Lebesgue’s does not, and in some cases the proofs seem to require at least some measure-theoretic ideas. The automatic integrability of the limit function was important in the uniqueness question at the beginning of this section.

As a final remark let us note that Lebesgue in [15] constructs an example of a measurable function that is not Riemann integrable but that is representable everywhere by its Fourier series. Thus, one can see that the class of functions representable by trigonometric series is quite large—beyond even the minimal regularity required by Riemann integrability. His construction is roughly as follows. He lets  $E$  be a generalized Cantor set in  $[-\pi, \pi]$ , i.e., a closed, nowhere dense set of positive measure. On each subinterval of the complement of  $E$  he constructs a differentiable function that is infinitely oscillating and has no limit at each endpoint of the interval. His function  $f$  is defined on  $[-\pi, \pi]$  as zero at each point of  $E$  and the appropriate oscillating function on each subinterval of the complement of  $E$ . Non-Riemann

integrability follows from the discontinuity of  $f$  at each point of  $E$ , a set of positive measure. The oscillating functions are chosen in such a way that  $f$  is Lebesgue integrable. The Fourier series of  $f$  (with coefficients necessarily defined as Lebesgue integrals) converges to  $f$  on each complementary interval of  $E$  since  $f$  is differentiable there, and converges also on  $E$  by a computation dependent on the explicit choice of the oscillating functions.

Lebesgue relates this construction to the question of integral representation of the coefficients in a trigonometric expansion of an arbitrary function. This is the major theorem at the beginning of this section: "I have thus demonstrated that there exist non [Riemann] integrable functions representable trigonometrically; the calculation of the coefficients done at the beginning is therefore not without object." Clearly if all functions representable trigonometrically were Riemann integrable there would be no reason to concern oneself with a more general integration theory, and one could hope that theorems about the Riemann integral alone could deal with this coefficient problem.

There is much more that can be said using the Lebesgue theory; in particular the whole topic of "convergence in square mean" is a triumph for the theory, identifying the class of square integrable Lebesgue measurable functions as precisely that class that can be represented by Fourier series with square-summable coefficients, provided that convergence is considered "in the mean". We have chosen to emphasize the contributions of Lebesgue integration given here because of the impression that mean-square convergence, although in some sense the "correct" way of considering convergence of Fourier series, is perhaps not as immediately accessible as the pointwise kind, with which students are already familiar.

## Conclusion

I hope that this tour through the theories of integration has provided an understanding of the need for refinement and development of the idea of the integral in the context of trigonometric series. Perhaps it is possible this way to see that the new theories did not develop in a vacuum—that they were responses to questions posed by previous attempts and in all cases overcame some of the inadequacies of their predecessors. Of course integration theory can be applied to many other problems, and Lebesgue's approach is capable of vast generalization into settings in which these series topics do not even make sense. But a focus on trigonometric series can provide a way of grounding the subject in a set of issues easily accessible to students of analysis.

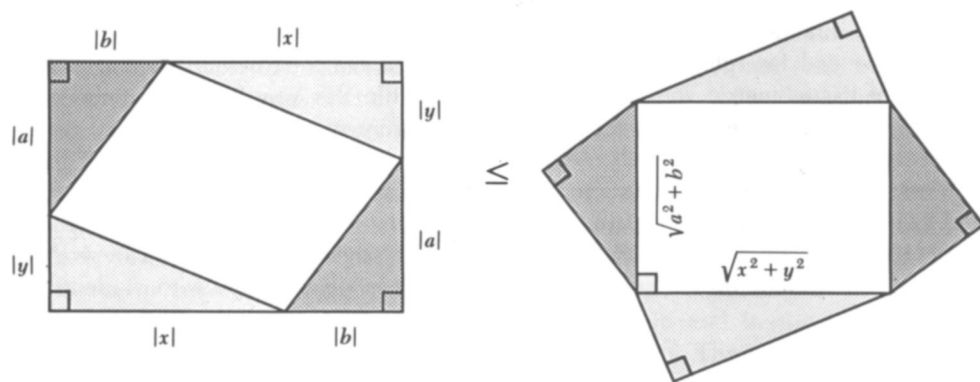
## REFERENCES

1. J. Marshall Ash, Uniqueness of representation by trigonometric series, *Amer. Math. Monthly*, 96 (1989), 873–885.
2. Ossian Bonnet, *Mémoires Couronnés et Mémoires des Savants Étrangers publiés par l'Académie Royale des Sciences, des Lettres et des Beaux-Arts de Belgique*, 23 (1850), p. 8.
3. Umberto Bottazzini, *The Higher Calculus: A History of Real and Complex Analysis from Euler to Weierstrass*, Springer-Verlag New York Inc., New York, 1986.
4. Georg Cantor, Über einen die trigonometrischen Reihen betreffenden Lehrsatz, *Journal für die reine und angewandte Mathematik*, 72 (1870), 130–138.
5. H. S. Carslaw, *An Introduction to the Theory of Fourier's Series and Integrals*, Dover Publications, Mineola, NY, 1950.
6. Augustin-Louis Cauchy, *Résumé des Leçons données à l'École Royale Polytechnique sur le calcul infinitésimal*, Paris, de Bure, 1823; see *Oeuvres complètes*, Vol. 4, ser. 2, pp. 5–26.

7. Augustin-Louis Cauchy, *Mémoire sur les développements des fonctions en séries périodiques*. Paris, Mémoires de l'Académie des Sciences, 6, 1827, pp. 603ff., Oeuvres complètes (1), 2, (1908), pp. 12–19.
8. Joseph W. Dauben, Georg Cantor and the origins of transfinite set theory, *Scientific American*, June 1983, pp. 122–131.
9. Johan P. G. L. Dirichlet, 1829, Sur la convergence des séries trigonométriques qui servent à représenter une fonction arbitraire entre des limites données; *Journal für die reine und angewandte Mathematik*, 4 (1829), 157–169.
10. Leopold Fejér, Untersuchungen über Fouriersche Reihen, *Mathematische Annalen* 58 (1904), 51–69.
11. J. Grabiner, *The Origins of Cauchy's Rigorous Calculus*, MIT Press, Cambridge, MA, 1981.
12. I. Grattan-Guinness, *The Development of the Foundations of Mathematical Analysis from Euler to Riemann*, MIT Press, Cambridge, MA, 1970.
13. Joseph Fourier, *La Théorie Analytique de la Chaleur*. Paris, Didot 1822.
14. Thomas Hawkins, *Lebesgue's Theory of Integration*, Chelsea Publishing Co., New York 1979.
15. Henri Lebesgue, Sur les séries trigonométriques, Paris, *Annales Scientifiques de l'École Normale Supérieure* (3) 20, (1903), 453–485.
16. Henri Lebesgue, *Leçons sur les séries trigonométriques*, Gauthier-Villars, Paris, 1906.
17. J. Lewin and M. Lewin, *An Introduction to Mathematical Analysis*, Random House, New York, 1988.
18. Benoit Mandelbrot, *The Fractal Geometry of Nature*, W. H. Freeman and Co., San Francisco, 1982.
19. Bernard Riemann, Über die Darstellbarkeit einer Function durch eine trigonometrische Reihe, in *Gesammelte Mathematische Werke und Wissenschaftlicher Nachlass* pp. 227–265 (first ed. 1876; 1902 ed. reprinted by Dover Publications, Mineola, NY, 1953, pp. 227–271. Citations are from the Dover reprint.
20. Georgi P. Tolstov, *Fourier Series*, Dover Publications, Mineola, NY, 1962.
21. E. B. Van Vleck, The influence of Fourier's series upon the development of mathematics, *Science*, 39, No. 995, 113–124.

## Proof without Words: The Cauchy-Schwarz Inequality

$$|\langle a, b \rangle \cdot \langle x, y \rangle| \leq \|\langle a, b \rangle\| \|\langle x, y \rangle\|$$



$$(|a| + |y|)(|b| + |x|) \leq 2\left(\frac{1}{2}|a||b| + \frac{1}{2}|x||y|\right) + \sqrt{a^2 + b^2} \sqrt{x^2 + y^2}$$

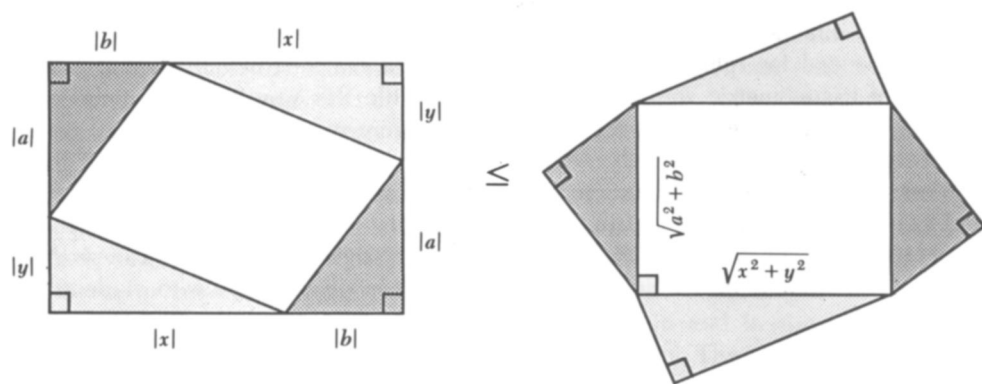
$$\therefore |ax + by| \leq |a||x| + |b||y| \leq \sqrt{a^2 + b^2} \sqrt{x^2 + y^2}$$

—ROGER B. NELSEN  
LEWIS AND CLARK COLLEGE  
PORTLAND, OR 97219

7. Augustin-Louis Cauchy, *Mémoire sur les développements des fonctions en séries périodiques*. Paris, Mémoires de l'Académie des Sciences, 6, 1827, pp. 603ff., Oeuvres complètes (1), 2, (1908), pp. 12–19.
8. Joseph W. Dauben, Georg Cantor and the origins of transfinite set theory, *Scientific American*, June 1983, pp. 122–131.
9. Johan P. G. L. Dirichlet, 1829, Sur la convergence des séries trigonométriques qui servent à représenter une fonction arbitraire entre des limites données; *Journal für die reine und angewandte Mathematik*, 4 (1829), 157–169.
10. Leopold Fejér, Untersuchungen über Fouriersche Reihen, *Mathematische Annalen* 58 (1904), 51–69.
11. J. Grabiner, *The Origins of Cauchy's Rigorous Calculus*, MIT Press, Cambridge, MA, 1981.
12. I. Grattan-Guinness, *The Development of the Foundations of Mathematical Analysis from Euler to Riemann*, MIT Press, Cambridge, MA, 1970.
13. Joseph Fourier, *La Théorie Analytique de la Chaleur*. Paris, Didot 1822.
14. Thomas Hawkins, *Lebesgue's Theory of Integration*, Chelsea Publishing Co., New York 1979.
15. Henri Lebesgue, Sur les séries trigonométriques, Paris, *Annales Scientifiques de l'École Normale Supérieure* (3) 20, (1903), 453–485.
16. Henri Lebesgue, *Leçons sur les séries trigonométriques*, Gauthier-Villars, Paris, 1906.
17. J. Lewin and M. Lewin, *An Introduction to Mathematical Analysis*, Random House, New York, 1988.
18. Benoit Mandelbrot, *The Fractal Geometry of Nature*, W. H. Freeman and Co., San Francisco, 1982.
19. Bernard Riemann, Über die Darstellbarkeit einer Function durch eine trigonometrische Reihe, in *Gesammelte Mathematische Werke und Wissenschaftlicher Nachlass* pp. 227–265 (first ed. 1876; 1902 ed. reprinted by Dover Publications, Mineola, NY, 1953, pp. 227–271. Citations are from the Dover reprint.
20. Georgi P. Tolstov, *Fourier Series*, Dover Publications, Mineola, NY, 1962.
21. E. B. Van Vleck, The influence of Fourier's series upon the development of mathematics, *Science*, 39, No. 995, 113–124.

## Proof without Words: The Cauchy-Schwarz Inequality

$$|\langle a, b \rangle \cdot \langle x, y \rangle| \leq \|\langle a, b \rangle\| \|\langle x, y \rangle\|$$



$$(|a| + |y|)(|b| + |x|) \leq 2\left(\frac{1}{2}|a||b| + \frac{1}{2}|x||y|\right) + \sqrt{a^2 + b^2} \sqrt{x^2 + y^2}$$

$$\therefore |ax + by| \leq |a||x| + |b||y| \leq \sqrt{a^2 + b^2} \sqrt{x^2 + y^2}$$

—ROGER B. NELSEN  
LEWIS AND CLARK COLLEGE  
PORTLAND, OR 97219

## Mathematical Certificates

JOHN HIGGINS  
DOUGLAS CAMPBELL  
Brigham Young University  
Provo, UT 84602

Some domains of knowledge permit easy verification but difficult creation. Consider the primes. Suppose that the RSA public key encryption system [7], that requires large (hundred-digit) primes, becomes a household item. The generation of hundred-digit primes could become a cottage industry.

*Jill:* But imagine a lawyer trying to convince a judge that the hundred-digit number on a crumpled piece of paper is indeed prime. It does not suffice to point out that the number does not end in 0, 2, 4, 6, 8, or 5.

*John:* Let  $N$  be the hundred-digit number. Let the lawyer introduce the results of the Sieve of Eratosthenes up to  $10^{100}$ . The judge can verify that  $N$  appears.

*Jill:* Won't do. A naive implementation of the Sieve of Eratosthenes up to  $N$ , requires on the order of  $N$  additions. A year has about  $10^8$  seconds. On a computer doing 100 billion additions per second, your procedure would take the lawyer about  $10^{81}$  years.

*John:* I see the problem. The Sieve of Eratosthenes produces all primes *and* all composite numbers up to  $N$ . But the lawyer is only interested in showing that her client's  $N$  is prime.

Here's my second scheme.  $N$  is a prime if, and only if, it has no divisors less than or equal to  $N^{1/2}$ . The lawyer computes the remainder of  $N$  divided by every number less than  $N^{1/2}$ . If the judge sees that they are all nonzero, the judge can see that the number is prime.

*Jill:* Oh that's much better. There are at most  $10^{50}$  numbers less than the square root of a hundred-digit prime. On a computer doing 100 billion divisions per second, your second procedure takes only  $10^{31}$  years!

*John:* Well, the lawyer should persuade her client to reproduce in court the sequence of steps that created the hundred-digit number  $N$ .

*Jill:* There are four objections to your third scheme.

- (1) It could be time consuming; it uses the *creation* procedure for *verification*.
- (2) The primer may claim a noncomputational method for generating large primes, such as maintaining a ritually pure food diet, observing various cultic rites, followed by fasting and prayer.
- (3) The primer may claim that his prime creation algorithm is a trade secret, and he will suffer irreparable loss by having to disclose it.

However, he is willing to provide a *certificate* that the number is prime. The certificate separates *verification* from *creation*. Let me explain.

Recall that every integer  $N$  has a *short* decimal description: its  $\log_{10} N$  decimal digits. A number's *description* is the logarithm of its *magnitude*.

A *verification procedure* for property  $P$  is an algorithm that takes the *description* of  $N$ , performs various basic computational steps, and outputs yes if, and only if,  $N$  has property  $P$ . For example, the Sieve of Eratosthenes is a verification procedure for



primality that uses additions. The square root algorithm is a verification procedure for primality that uses divisions. Unfortunately, these are also time-consuming creation procedures.

In any verification procedure, the maximum number of basic computational steps taken by the algorithm is a function of the input description length. (Concentrating on basic computational steps lets us ignore details of array manipulation, indices, questions of efficient implementation of basic computational steps, and relative time differences for basic computational steps.) We focus on the *number* of basic computational steps. Computational complexity defines a verification procedure for property  $P$  of an integer  $N$  to be *short* if the maximum number of basic computational steps is bounded by a power of the description of  $N$ .

*John:* Let's see whether the Sieve of Eratosthenes is short. Since it marks every integer from 1 to  $N$  as prime or composite, it must take at least  $N$  computational steps. But  $N$  is not bounded by any power of  $\log_{10} N$ . Therefore, the Sieve of Eratosthenes is not a short verification procedure.

Let me now check the naive square-root procedure. The naive square-root procedure divides every number less than  $N^{1/2}$  into  $N$ . Therefore, it requires at least  $N^{1/2}$  computational steps. But  $N^{1/2}$  is not bounded by any power of  $\log_{10} N$ . Therefore, it also is not a short verification procedure.

*Jill:* Quite right. Even if we had hardware, tailored to do 100 billion basic computation steps every second, for each of the  $10^8$  seconds a year, for each of the 10 billion years of the known universe, we get a paltry  $10^{29}$ , nowhere near the  $10^{50}$  basic computations needed by the square root algorithm for the primality of a random hundred-digit number.

Let me introduce the notion of a *certificate* that an integer  $N$  has property  $P$ . A *certificate* that the integer  $N$  has property  $P$  is a set of integers  $I$  and an algorithm whose successful termination on the set of integers in  $I$  guarantees that  $N$  has property  $P$ . We say that a property  $P$  has *short* certificates if, and only if, every number with property  $P$  has a certificate for which the number of the algorithm's computational steps is bounded by a power of  $\log_{10} N$ , the length of the description of  $N$ .

One certificate for the primality of an odd number  $N$  consists of  $I$ , all *odd integers* from 3 to  $N - 2$ , together with the algorithm that checks whether each integer in  $I$  fails to divide  $N$ . The algorithm requires  $N/2$  computational steps. Since  $N/2$  is not bounded by a power of  $\log_{10} N$ , this certification procedure is not short.

Another certificate for the primality of an odd number  $N$  consists of  $I$ , all *primes* from 3 to  $N^{1/2}$ , and the algorithm that checks that each element in  $I$  fails to divide  $N$ . This is more parsimonious than the naive square root algorithm. By the prime number theorem [3, p. 366], the cardinality of  $I$  is essentially  $(N^{1/2}/\log_{10} N)$ , a quantity that again is not bounded by any power of  $\log_{10} N$ . Thus, this parsimonious square root procedure is still not short.

However, before you lose hope, the primes do have short certificates. In fact, for a prime  $N$ , I claim that the cardinality of  $I$  in a certificate for  $N$  is bounded by a *linear* function of  $\log_{10} N$ , the input description of the number  $N$ . Thus, the certificate for a hundred-digit prime is not more than forty times the length of any certificate for any prime up to 999 (three-digit prime).

*John:* So you have a fast method for primes?

*Jill:* Yes. I can verify that a number  $N$  is prime by exponentiations and divisions, the number of which is bounded by a multiple of  $\log_{10} N$ . By your face, I see that I need to begin with a theoretical explanation. Please humor me by giving Fermat's theorem in a computational form.

*John:* From a computational standpoint, Fermat's theorem [2, p. 39] states

*If  $p$  is prime and if  $b$  is relatively prime to  $p$ , then  $b^{p-1} = 1 \pmod{p}$ .*

For example, since 11 is prime and 2 is relatively prime to 11,  $2^{10}$  must be 1 modulo 11. We verify this by noting that  $2^{10} = 1024 = 11 * 93 + 1$ .

But, the converse is not always true. There are  $b$ 's that are relatively prime to  $s$  and for which  $b^{s-1}$  is 1 mod  $s$ , but for which  $s$  is not prime. For example, 15 and 4 are relatively prime; let  $b = 4$  and  $s = 15$ , then  $4^{14} = (4^2)^7 = (1)^7 = 1 \pmod{15}$ , but 15 is not prime.

*Jill:* But a true converse is given by Lucas' theorem [4].

*Let  $p$  and  $b$  be positive integers. Suppose*

1.  $b^{p-1} = 1 \pmod{p}$ ; and
2. *for each prime divisor  $q$  of  $p - 1$ ,  $b^{(p-1)/q} \neq 1 \pmod{p}$ . Then  $p$  is prime.*

Let me show you how Lucas' theorem lets me write a certificate for the primality of 211. A certificate of primality for 211 consists of the special number 2, and the four primes 2, 3, 5, and 7. It is written (2; 2, 3, 5, 7). The execution of the algorithm consists of two parts:

- (a) Verifying that 2, 3, 5, and 7 are the sole prime divisors of  $211 - 1 = 210 = 2 * 3 * 5 * 7$ ;
- (b) verifying the five computational steps:

$$\begin{array}{lll}
 (1) & 2^{210} & = 1 \pmod{211}, \\
 (2a) & 2^{105} & = 210 \pmod{211}, \\
 (2b) & 2^{70} & = 196 \pmod{211}, \\
 (2c) & 2^{42} & = 107 \pmod{211}, \\
 (2d) & 2^{30} & = 171 \pmod{211}.
 \end{array}$$

*John:* Wait! I have four objections.

- (1) The five basic computational steps (1), (2a), (2b), (2c), and (2d), involve exponentiation to a large power and are not cheap.
- (2) You did not explain where you got  $b$  and the factors for  $p - 1$ .
- (3) Your certificate that 211 is prime depends on knowledge that 2, 3, 5, and 7 (or in general, some smaller numbers) are prime. The method therefore has hidden costs.
- (4) You haven't shown the existence of a constant  $M$  such that the length of the certificate for primality of  $p$  is bounded by  $M * \log_{10} p$ .

*Jill:* Patience. I wanted you to understand the method before I went into the details. The use of Lucas' theorem to provide short certificates for primality is due to Vaughan R. Pratt [6]. Table 1 (see p. 24) gives certificates for all primes from 3 to 211.

Examine, say, the certificate for 107: (2; 2, 53). Since  $2 * 53 = 107 - 1$ , we require the knowledge that 2 and 53 are prime. The certificate for 53 is (2; 2, 2, 13). Since  $2 * 2 * 13 = 53 - 1$ , we require the additional knowledge that 13 is prime. The certificate for 13 is (2; 2, 2, 3). Since  $2 * 2 * 3 = 13 - 1$ , we require the additional knowledge that 3 is prime. The certificate for 3 is (2; 2). Since  $2 = 3 - 1$ , we see that the proof of the primality of 107 reduces to the fact that 2 itself is prime. All prime certificates reduce to 2 being prime.

The prime certificate method is recursive. It's like induction: to prove that the  $n$ -th prime has a certificate, we make use of the fact that all primes from the first prime to the  $(n - 1)$ st prime have certificates.

TABLE 1: CERTIFICATES OF PRIMALITY

$p$	$b$	Certificate							
3	2	2							
5	2	2	2						
7	3	2	3						
11	2	2	5						
13	2	2	2	3					
17	3	2	2	2	2				
19	2	2	3	3					
23	5	2	11						
29	2	2	2	7					
31	3	2	3	5					
37	2	2	2	3	3				
41	6	2	2	2	5				
43	3	2	3	7					
47	5	2	23						
53	2	2	2	13					
59	2	2	2	29					
61	2	2	2	3	5				
67	2	2	3	11					
71	7	2	5	7					
73	5	2	2	2	3	3			
79	3	2	3	13					
83	2	2	41						
89	3	2	2	2	11				
97	5	2	2	2	2	2	3		
101	2	2	2	5	5				
103	5	2	3	17					
107	2	2	53						
109	6	2	2	3	3	3			
113	3	2	2	2	2	7			
127	3	2	3	3	7				
131	2	2	5	13					
137	3	2	2	2	17				
139	2	2	3	23					
149	2	2	2	37					
151	6	2	3	5	5				
157	5	2	2	3	13				
163	2	2	3	3	3	3			
167	5	2	83						
173	2	2	2	43					
179	2	2	89						
181	2	2	2	3	3	5			
191	19	2	5	19					
193	5	2	2	2	2	2	2	3	
197	2	2	2	7	7				
199	3	2	3	3	11				
211	2	2	3	5	7				

*John:* But you just do three things before you invoke the magic word ‘recursion.’

(1) You must show that all primes have certificates under this scheme. (2) You must show that no composite number has a certificate under this scheme. (3) You must show that the certificate is still short even if *all* recursively defined steps are taken into account.

*Jill:* First I show that every prime has at least one certificate. If  $p$  is a prime, then there exists a certificate  $(b; p_1, p_2, \dots, p_n)$ , where  $b$  is a positive integer, and  $p_i$ ,  $0 < i \leq n$ , is a nondecreasing sequence of positive integers such that

- (1)  $p_1 * \dots * p_n = p - 1$ .
- (2)  $b^{p-1} = 1 \pmod{p}$ ;
- (3) for each  $i$ ,  $0 < i \leq n$ ,  $b^{(p-1)/p_i} \neq 1 \pmod{p}$ ;
- (4) for each  $i$ ,  $0 < i \leq n$ ,  $p_i$  is a prime.

*John:* I’ll grant that the existence of  $(b; p_1, \dots, p_n)$  satisfying 1, 2, 3, and 4, together with Lucas’ theorem, constitutes a certificate for  $p$  being prime. However, maybe you stopped Table 1 when the certificates stopped!

*Jill:* I’m hurt. What do you know about the existence of primitive roots for primes?

*John:* Every prime has a primitive root [3, p. 20, Theorem C], that is, for every prime  $p$  there exists a positive number  $b$  such that

- (I) for  $1 \leq j \leq p - 2$ ,  $b^j \neq 1 \pmod{p}$ ,
- (II) for  $j = p - 1$ ,  $b^j = 1 \pmod{p}$ .

*Jill:* Thank you. Here’s how to show each prime has a certificate. Let  $b$  be a primitive root for  $p$ . You’ve already guaranteed it exists. Let  $p_1, \dots, p_n$  be the primes in the prime factorization of  $p - 1$ . They certainly exist. Hence  $(b; p_1, \dots, p_n)$  exists. Since  $p_1 * \dots * p_n = p - 1$ , we have (1). Since  $p_1, \dots, p_n$  are primes we have (4). By (I) and (II), both (2) and (3) hold.

*John:* Tilt. I see the pea under the walnut shell. You have hidden all the work! I mean of course in that sense certificates exist. But how do you propose to *find* a primitive root? How do you propose to *factor*  $p - 1$ , when  $p$  is a hundred-digit number?

*Jill:* *Don’t change the game John.* The client doesn’t have to tell you where he got the certificate. He only has to prove a way for you to verify his claim that he has a prime. And, unlike your methods, *you can verify* a certificate in a few moments instead of  $10^{29}$  years.

*John:* But you also must show that no composite number has a certificate.

*Jill:* But John, that is precisely what Lucas’ Theorem proves: If a number  $p$  has a certificate,  $(b; p_1, p_2, \dots, p_n)$ , then  $p$  must be prime.

*John:* What about my other two objections:

- (1) Perhaps the length of the extended nonrecursive certificate is not bounded by a polynomial function of the input length;
- (2) Perhaps the time it takes to check all the multiplications and exponentiations is not bounded by a polynomial function of the input length.

*Jill:* Let me relieve your mind of both anxieties by bounding the number of exponentiations and multiplications when all recursive calls are taken into account. For simplicity, let’s count raising a number to a power and then taking its residue as a single *Basic Operation*. Similarly, let’s count a series of multiplications as a single *Basic Operation*.

Let  $\log_2 p$  be the length of  $p$  written in binary rather than decimal notation. Let  $f(p)$  be the number of Basic Operations needed to verify that  $p$  is a prime taking recursion into account. I will prove by induction that

$$f(p) \leq -1 + 4 * (\log_2 p).$$

Recall that Lucas' theorem requires four different types of checks on a certificate  $(b; p_1, p_2, \dots, p_n)$  for  $p$ :

- (1)  $p - 1 = p_1 * p_2 * \dots * p_n$ ;
- (2)  $b^{p-1} = 1 \pmod{p}$ ;
- (3) for each  $i$ ,  $0 < i \leq n$ ,  $b^{(p-1)/p_i} \neq 1 \pmod{p}$ ;
- (4) for each  $i$ ,  $1 < i \leq n$ ,  $p_i$  is a prime.

(Since the  $p_i$  are nondecreasing and  $p$  is odd,  $p_1$  is always two; since  $p_1 = 2$  is prime, we don't check that  $p_1$  is prime). To check (1) involves a series of multiplications, and is therefore one Basic Operation; to check (2) requires raising a number to a power and computing the residue, a Basic Operation; to check (3) requires  $n$  Basic Operations; to check (4) requires recursive calls to the certification procedure for  $p_2$  through  $p_n$ .

$$f(p) = 1 + 1 + n + f(p_2) + \dots + f(p_n). \quad (1)$$

Let's first check that the induction basis holds for  $p = 2$ . Since no multiplications nor exponentiations are needed to verify that two is prime,  $f(2) = 0$ . But  $-1 + 4 * \log_2 2 = 3$ . Since  $0 \leq 3$ , the basis  $f(p) \leq -1 + 4 * \log_2 p$  holds for  $p = 2$ .

Assume by induction that the bound holds for all primes less than  $p$ . Equation (1) becomes

$$\begin{aligned} f(p) &= 2 + n + f(p_2) + \dots + f(p_n) \\ &\leq 2 + n + (-1 + 4 * \log_2 p_2) + \dots + (-1 + 4 * \log_2 p_n) \\ &= 3 + 4 * (\log_2 p_2 + \dots + \log_2 p_n) \\ &= 3 + 4 * \log_2 (p_2 * \dots * p_n) \\ &= 3 + 4 * \log_2 ((p - 1) / 2) \\ &= -1 + 4 * \log_2 (p - 1) \\ &\leq -1 + 4 * \log_2 p \end{aligned}$$

as claimed. Thus, the nonrecursive, expanded certificate is at most a multiple of the length of the prime.

*John:* Neat! It permits short pedigrees for primes; it says nothing about how the certificate is obtained. Finding the certificate is exclusively the problem of the client. But so far I've only seen certificates for small primes.

*Jill:* Let me give you a certificate for the thirty-digit number

$$N = 909,090,909,090,909,090,909,091.$$

It is: (3; 2, 3, 3, 3, 5, 7, 13, 31, 37, 41, 211, 241, 271, 2161, 9091, 2906161).

If you multiply the sixteen numbers to the right of the semicolon you will obtain  $N - 1$ . On a standard personal computer with a straightforward string exponentiation procedure, you may check in a couple of seconds that  $3^{N-1} \pmod{N}$  is 1.

*John:* So far, so good. But that does not prove that three is a primitive root. You must now show for each  $j$  between 0 and  $N - 1$  that  $3^j \pmod{N}$  is not 1.

*Jill:* No. Although the proof for theorem 2 shows that there always is a certificate with  $b$  a primitive root, there may be other certificates for which  $b$  is not a primitive root. We do not have to prove that the  $b$  in a certificate is a primitive root!

I have shown how to do the first two checks for  $N$ . Using the same string exponentiation procedure I can check the third step for each of the distinct numbers to the right of the semicolon in  $N$ 's short certificate. To complete the proof, I need only demonstrate that each of the numbers to the right of the semicolon in my short certificate is a prime.

*John:* Look, I don't have all day. Clearly, 2, 3, 5, 7, 13, 31, 37, and 41 are prime. I can look up 211, 241, 271, 2161, and 9091 in the standard CRC engineering tables [1]. But since you won't find 2906161 in any standard table of primes, I want to see its certificate.

*Jill:* A certificate for 2906161 is (2; 2, 2, 2, 2, 3, 5, 12109). We check that  $2 * 2 * 2 * 2 * 3 * 5 * 12109 = 2906160$ . We then use the string exponentiation routine to verify that

$$\begin{array}{ll} 2^{2906160} & \text{mod } 2906161 = 1, \\ 2^{2906160/2} & \text{mod } 2906161 = 1453080, \\ 2^{2905160/3} & \text{mod } 2906161 = 968720, \end{array}$$

and...

*John:* OK, I get the idea. But what about 12109?

*Jill:* Its certificate is (6; 2, 2, 3, 1009) and we check that  $2 * 2 * 3 * 1009 = 12108$  and then we...or 12109 is right here in Lehmer's Tables [5]. We now agree that all factors of 2906160 are prime. We need to give an upper bound for the number of steps that it would take if we recursively did every certificate all the way down to 2.

*John:* Since  $-1 + 4 * \log_2 N$  is 419 to the nearest larger whole number, this means that we could outline a complete verification of  $N$  in no more than 419 steps.

*Jill:* When a problem has an algorithm that terminates in time bounded by a polynomial of its input length, the problem is said to belong to the *class P* (for polynomial time). For example, since there is an algorithm that determines in  $O(n^3)$  steps whether a graph of  $n$  vertices is connected, the problem of determining connectivity for a graph belongs to  $P$ . Problems in  $P$  have algorithms that *create the answer* from scratch in polynomial time. A problem is in NP if we can *verify the answer* in polynomial time. The answer can come from such non-deterministic mechanisms as guessing, inspiration, intuition, or deciphering background radiation from Alpha Centauri. We have just seen that the problem 'Is  $x$  prime?' belongs to the class NP. Namely, if we guess a certificate, then we can verify the guess in polynomial (linear) time of the input length.

*John:* What about composite numbers?

*Jill:* They are examples of the class CO-NP. A problem belongs to CO-NP if we can verify in polynomial time that something does *not* have a property. For example, it is easy to verify in polynomial time that  $x$  is not a prime ( $x$  is composite): Simply show that the product of two nontrivial integers is  $x$ , the number in question. Since we can multiply two  $N$ -digit factors in time bounded by  $N^2$ , we can verify in polynomial time that a number  $x$  is not a prime. Thus, the problem 'Is  $x$  a prime?' is in both NP and CO-NP.

*John:* Wait! Doesn't this have something to do with NP-completeness?

*Jill:* Yes. Since the early 1970s computer scientists have discovered hundreds of problems that are called *NP-complete*. The NP-complete problems are the *hardest* NP problems. The shortest traveling salesperson tour, the chromatic number of a graph, and the satisfiability of propositional formulas are examples of these NP-complete problems. If a single NP-complete problem can be shown to belong to  $P$ , then in fact all NP problems belong to  $P$ . It is widely conjectured that no NP-complete problem is in  $P$ —that is, it is conjectured that  $P \neq \text{NP}$ . If NP equals  $P$ , then CO-NP must equal  $P$ . But it is possible that  $\text{NP} \neq P$ , while  $\text{NP} = \text{CO-NP}$ . In this direction,  $\text{CO-NP} = \text{NP}$  if and only if some NP-complete problem is in  $P$ . It is conjectured that  $\text{NP} \neq \text{CO-NP}$ , that is, there are no NP-complete problems that are both in NP and in CO-NP. The primes are the first example of a problem that is not known to be in  $P$ , but for which there are polynomial procedures for both verification and falsification.

## REFERENCES

1. C. Hodgman, *CRC Standard Mathematical Tables*, Chemical Rubber Publishing, Cleveland, 1959.
2. D. E. Knuth, *The Art of Computer Programming, Fundamental Algorithms*, Vol. 1, 2nd edition, Addison-Wesley Publishing Co., Reading, MA, 1973.
3. D. E. Knuth, *The Art of Computer Programming, Seminumerical Algorithms*, Vol. 2, 2nd edition, Addison-Wesley Publishing Co., Reading, MA, 1981.
4. D. H. Lehmer, *Tests for primality by the converse of Fermat's Theorem*, Bulletin Amer. Math. Society, 33 (1927), 327–340.
5. D. N. Lehmer, *List of Prime Numbers from 1 to 10,006,721*, Hafner, New York, 1956.
6. V. R. Pratt, *Every prime has a succinct certificate*, SIAM J. Computation, 4 (1975), 214–220.
7. R. L. Rivest, A. Shamir, and L. Adleman, *A method for obtaining digital signatures and public-key cryptosystems*, Communications of the ACM, 21, (1978)

---

## Pictures, Projections, and Proverbs

RICHARD L. FRANCIS

Southeast Missouri State University  
Cape Girardeau, MO 63701

During the warm, humid days of a long past summer, I was, as a graduate student at the University of Missouri, caught up in the exciting world of projective geometry. An out-of-the-ordinary course and an inspiring professor were, as I now recall, all the ingredients for a challenging summer. Rarely did this professor, who loved to quote famous sayings, send his students to the blackboard without asking them to draw the appropriate geometric figure as well. “Remember,” he said, “a picture is worth a thousand words.”

Not only was the course enlightening, it was at times overwhelming. In particular, the professor had just announced an assignment of a special paper on the subject of finite geometries; its length was to be *two thousand words*. Two thousand words seemed to me excessive for such an unfamiliar subject, and, inspired by his oft-quoted proverb as to the word value of a single drawing, I asked, not all so seriously, if I might submit two pictures (just two). Needless to say, the words of the proverb failed to cover the situation at hand. Accordingly, much of my summer's hard work was taken up geometrizing, not in pictures, but in words (many words).

*John:* Wait! Doesn't this have something to do with NP-completeness?

*Jill:* Yes. Since the early 1970s computer scientists have discovered hundreds of problems that are called *NP-complete*. The NP-complete problems are the *hardest* NP problems. The shortest traveling salesperson tour, the chromatic number of a graph, and the satisfiability of propositional formulas are examples of these NP-complete problems. If a single NP-complete problem can be shown to belong to  $P$ , then in fact all NP problems belong to  $P$ . It is widely conjectured that no NP-complete problem is in  $P$ —that is, it is conjectured that  $P \neq \text{NP}$ . If NP equals  $P$ , then CO-NP must equal  $P$ . But it is possible that  $\text{NP} \neq P$ , while  $\text{NP} = \text{CO-NP}$ . In this direction,  $\text{CO-NP} = \text{NP}$  if and only if some NP-complete problem is in  $P$ . It is conjectured that  $\text{NP} \neq \text{CO-NP}$ , that is, there are no NP-complete problems that are both in NP and in CO-NP. The primes are the first example of a problem that is not known to be in  $P$ , but for which there are polynomial procedures for both verification and falsification.

## REFERENCES

1. C. Hodgman, *CRC Standard Mathematical Tables*, Chemical Rubber Publishing, Cleveland, 1959.
2. D. E. Knuth, *The Art of Computer Programming, Fundamental Algorithms*, Vol. 1, 2nd edition, Addison-Wesley Publishing Co., Reading, MA, 1973.
3. D. E. Knuth, *The Art of Computer Programming, Seminumerical Algorithms*, Vol. 2, 2nd edition, Addison-Wesley Publishing Co., Reading, MA, 1981.
4. D. H. Lehmer, *Tests for primality by the converse of Fermat's Theorem*, Bulletin Amer. Math. Society, 33 (1927), 327–340.
5. D. N. Lehmer, *List of Prime Numbers from 1 to 10,006,721*, Hafner, New York, 1956.
6. V. R. Pratt, *Every prime has a succinct certificate*, SIAM J. Computation, 4 (1975), 214–220.
7. R. L. Rivest, A. Shamir, and L. Adleman, *A method for obtaining digital signatures and public-key cryptosystems*, Communications of the ACM, 21, (1978)

## Pictures, Projections, and Proverbs

RICHARD L. FRANCIS

Southeast Missouri State University  
Cape Girardeau, MO 63701

During the warm, humid days of a long past summer, I was, as a graduate student at the University of Missouri, caught up in the exciting world of projective geometry. An out-of-the-ordinary course and an inspiring professor were, as I now recall, all the ingredients for a challenging summer. Rarely did this professor, who loved to quote famous sayings, send his students to the blackboard without asking them to draw the appropriate geometric figure as well. "Remember," he said, "a picture is worth a thousand words."

Not only was the course enlightening, it was at times overwhelming. In particular, the professor had just announced an assignment of a special paper on the subject of finite geometries; its length was to be *two thousand words*. Two thousand words seemed to me excessive for such an unfamiliar subject, and, inspired by his oft-quoted proverb as to the word value of a single drawing, I asked, not all so seriously, if I might submit two pictures (just two). Needless to say, the words of the proverb failed to cover the situation at hand. Accordingly, much of my summer's hard work was taken up geometrizing, not in pictures, but in words (many words).



---

# NOTES

---

## An Advanced Calculus Approach to Finding the Fermat Point

MOWAFFAQ HAJJA  
Yarmouk University  
Irbid, Jordan

Steiner's problem, or Fermat's problem to Torricelli as it is sometimes called, asks for the location of the point in the plane of a given triangle whose distances from the vertices have a minimum sum. Several noncalculus solutions can be found in [2, pp. 24–34], [5, pp. 156–162], [1, pp. 354–361], and in many other books. However, one wonders why this beautiful extremum problem is not usually presented to students of advanced calculus. Possibly the negative speculation of D. C. Kay in his well-known *College Geometry* [3, p. 271] has tended to direct instructors away from the problem: "Any attempt to solve this by means of calculus would most probably end in considerable frustration." In any case, the following simple solution using calculus shows that it is well within the reach of a student of advanced calculus. It also extends to a *correct* solution of the *complementary* problem discussed by R. Courant and H. R. Robbins in their admirable *What is Mathematics?* [1, Chapter VII, §5.3, p. 358].

Let  $P_1 = (x_1, y_1)$ ,  $P_2 = (x_2, y_2)$ , and  $P_3 = (x_3, y_3)$  be the vertices of the given triangle, and for any  $P = (x, y)$  in the plane, let  $r_i$ ,  $\alpha_i$  and  $A_i$ ,  $i = 1, 2, 3$ , be assigned as in FIGURE 1. Thus

$$r_1 = PP_1 (=d(P, P_1)), \quad \alpha_1 = \angle P_2PP_3, \quad A_1 = \text{area of } \Delta P_2PP_3$$

and  $\alpha_i$  is not defined if  $P$  is a vertex other than  $P_i$ . Then the function that we want to minimize is

$$f(P) = f(x, y) = r_1 + r_2 + r_3 = \sum_{i=1}^3 \left[ (x - x_i)^2 + (y - y_i)^2 \right]^{1/2}$$

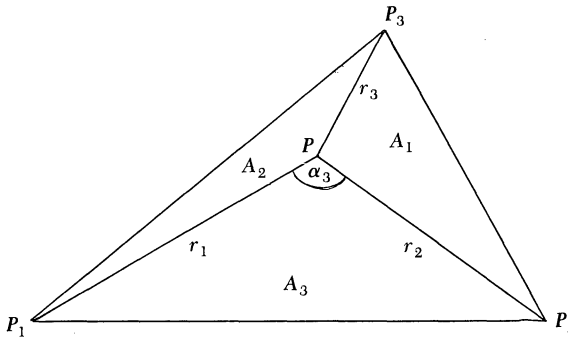


FIGURE 1

As noted in [4, p. 107], if  $P$  is properly outside the triangle then there is a line  $L$  that strictly separates  $P$  from the triangle (FIGURE 2). If  $Q$  is the foot of the line from  $P$  perpendicular to  $L$ , then for  $i = 1, 2, 3$ ,  $QP_i < PP_i$  and therefore  $f(Q) < f(P)$ . This shows that  $f$  has a minimum and that the minimum is attained at some point in or on the triangle.

Also, the only points in the plane at which  $\partial f/\partial x$  and  $\partial f/\partial y$  do not exist are the vertices of the triangle. Thus at every *critical* point  $P = (x, y)$  other than these, we have

$$\partial f/\partial x = 0 \quad \text{and} \quad \partial f/\partial y = 0.$$

These equations directly yield

$$(x - x_1, x - x_2, x - x_3) \cdot (1/r_1, 1/r_2, 1/r_3) = 0$$

$$(y - y_1, y - y_2, y - y_3) \cdot (1/r_1, 1/r_2, 1/r_3) = 0$$

where “ $\cdot$ ” represents the dot product of vectors (or points) in  $\mathbb{R}^3$ . Letting

$$\xi = (x - x_1, x - x_2, x - x_3), \quad \eta = (y - y_1, y - y_2, y - y_3), \quad \text{and}$$

$$\rho = (1/r_1, 1/r_2, 1/r_3),$$

these equations are simply

$$\xi \cdot \rho = \eta \cdot \rho = 0.$$

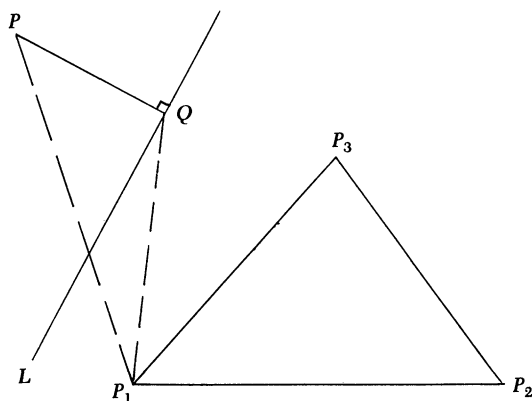


FIGURE 2

Since  $\xi$  and  $\eta$  are linearly independent (because  $P_1$ ,  $P_2$ , and  $P_3$  are not collinear), then  $\rho$ , being perpendicular to both  $\xi$  and  $\eta$ , is parallel to their cross product

$$\begin{aligned} \xi \times \eta &= \begin{bmatrix} i & j & k \\ x - x_1 & x - x_2 & x - x_3 \\ y - y_1 & y - y_2 & y - y_3 \end{bmatrix} \\ &= \left( \begin{vmatrix} x - x_2 & x - x_3 \\ y - y_2 & y - y_3 \end{vmatrix}, \begin{vmatrix} x - x_3 & x - x_1 \\ y - y_3 & y - y_1 \end{vmatrix}, \begin{vmatrix} x - x_1 & x - x_2 \\ y - y_1 & y - y_2 \end{vmatrix} \right) \\ &= (\mp 2A_1, \mp 2A_2, \mp 2A_3) \\ &= (\mp (r_2 r_3 \sin \alpha_1, r_3 r_1 \sin \alpha_2, r_1 r_2 \sin \alpha_3)), \end{aligned}$$

because all the coordinates of the parallel vector  $\rho$  have the same (positive) sign.

Therefore, for some  $\lambda$ ,

$$\rho = (1/r_1, 1/r_2, 1/r_3) = \lambda(r_2 r_3 \sin \alpha_1, r_3 r_1 \sin \alpha_2, r_1 r_2 \sin \alpha_3).$$

Hence

$$\sin \alpha_1 = \sin \alpha_2 = \sin \alpha_3 \quad (= 1/(\lambda r_1 r_2 r_3)). \quad (1)$$

Therefore, a nonvertex that minimizes  $f$  must lie inside the triangle and must have  $\alpha_1 = \alpha_2 = \alpha_3 = 120^\circ$  (for if two  $\alpha_i$ 's are supplementary for an inner point, then the third would be  $180^\circ$ , contradicting (1)). If none of the angles of our triangle is  $\geq 120^\circ$ , then there is a unique such point, called the (*interior*) *Fermat point* whose easy construction is given in the first figure of [2, p. 25] and at which  $f$  can be shown to attain its minimum. Otherwise, the minimum of  $f$  is attained at a vertex of the triangle and it is easy to see that it is attained at the vertex holding the largest angle.

It is interesting to note that if a point other than a vertex minimizes the similar expression

$$g(x, y) = r_1 + r_2 - r_3, \quad (2)$$

then the *same* equation (1) must be satisfied at that point. This is obtained by carrying out the same computations above with  $\rho$  replaced by  $\tilde{\rho} = (1/r_1, 1/r_2, -1/r_3)$ . However, a point  $P$  that minimizes  $g$  cannot lie inside the triangle since the intersection (FIGURE 3)  $Q$  of  $P_3P$  and  $P_1P_2$  gives

$$g(Q) = QP_1 + QP_2 - QP_3 = P_1P_2 - QP_3 < PP_1 + PP_2 - PP_3 = g(P).$$

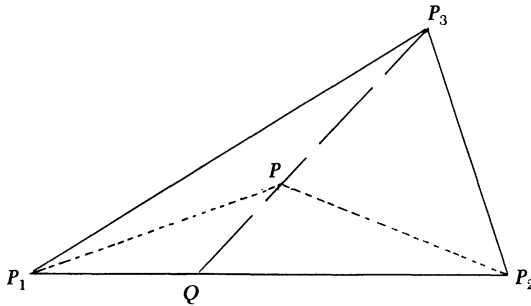


FIGURE 3

Thus a nonvertex point that minimizes  $g$  cannot lie inside the triangle and, having to satisfy (1), it must have  $(\alpha_1, \alpha_2, \alpha_3) = (120^\circ, 60^\circ, 60^\circ)$  in some order. Calling such points *exterior Fermat points*, we proceed to examine their existence and uniqueness and their significance in minimizing such expressions as  $g$ . We show in particular that, in contrast to the impression and statements given in [1], *every* triangle has an external Fermat point unless exactly one of its angles is  $60^\circ$ , and that an exterior Fermat point does minimize  $g$  (or a similar expression) precisely when two angles of the triangle are  $\geq 60^\circ$  each (*and not when one of the angles is  $> 120^\circ$* ).

Suppose  $P_0$  is an exterior Fermat point of  $\Delta P_1P_2P_3$ . For simplicity, let  $\alpha_3$  be the largest among  $\alpha_1, \alpha_2, \alpha_3$ . Then  $\alpha_3 (= \angle P_1P_0P_2) = 120^\circ$ ,  $\alpha_1 = \alpha_2 = 60^\circ$  and  $P_3$  lies on the bisector  $L$  of  $\angle P_1P_0P_2$  (FIGURE 4). Let  $L$  intersect the circle circumscribing  $P_1P_0P_2$  at  $X$  and  $P_1P_2$  at  $O$ . Applying Ptolemy's Theorem [3, p. 8] to the cyclic quadrilateral  $P_1P_0P_2X$  and using the fact that  $\Delta P_1P_2X$  is equilateral, one concludes that

$$P_0X = P_0P_1 + P_0P_2. \quad (3)$$

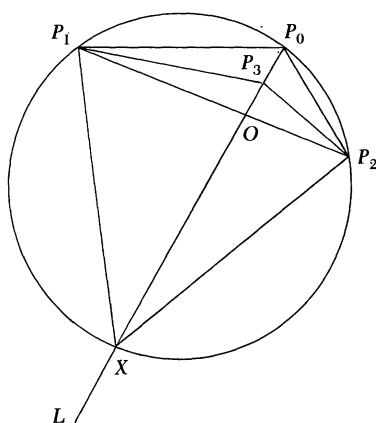


FIGURE 4

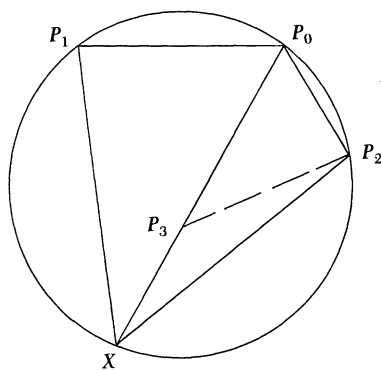


FIGURE 5

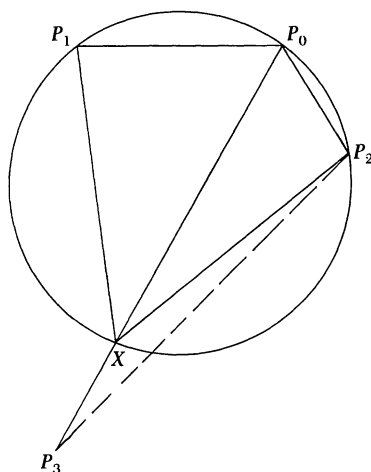


FIGURE 6

We now use this to show that if  $P_3$  lies on the segment  $P_0X$  (and in particular on  $P_0O$ ) (FIGURE 5) then, contrary to the claim made in [1],  $P_0$  does *not* minimize the expression  $g(P)$  of (2); while if  $P_3$  does *not* lie on  $P_0X$  (FIGURE 6), then  $P_0$  *does* minimize  $g$ .

Suppose that  $P_3$  lies on the (open) segment  $P_0X$  (FIGURE 5). (Note that this includes the possibility that  $P_3$  lies on  $P_0O$ , which corresponds to the (only) configuration considered in [1]). Then

$$\begin{aligned} g(P_0) &= P_0P_1 + P_0P_2 - P_0P_3 = P_0X - P_0P_3 = P_3X \\ g(P_2) &= P_2P_1 + 0 - P_2P_3 = P_2X - P_2P_3. \end{aligned}$$

By the triangle inequality, we have  $P_2X < P_2P_3 + P_3X$ . Therefore  $g(P_2) < g(P_0)$  and hence  $g$  does not take its minimum at  $P_0$ . In fact  $g$  is minimum at the vertex with the smallest angle. Also, to remove all possible remaining doubts concerning the validity of the statement about  $P_0$  made in [1], one can similarly show that  $P_0$  does not minimize any of the similar expressions

$$h(P) = r_2 + r_3 - r_1, \quad k(P) = r_3 + r_1 - r_2.$$

Now suppose that  $P_3$  does not lie on the segment  $P_0X$  (FIGURE 6). Then

$$g(P_0) = P_0P_1 + P_0P_2 - P_0P_3 = P_0X - P_0P_3 = P_3X.$$

Since

$$g(P_2) = P_2P_1 + 0 - P_2P_3 = P_2X - P_2P_3 > -P_3X = g(P_0),$$

$$g(P_1) = 0 + P_1P_2 - P_1P_3 = P_1X - P_1P_3 > -P_3X = g(P_0),$$

and since  $g(P_3) > 0$ , then  $g$  takes its minimum at  $P_0$ . It is also easy to see that neither  $h$  nor  $k$  is minimized at  $P_0$ .

This completes the description of the role played by an external Fermat point  $P_0$  of a triangle  $P_1P_2P_3$  in terms of how a specific vertex (say  $P_3$ ) is located relative to  $P_0, P_1, P_2$ . The next theorem describes how to locate the Fermat points of a given triangle. Its first part is evident while the last part follows from examining how  $P_0, P_1, P_2, P_3$  in FIGURES 4 and 5 are inter-located and from the previous discussion of interior Fermat points.

**THEOREM.** *Let  $P_1P_2P_3$  be a triangle none of whose angles is  $60^\circ$  or  $120^\circ$ . Let any (of the two possible) equilateral triangles, say  $P_1P_2X$ , be drawn on one of the sides. If the straight line passing through  $X$  and  $P_3$  intersects the circle circumscribing  $P_1P_2X$  at a point  $F$ , then  $F$  is a Fermat point. Moreover, all Fermat points are obtained this way using for the side  $P_1P_2$  the (unique) side whose two angles  $P_1$  and  $P_2$  are  $> 60^\circ$  each (FIGURE 7) or  $< 60^\circ$  each (FIGURE 8) and using both equilateral triangles  $P_1P_2X$  and  $P_1P_2X'$  that can be drawn on  $P_1P_2$ . Consequently, the triangle has exactly two Fermat points, one on each of the two small arcs  $P_1P_2$ .*

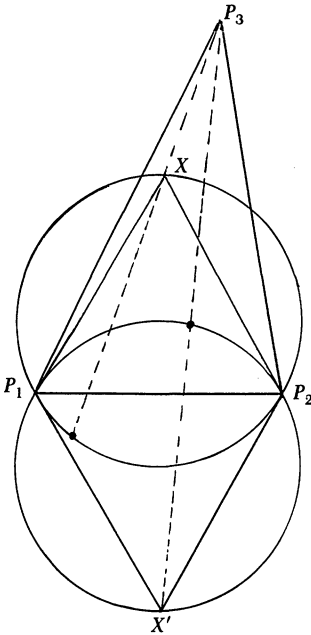


FIGURE 7

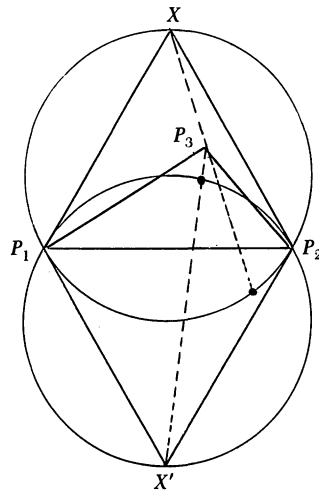


FIGURE 8

The artificial restriction on the angles of the triangle can be relaxed if we agree to call a vertex holding  $60^\circ$  or  $120^\circ$  a (*boundary*) Fermat point. Let us further agree to call a Fermat point *significant* if it minimizes any of the expressions  $f, g, h, k$  described above. With this terminology, the following theorem is quite evident.

**THEOREM.** *If a triangle is nonequilateral, then it has exactly two Fermat points. At most, one of them is boundary (and that happens precisely at the vertex, if any, at which the angle is  $60^\circ$  or  $120^\circ$ ) and, at most, one of them is interior (and that happens precisely when none of the angles is  $> 120^\circ$ ). Every interior or boundary Fermat point is significant while an exterior Fermat point is significant if, and only if, two angles of the triangle are  $> 60^\circ$  each. If the triangle is equilateral, then (by (3)) every point on its circumscribing circle is a significant Fermat point.*

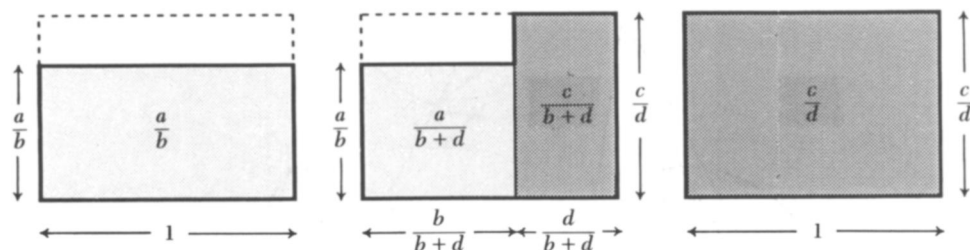
**Acknowledgment.** The author thanks Ross Honsberger for his invaluable suggestions.

## REFERENCES

1. Courant, Richard, and Herbert R. Robbins, *What is Mathematics?*, 4th edition, Oxford University Press, New York, 1978.
2. Honsberger, Ross, *Mathematical Gems I*, Dolciani Mathematical Expositions, No. 1, MAA, Washington, DC, 1973.
3. Kay, D. C., *College Geometry*, Holt, Rinehart, and Winston, New York, 1969.
4. Melzak, Z. A., *Invitation to Geometry*, John Wiley & Sons, Inc., New York, 1983.
5. Niven, Ivan, *Maxima and Minima without Calculus*, Dolciani Mathematical Expositions, No. 6, MAA, Washington, DC, 1981.

## Proof without Words: Regle des Nombres Moyens

[Nicolas Chuquet, *Le Triparty en la Science des Nombres*, 1484]



$$a, b, c, d > 0; \quad \frac{a}{b} < \frac{c}{d} \quad \rightarrow \quad \frac{a}{b} < \frac{a+c}{b+d} < \frac{c}{d}.$$

$$\frac{a}{b} < \frac{a}{b+d} + \frac{c}{b+d} < \frac{c}{d}$$

**THEOREM.** *If a triangle is nonequilateral, then it has exactly two Fermat points. At most, one of them is boundary (and that happens precisely at the vertex, if any, at which the angle is  $60^\circ$  or  $120^\circ$ ) and, at most, one of them is interior (and that happens precisely when none of the angles is  $> 120^\circ$ ). Every interior or boundary Fermat point is significant while an exterior Fermat point is significant if, and only if, two angles of the triangle are  $> 60^\circ$  each. If the triangle is equilateral, then (by (3)) every point on its circumscribing circle is a significant Fermat point.*

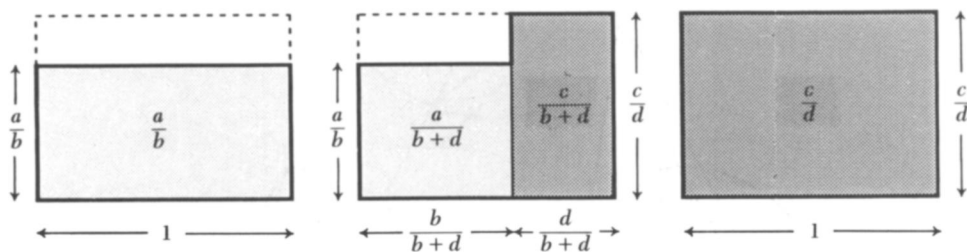
**Acknowledgment.** The author thanks Ross Honsberger for his invaluable suggestions.

## REFERENCES

1. Courant, Richard, and Herbert R. Robbins, *What is Mathematics?*, 4th edition, Oxford University Press, New York, 1978.
2. Honsberger, Ross, *Mathematical Gems I*, Dolciani Mathematical Expositions, No. 1, MAA, Washington, DC, 1973.
3. Kay, D. C., *College Geometry*, Holt, Rinehart, and Winston, New York, 1969.
4. Melzak, Z. A., *Invitation to Geometry*, John Wiley & Sons, Inc., New York, 1983.
5. Niven, Ivan, *Maxima and Minima without Calculus*, Dolciani Mathematical Expositions, No. 6, MAA, Washington, DC, 1981.

## Proof without Words: Regle des Nombres Moyens

[Nicolas Chuquet, *Le Triparty en la Science des Nombres*, 1484]



$$a, b, c, d > 0; \quad \frac{a}{b} < \frac{c}{d} \quad \rightarrow \quad \frac{a}{b} < \frac{a+c}{b+d} < \frac{c}{d}.$$

$$\frac{a}{b} < \frac{a}{b+d} + \frac{c}{b+d} < \frac{c}{d}$$

# Remarks on Problem B-3 on the 1990 William Lowell Putnam Mathematical Competition\*

JIUQIANG LIU  
ALLEN J. SCHWENK  
Western Michigan University  
Kalamazoo, MI 49008

## Problem B-3 [1]

Let  $S$  be a set of  $2 \times 2$  integer matrices whose entries  $a_{ij}$

- (1) are all squares of integers, and,
- (2) satisfy  $a_{ij} \leq 200$ .

Show that if  $S$  has more than 50,387 ( $= 15^4 - 15^2 - 15 + 2$ ) elements, then it has two elements that commute.

## Solution

Since the number 50,387 is far from the best possible, there are many potential solutions. For example, consider the three subsets

$$\begin{aligned} M_1 &= \left\{ \begin{pmatrix} a^2 & b^2 \\ b^2 & a^2 \end{pmatrix} : a \neq b \right\} \\ M_2 &= \left\{ \begin{pmatrix} a^2 & b^2 \\ 0 & a^2 \end{pmatrix} : b \neq 0 \right\} \\ M_3 &= \left\{ \begin{pmatrix} a^2 & 0 \\ c^2 & a^2 \end{pmatrix} : c \neq 0 \right\}. \end{aligned}$$

It is easy to verify that these three sets are disjoint and any two matrices in the same  $M_i$  must commute. Since each set has order  $15 \cdot 14 = 210$ , any set with more than  $49,998 = 15^4 - 3 \cdot 210 + 3$  elements must contain two from the same  $M_i$ , that is, a commuting pair. In fact, any two of these three subsets can be used to establish a solution within the number allowed in the problem.

Let  $U$  be the collection of  $50,625 = 15^4$  matrices of the form  $\begin{pmatrix} a^2 & b^2 \\ c^2 & d^2 \end{pmatrix}$  that satisfy conditions (1) and (2). While it is beyond what students could reasonably be expected to solve during the exam, clearly the elegant question is: What is the maximum order of a subset  $S \subseteq U$  in which no two elements commute? We shall show that this order is 32,390, so that any subset of 32,391 matrices in  $U$  must include a commuting pair.

The method we use can be applied to numerous variations of the problem posed. For our purposes, the requirement of perfect square entries creates an extra level of difficulty. Consequently, let us first solve a variation that seems more natural. We shall then return to the original question of finding the largest totally noncommuting set  $S$  when the entries must be perfect squares. Specifically, we ask what is the order of a maximum noncommuting set among all  $2 \times 2$  matrices  $\begin{pmatrix} a & b \\ c & d \end{pmatrix}$  with nonnegative integer entries bounded by  $0 \leq a, b, c, d \leq n$ . In this case our universal set  $U$  has  $(n+1)^4$  matrices in it.

---

\*Research supported in part by the Office of Naval Research, Contract N001491-J-1364.



We begin by partitioning  $U$  into four subsets depending on the zero/nonzero character of the entries  $b$  and  $c$ . Namely,  $U = S_{00} \cup S_{b0} \cup S_{0c} \cup S_{bc}$ . In each  $S_{ij}$ , if the first subscript is 0, then  $b = 0$ , while, if it is  $b$ , then  $b \neq 0$ . Similarly, the second subscript determines the character of  $c$ . Next, we further partition  $S_{00}$  into  $Z \cup S'_{00}$  where  $Z$  consists of all the scalars  $aI$  and  $S'_{00}$  has all the matrices of  $S_{00}$  for which  $a \neq d$ . We shall examine a typical matrix from each of the five  $S_{ij}$  to see which other matrices commute with it.

First, the matrices in  $Z$  commute with everything, so we dare not include any in a noncommuting set  $S$ .

Next, suppose a matrix  $\begin{pmatrix} a & 0 \\ 0 & d \end{pmatrix} \in S'_{00}$  commutes with an arbitrary matrix

$$X = \begin{pmatrix} w & x \\ y & z \end{pmatrix}.$$

Since  $a \neq d$ , we must have  $x = y = 0$ . But then, for all choices of  $w$  and  $z$ , the matrices commute. Thus, each matrix in  $S'_{00}$  commutes precisely with the matrices in  $S_{00}$ . Consequently, at most one matrix from  $S'_{00}$  may be included in  $S$ .

Now any upper triangular matrix  $\begin{pmatrix} a & b \\ 0 & d \end{pmatrix} \in S_{b0}$  commuting with

$$X = \begin{pmatrix} w & x \\ y & z \end{pmatrix}$$

immediately forces  $y = 0$ , so we have the form  $X = \begin{pmatrix} w & x \\ 0 & z \end{pmatrix}$ . Next we find that  $x(a - d) = b(w - z)$ . Now if  $x = 0$ , we are forced to have  $w = z$  and  $X \in Z$ . Otherwise,  $X \in S_{b0}$ , but even stronger, we must have

$$\frac{a - d}{b} = \frac{w - z}{x}.$$

That is, two matrices in  $S_{b0}$  commute if, and only if, they have the same value for

$$t = \frac{a - d}{b}.$$

Consequently, the number of matrices from  $S_{b0}$  that can be included in a noncommuting set  $S$  is precisely the number of different values that can occur for  $t$ . Obviously, there is a single value of  $t = 0$  possible, and just as many positive values as negative ones. We use inclusion and exclusion to count the positive values,  $N^+(t)$ . Our initial estimate is that there are  $n$  possible numerators and  $n$  possible denominators, so  $N^+(t) = n^2$ , but we have overcounted values of  $t$  that can be produced in more than one way. Whenever a prime  $p$  divides both numerator and denominator, the ratio  $t$  has been counted twice, so we should subtract the number of ways this can happen, specifically  $\lfloor n/p \rfloor^2$ . But now if a fraction has two prime cancellable factors, say  $p$  and  $q$ , we have taken it out of our counting formula twice, so we need to reinsert it by adding  $\lfloor n/pq \rfloor^2$ . Continuing in this way, we subtract terms with three prime factors, and so on. The final result is the formula

$$N^+(t) = n^2 - \sum_{p \leq n \text{ is prime}} \left\lfloor \frac{n}{p} \right\rfloor^2 + \sum_{p, q \leq n \text{ are prime}} \left\lfloor \frac{n}{pq} \right\rfloor^2 - \cdots.$$

Since there are equally many negative values for  $t$ , plus the single value  $t = 0$ , we conclude that we may select  $1 + 2N^+(t)$  matrices from  $S_{b0}$ . The analysis for  $S_{0c}$  is precisely the same, except we must define  $t = (a - d/c)$ . Therefore,  $S_{0c}$  can also contribute at most  $1 + 2N^+(t)$  elements to the noncommuting set  $S$ .

Finally, it remains to analyze  $S_{bc}$  where both  $b$  and  $c$  must be nonzero. In this case, we define two quantities,

$$t = \frac{a-d}{b} \text{ and } s = \frac{c}{b}.$$

First, if  $\begin{pmatrix} a & b \\ c & d \end{pmatrix} \in S_{bc}$  commutes with any  $X = \begin{pmatrix} w & x \\ y & z \end{pmatrix}$ , we find that  $by = cx$ . Now  $x = 0$  forces  $y = 0$  and subsequently  $w = z$  so that  $X \in Z$ . Otherwise,  $x \neq 0$  forces  $y \neq 0$  and  $X \in S_{bc}$ . But, we also have

$$s = \frac{c}{b} = \frac{y}{x}$$

and

$$t = \frac{a-d}{b} = \frac{w-z}{x}.$$

That is, two matrices from  $S_{bc}$  commute if, and only if, they have precisely the same pair  $(s, t)$  of parameters. Again we can use inclusion and exclusion to count these pairs. By definition,  $s$  is positive, but let us first count the number of pairs  $N(s, t^+)$  with positive  $t$ 's. Our original crude count is  $n$  choices for  $b$ , likewise  $n$  choices for  $c$ , and  $n$  possible positive numerators  $a - d$ . This gives  $n^3$  positive pairs. Again we must adjust the count for fractions that can be reduced, but we shall only reduce a pair by a factor of  $p$  when both  $s$  and  $t$  admit this cancellation. Analogous to the analysis in the  $S_{b0}$  case, we find

$$N(s, t^+) = n^3 - \sum_{p \leq n \text{ is prime}} \left\lfloor \frac{n}{p} \right\rfloor^3 + \sum_{p, q \leq n \text{ are prime}} \left\lfloor \frac{n}{pq} \right\rfloor^3 - \cdots.$$

Now the same formula counts the number of choices with negative values for  $t$ , but in this case we have to be more careful analyzing the choices for  $N(s, 0)$ . The value  $t = 0$  forces  $a = d$ , but we must still count how many fractions we can form for  $s = (c/b)$ . But this is precisely the same problem we solved above when we counted  $N^+(t)$ . That is,  $N(s, 0) = N^+(t)$ . Therefore, the set  $S_{bc}$  can contribute  $2N(s, t^+) + N^+(t)$  to the set  $S$ .

In conclusion, we have partitioned  $U$  into five subsets. The table below summarizes what we have determined about these subsets.

Subset	Defining conditions	Order	Maximum number contributed to S
$Z$	$b = c = 0; a = d$	$n + 1$	0
$S'_{00}$	$b = c = 0; a \neq d$	$(n + 1)n$	1
$S_{b0}$	$b \neq 0; c = 0$	$(n + 1)^2 n$	$1 + 2N^+(t)$
$S_{0c}$	$b = 0; c \neq 0$	$(n + 1)^2 n$	$1 + 2N^+(t)$
$S_{bc}$	$b \neq 0; c \neq 0$	$(n + 1)^2 n^2$	$2N(s, t^+) + N^+(t)$
Totals		$(n + 1)^4$	$2N(s, t^+) + 5N^+(t) + 3$

The maximum order of  $S$  is  $2N(s, t^+) + 5N^+(t) + 3$ . For large  $n$  we can estimate  $|S|$  by

$$2n^3 \prod_{p \text{ prime}} \left(1 - \frac{1}{p^3}\right) + 5n^2 \prod_{p \text{ prime}} \left(1 - \frac{1}{p^2}\right) + 3.$$

The infinite products converge to approximately 0.8319 and 0.6079, so that  $|S| \approx 1.6638n^3 + 3.0395n^2 + 3$ .

Now let us return to the problem in which the entries are perfect squares between 0 and  $n^2$ . The partition into sets  $Z$ ,  $S'_{00}$ ,  $S_{b0}$ ,  $S_{0c}$ , and  $S_{bc}$  proceeds just as above. The only difference is that now the formula for  $t$  is  $((a^2 - d^2)/b^2)$  and the formula for  $s$  is  $c^2/b^2$ . Moreover, when we attempt to count the possible fractions for  $t$  by inclusion and exclusion, we must contend with the complexity of determining how many numerators are possible for  $t$ . Previously, there were simply  $n$  possible numerators. But now we need to know how many values between 1 and  $n^2$  can be written as  $a^2 - d^2$  with  $a$  and  $d$  both between 0 and  $n$ . We shall denote this set of possible positive numerators as  $\text{Num}^+$ . Next, we write  $N((a^2 - d^2)/b^2)$  for the number of times

$$\frac{a^2 - d^2}{b^2} \in \text{Num}^+.$$

We find the critical formulas are

$$\begin{aligned} N^+(t) &= n \cdot |\text{Num}^+| - \sum_{p \leq n \text{ is prime}} \left\lfloor \frac{n}{p} \right\rfloor \cdot N\left(\frac{a^2 - d^2}{p^2}\right) \\ &\quad + \sum_{p, q \leq n \text{ are prime}} \left\lfloor \frac{n}{pq} \right\rfloor \cdot N\left(\frac{a^2 - d^2}{p^2 q^2}\right) + \cdots \\ N(s, t^+) &= n^2 \cdot |\text{Num}^+| - \sum_{p \leq n \text{ is prime}} \left\lfloor \frac{n}{p} \right\rfloor^2 \cdot N\left(\frac{a^2 - d^2}{p^2}\right) \\ &\quad + \sum_{p, q \leq n \text{ are prime}} \left\lfloor \frac{n}{pq} \right\rfloor^2 \cdot N\left(\frac{a^2 - d^2}{p^2 q^2}\right) - \cdots \\ N(s, 0) &= n^2 - \sum_{p \leq n \text{ is prime}} \left\lfloor \frac{n}{p} \right\rfloor^2 + \sum_{p, q \leq n \text{ are prime}} \left\lfloor \frac{n}{pq} \right\rfloor^2 - \cdots \end{aligned}$$

For arbitrary  $n$ , we do not know how to determine the set  $\text{Num}^+$  or the values  $N((a^2 - d^2)/b^2)$ , so this is the best we can do. However, to analyze the original problem in which  $n = 14$ , we can proceed by detailed inspection to find that there are 81 possible positive numerators of the form  $a^2 - d^2$ ; specifically,

$$\begin{aligned} \text{Num}^+ = \{ &1, 3, 4, 5, 7, 8, 9, 11, 12, 13, 15, 16, 17, 19, 20, 21, 23, 24, 25, 27, 28, 32, 33, 35, \\ &36, 39, 40, 44, 45, 48, 49, 51, 52, 55, 56, 57, 60, 63, 64, 65, 69, 72, 75, 77, 80, \\ &81, 84, 85, 88, 91, 95, 96, 99, 100, 105, 108, 112, 115, 117, 119, 120, 121, 128, \\ &132, 133, 135, 140, 143, 144, 147, 153, 160, 165, 168, 169, 171, 180, 187, 192, \\ &195, 196\}. \end{aligned}$$

Roughly then,  $N^+(t) = 14 \cdot 81$ , except that certain values of  $t$  will be counted more than once, whenever it is possible to cancel a prime factor  $p^2$  from both the denominator  $b^2$  and the numerator  $a^2 - d^2$ . But to cancel  $p^2$  from the numerator, it is not enough that  $(a^2 - d^2)/p^2$  be an integer. In addition we must have  $(a^2 - d^2)/p^2 \in \text{Num}^+$ . Thus, for  $p = 2$ , we find that 26 numerators, namely

$$\{4, 12, 16, 20, 28, 32, 36, 44, 48, 52, 60, 64, 80, 84, 96, 100, 108, 112, 128, 132, \\ 140, 144, 160, 180, 192, 196\},$$

allow cancellation, but eight others,  $\{8, 24, 40, 56, 72, 88, 120, 168\}$ , do not because dividing by 4 leaves a quotient not in  $\text{Num}^+$ . Together with the seven possible even denominators, these 26 numerators lead us to subtract  $7 \cdot 26$  from our original count.

Similarly, for  $p = 3$  we must subtract  $4 \cdot 15$ ; for  $p = 5$  we must subtract  $2 \cdot 3$ ; for  $p = 7$  we also subtract  $2 \cdot 3$ ; for  $p = 11$  we must subtract  $1 \cdot 1$ ; and finally for  $p = 13$  we must subtract  $1 \cdot 1$ . But now for  $6 = 2 \cdot 3$  we have subtracted twice, so we must add a correction whenever  $6^2$  divides both numerator and denominator. There are two denominators, specifically  $6^2$  and  $12^2$ , and there are four numerators, 36, 108, 144, 180, but not 72 because  $72/36 = 2 \notin \text{Num}^+$ . Similarly, we add  $1 \cdot 1$  for  $10 = 2 \cdot 5$  and  $1 \cdot 1$  for  $14 = 2 \cdot 7$ . Thus, the accurate count for positive values of  $t$  is

$$N^+(t) = 14 \cdot 81 - 7 \cdot 26 - 4 \cdot 15 - 2 \cdot 3 - 2 \cdot 3 - 1 \cdot 1 - 1 \cdot 1 + 2 \cdot 4 + 1 \cdot 1 + 1 \cdot 1 = 888.$$

The maximum number of elements from  $S_{b_0}$  in  $S = 1 + 2N^+(t) = 1777$ .

The analysis for  $S_{0c}$  is precisely the same, except we must define

$$t = \frac{a^2 - d^2}{c^2}.$$

Therefore,  $S_{0c}$  can also contribute at most 1777 elements to  $S$ .

Finally, it remains to analyze  $S_{bc}$  where both  $b$  and  $c$  must be nonzero. In this case, we define the quantities,

$$t = \frac{a^2 - d^2}{b^2} \text{ and } s = \frac{c^2}{b^2}.$$

As above, two matrices from  $S_{bc}$  commute if and only if they have precisely the same pair  $(s, t)$  of parameters. Again we can use inclusion and exclusion to count these pairs. By definition,  $s$  is positive, but let us first count pairs  $N(s, t^+)$  with positive  $t$ 's. Our original crude count is 14 choices for  $b$ , 14 choices for  $c$ , and 81 possible numerators  $a^2 - d^2$ . This gives  $14^2 \cdot 81$  positive pairs. Again we must adjust the count for fractions that can be reduced, but we shall only reduce a pair by a factor of  $p^2$  when both  $s$  and  $t$  admit this cancellation. Analogous to the analysis in the  $S_{b_0}$  case, we find

$$\begin{aligned} N(s, t^+) &= 14^2 \cdot 81 - 7^2 \cdot 26 - 4^2 \cdot 15 - 2^2 \cdot 3 - 2^2 \cdot 3 - 1^2 \cdot 1 - 1^2 \cdot 1 \\ &\quad + 2^2 \cdot 4 + 1^2 \cdot 1 + 1^2 \cdot 1 = 14,354. \end{aligned}$$

Of course, this also means there are 14,354 pairs with negative  $t$  values. When  $t = 0$ , we must still analyze how many different values  $s$  can have. Again the answer is obtained by inclusion and exclusion.

$$N(s, 0) = 14^2 - 7^2 - 4^2 - 2^2 - 2^2 - 1^2 - 1^2 + 2^2 + 1^2 + 1^2 = 127.$$

Altogether,  $S_{bc}$  can contribute  $127 + 2 \cdot 14,354 = 28,835$  elements to a noncommuting set  $S$ .

In conclusion, we have partitioned  $U$  into five subsets. The table below summarizes what we have determined about these subsets.

Subset	Defining conditions	Maximum number contributed		
		Order	Order	to $S$
$Z$	$b = c = 0; a = d$	$15 =$	$15$	$0$
$S'_{00}$	$b = c = 0; a \neq d$	$15 \cdot 14 =$	$210$	$1$
$S_{b_0}$	$b \neq 0; c = 0$	$15^2 \cdot 14 =$	$3150$	$1777$
$S_{0c}$	$b = 0; c \neq 0$	$15^2 \cdot 14 =$	$3150$	$1777$
$S_{bc}$	$b \neq 0; c \neq 0$	$15^2 \cdot 14^2 =$	$44,100$	$28,835$
Totals		$15^4 =$	$50,625$	$32,390$

While this completes our analysis, we should mention that our three uses of inclusion and exclusion in this original version of the problem, specifically to count  $N^+(t)$ ,  $N(s, t^+)$ , and  $N(s, 0)$ , have all been confirmed by a computer program that determined the same information by a direct construction of the numbers in question and then deletion of any instance that reproduced any value found earlier.

## REFERENCE

1. *51st Annual William Lowell Putnam Mathematical Competition*, this MAGAZINE 64 (1991), 143.

# The Hiding Path

HERBERT R. BAILEY

Rose-Hulman Institute of Technology  
Terre Haute, IN 47803

The classic pursuit path problem involves a fox chasing a rabbit with both traveling at constant speed. The rabbit path is a straight line and the fox always runs toward the current position of the rabbit. Finding the differential equation describing the path of the fox provides an interesting application of arc length, implicit differentiation, and chain rule.

This problem is used as an example in many differential equations text books. Leonardo da Vinci was apparently one of the first to pose the problem. References to early solutions and discussion of more general pursuit problems are given in [2]. Pursuit paths are currently of great interest in modern warfare.

In this paper we consider the more humane problem of finding a hiding path for the rabbit. The derivation and numerical solution of the differential equation can be done in the first year calculus course. The analysis of the qualitative features of the solution is a nice example for a differential equations course. An interesting feature of the problem is that the differential equation can involve a discontinuous function depending on path choices of the rabbit during flight. This leads to hiding paths that include points where the derivative is discontinuous.

**The problem** A fox is initially  $b$  units due South of a tree and a rabbit is initially  $c$  units due North of the tree. The fox runs due West along a straight line path with constant speed  $V_f$ . The rabbit also runs at constant speed  $V_r$  and must remain hidden from the sight of the fox by keeping fox, tree and rabbit on a straight line.

The situation is shown in FIGURE 1 with the tree at the origin and the positive  $y$ -axis to the North. If the rabbit runs due East, as shown in FIGURE 1, then by similar triangles we have

$$\frac{V_r t}{c} = \frac{V_f t}{b},$$

where  $t$  is time. Thus, in this case, the rabbit speed must be

$$V_r = \frac{cV_f}{b}.$$

While this completes our analysis, we should mention that our three uses of inclusion and exclusion in this original version of the problem, specifically to count  $N^+(t)$ ,  $N(s, t^+)$ , and  $N(s, 0)$ , have all been confirmed by a computer program that determined the same information by a direct construction of the numbers in question and then deletion of any instance that reproduced any value found earlier.

## REFERENCE

1. *51st Annual William Lowell Putnam Mathematical Competition*, this MAGAZINE 64 (1991), 143.

# The Hiding Path

HERBERT R. BAILEY

Rose-Hulman Institute of Technology  
Terre Haute, IN 47803

The classic pursuit path problem involves a fox chasing a rabbit with both traveling at constant speed. The rabbit path is a straight line and the fox always runs toward the current position of the rabbit. Finding the differential equation describing the path of the fox provides an interesting application of arc length, implicit differentiation, and chain rule.

This problem is used as an example in many differential equations text books. Leonardo da Vinci was apparently one of the first to pose the problem. References to early solutions and discussion of more general pursuit problems are given in [2]. Pursuit paths are currently of great interest in modern warfare.

In this paper we consider the more humane problem of finding a hiding path for the rabbit. The derivation and numerical solution of the differential equation can be done in the first year calculus course. The analysis of the qualitative features of the solution is a nice example for a differential equations course. An interesting feature of the problem is that the differential equation can involve a discontinuous function depending on path choices of the rabbit during flight. This leads to hiding paths that include points where the derivative is discontinuous.

**The problem** A fox is initially  $b$  units due South of a tree and a rabbit is initially  $c$  units due North of the tree. The fox runs due West along a straight line path with constant speed  $V_f$ . The rabbit also runs at constant speed  $V_r$  and must remain hidden from the sight of the fox by keeping fox, tree and rabbit on a straight line.

The situation is shown in FIGURE 1 with the tree at the origin and the positive  $y$ -axis to the North. If the rabbit runs due East, as shown in FIGURE 1, then by similar triangles we have

$$\frac{V_r t}{c} = \frac{V_f t}{b},$$

where  $t$  is time. Thus, in this case, the rabbit speed must be

$$V_r = \frac{cV_f}{b}.$$

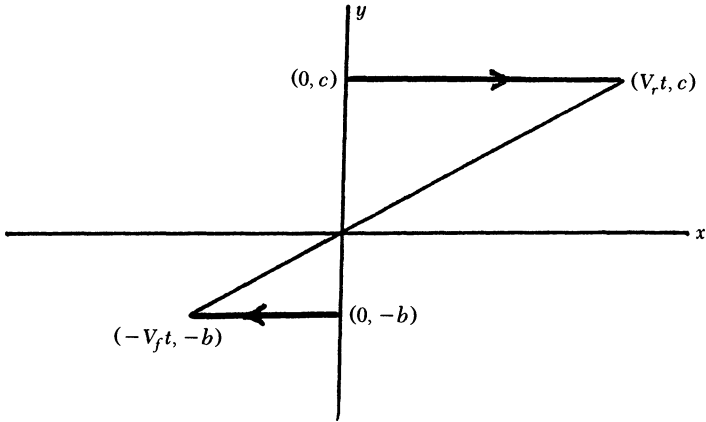


FIGURE 1  
Rabbit runs due east.

We will show that if  $V_r < cV_f/b$ , then the rabbit is not fast enough to remain hidden no matter what path he chooses. If  $V_r > cV_f/b$ , then the rabbit can still remain hidden if he chooses the proper path. The problem is to find this hiding path.

**The path equations** Let  $x(t), y(t)$  be the coordinates of the rabbit as he runs along a hiding path as shown in FIGURE 2. The coordinates of the fox are  $(-V_f t, -b)$ , since she runs in the negative  $x$  direction starting at  $(0, -b)$ .

If  $s(t)$  is the distance traveled by the rabbit along this path, then  $ds/dt = V_r$  and

$$V_r^2 = \left(\frac{ds}{dt}\right)^2 = \left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2. \tag{1}$$

Also, from the similar triangles in FIGURE 2, we have

$$x = \frac{V_f t}{b} y. \tag{2}$$

Differentiating with respect to time gives

$$\frac{dx}{dt} = \frac{V_f}{b} \left( t \frac{dy}{dt} + y \right).$$

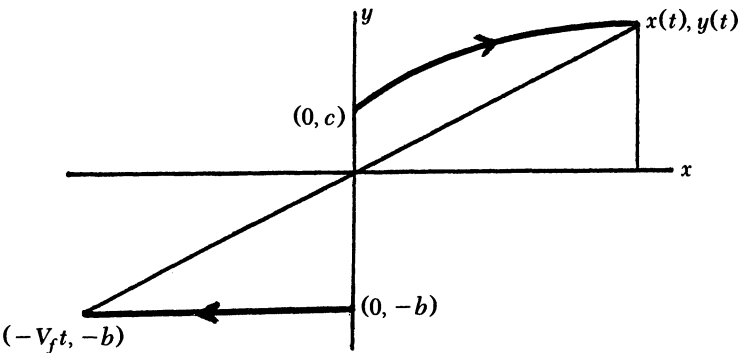


FIGURE 2  
Rabbit runs on curved path.

Substituting the above expression for  $dx/dt$  into equation (1) and combining terms gives

$$\left[1 + \left(\frac{V_f t}{b}\right)^2\right] \left(\frac{dy}{dt}\right)^2 + 2\left(\frac{V_f}{b}\right)^2 y t \frac{dy}{dt} + \left(\frac{V_f}{b}\right)^2 y^2 = V_r^2. \quad (3)$$

If we introduce the dimensionless variables

$$X = \frac{V_f}{bV_r} x, \quad Y = \frac{V_f}{bV_r} y, \quad \text{and} \quad T = \frac{V_f t}{b} \quad (4)$$

and note that

$$\frac{dy}{dt} = \frac{d(bV_r Y/V_f)}{dT} \frac{dT}{dt} = V_r \frac{dY}{dT},$$

then equation (3) becomes

$$(1 + T^2) \left(\frac{dY}{dT}\right)^2 + 2YT \frac{dY}{dT} + Y^2 - 1 = 0.$$

Solving this quadratic for  $dY/dT$  gives

$$\frac{dY}{dT} = \frac{-YT \pm \sqrt{(YT)^2 - (1 + T^2)(Y^2 - 1)}}{1 + T^2} = \frac{-YT \pm \sqrt{1 + T^2 - Y^2}}{1 + T^2}. \quad (5)$$

Combining equations (2) and (4) gives the following equation relating  $X$ ,  $Y$ , and  $T$ .

$$X = \frac{V_f}{bV_r} x = \frac{V_f}{bV_r} \frac{V_f t}{b} y = YT. \quad (6)$$

Differentiating the equation  $X = YT$  with respect to  $T$  and combining with (5) gives

$$\frac{dX}{dT} = \frac{Y \pm T\sqrt{Y^2 T^2 - (1 + T^2)(Y^2 - 1)}}{1 + T^2} = \frac{Y \pm T\sqrt{1 + T^2 - Y^2}}{1 + T^2}. \quad (7)$$

**Solutions** Explicit solutions of equations (5) and (7) were not found. In this section we give some qualitative results and also graphs of numerical solutions.

The rabbit path will depend on the initial value  $Y_0$  of  $Y$ , where  $Y_0 = Y(0) = y(0)V_f/(V_r b) = cV_f/(V_r b)$ . Since  $c \geq 0$  and  $b \geq 0$ , we have  $Y_0 \geq 0$ . If  $Y_0 > 1$ , then the expression under the radical in equations (5) and (7) will be negative at  $T = 0$ . Thus  $Y_0$  is restricted to the interval

$$0 \leq Y_0 \leq 1.$$

Note that the restriction  $Y_0 \leq 1$  is equivalent to  $V_r \geq cV_f/b$ , thus if  $V_r < cV_f/b$ , the rabbit is not fast enough to hide.

The rabbit path will also depend on whether we choose the plus signs before the radicals in equations (5) and (7) (the ‘plus’ case) or the minus signs (the ‘minus’ case). In the ‘plus’ case we note that the constant solution  $X = T$ ,  $Y = 1$  satisfies equations (5) and (7) with the initial condition  $Y_0 = 1$ . This solution corresponds to the constant solution  $y = c$  as shown in FIGURE 1.

By the uniqueness theorem for ordinary differential equations (e.g. [1]), we know that no solution can cross this constant solution. Hence, in the ‘plus’ case,



$Y(T) < 1$  since  $Y_0 < 1$ . Also if  $0 \leq Y(T) < 1$  then  $Y^2 - 1 < 0$  and  $\sqrt{(YT)^2 - (1 + T^2)(Y^2 - 1)} > YT$ . Thus, by equation (5),  $dY/dT > 0$  and  $Y$  is an increasing function of  $T$ . Hence  $Y(T)$  will approach but never reach  $Y = 1$ .

In the 'minus' case we see from equation (5) that  $dY/dT < 0$  and thus  $Y$  is a decreasing function of  $T$ .

Turning to equation (7) we see that  $dX/dT$  is always positive in the 'plus' case. For the 'minus' case,  $dX/dT$  can be zero and this happens when

$$Y - T\sqrt{1 + T^2 - Y^2} = 0.$$

Solving the above equation gives  $T = Y$ . Thus, in the 'minus' case,  $dX/dT > 0$  for  $T < Y$  and  $dX/dT < 0$  for  $T > Y$ . Also for  $Y = T$  and  $X = TY = T^2$ , the path has a vertical tangent line.

The above qualitative results are summarized in graphical form in FIGURE 3. Note that the rabbit must hit the tree in the 'minus' case since  $Y$  is always decreasing and when  $Y$  reaches zero,  $X$  must also be zero since  $X = YT$ . Three path continuations are possible when the rabbit hits the tree. (1) He could be stunned and fall hidden behind the tree. (2) He could pass through the tree (origin) into the third quadrant and remain in line with fox and tree but between them. (3) He could switch to the 'plus' case and remain in the first quadrant on a rising curve. In continuation 3, he would remain hidden since both the 'plus' and 'minus' cases are hiding paths.

Equation (5) is easily solved numerically using Euler's method and a PC or any programmable calculator. The solutions shown in FIGURES 4 and 5 are the graphical output for some numerical solutions. Euler's method simply replaces the derivative with a difference quotient, and the solution of equation (5) is approximated by

$$Y(T + \Delta T) = Y(T) + M\Delta T,$$

with

$$M = \frac{dY}{dT} = \frac{-YT \pm \sqrt{1 + T^2 - Y^2}}{1 + T^2}.$$

$Y(T)$  is then calculated successively after choosing sufficiently small  $\Delta T$  and starting with  $Y = Y_0$  at  $T = 0$ . The  $X$  values are calculated from the equation  $X = YT$ .

FIGURE 4 is a graph of the numerical results for the  $X, Y$  coordinates for two hiding paths. The upper solid curve is for the 'plus' case with  $Y_0 = .6$ . The solid part of the lower curve is for the 'minus' case with  $Y_0 = .6$ . For the dashed part of the lower curve, we have assumed that the rabbit changed suddenly from 'minus' to 'plus' when he hit the tree.

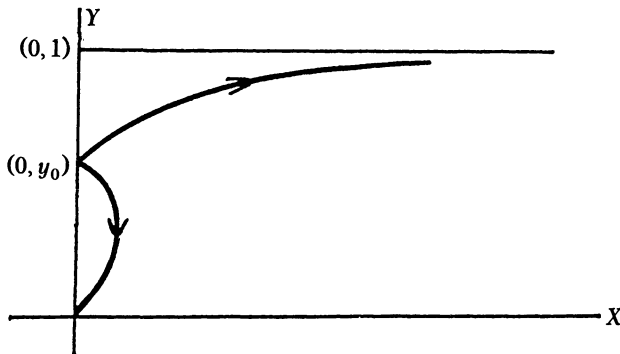


FIGURE 3  
Qualitative hiding paths for the 'plus' and 'minus' cases.

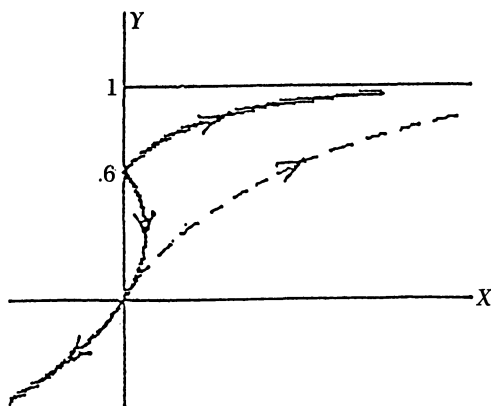


FIGURE 4

Numerical solutions for the 'plus' and 'minus' cases.

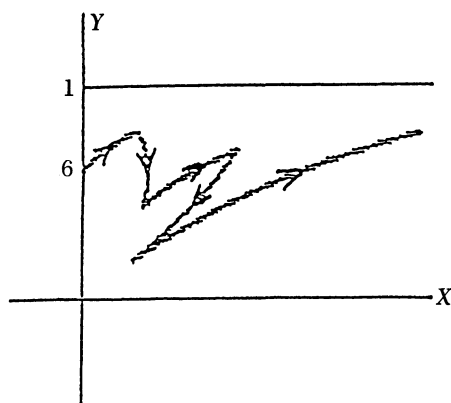


FIGURE 5

Numerical solutions with random switching between 'plus' and 'minus'.

The rabbit could make these sudden switches between 'plus' and 'minus' at any time along his path. Numerical results for an example with this 'random' switching are shown in FIGURE 5.

Polar coordinates  $(r, \theta)$  are more natural for the hiding path problem and lead to the simpler differential equation

$$\frac{dR}{d\theta} = \pm \sqrt{\csc^4 \theta - R^2},$$

where  $R = (rV_f)/(bV_r)$ . However the qualitative features of the solution are not as apparent from this equation as they were from the parametric equations (5) and (7). The derivation of this polar equation provides an application of the arc length formula  $(ds/d\theta)^2 = r^2 + (dr/d\theta)^2$  in polar coordinates.

**Acknowledgement.** This problem was suggested to the author by Mark Roseberry.

## REFERENCES

1. W. E. Boyce and R. C. DiPrima, *Elementary Differential Equations and Boundary Value Problems*, 3rd edition, John Wiley & Sons, Inc., New York, 1977, pp. 70–71.
2. Morley, F. V., A curve of pursuit, *Amer. Math. Monthly* 28 (1921), 54–61.

# How Few Transpositions Suffice? ... You Already Know!

JOHN O. KILTINEN

Northern Michigan University  
Marquette, MI 49855

It is well known that every permutation  $\alpha$  in the set  $S_n$  of all permutations of the set  $\{1, 2, \dots, n\}$ , can be expressed as a product (composition, really) of transpositions, or two-cycles. However the question of how one can do this in a minimal way is generally overlooked. Several hours of searching in the library uncovered only one rather old book [3, p. 15, exercise 15] dating back to 1937 and reissued in 1956 in which this issue was addressed.

One of the students (Stacy Schwenke) in my abstract algebra class raised this question recently. Once asked, the question seemed very natural, but I had to admit that I had never thought about it. It did not take too long to work out a proof that the standard approach that is taken to expressing a permutation in terms of transpositions does indeed give a minimal expression.

The proof was as interesting as the result. One can devise various proofs. (Indeed one of the referees offered an interesting one using some basic results in graph theory.) The virtue of the proof given here is that it uses the elegant but simple technique that is currently favored for proving the invariance of parity theorem for permutations. This result shows students that the method, typically seen only once, is useful for other things—an observation worth sharing with them.

The standard way of expressing a permutation  $\alpha$  as a product of transpositions is first to express  $\alpha$  as a product  $\sigma_1 \circ \sigma_2 \circ \dots \circ \sigma_r$  of disjoint cycles, and then to express each of the cycles as a product of transpositions using the formula

$$(a_1, a_2, \dots, a_k) = (a_1, a_k)(a_1, a_{k-1}) \cdots (a_1, a_3)(a_1, a_2). \quad (1)$$

(Note that here we evaluate permutations from right to left. See most any introductory text such as [7, pp. 81–86] for the necessary background.) This produces an expression for  $\alpha$  as a product of

$$\sum_{i=1}^r (k_i - 1) = (k_1 + k_2 + \dots + k_r) - r \quad (2)$$

transpositions, where  $k_i$  is the length of  $\sigma_i$ . Since, as is well known, permutations have an essentially unique cycle structure, this number, the sum of the lengths of  $\alpha$ 's disjoint cycles, minus the number of cycles, is uniquely determined by  $\alpha$ . We will show that it is the minimal number we seek.

We will do so by demonstrating that *any* minimal expression for  $\alpha$  as a product of transpositions can be converted into the standard expression described above *without changing the number of transpositions*. This of course shows that the standard expression is also minimal.

Our technique, as mentioned above, is the one most often used to prove the parity theorem for  $S_n$ . See, for example, [1–2, 4–8]. In each of these sources, it is first shown that if the identity function  $\varepsilon$  is represented by  $m$  transpositions, then  $m$  must be even. This is done by showing that any such representation can be shortened by two transpositions. If the number were odd to begin with, this would mean that  $\varepsilon$  could eventually be expressed as a single transposition, which is clearly impossible. The

method involves pushing the rightmost occurrence of one of the numbers involved in the expression to the left until an identical pair of adjacent transpositions occurs. When this happens, and one can argue that it must, the pair can be dropped since any transposition combined with itself yields the identity.

We will use this same pushing method, but not to get a reduction. Instead, we will convert a given product to a desired one while keeping the length the same. Let us begin by taking any  $\alpha$  in  $S_n$ , and letting  $\alpha = \sigma_1 \circ \sigma_2 \circ \cdots \circ \sigma_r$  be a decomposition of  $\alpha$  into disjoint cycles. Let us single out the cycle  $\sigma_1$  and denote it by  $\sigma_1 = (a_1, a_2, \dots, a_k)$ . (We are dropping the subscript from  $k$ , the length of the cycle to simplify the notation a bit.) Now let  $\tau_1 \circ \tau_2 \circ \cdots \circ \tau_s$  be an expression of minimal length for  $\alpha$  as a product of transpositions.

In the expression  $\alpha = \tau_1 \circ \tau_2 \circ \cdots \circ \tau_s$ , let us find the rightmost occurrence of  $a_k$ . It must be there since  $\alpha$  moves  $a_k$ . Let us say that it is in the transposition  $\tau_j = (b, a_k) = (a_k, b)$ . Now consider the transposition just to the left of  $\tau_j$ , assuming  $j \neq 1$ . One cannot have  $\tau_{j-1} = \tau_j$ ; otherwise the product of these two transpositions would equal  $\varepsilon$ , the identity, and the expression for  $\alpha$  could be shortened by two, contradicting our minimality hypothesis. Thus, the two entries in  $\tau_{j-1}$  can overlap with the two in  $\tau_j$  by at most one. If  $\tau_j$  and  $\tau_{j-1}$  are disjoint, then they commute with each other. We commute them, and thereby move the rightmost occurrence of  $a_k$  one transposition to the left.

There are two other possibilities to consider. We might have  $\tau_{j-1}$  of the form  $(a_k, c)$ , where  $c \neq b$ . We observe that  $\tau_{j-1} \circ \tau_j = (a_k, c)(a_k, b) = (a_k, b)(b, c)$ . By replacing  $\tau_{j-1} \circ \tau_j$  by the pair on the right, we have again moved the rightmost occurrence of  $a_k$  one spot to the left. Finally, we may have  $\tau_{j-1}$  of the form  $(b, c)$ , where  $c \neq a_k$ . In this case,  $\tau_{j-1} \circ \tau_j = (b, c)(a_k, b) = (a_k, c)(b, c)$ , and we again replace  $\tau_{j-1} \circ \tau_j$  by the pair on the right.

In each case, we see that we have been able to move the rightmost occurrence of  $a_k$  one place to the left. By repeated applications of this technique, we can move the rightmost occurrence of  $a_k$  so that it is in the leftmost transposition—and we have not changed the number of transpositions. Since  $a_k$  maps to  $a_1$  under  $\alpha$ , it must be the case that this leftmost transposition is of the form  $(a_1, a_k)$ .

We now find the rightmost occurrence of  $a_{k-1}$ . Using the same method, we again alter the expression to put the leftmost (and only) occurrence of  $a_{k-1}$  into the transposition which is second from the left. Since  $\alpha$  maps  $a_{k-1}$  to  $a_k$ , we see that this transposition must be  $(a_1, a_{k-1})$ . Repeating this procedure for  $a_{k-2}, \dots, a_3, a_2$ , we see that we can convert our expression for  $\alpha$  in terms of transpositions into the one of the form

$$\alpha = (a_1, a_k)(a_1, a_{k-1}) \cdots (a_1, a_3)(a_1, a_2)\tau_k^* \circ \tau_{k+1}^* \circ \cdots \circ \tau_s^*,$$

where the  $\tau_j^*$ 's are the transpositions obtained by the necessary changes in the  $\tau_j$ 's. Moreover, *we have not changed the length of the expression*. If we apply this same procedure successively to the cycles  $\sigma_2, \sigma_3, \dots, \sigma_r$ , we see that we can convert our given minimal expression for  $\alpha$  in terms of transpositions into one that contains the standard one. If there were any of the  $\tau_j^*$ 's left from the original expression, then the standard expression would be shorter than the minimal one, a contradiction. Thus, the standard one is minimal, and the formula in (2) gives the minimal length.

## REFERENCES

1. John A. Beachy and William D. Blair, *Abstract Algebra with a Concrete Introduction*, Prentice-Hall, Englewood Cliffs, NJ, 1990.

2. David C. Buchthal and Douglas E. Cameron, *Modern Abstract Algebra*, PWS Publishers, Boston, MA, 1987.
3. Robert D. Carmichael, *Introduction to the Theory of Groups of Finite Order*, Dover Publications, Mineola, NY, 1956.
4. Richard A. Dean, *Classical Abstract Algebra*, Harper & Row, New York, 1990.
5. Joseph A. Gallian, *Contemporary Abstract Algebra*, 2nd edition, D. C. Heath, Lexington, MA, 1990.
6. Neal H. McCoy and Gerald J. Janusz, *Introduction to Modern Algebra*, 4th edition, Allyn & Bacon, Boston, 1987.
7. Charles C. Pinter, *A Book of Abstract Algebra*, 2nd edition, McGraw Hill, New York, 1990.
8. Elbert A. Walker, *Introduction to Abstract Algebra*, Random House/Birkhäuser, New York, 1987.

**Added in proof.** A later literature search uncovered the following paper in which the result given here is noted: W. Feit, R. Lyndon, and L. Scott, A remark about permutations, *J. Combinatorial Theory* (A) 18 (1975), 234–235. The following paper explores generalizations of our result: M. Herzog and K. B. Reed, Representation of permutations as products of cycles of fixed length, *J. Austral. Math. Soc.* 22(Series A) (1976), 321–331.

## A Note on the Gaussian Integral

CONSTANTINE GEORGAKIS

DePaul University  
Chicago, IL 60614

The elementary derivation given below for the Gaussian integral

$$I = \int_0^{\infty} e^{-x^2} dx = \frac{\sqrt{\pi}}{2}$$

uses integration in Cartesian coordinates and a dilation kind of change of variable. It is a better alternative to the usual method of reduction to polar coordinates that is found in texts on advanced calculus or probability and statistics.

Let  $y = xs$ ,  $dy = xds$ , then

$$\begin{aligned} I^2 &= \int_0^{\infty} \left( \int_0^{\infty} e^{-(x^2+y^2)} dy \right) dx = \int_0^{\infty} \left( \int_0^{\infty} e^{-x^2(1+s^2)} x ds \right) dx \\ &= \int_0^{\infty} \left( \int_0^{\infty} e^{-x^2(1+s^2)} x dx \right) ds \\ &= \int_0^{\infty} \left[ \frac{1}{-2(1+s^2)} e^{-x^2(1+s^2)} \right]_0^{\infty} ds = \frac{1}{2} \int_0^{\infty} \frac{ds}{(1+s^2)} \\ &= \frac{1}{2} \arctan s \Big|_0^{\infty} = \frac{\pi}{4}. \end{aligned}$$

The idea utilized in this derivation is implicit in a similar method used for establishing the functional equation:  $\Gamma(a)\Gamma(b) = \Gamma(a+b)B(a,b)$  for the gamma and beta function, in the special case where  $a = b = \frac{1}{2}$ .

2. David C. Buchthal and Douglas E. Cameron, *Modern Abstract Algebra*, PWS Publishers, Boston, MA, 1987.
3. Robert D. Carmichael, *Introduction to the Theory of Groups of Finite Order*, Dover Publications, Mineola, NY, 1956.
4. Richard A. Dean, *Classical Abstract Algebra*, Harper & Row, New York, 1990.
5. Joseph A. Gallian, *Contemporary Abstract Algebra*, 2nd edition, D. C. Heath, Lexington, MA, 1990.
6. Neal H. McCoy and Gerald J. Janusz, *Introduction to Modern Algebra*, 4th edition, Allyn & Bacon, Boston, 1987.
7. Charles C. Pinter, *A Book of Abstract Algebra*, 2nd edition, McGraw Hill, New York, 1990.
8. Elbert A. Walker, *Introduction to Abstract Algebra*, Random House/Birkhäuser, New York, 1987.

**Added in proof.** A later literature search uncovered the following paper in which the result given here is noted: W. Feit, R. Lyndon, and L. Scott, A remark about permutations, *J. Combinatorial Theory (A)* 18 (1975), 234–235. The following paper explores generalizations of our result: M. Herzog and K. B. Reed, Representation of permutations as products of cycles of fixed length, *J. Austral. Math. Soc.* 22(Series A) (1976), 321–331.

## A Note on the Gaussian Integral

CONSTANTINE GEORGAKIS  
DePaul University  
Chicago, IL 60614

The elementary derivation given below for the Gaussian integral

$$I = \int_0^{\infty} e^{-x^2} dx = \frac{\sqrt{\pi}}{2}$$

uses integration in Cartesian coordinates and a dilation kind of change of variable. It is a better alternative to the usual method of reduction to polar coordinates that is found in texts on advanced calculus or probability and statistics.

Let  $y = xs$ ,  $dy = xds$ , then

$$\begin{aligned} I^2 &= \int_0^{\infty} \left( \int_0^{\infty} e^{-(x^2+y^2)} dy \right) dx = \int_0^{\infty} \left( \int_0^{\infty} e^{-x^2(1+s^2)} x ds \right) dx \\ &= \int_0^{\infty} \left( \int_0^{\infty} e^{-x^2(1+s^2)} x dx \right) ds \\ &= \int_0^{\infty} \left[ \frac{1}{-2(1+s^2)} e^{-x^2(1+s^2)} \right]_0^{\infty} ds = \frac{1}{2} \int_0^{\infty} \frac{ds}{(1+s^2)} \\ &= \frac{1}{2} \arctan s \Big|_0^{\infty} = \frac{\pi}{4}. \end{aligned}$$

The idea utilized in this derivation is implicit in a similar method used for establishing the functional equation:  $\Gamma(a)\Gamma(b) = \Gamma(a+b)B(a, b)$  for the gamma and beta function, in the special case where  $a = b = \frac{1}{2}$ .

# Lemniscates and Osculatory Interpolation

DONALD TEETS

South Dakota School of Mines and Technology  
Rapid City, SD 57701

PATRICK LANG

Idaho State University  
Pocatello, ID 83209

The “circle of convergence” idea for a complex Taylor series is well known to those who have studied complex variables. Yet few recognize it as a special case of a more general result on the convergence of sequences of interpolating polynomials in the complex plane.

An interpolating polynomial approximates a given function by taking on the same values as the function at a prescribed set of points. A basic result of interpolation theory is that for a function  $f$  defined at  $n + 1$  distinct complex points  $z_0, \dots, z_n$ , there is a unique polynomial  $p_n(z)$  of degree at most  $n$  such that  $p_n(z_j) = f(z_j)$ ,  $j = 0, 1, \dots, n$  [4, p. 38]. If the interpolation points  $z_0, \dots, z_n$  are not distinct, i.e., some point  $z_k$  is repeated  $m$  times, the condition

$$p_n(z_k) = f(z_k)$$

is supplemented by the derivative conditions

$$p'_n(z_k) = f'(z_k), p''_n(z_k) = f''(z_k), \dots, p_n^{(m-1)}(z_k) = f^{(m-1)}(z_k),$$

so that, again, a total of  $n + 1$  conditions are imposed to uniquely determine the  $n + 1$  coefficients of  $p_n(z)$  [4, pp. 52–53]. For example, the interpolating polynomial for  $f$  at the five points  $z_0, z_0, z_0, z_1, z_1$  (where  $z_0 \neq z_1$ ) is the unique 4th-degree polynomial  $p_4(z)$  satisfying

$$\begin{aligned} p_4(z_0) &= f(z_0), & p'_4(z_0) &= f'(z_0), & p''_4(z_0) &= f''(z_0), \\ p_4(z_1) &= f(z_1), & p'_4(z_1) &= f'(z_1). \end{aligned}$$

(Here it is assumed that  $f$  and its derivatives exist at the appropriate points.) Polynomials constructed in this manner are called *osculatory* (or *Hermite*) interpolating polynomials, or, in the special case in which the interpolation points  $z_0, \dots, z_n$  are distinct, *Lagrange* interpolating polynomials.

As the number of interpolation points  $n$  increases, a sequence of polynomials  $\{p_n\}_{n=0}^\infty$  is generated, with  $p_n$  matching the interpolated function  $f$  at one more point (or in one higher derivative) than  $p_{n-1}$ . Thus it is natural to ask, “For which  $z$ -values does  $\{p_n(z)\}$  converge to  $f(z)$ ?” This question leads to some surprising results, perhaps the best known of which is an example due to Runge showing that  $\{p_n(z)\}$  need not converge to  $f(z)$  even under fairly strong hypotheses [1, p. 78, 4, p. 51]. This question also motivates the following development.

Let  $\xi_1, \xi_2, \dots, \xi_k$  be  $k$  complex numbers, not necessarily distinct, and let  $\rho$  be a positive real number. The set of points  $L_\rho$  in the complex plane  $\mathbb{C}$  satisfying

$$|(z - \xi_1)(z - \xi_2) \cdots (z - \xi_k)| = \rho^k \tag{1}$$

is called a *lemniscate* with *foci*  $\xi_1, \xi_2, \dots, \xi_k$  and *radius*  $\rho$ . For fixed  $k$  and foci, the

lemniscates corresponding to different values of  $\rho$  make up a *confocal family* of lemniscates.

The geometry of lemniscates clearly depends on their foci and radii. In the case that  $\xi_1, \xi_2, \dots, \xi_k$  are distinct and  $\rho$  is sufficiently small, the lemniscate  $L_\rho$  is made up of  $k$  simple closed curves, each containing one of the foci  $\xi_i$  in its interior. As  $\rho$  increases, the curves increase in size until two or more of them merge into a single larger curve. As  $\rho$  increases still further, this merging of curves continues until all foci are contained in the interior of a single simple closed curve. Finally, as  $\rho \rightarrow \infty$ , this single curve becomes arbitrarily close in shape to a circle (see FIGURE 1). A more complete description of these facts, together with their proofs can be found in [1, pp. 87–89].

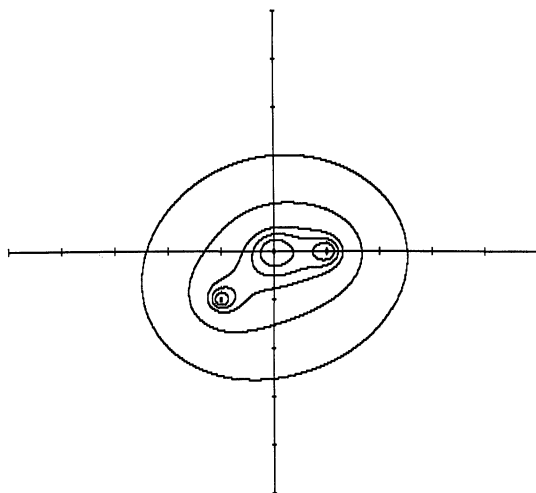


FIGURE 1

The confocal family of lemniscates  $|z(z-1)(z-(-1-i))| = \rho^3$ .

The following theorem, whose proof can be found in [1, p. 90–92], relates lemniscates to polynomial interpolation in the special case that the complex interpolation points  $z_0, z_1, \dots$  approach  $k$  limit points  $\xi_1, \xi_2, \dots, \xi_k$  cyclically, i.e.,

$$\lim_{i \rightarrow \infty} z_{ik+j} = \xi_j, \quad j = 1, 2, \dots, k.$$

**THEOREM 1.** *Let  $L_\rho$  be as in (1) with interior*

$$I_\rho = \{z \mid |(z - \xi_1)(z - \xi_2) \cdots (z - \xi_k)| < \rho^k\}.$$

*Suppose that  $z_0, z_1, \dots$  lie in  $I_\rho$  and approach  $\xi_1, \xi_2, \dots, \xi_k$  cyclically and suppose that  $f$  is analytic in  $I_\rho$ , but not in any  $I_{\rho_1}$ , with  $\rho_1 > \rho$ . If  $p_n$  is the polynomial of degree at most  $n$  that agrees with  $f$  at  $z_0, z_1, \dots, z_n$ , then  $p_n(z) \rightarrow f(z)$  as  $n \rightarrow \infty$  for all  $z \in I_\rho$ , with the convergence being uniform in any closed subset of  $I_\rho$ . For  $z$  outside  $I_\rho$ , the sequence  $\{p_n(z)\}_{n=0}^\infty$  diverges.*

It should be noted that this theorem is but a special case of a much more general theory that is fully developed in [5]. The following examples illustrate Theorem 1.

**Example 1.** If  $k = 1$ , the confocal lemniscates  $L_\rho$  are concentric circles centered at  $\xi_1$ , and the interpolation points  $z_0, z_1, \dots$  converge to  $\xi_1$ . The polynomials  $\{p_n\}$  converge to  $f$  inside the largest circle  $|z - \xi_1| = \rho$  in which  $f$  is analytic (uniformly in any disk  $|z - \xi_1| \leq \rho_1 < \rho$ ), and diverge outside this circle. If  $z_0 = z_1 = \dots$ , the interpolating polynomials  $\{p_n\}$  are partial sums of the Taylor series of  $f$  centered at  $\xi_1$ , and the familiar convergence theory for Taylor series results.



*Example 2.* Let

$$f(z) = \frac{e^z \sin z}{(z-2)(z-3)}.$$

Let  $\{z_j\}$  be the sequence

$$1, 1+2i, 4, 3+2i, 1, 1+2i, 4, 3+2i, \dots,$$

and for  $n = 0, 1, \dots$ , let  $p_n(z)$  interpolate  $f(z)$  at  $z_0, z_1, \dots, z_n$ . For  $\rho = 10^{1/4}$ , the lemniscate  $L_\rho$  with foci at 1,  $1+2i$ , 4, and  $3+2i$  is given by

$$|(z-1)(z-(1+2i))(z-4)(z-(3+2i))| = 10$$

(see FIGURE 2). This lemniscate passes through the pole of  $f$  at  $z = 2$ , but contains no singularities of  $f$  in its interior  $I_\rho$ . Thus  $p_n(z) \rightarrow f(z)$  as  $n \rightarrow \infty$  for  $z \in I_\rho$ , but  $\{p_n(z)\}$  diverges for all  $z \notin \bar{I}_\rho$ . The convergence is uniform on closed subsets of  $I_\rho$ .

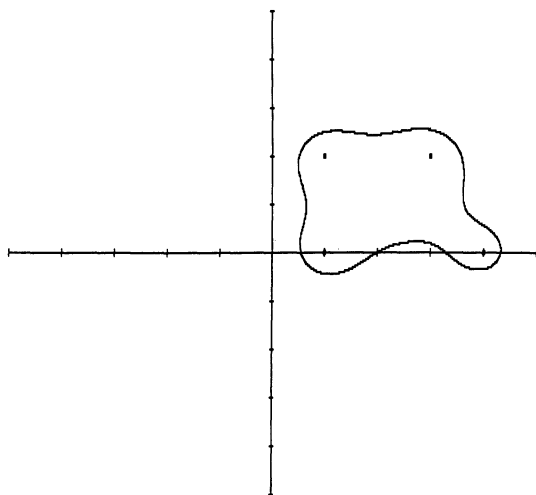


FIGURE 2

The lemniscate  $|(z-1)(z-(1+2i))(z-4)(z-(3+2i))| = 10$ .

*Example 3.* Let

$$f(z) = \frac{1}{z-1},$$

and for  $j = 0, 1, \dots$ , set

$$z_j = \begin{cases} 0, & j \text{ even} \\ 3, & j \text{ odd.} \end{cases}$$

Let  $p_n(z)$  interpolate  $f(z)$  at  $z_0, z_1, \dots, z_n$ ,  $n = 0, 1, \dots$ . Since  $f$  is analytic in the interior of the lemniscate  $|z(z-3)| = 2$  (shown in FIGURE 3), but not in the interior of any larger lemniscate with foci at 0 and 3, the polynomials  $\{p_n(z)\}$  converge to  $f(z)$  if, and only if,  $|z(z-3)| < 2$ .

Examples 1, 2, and 3 suggest that the maximal convergence regions for certain sequences of interpolating polynomials can have rather arbitrary shapes. The purpose of this paper is to make this idea precise. This is accomplished in Theorem 3. Its proof depends on the following two preliminary results.

**THEOREM 2.** (Hilbert, 1897) *If  $C$  and  $C'$  are two simple closed curves,  $C'$  interior to  $C$ , then there exists a lemniscate  $L$  that is interior to  $C$  and contains  $C'$  in its interior.*

*Proof.* See [3, v.II pp. 294–295].

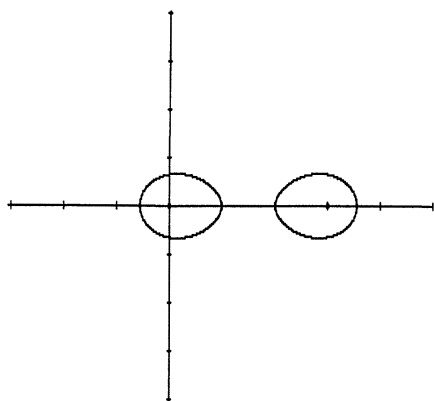


FIGURE 3

Maximal convergence region for Example 3.

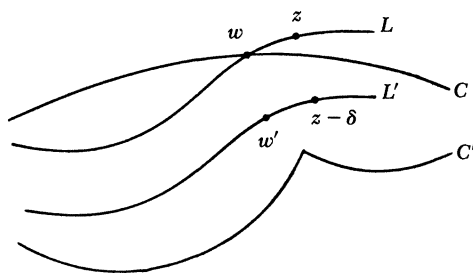


FIGURE 4

In the following corollary, the *distance* between a point  $z$  and a curve  $C$  is defined by

$$d(z, C) = \min_{w \in C} |z - w|.$$

Note that for any points  $z$  and  $\xi$  and any curve  $C$ , this distance satisfies

$$d(z, C) \leq d(\xi, C) + |z - \xi|$$

[2, pp. 184–185].

**COROLLARY 1.** *Let  $C$  be a simple closed curve, and let  $w$  be a point of  $C$ . Then for any  $\varepsilon > 0$ , there exists a lemniscate  $L$  such that  $w \in L$  and  $d(z, C) < \varepsilon$  for all  $z \in L$ .*

*Proof.* Let  $\varepsilon > 0$  be given. Consider an open cover of  $C$  by circular disks of radius at most  $r < \varepsilon/4$ . Since  $C$  is compact, there exists a finite subcover. Using the circular arcs that bound the disks of the finite subcover, we construct a simple closed curve  $C'$  such that i)  $C'$  is the union of a finite number of circular arcs that turn their concave sides toward  $C$ ; and ii)  $C'$  lies in the interior of  $C$ . By Theorem 2, there is a lemniscate  $L'$ , given by

$$|(z - \xi_1)(z - \xi_2) \cdots (z - \xi_k)| = \rho^k, \quad (2)$$

which is interior to  $C$  and contains  $C'$  in its interior (see FIGURE 4).

Choose  $w' \in L'$  such that  $|w - w'| < \varepsilon/2$ , and define  $\delta = w - w'$ . Consider the lemniscate  $L$  given by

$$|(z - \xi_1 - \delta)(z - \xi_2 - \delta) \cdots (z - \xi_k - \delta)| = \rho^k. \quad (3)$$

Since  $w' \in L'$ ,  $w'$  satisfies (2). Therefore,  $w = w' + \delta$  satisfies (3), so that  $w \in L$ .

Let  $z$  be a point of  $L$ . Then  $z - \delta \in L'$ , so that

$$\begin{aligned} d(z, C) &\leq d(z - \delta, C) + |z - (z - \delta)| \\ &\leq 2r + |\delta| \\ &< \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

Hence, each point of the lemniscate  $L$  lies within  $\varepsilon$  of the curve  $C$ .

For any simple closed curve  $C$  and any  $\varepsilon > 0$ , the following notation is introduced:  $I_\varepsilon(C)$  denotes the set of all points  $z$  in the interior of  $C$  satisfying  $d(z, C) \geq \varepsilon$  and  $E_\varepsilon(C)$  denotes the set of all points  $z$  in the exterior of  $C$  for which  $d(z, C) > \varepsilon$ . In what follows,  $\varepsilon$  is assumed to be small enough to ensure that  $I_\varepsilon(C)$  is nonempty.

The next theorem is the main result of this paper. It shows how remarkably irregular the maximal convergence regions for sequences of interpolating polynomials can be.

**THEOREM 3.** *Suppose the singularities of  $f$  are all isolated singularities. Let  $C$  be any simple closed curve passing through a singularity  $w$  of  $f$  but containing no singularities in its interior. Then for any  $\varepsilon > 0$ , there exists a sequence of points  $\{z_i\}_{i=0}^\infty$  with the following property: If  $p_n(z)$  is the polynomial of degree at most  $n$  that agrees with  $f(z)$  at  $z_0, z_1, \dots, z_n$ , then the sequence  $\{p_n(z)\}_{n=0}^\infty$  converges uniformly to  $f$  on  $I_\varepsilon(C)$  and diverges at each point of  $E_\varepsilon(C)$ .*

*Proof.* Let  $\varepsilon > 0$  be given, and define  $I_\varepsilon(C)$  and  $E_\varepsilon(C)$  as above. Choose  $\varepsilon' < \varepsilon$  and less than the minimum distance between  $C$  and any singularity of  $f$  exterior to  $C$ . By Corollary 1, there exists a lemniscate  $L$  given by

$$|(z - \xi_1)(z - \xi_2) \cdots (z - \xi_k)| = \rho^k$$

such that  $w$  lies on  $L$  and such that  $d(z, C) < \varepsilon'$  for each  $z \in L$ . Note that the interior of  $L$  contains no singularities of  $f$ , but does contain the set  $I_\varepsilon(C)$ .

Choose the interpolation points  $z_0, z_1, \dots$  as follows:

$$\begin{aligned} z_0 &= \xi_1 \\ z_{mk+1} &= \xi_1 \\ z_{mk+2} &= \xi_2 \\ &\vdots \\ z_{mk+k} &= \xi_k, \quad m = 0, 1, \dots \end{aligned}$$

By Theorem 1,  $p_n \rightarrow f$  as  $n \rightarrow \infty$  uniformly on  $I_\varepsilon(C)$  and  $\{p_n(z)\}$  diverges at each point  $z \in E_\varepsilon(C)$ .

With Theorem 3, it has been shown that the convergence region for a sequence of interpolating polynomials can be quite arbitrary indeed.

## REFERENCES

1. P. J. Davis, *Interpolation and Approximation*, Dover Publications, Inc., Mineola, NY, 1975.
2. J. Dugundji, *Topology*, Allyn and Bacon, Boston, 1956.
3. E. Hille, *Analytic Function Theory*, Ginn and Company, Boston, 1959.
4. J. Stoer and R. Bulirsch, *Introduction to Numerical Analysis*, Springer-Verlag New York, Inc., New York, 1980.
5. D. A. Teets, *Convergence Regions for Sequences of Complex Interpolating Functions*, D. A. Thesis, Idaho State University, May 1989.

# Note on the Evaluation of $\int_0^x \frac{1}{1+t^{2^n}} dt$

M. A. GOPALAN  
V. RAVICHANDRAN  
National College  
Trichy-620 001, India

In [1], the integral

$$I = \int_0^x \frac{1}{1+t^m} dt$$

has been evaluated for  $m = 1, 2, 3, 4, 5, 6, 8$ , and 10. In an attempt to evaluate  $I$  for all values of  $m$ , we find that it is possible to obtain the general value of  $I$  when  $m = 2^n (n \geq 2)$ . The steps involved are as follows:

*Step I:* By factoring  $(t^{2^n} + 1)$ , we get  $2^{n-1}$  distinct quadratic factors; of these,  $2^{n-2}$  quadratic factors are of the form  $(t^2 + at + 1)$ , and the other  $2^{n-2}$  quadratic factors are of the form  $(t^2 - at + 1)$ , where  $a$  is one of  $2^{n-2}$  numbers given by

$$\sqrt{2 \pm \sqrt{2 \pm \sqrt{2 \pm \cdots \pm \sqrt{2}}}}. \quad (1)$$

Note that the number of 2's in  $a$  is  $(n-1)$ .

*Step II:* On separating  $\frac{1}{1+t^{2^n}}$  into partial fractions, we have

$$\frac{1}{1+t^{2^n}} = \sum_{k=1}^{2^{n-2}} \left[ \frac{A_k t + B_k}{t^2 + a_k t + 1} + \frac{C_k t + D_k}{t^2 - a_k t + 1} \right]. \quad (2)$$

Multiply equation (2) by  $(1+t^{2^n})$  on both sides and take the limit as  $t$  tends to  $\alpha$ , where

$$\alpha = \left[ -a_k + i\sqrt{(4 - a_k^2)} \right] / 2, \text{ to get}$$

$$\lim_{t \rightarrow \alpha} \left[ \frac{(A_k t + B_k)(t^{2^n} + 1)}{(t^2 + a_k t + 1)} \right] = 1.$$

Note that  $1 + \alpha^{2^n} = 0$ , and  $\alpha^2 + a_k \alpha + 1 = 0$ .

Evaluating the limit and rearranging the result, we see that

$$2 + a_k \alpha - 2^n (A_k \alpha + B_k) = 0. \quad (3)$$

Equating the real and imaginary parts of (3) to zero and solving the resulting equations, it is seen that

$$A_k = a_k 2^{-n}, \text{ and } B_k = 2^{1-n}. \quad (4)$$

Similarly

$$C_k = -a_k 2^n, \text{ and } D_k = 2^{1-n}. \quad (5)$$

Substituting (4) and (5) in (2) and integrating both sides with respect to  $t$  between the limits zero and  $x$ , we get

$$\int_0^x \frac{1}{1+t^{2^n}} dt = \sum_{k=1}^{2^n-2} \frac{a_k}{2^{n+1}} \left[ \log \left( \frac{x^2 + a_k x + 1}{x^2 - a_k x + 1} \right) + 2 \tan^{-1} \left( \frac{a_k x}{1-x^2} \right) \right]. \quad (6)$$

### Examples

(1) For  $n = 2$ , the number of 2's in  $a$  is 1 and  $a_1 = \sqrt{2}$ .

From (6), we get

$$\int_0^x \frac{1}{1+t^4} dt = \frac{1}{2\sqrt{2}} \left[ \log \left( \frac{x^2 + \sqrt{2}x + 1}{x^2 - \sqrt{2}x + 1} \right) + 2 \tan^{-1} \left( \frac{\sqrt{2}x}{1-x^2} \right) \right].$$

(2) For  $n = 3$ , the number of 2's in  $a$  is 2 and therefore

$$a_1 = \sqrt{2 + \sqrt{2}}, \text{ and } a_2 = \sqrt{2 - \sqrt{2}}.$$

(3) For  $n = 4$ , the number of 2's in  $a$  is 3 and

$$\begin{aligned} a_1 &= \sqrt{2 + \sqrt{2 + \sqrt{2}}}, \quad a_2 = \sqrt{2 + \sqrt{2 - \sqrt{2}}}, \\ a_3 &= \sqrt{2 - \sqrt{2 + \sqrt{2}}}, \quad \text{and } a_4 = \sqrt{2 - \sqrt{2 - \sqrt{2}}}. \end{aligned}$$

It is seen that Examples (1) and (2) are in agreement with [1].

### REFERENCE

1. *Notebooks of Srinivasa Ramanujan*, Vol. 2, Tata Institute of Fundamental Research, Bombay, 1957.

# How Expected is the Unexpected Hanging?

DEAN CLARK  
University of Rhode Island  
Kingston, RI 02881

This is about a famous paradox introduced by philosophers in the 1940s. It was casually dismissed by Quine [1] in 1953. He cited five papers that preceded his in the prestigious British journal *Mind*. Martin Gardner wrote about unexpected quizzes, unexpected eggs, unexpected tigers, etc. in the 60s and 70s [2], [3], and his bibliography in [2] listed 23 references. Twenty years after Quine the paradox was dismissed again, in *Mind*, by A. J. Ayer [4]. In the 80s the paradox of the unexpected hanging was attacked with the methods of formal logic [5] by Professor R. M. Sainsbury, the current editor of *Mind*. It also reached the popular book stores [6]. Now we are into the 90s and Martin Gardner's book [2] has been reissued by the University of Chicago Press in a new edition [8] listing 57 references on the paradox!

After a half-century of effort and nearly 60 papers it seems unlikely that any single article will resolve it. In fact, it is unclear which domain of human intelligence should take custody of it. Is it a problem in pure logic? semantics? psychology? probability theory? Is it a problem without a solution or with multiple solutions? One thing is clear: The paradox of the unexpected hanging has rarely been approached as a *mathematical* problem. Questions like "What is the numerical probability of the unexpected hanging?" or "What are necessary and sufficient conditions for the unexpected hanging?" have not been posed. That is what we do here.

**Problem A** On Sunday evening a judge tells a condemned prisoner that he will be awakened and hanged on the morning of one of the following five days. The judge says that it will happen unexpectedly, i.e., the prisoner will be uncertain about when the hanging will occur until the moment the attendants arrive. But the prisoner's attorney convinces him that *no such hanging is possible*. The first step in the attorney's argument is to eliminate Friday as execution day: If the judge sets Friday as the morning of the hanging, the prisoner will know it on Thursday because he is still alive and realizes that tomorrow is the final day of the execution period. Friday is ruled out. But then, by the process of elimination, so are Thursday, Wednesday, Tuesday, and Monday. Of course, on Wednesday, the prisoner is hauled out of bed, much to his surprise, and hanged.

The unexpected hanging is impossible: *true* by logical argument, *false* by counterexample. A claim that is both true and false is unacceptable, yet two things are clear: (1) the Wednesday hanging *was* unexpected; (2) Friday *is* ruled out (and thus so are Thursday, Wednesday, etc.)!

Problems B and C, below, are related to Problem A. The reader will see that Problem B has the same underlying structure as Problem A, even though words like "unexpected", "uncertain", etc. are not mentioned.

**Problem B** Suppose I tell you that I am thinking of either the name "Paul" or "Paula", and that each name has *exactly one white letter*. The rest of the letters are black. I will now expose the letters one at a time.

P

PA

Consider the statement  $S =$  “Five exposures are necessary to identify the name correctly.” Now,  $S$  is surely false (because the first four letters could be black, implying that the name is Paula; that’s *four* exposures).

**PAU**

But,  $S$  is surely true. Look:

**PAUL**.

You cannot possibly deduce the name. Will the next exposure be **A** or will it be blank? So five exposures *are* necessary!

With the question posed at the end of Problem B, the element of uncertainty becomes evident, and so does the correspondence between the position of the white letter and the day of the hanging. In both Problems A and B there is precisely one instance, among a finite number of possibilities, in which a deduction can be made.

Now we open an infinite Pandora’s box of impossibilities:

**Problem C [6]** Two mathematicians on a train, X and Y, each think of a positive integer and whisper it to a fellow rider Z. Z gets up and announces, “This is my stop. You have each thought of a different number, and neither of you can deduce whose number is larger. You may not speak or write messages to one another.” X and Y continue their travel in silence. X, whose number was 62, thinks, “Obviously Y didn’t choose 1. If he did, he’d know that my number was the larger, just from Z’s statement that we chose different numbers. Just as obviously, Y knows that I didn’t choose 1. Our choices *must* have come from the subset  $\{2, 3, \dots\}$ . But that would mean that neither of us chose 2. By mathematical induction, all positive integers are eliminated—which is completely absurd!”

We use the standard apparatus of probability theory, i.e., sample space  $\Omega$  with sample points  $\omega \in \Omega$ , a collection of subsets of  $\Omega$  called *events*, and probability measures  $P_n$ . I claim that Problems A and B are excellent examples of confusion between  $\Omega$  and  $\omega$ , between w.p.1 (with probability one) and probabilities that approach 1 in the limit, between *statements* and *truth-valued random variables*, between *deductive* and *inductive* reasoning.

We can deal with Problems A and B simultaneously without loss of generality. Both involve statements that seem to be both true *and* false. Consider  $S$  of Problem B, which, according to elementary logic, is not a statement at all (at least, it is no more a statement than “ $x + 3 > 5$ ”). Instead,  $S$  is a random variable defined on a sample space  $\Omega$  whose nine possible outcomes are

**PAUL, PAUL, PAUL, PAUL,**

**PAULA, PAULA, PAULA, PAULA, PAULA**

$S$  has two possible values: 0 (= false) and 1 (= true). We saw that  $S(\text{PAULA}) = 0$  and all the other sample points have  $S$ -value of 1. The “necessary” in  $S$  does not mean “logically necessary”, i.e., already implicit in a set of valid premises. It does not even mean “probably necessary” in the colloquial sense until more is known about the probability distribution on  $\Omega$ . If a tenth, all-black point **PAUL** were added to  $\Omega$  (making it a new sample space  $\Omega'$ ),  $S$  would be true for every sample point, in other words w.p.1.  $S$  would then be indistinguishable from a true statement. Thus a small change in  $\Omega$  results in a radical change in  $S$ .

Any treatment that purports to resolve the paradox of the unexpected hanging

ought to address these issues (since the above problem of the unexpected white letter is just a thinly disguised version of the original). Instead, several writers have claimed that the prisoner can be alive on Thursday night and be *unable* to deduce that he will hang the next day. Given the knowledge that the unknown name was chosen from  $\Omega$  these same writers would presumably be *unable* to identify PAULA after four exposures!

Let a sample space for Problem A be the set of nine points given above. Let us agree that the method of deciding the day of the hanging is to choose a point from  $\Omega$  and expose one letter each morning. A “white letter day” is the prisoner’s unlucky day. We imagine a bag filled with many slips of paper, each slip with one of the nine names on it, into which the judge can reach to make his decision. In fact, we imagine a sequence of  $n$  bags, in which the relative proportion of slips with the name PAULA gets ever smaller. Each bag gives a probability measure  $P_n$  on  $\Omega$ , and these measures converge (in distribution) to a measure  $P$  on  $\Omega$  with  $P(\text{PAULA}) = 0$ .

Of course, this is an imperfect model of how the judge makes his decision, but it, or some variant of it, is an orthodox way to view the problem. After all, think of coin tossing and how we come to grips with describing numerically the uncertainty of the outcome. We have even less understanding of the psychic path that leads to the judge’s decision. All we have is the *assertion* that the hanging *will* occur (covered by the above model) and that it *will* be unexpected (under conditions yet to be discovered within the model).

Here is the judge’s claim: “For some day  $k$  which I shall choose ( $k$  from Monday (1) to Friday (5)), the hanging on day  $k$  and its unexpectedness are both true”. The judge does not say that the unexpectedness is true w.p.1, i.e., for *any* or *every*  $k = 1, 2, \dots, 5$ . He *cannot* say that, since the hanging on day 5 and its unexpectedness are contradictory and imply that  $k$  cannot equal 5. The judge’s “statement”, like  $S$  in Problem B, is a random variable. It does not have deterministic truth value. It only makes sense to talk about  $p_n = P_n(k)$ : the unexpected hanging occurs on day  $k$ . The valid conclusion is that  $p_n < 1$ . For instance, if the nine sample points are equally probable,  $p_n = 8/9$ . More ominously,  $p_n \rightarrow 1$  as  $P_n(\text{PAULA}) \rightarrow 0$ .

To sum up, if the prisoner is alive on *Wednesday* night he is entitled to think: “Either tomorrow I hang, or else the judge’s claim about the unexpectedness is false.” And if he survives to *Thursday* night, the second half of the preceding sentence is all that is left. Given the way the judge makes his decision, this last likelihood is remote in the extreme. Ironically, the sought for condition that ensures the virtual certainty of the unexpected hanging is nothing more than the infinitesimally small but *positive* probability that the judge’s claim is *false*!

The author of [6] attributes Problem C to Martin Hollis, although no reference is given. It is suggested that Problem C is a generalization or infinite extension of the paradox of the unexpected hanging. In fact, it is a different and cruder problem.

In Problem A the truth of the judge’s edict was truth in an “inductive” sense. That is, the unexpectedness of the hanging was not a logically necessary conclusion derived from primitive axioms, i.e., not a *deduction*. The judge could no more prove that the hanging would be unexpected than he could prove that it would rain on the day of the hanging. However, the judge *could* be infinitely more *confident* in the unexpectedness of the hanging than he could be of the rain since he controls the probabilities—the *inductive* strength of the claim—with exact numerical precision. For an excellent discussion of these things see [7].

Mr. Z’s claim could also be true in the course of events. X and Y may still be cogitating on that train, but it’s unlikely. Both would have realized that their “mathematical induction” was the very nonsense that invalidated the premise of



neither being able to deduce who had the larger number. They *can* deduce who has the larger number but *how*? Actually, it is not up to us to come up with a practical method for doing this. The procedure we are about to describe requires the kind of spontaneous and intelligent coordination sometimes exhibited by identical twins. But now that we sense the danger of mixing inductive with deductive reasoning, we assert our right as mathematicians—deductive reasoners—to simply provide a contradiction-free procedure for either X or Y to deduce who has the larger number without learning the value of that number.

Is it unthinkable that (having realized at once how to solve the problem) at a given signal X and Y could use the second hand of their watches as a counter and give a second signal when the one with the smaller number reaches it? The second signal could be the statement “I’ve just deduced who has the larger number—you do!” Unthinkable? Certainly not. Mathematicians accustomed to looking for counterexamples and not terribly worried about the practicalities would regard *S* in the “Paul-Paula” experiment as false as long as  $P(\text{PAULA}) > 0$ . This is the same opinion that X and Y have about Z’s assertion as they get off the train.

## REFERENCES

1. Willard van O. Quine, On a so-called paradox, *Mind* 62 (1953), 65–67.
  2. Martin Gardner, *The Unexpected Hanging and Other Mathematical Diversions*, Simon and Schuster, New York, 1969.
  3. ———, *Aha! Gotcha*, W. H. Freeman and Company, San Francisco, 1982.
  4. A. J. Ayer, On a supposed antinomy, *Mind* 82 (1973), 125–126.
  5. R. M. Sainsbury, *Paradoxes*, Cambridge University Press, Cambridge, 1988.
  6. William Poundstone, *Labyrinths of Reason*, Anchor Press Doubleday, New York, 1988.
  7. Brian Skyrms, *Choice and Chance: an introduction to inductive logic*, Wadsworth, Belmont, CA, 1986.
  8. Re-publication of [2] with a new afterword, University of Chicago Press, Chicago, 1991.
-

# Sums of Like Powers of Multivariate Linear Forms

ISMOR FISCHER  
Naval Postgraduate School  
Monterey, CA 93943

**Introduction** The subject of this note began with a very elementary observation, starting with the algebraic identity

$$(x + y)^2 - (x - y)^2 = 4xy. \quad (1)$$

The left-hand side is clearly a polynomial of degree at most 2 in two variables. However it has special structure in that its powers and arrangement of elementary operations are “just right” to cancel all but the uniformly mixed cross term. By formally replacing  $x$  with  $ax$  and  $y$  with  $by$ , we obtain

$$(ax + by)^2 - (ax - by)^2 = 4abxy. \quad (2)$$

If we interpret  $x$  and  $y$  as fixed formal parameters, then Equation (2) says that we may interchange the coefficients  $a$  and  $b$  on the left-hand side and yet, by virtue of the symmetry on the right-hand side, leave the value unchanged. Hence,

$$(ax + by)^2 - (ax - by)^2 = (bx + ay)^2 - (bx - ay)^2,$$

or

$$A^2 + B^2 = C^2 + D^2, \quad (3)$$

with  $A = ax + by$ ,  $B = bx - ay$ ,  $C = bx + ay$ ,  $D = ax - by$ . Trivial solutions will result, if and only if, any of  $a, b$ , or  $a \pm b = 0$ , with similar conditions on  $x$  and  $y$ . Hence, with appropriate conditions on  $a, b, x, y$  we have succeeded in generating integers that are expressible as a sum of two squares in (at least) two distinct ways, and actually one can prove that the general solution of Equation (3) in integers must essentially be of the form derived above; see [2].

Can the special structure of Equation (1) be extended to higher powers? Specifically, is it possible to find a sequence of sums and differences (hereafter referred to simply as “operations”) of like powers of linear forms in  $n$  variables to produce exactly the right amount of cancellation needed to leave only the uniform cross term  $C_n x_1 \dots x_n$ ? If so, what implications does this have for characterizing those integers that are expressible as sums of like powers in more than one way?

**Procedure** We start with some notation and definitions for one method of answering some of these questions.

Let  $S = \{2, 3, \dots, n\}$ ,  $n > 1$ . For each  $m$  ( $0 \leq m \leq n - 1$ ), there are exactly  $N = \binom{n-1}{m}$  subsets of  $S$  of size  $m$ . Order these subsets in some fashion, calling the  $k$ th subset  $S_k^m$  for  $k = 1, 2, \dots, N$ .

Now for each  $i \in S \cup \{1\}$  we define a mapping  $\sigma$  from the power set  $\mathcal{P}(S)$  into the set of linear forms in  $n$  variables by

$$\sigma(S_k^m) = \sum_{i \notin S_k^m} x_i - \sum_{i \in S_k^m} x_i.$$

Let

$$P_n(x) = \sum_{m=0}^{n-1} \sum_{k=1}^N (-1)^m [\sigma(S_k^m)]^n, \quad (4)$$

where  $x = (x_1, x_2, \dots, x_n)$ .

THEOREM.

$$P_n(x) = 2^{n-1} n! \prod_{i=1}^n x_i. \quad (5)$$

In particular,

$$P_n(a \cdot x) = P_n(\pi(a) \cdot x), \quad (6)$$

where  $\pi(a)$  is any nontrivial permutation on  $a = (a_1, \dots, a_n)$ , and  $\cdot$  denotes the standard Euclidean dot product.

Before proving this result, let us consider an example. Take  $n = 4$ . We partition the power set  $\mathcal{P}(S)$  of  $S = \{2, 3, 4\}$  into eight subsets  $S_k^m$  where  $m$  denotes the cardinality of the set. The sequence of signs within each linear form depends on  $S_k^m$ ; the operation assigned to it in the full sum  $P_n(x)$  (either  $+$  or  $-$ ) depends on  $m$ . For  $x = (x_1, x_2, x_3, x_4)$ ,

$$\begin{array}{lll} P_4(x) = + (x_1 + x_2 + x_3 + x_4)^4 & S_1^0 = \phi & m = 0 \\ \begin{array}{l} - (x_1 - x_2 + x_3 + x_4)^4 \\ - (x_1 + x_2 - x_3 + x_4)^4 \\ - (x_1 + x_2 + x_3 - x_4)^4 \end{array} & \begin{array}{l} S_1^1 = \{2\} \\ S_2^1 = \{3\} \\ S_3^1 = \{4\} \end{array} & \left. \vphantom{\begin{array}{l} - (x_1 - x_2 + x_3 + x_4)^4 \\ - (x_1 + x_2 - x_3 + x_4)^4 \\ - (x_1 + x_2 + x_3 - x_4)^4 \end{array}} \right\} m = 1 \\ \begin{array}{l} + (x_1 - x_2 - x_3 + x_4)^4 \\ + (x_1 - x_2 + x_3 - x_4)^4 \\ + (x_1 + x_2 - x_3 - x_4)^4 \end{array} & \begin{array}{l} S_1^2 = \{2, 3\} \\ S_2^2 = \{2, 4\} \\ S_3^2 = \{3, 4\} \end{array} & \left. \vphantom{\begin{array}{l} + (x_1 - x_2 - x_3 + x_4)^4 \\ + (x_1 - x_2 + x_3 - x_4)^4 \\ + (x_1 + x_2 - x_3 - x_4)^4 \end{array}} \right\} m = 2 \\ - (x_1 - x_2 - x_3 - x_4)^4 & S_1^3 = \{2, 3, 4\} & m = 3 \end{array}$$

$$= 2^3 4! x_1 x_2 x_3 x_4.$$

*Proof.* Since  $P_n(x)$  has degree  $n$ , it suffices to show that  $x_i$  divides  $P_n(x)$  for all  $i \in S \cup \{1\}$ . This will be accomplished by demonstrating that  $P_n(x) = 0$  if  $x_i = 0$ . We distinguish two cases.

First let  $i \in S$  be fixed but arbitrary and suppose that the coefficient of  $x_i$  is  $+1$  for some  $\sigma(S_{k_1}^m)$ ,  $0 \leq m \leq n-2$ . (The argument for  $-1$  is similar.) By construction, there exists exactly one  $k_2$  such that the coefficient of  $x_i$  in  $\sigma(S_{k_2}^{m+1})$  is  $-1$ , while all other coefficients are precisely the same as those in  $\sigma(S_{k_1}^m)$ . Hence

$$\sigma(S_{k_1}^m)|_{x_i=0} = \sigma(S_{k_2}^{m+1})|_{x_i=0},$$

and of course equality is preserved when both sides are raised to the  $n$ th power. Since from Equation (4) the operation corresponding to the left-hand term in  $P_n(x)$  is  $(-1)^m$  while that of the right is  $(-1)^{m+1}$ , it follows that  $P_n(x)|_{x_i=0} = 0$  for  $i \in S$ .

Now let  $i = 1$ . By construction (and via reasoning similar to above), for each  $k_1$  there exists exactly one  $k_2$  such that

$$\sigma(S_{k_1}^m)|_{x_1=0} = -\sigma(S_{k_2}^{n-1-m})|_{x_1=0}.$$

Since the operation corresponding to the left-hand term in  $P_n(x)$  is  $(-1)^m$  while that

of the right is  $(-1)^n(-1)^{n-1-m} = (-1)^{m+1}$ , we have  $P_n(x)|_{x_1=0} = 0$ . Therefore  $P_n(x) = C_n \prod_{i=1}^n x_i$  for some constant  $C_n$ .

To pin down the value of  $C_n$  for completeness (and to show that it is not identically zero), observe that from the Multinomial Theorem

$$\left( \sum_{i=1}^p x_i \right)^n = \sum \frac{n!}{j_1! \cdots j_p!} x_1^{j_1} \cdots x_p^{j_p}$$

(where the sum is taken over all nonnegative integers  $j_1 + \cdots + j_p = n$ ) for  $p = n$ , the coefficient of  $x_1 \cdots x_n$  is  $n!$ . Hence the coefficient of  $x_1 \cdots x_n$  in  $(-1)^m [\sigma(S_k^m)]^n$  is  $(-1)^m \cdot (-1)^m n! = n!$ . Since there are  $|\mathcal{P}(S)| = 2^{n-1}$  of these terms, the result follows.

An interesting consequence is the summation identity

$$\sum_{m=0}^{n-1} (-1)^m \binom{n-1}{m} (n-2m)^n = 2^{n-1} n!,$$

obtained by choosing all  $x_i = 1$  in Equations (4) and (5).

By replacing  $x_i$  with  $a_i x_i$  in a manner analogous to the  $n = 2$  case, this procedure gives rise to a way of generating integers that are expressible as a sum of  $2^{n-1}$  powers of  $n$  in at least two ways. Note that one may obtain similar relations for lower powers simply by taking successive partial derivatives of Equations (5) or (6) with respect to  $x_1$ .

We end by mentioning the possible existence of a converse to the above, analogous to the  $n = 2$  case, and pose it as an open problem.

**Conjecture** Under suitable necessary and sufficient conditions on  $n$ -tuples  $a$  and  $x$  (find them), the general solution of

$$\sum_{i=1}^{2^{n-1}} A_i^n = \sum_{i=1}^{2^{n-1}} B_i^n$$

in integers is characterized by (a slight algebraic reformulation of) the  $n! - 1$  relations given by Equation (6).

The interested reader is referred to [1] for related material.

**Acknowledgment.** The author gratefully wishes to acknowledge the Naval Postgraduate School for its support during this research, as well as the referees for their many helpful suggestions and comments.

## REFERENCES

1. Aissen, M., Some identities involving partitions, *Annals New York Academy of Sciences*, 175 (1970), 7-22.
2. Dickson, L. E., *Introduction to the Theory of Numbers*, Dover Publications, Mineola, NY, 1957 (originally by University of Chicago Press, 1929).

# Smith Numbers

UNDERWOOD DUDLEY

DePauw University  
Greencastle, IN 46135

Were a reporter to ask any of us about a piece of mathematical research appearing in a journal, we could answer five of the traditional six questions instantly: The answers to “Who?,” “Where?,” and “When?” are obvious, and while the answer to “What?” might be harder to make clear to a reporter, it too would be obvious to any of us expert in the field of the research. “How?” would have to be replied to with some appeal to the mysteries of the creative process and the unfathomable sources of inspiration, but that too we could do instantly. On the other hand, “Why?” could not be responded to as quickly or, in some cases, at all, because it is not a question that we ask ourselves very often. We ought to ask it more often. It is a good thing to know why we do what we do.

The purpose of this note is to give a survey of Smith numbers, a recently created and thus far very small subfield of number theory, and to attempt to answer the last question for it.

The origin of Smith numbers is a paper by Albert Wilansky that appeared in the January 1982 issue of the *Two-Year College Mathematics Journal* [10]. Here it is:

A *Smith number* is a composite number the sum of whose digits is the sum of all the digits of all its prime factors. The (rather startling) reason for the name is mentioned below.

*Examples.*  $9985 = 5 \times 1997$ ;  $9 + 9 + 8 + 5 = 5 + 1 + 9 + 9 + 7$ ;  $6036 = 2 \times 2 \times 3 \times 503$ ;  $6 + 0 + 3 + 6 = 2 + 2 + 3 + 5 + 0 + 3$ .

The number of Smith numbers between  $n$  thousand and  $n$  thousand + 999 for  $n = 0, 1, 2, \dots, 9$ , is, respectively, 47, 32, 42, 28, 33, 32, 32, 37, 37, 40.

I wonder whether there are infinitely many Smith numbers.

The largest Smith number known is due to my brother-in-law H. Smith who is not a mathematician. It is his telephone number: 4937775!

The paper is a gem. Professor Wilansky is to be congratulated for sending it in, and Donald Albers, the editor of the *Journal* at the time, is to be congratulated for printing it. It is short, comprehensible, amusing, unpretentious, and admirably serves its purpose of mathematically stimulating its readers: They can verify that 4937775 is indeed a Smith number ( $4937775 = 3 \cdot 5^2 \cdot 65837$ , and  $4 + 9 + 3 + 7 + 7 + 7 + 5 = 3 + 5 + 5 + 6 + 5 + 8 + 3 + 7$ ), they can search for the first Smith number greater than 10000 ( $10086 = 2 \cdot 3 \cdot 41^2$ ), or, for a greater challenge, find a Smith number larger than the largest one then known (the next one is 4937818 =  $2 \cdot 2468909$ ), and they can even think about how to go about showing that there are infinitely many of them.

One of the characteristics of gems is that they sparkle and hence attract attention, and Professor Wilansky's was no exception. In 1983, a paper appeared in this MAGAZINE [9], giving a larger Smith number. The authors' discovery was that if  $p$  is a prime whose digits are all 1s, then  $3304p$  is a Smith number. After that clever discovery was made, proving it is easy: If  $p = 111\dots 1$  ( $n$  1s), then  $3304p$  is, on the one hand,  $3671112\dots 10744$  ( $n - 4$  1s) with digit sum

$$3 + 6 + 7 + (n - 4) + 7 + 4 + 4 = 27 + n,$$

and on the other,  $2^3 \cdot 7 \cdot 59 \cdot 111 \dots 1$  ( $n$  1s), with digit sum

$$2 + 2 + 2 + 7 + 5 + 9 + n = 27 + n.$$

However, in the proof, the authors could not keep notation from creeping in:

Let  $S(n)$  denote the sum of digits of  $n$  and let  $S_p(n)$  denote the sum of the digits of a prime factorization of  $n$ . Thus, a composite number  $n$  such that  $S(n) = S_p(n)$  is a Smith number.

They went on to note that 3304 is just one of several integers that could serve the same purpose. Moreover, if  $S(n) - S_p(n)$  is a multiple of 7, then multiplying  $n$  by a power of 10 will produce a Smith number, as from  $69 = 3 \cdot 23$ ,  $6 + 9 = 15$ ,  $3 + 2 + 3 = 8$  we get  $690 = 2 \cdot 3 \cdot 5 \cdot 23$ ,  $6 + 9 + 0 = 15$ ,  $2 + 3 + 5 + 2 + 3 = 15$ . Similarly, from the prime  $p = 111 \dots 1$  (317 ones) comes the 362-digit Smith number  $2p \cdot 10^{45}$ .

The paper is short (though longer than Wilansky's ur-paper), comprehensible (though you need a pencil and paper to verify that one of the authors' "it is easy to verify" is correct), amusing (though not as likely to raise as broad a smile as the original paper), and capable of moving readers to mathematical action. Thus, it was well worth publishing.

So was the paper that appeared in 1987 [6] that removed Professor Wilansky's wonderment about the number of Smith numbers by showing that there were indeed infinitely many of them. Its author was inspired by the preceding paper, since his method was to show that from any number  $111 \dots 1$  a Smith number could be constructed by multiplying it by an appropriate integer. It is always good to settle outstanding problems.

After that, with the outstanding problem about them settled, Smith numbers might have been left alone, but it was not to be. In a 1986 issue of this MAGAZINE another method for generating Smith numbers was presented [11], leading to Smith numbers such as

$$5 \cdot 1110110110111 \cdot (2 \cdot 5)^5 = 555055055055500000$$

and to one with 2,592,699 digits,

$$18 \times 111 \dots 1 \text{ (1031 digits)} \cdot (138 \cdot 10^{4071} + 1)^{480} \cdot 10^{636560}.$$

The shortness is going away, as is the unpretentiousness, the comprehensibility is on the decline, and the amusement, though still present—who could not be amused by the grotesque integer above?—is being accompanied by thoughts of "Why is this being done?"

The *Abstracts* of papers presented to the American Mathematical Society in 1987 contained [4]

Let  $s(n)$  denote the digital sum of  $n$  and  $s_p(n)$  denote the sum of the digital sums of the prime factors of  $n$ . An integer  $n$  is called a "Smith number" if  $s(n) = s_p(n)$ . We investigate the digital sums of powers of certain integers. In particular, we consider conditions in order that a power of an integer is a Smith number.

Smith numbers enjoyed a banner year in 1987 in the *Journal of Recreational Mathematics*, with three papers making their appearance [7, 12, 13]. There we find palindromic Smith numbers, as 12345554321, the definition of *Smith Brothers* (consecutive Smith numbers), as 728 and 729, the fact that there are more Smith numbers congruent to 4 (mod 9) than to any other residue, and another gigantic Smith number

of the same form as the 2 million digit example above but with 10,694,685 digits. This is all fun, and the *JoRM* was the proper place for it, since Smith numbers are really recreation and not serious mathematics. One reason they are not serious is that all Smith numbers depend on how integers are represented with digits in the base ten. Thus, those who investigate Smith numbers are not trying to penetrate deep into the secrets of integers; they are instead observing mere accidents of their representation in an arbitrary system. But there is nothing wrong with that as recreation.

However, 1987 also saw an ominous sign: the first generalization of Smith numbers. Generalization is usually a good thing, providing as it does opportunities to see results in wider settings and thus understand them better, but not everything benefits from being made more general. The *Fibonacci Quarterly* [6] is where the *k-Smith number* first made its appearance. The definition is obvious once it is thought of. Two examples of 2-Smith numbers are  $42 = 2 \cdot 3 \cdot 7$ ;  $2(4 + 2) = 2 + 3 + 7$ , and  $32 = 2^5$ ;  $2(3 + 2) = 2 + 2 + 2 + 2 + 2$ . It opens up a whole new field! After *k-Smith numbers* are exhausted, we can consider  $(r, s)$ -Smith numbers, where  $r$  times the sum of the digits of an integer equals  $s$  times the sum of the digits of its prime divisors. Then can come  $(r, s, b)$ -Smith numbers, where we write numbers in the base  $b$ , and after that we can start all over with *Smith numbers of the second kind*, where we count only *distinct* prime divisors. We can start the *Journal of Smith Numbers and Their Applications*. *Lifetimes* of mathematical effort can be spent on Smith numbers.

None of this has happened yet, but it may. In 1988 in the *College Mathematics Journal* there was a review of an article on Smith numbers [2] that concluded

The vistas for research on Smith numbers and their generalizations are opening up.

Irony may have been intended, but a danger with irony is that some readers will not catch it.

In 1989 a Smith number paper appeared in a Dutch mathematics journal [8]. I was greatly pleased to read its review in *Mathematical Reviews* [5], because it concluded

The authors give a history of the problem and prove two propositions that produce some pretty large Smith numbers. There are [sic] also a list of serious number theory papers, by Lucas, Kummer and others that mention digits (usually to a prime base). But the reviewer is not convinced thereby that Smith numbers are not a rathole down which valuable mathematical effort is being poured.

*Mathematical Reviews* seldom contains such expression of opinion.

Smith numbers returned to the *Journal of Recreational Mathematics* in 1990 [1] in a paper that introduced *Jones numbers* (they are next door to Smith numbers, you see). That is appropriate whimsy, but the rest of the article sounds, if not dead serious, at least as if the author had his tongue firmly out of his cheek.

It remains to be seen what will become of the Smith number phenomenon, but the reason for it is clear, I think. It is that mathematics is hard. In other sciences, the biological sciences especially, I think it is accurate to say that there are many more problems to work on than there are people to work on them. There are millions of undiscovered species to describe. The human body contains thousands of unexplained mysteries. Any person with the training that goes with the attainment of the Ph.D. degree should be able to find something to work on that is new and potentially useful. Mathematicians sometimes sourly observe that it is wonderful to be a chemist: All you have to do is boil something new and you have a paper. And to be a sociologist! You

send out a survey, tabulate its results, and then you have a *book*. In mathematics it is different. There are hardly any small but useful projects that can be done by people with sufficient training but limited talent. Hence the proliferation of Smith numbers. They serve the purpose of maintaining a measure of mathematical vitality in those who read and write about them, but they have the potential to mislead about the purpose of the mathematical enterprise.

In *A Mathematician's Apology* [3], G. H. Hardy wrote

It is a melancholy experience for a professional mathematician to find himself writing about mathematics. The function of a mathematician is to do something...

It is also a melancholy experience for a professional mathematician not to be able to do something of significance. The answer, then, to "Why Smith numbers?" is, "We have to do *something*."

We will not go into the question of how much of modern mathematics is analogous to Smith numbers.

## REFERENCES

1. Bishop, Robert L., How to generate all types of Smith numbers, *J. of Recreational Math.* 22 (1990), 262–270.
  2. Bushaw, Donald L., Review of "Smith numbers", *College Math. J.* 19 (1988), 375.
  3. Hardy, G. H., *A Mathematician's Apology*, Cambridge University Press, London, 1967.
  4. Kennedy, Robert E. and Curtis N. Cooper, Abstract 831-11-298, *Abstracts of Papers Presented to the Amer. Math. Society* 8 (1987), 18.
  5. Linderholm, Carl, Review 90g:11014, *Mathematical Reviews* (1990), 3799.
  6. McDaniel, Wayne, The existence of infinitely many  $k$ -Smith numbers, *Fibonacci Quarterly* 25 (1987), 76–80.
  7. McDaniel, Wayne L., Palindromic Smith numbers, *J. of Recreational Math.* 19 (1987), 34–37.
  8. McDaniel, Wayne L. and Samuel Yates, The sum of digits function and its application to a generalization of the Smith number problem, *Nieuw Archief voor Wiskunde* (4) 7 (1989), 39–51.
  9. Oltikar, Sham and Keith Wayland, Construction of Smith numbers, this *MAGAZINE* 56 (1983), 36–37.
  10. Wilansky, A., Smith numbers, *Two-Year College Math. J.* 13 (1982), 21.
  11. Yates, Samuel, Special sets of Smith numbers, this *MAGAZINE* 59 (1986), 293–296.
  12. Yates, Samuel, Smith numbers congruent to 4 (mod 9), *J. of Recreational Math.* 19 (1987), 139–141.
  13. Yates, Samuel, How odd the Smiths are, *J. of Recreational Math.* 19 (1987), 268–274.
-



---

# PROBLEMS

---

LOREN C. LARSON, *editor*  
St. Olaf College

GEORGE GILBERT, *associate editor*  
Texas Christian University

## Proposals

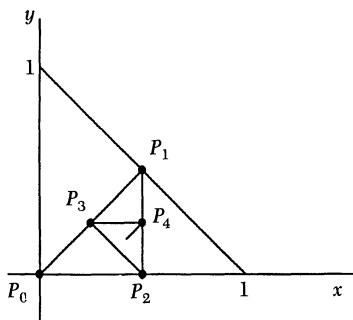
*To be considered for publication, solutions should be received by July 1, 1994.*

**1438.** *Proposed by David M. Bloom, Brooklyn College of CUNY, Brooklyn, New York.*

Let  $p$  be a prime with  $p > 5$ , and let  $S = \{p - n^2 : n \in \mathbb{Z}^+, n^2 < p\}$ . (For example, if  $p = 31$  then  $S = \{6, 15, 22, 27, 30\}$ .) Prove that  $S$  contains two elements  $a, b$  such that  $1 < a < b$  and  $a$  divides  $b$ .

**1439.** *Proposed by Charles Vanden Eynden, Illinois State University, Normal, Illinois.*

All the lines in the sketch have slope 0, 1, or  $-1$ , or are vertical. What point do the points  $P_n$  approach?



---

ASSISTANT EDITORS: CLIFTON CORZAT, BRUCE HANSON, RICHARD KLEBER, KAY SMITH, and THEODORE VESSEY, *St. Olaf College* and MARK KRUSEMEYER, *Carleton College*. We invite readers to submit problems believed to be new and appealing to students and teachers of advanced undergraduate mathematics. Proposals should be accompanied by solutions, if at all possible, and by any other information that will assist the editors and referees. A problem submitted as a Quickie should have an unexpected, succinct solution. An asterisk (\*) next to a problem number indicates that neither the proposer nor the editors supplied a solution.

Solutions should be written in a style appropriate for *Mathematics Magazine*. Each solution should begin on a separate sheet containing the solver's name and full address.

Solutions and new proposals should be mailed in duplicate to Loren Larson, Department of Mathematics, St. Olaf College, 1520 St. Olaf Ave., Northfield, MN 55057-1098 or mailed electronically via fax: (507) 663-3549 or e-mail: [larson@stolaf.edu](mailto:larson@stolaf.edu).

**1440.** *Proposed by Bruce Reznick, University of Illinois, Urbana, Illinois.*

Let  $r \geq 2$  be an integer. We say that a set of positive integers  $A$  is  **$r$ -fold-free** if  $k \in A$  implies  $rk \notin A$ . Let  $f_r(n)$  denote the cardinality of the largest  $r$ -fold-free subset of  $\{1, 2, \dots, n\}$ . It is easy to calculate  $f_r(n)$  for small values of  $r$  and  $n$  and corresponding sets of maximal cardinality. For example,  $f_2(1) = 1, (\{1\})$ ;  $f_2(2) = 1, (\{1\} \text{ or } \{2\})$ ;  $f_2(3) = 2, (\{1, 3\} \text{ or } \{2, 3\})$ ;  $f_2(4) = 3, (\{1, 3, 4\})$ ;  $f_2(5) = 4, (\{1, 3, 4, 5\})$ .

Prove the following closed formula for  $f_r(n)$ . Let  $[a_m a_{m-1} \dots a_0]_r$  denote the base  $r$  representation of  $n$ . Then

$$f_r(n) = \left( \frac{1}{r+1} \right) \left( r \cdot n + \sum_{k=0}^m (-1)^k a_k \right).$$

For example,  $3000 = [101110111000]_2$ , so

$$f_2(3000) = (6000 - 1 + 1 - 1 - 1 + 1 - 1 - 1)/3 = 1999.$$

**1441.** *Proposed by S. C. Woon, Imperial College, London, United Kingdom.*

Let  $\pi(0) = 1$  and for each integer  $n \geq 1$ ,

$$\pi(n) = \sqrt{1 + \left( \sum_{k=0}^{n-1} \pi(k) \right)^2}.$$

Show that

$$\lim_{n \rightarrow \infty} \frac{2^{n+1}}{\pi(n)} = \pi.$$

**1442.** *Proposed by W. O. Egerland and C. E. Hansen, Aberdeen Proving Ground, Maryland.*

Prove that two ellipses with exactly one focus in common intersect in at most two points.

## Quickies

*Answers to the Quickies are on page 74.*

**Q814.** *Proposed by Paul Erdős, Hungarian Academy of Sciences, Budapest, Hungary.*

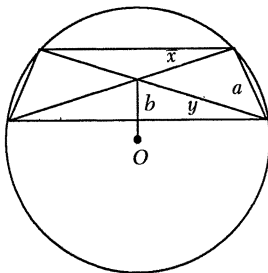
Let  $a_1 = 1, a_2 = 2, \dots$  be the sequence of positive integers of the form  $2^\alpha 3^\beta$ , where  $\alpha$  and  $\beta$  are nonnegative integers. Prove that every positive integer is expressible in the form  $a_{i_1} + a_{i_2} + \dots + a_{i_n}$ , where no summand is a multiple of any other.

**Q815.** *Proposed by Murray S. Klamkin, University of Alberta, Edmonton, Alberta, Canada.*

$A, B, C, D$  is a convex quadrilateral inscribed in the base of a right circular cone of vertex  $P$ . Show that for the pyramid  $PABCD$ , the sum of the dihedral angles with edges  $PA$  and  $PC$  equals the sum of the dihedral angles with edges  $PB$  and  $PD$ .

**Q816.** *Proposed by Michael Handelsman, Erasmus Hall High School, Brooklyn, New York.*

An isosceles trapezoid has a unit circumradius, a leg of length  $a$ , and diagonals that intersect at a distance  $b$  from the circumcenter. Also, the diagonals segment each other into lengths of  $x$  and  $y$ , with  $y > x$ . Show that  $y - x = ab$ .



## Solutions

### Union of lines

February 1993

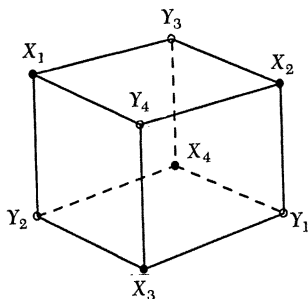
**1413.** *Proposed by Victor Klee, University of Washington, Seattle, Washington.*

For each subset  $X$  of  $\mathbf{R}^3$ , let  $\text{Con}_2(X)$  denote the union of  $X$  with all line segments joining pairs of points of  $X$ , and let  $\text{Aff}_2(X)$  denote the union of  $X$  with all lines determined by two points of  $X$ . Now suppose that  $X$  consists of the four vertices of a tetrahedron. Then  $\text{Con}_2(X)$  is the union of the six edges of the tetrahedron, and  $\text{Aff}_2(X)$  is the union of the six lines that are extensions of those edges. Also,  $\text{Con}_2(\text{Con}_2(X))$  is the entire tetrahedron. Give a geometric description of the set  $\text{Aff}_2(\text{Aff}_2(X))$ .

*Solution by the Con Amore Problem Group, Royal Danish School of Educational Studies, Copenhagen, Denmark.*

Notation: For arbitrary distinct points  $A$  and  $B$ , we use  $\lambda(A, B)$  for the line through  $A$  and  $B$ . For arbitrary noncollinear points  $A$ ,  $B$ , and  $C$ , we use  $\pi(A, B, C)$  for the plane through  $A$ ,  $B$ , and  $C$ . Using a linear transformation in  $\mathbf{R}^3$  if necessary, we may assume that  $X = \{X_1, X_2, X_3, X_4\}$  where  $X_1, X_2, X_3, X_4$  are vertices of a cube as in the following figure. Their opposite vertices will be called  $Y_1, Y_2, Y_3, Y_4$ , respectively. We will show that

$$\text{Aff}_2(\text{Aff}_2(X)) = \mathbf{R}^3 \setminus \{Y_1, Y_2, Y_3, Y_4\}.$$

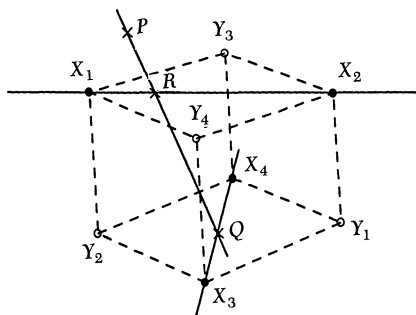
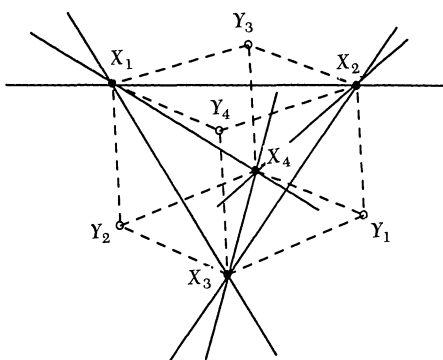


Letting our initial transformation go backwards, our  $Y_i$ 's will, of course, be replaced with vertices of a suitably defined parallelepiped. (The definition is that it is bounded by six planes, each of them through an edge of the tetrahedron  $X_1X_2X_3X_4$ , and parallel to the opposite edge.)

Let  $Z = \text{Aff}_2(X) = \bigcup_{\substack{i,j=1 \\ i \neq j}}^4 \lambda(X_i, X_j)$ , and  $W = \text{Aff}_2(\text{Aff}_2(X)) = \text{Aff}_2(Z)$ . We first note

that  $Y_1, Y_2, Y_3, Y_4 \notin W$ . For, as is clear from the following figure, a line through  $Y_4$  and a point of  $Z$  cannot meet  $Z$  in a second point, and similarly for  $Y_1, Y_2, Y_3$ .

On the other hand, consider a point  $P \neq Y_1, Y_2, Y_3, Y_4$ . We claim that  $P \in W$ . This is clear if  $P$  is  $X_1, X_2, X_3$ , or  $X_4$ . Otherwise  $P$  is either: not in  $\pi(X_1, X_2, X_3) \cup \pi(X_3, X_4, Y_1)$ , not in  $\pi(X_2, X_3, Y_4) \cup \pi(X_4, X_1, Y_2)$ , or not in  $\pi(X_1, X_3, Y_4) \cup \pi(X_2, X_4, Y_1)$ . Suppose, without loss of generality, that  $P \notin \pi(X_1, X_2, Y_3) \cup \pi(X_3, X_4, Y_1)$ .



Then  $\pi(X_1, X_2, P)$  meets  $\lambda(X_3, X_4)$  in a point  $Q$ , and  $\lambda(P, Q)$  meets  $\lambda(X_1, X_2)$  in a point  $R$ . Since  $Q, R \in Z$ ,  $P \in W$ , and this completes the proof.

Also solved by Manjul Bhargava (student), Momcilo Bjelica (Yugoslavia), Clayton Brooks, O. P. Lossers (The Netherlands), Howard Morris, Andreas Müller (Germany), Jack V. Wales, Jr., A. N't Woord (The Netherlands), John S. Sumner and Kevin L. Dove, and the proposer.

## Buffon's square

February 1993

**1414.** Proposed by Stan Wagon, Macalester College, St. Paul, Minnesota.

Suppose a square whose diagonal's length is  $\frac{5}{9}$  of an inch is thrown randomly (with uniform distribution) onto a flat surface ruled with parallel lines one inch apart. What is the probability that the square will touch one of the lines?

**I. Solution by Jerrold W. Grossman, Oakland University, Rochester, Michigan.**

Think of the parallel rules as being the lines  $y = n$  for all integers  $n$ , lying in the coordinate plane. We interpret the randomness assumption to imply that the smaller acute angle  $\theta$  made by a diagonal of the square with the horizontal has a uniform distribution on  $[0, \pi/4]$ . For given  $\theta$  it is easy to see that the vertical height of the square is  $\frac{5}{9} \cos \theta$ . Therefore the probability that a randomly thrown square yielding this value of  $\theta$  touches a line is  $\frac{5}{9} \cos \theta$ , and so the desired probability is the average of this value over the given interval:

$$\frac{4}{\pi} \int_0^{\pi/4} \frac{5}{9} \cos \theta \, d\theta = \frac{10\sqrt{2}}{9\pi} \approx 0.5002.$$

## II. Solution by David Callan, University of Wisconsin, Madison, Wisconsin.

This is an immediate consequence of a remarkable (and easily proved) generalization [1] of Buffon's classic needle problem: Suppose the plane is ruled with parallel lines distance  $w$  apart and suppose a convex polygon of perimeter  $l$  is tossed at random onto the plane. Then provided that the diameter of the polygon is smaller than  $w$ , so that it can intersect at most one ruled line, the probability that it does so is  $l/(\pi w)$ , independent of its shape! In fact, by taking limits, the result applies to any convex region of diameter  $< w$ . To recover Buffon's original result, consider the needle as a polygon of two sides.

## REFERENCE

1. J. V. Uspensky, *Introduction to Mathematical Probability*, 1937, McGraw Hill, pp. 251–253.

## III. Solution by Peter A. Rogerson, Center for Advanced Study in the Behavioral Sciences, Stanford, California.

The solution may be found using a result of Kendall and Moran [1]. They state that a line of length  $L$ , twisted into a definite shape, has an expected number of intersections with a set of unit-spaced lines equal to  $2L/\pi$ . (Also, see J. F. Ramely, "Buffon's noodle problem," *American Mathematical Monthly*, October, 1969, pp. 916–918.) For a square with diagonal equal to  $\frac{5}{9}$ ,

$$L = 20 \frac{\sqrt{2}}{18},$$

and the expected number of intersections is

$$40 \frac{\sqrt{2}}{18\pi}.$$

Call the desired probability  $p$ . When the square is thrown onto the lines it has two crossings with probability  $p$ , and zero crossings with probability  $1 - p$ . Thus

$$40 \frac{\sqrt{2}}{18\pi} = 0(1 - p) + 2p,$$

implying that

$$p = 20 \frac{\sqrt{2}}{18\pi} \cong .500176.$$

More generally, for squares with diagonal  $D < 1$  for unit spaced lines,

$$p = 2\sqrt{2} \frac{D}{\pi}.$$

## REFERENCE

1. Kendall, M. G. and Moran, P. A. P., *Geometrical Probability*, Hafner, New York, 1963.

Also solved by Robert A. Agnew, Michael Bertrand, Manjul Bhargava (student), Warren L. Bosch, Stephen D. Bronn, Con Amore Problem Group (Denmark), Bill Correll, Jr. (student), Richard Croutharmel, William E. Daniels, Robert L. Doucette, Mark Dugopolski, Milton P. Eisner, Herbert Gintis, Lee O. Hagglund, James Henderson, James C. Hickman, R. Daniel Hurwitz and David Vella, Hans Georg Killingbergtrø and Ivar Skau (Norway), Emil F. Knapp, Norman F. Lindquist, O. P. Lossers (The Netherlands), Andreas Müller (Germany), National Security Agency Problems Group, Greg Neumer, New

Mexico Tech Problem Solving Group, Larry Olson, Edward D. Onstott, F. C. Rembis, Michael Rosing (student), John P. Schneider (student), Harry Sedinger, Robert W. Sheets, SPOAK (Oklahoma State University), Daniel L. Stock, John S. Sumner and Kevin L. Dove, R. S. Tiberio, Trinity University Problem Group, Robert J. Wagner, Jack V. Wales, Jr., Rob Williams, A. N't Woord (The Netherlands), Harald Ziehms (Germany), Michael L. Zwilling, and the proposer.

### Determinant with binomial coefficient entries

February 1993

**1415.** Proposed by John Duncan and Mihalis Maliakas, University of Arkansas, Fayetteville, Arkansas.

Let  $m$  and  $n$  be integers satisfying  $m \geq 2n - 1 \geq 3$ , and let  $A(m, n)$  denote the  $n \times n$  matrix whose  $(i, j)$  entry is

$$\binom{m-j+1}{2n-2i} + \binom{m+j-2n}{2n-2i}.$$

Evaluate the determinant of  $A(m, n)$ .

*Solution by John S. Sumner, University of Tampa, Tampa, Florida.*

We prove by induction on  $n$  that the determinant is  $2^n$ . It is trivial to verify this for  $n = 2$ . Suppose  $n \geq 3$  and the determinant of  $A(m, n-1) = 2^{n-1}$ . Note that each entry in row  $n$  of  $A(m, n)$  equals 2. For  $j = 1, 2, \dots, n-1$ , subtract column  $j+1$  from column  $j$ . Call this new matrix  $A'(m, n)$ . Then  $\det A(m, n) = \det A'(m, n)$ . Let  $B$  be the  $(n-1) \times (n-1)$  upper left-hand submatrix of  $A'(m, n)$ . Then  $\det A'(m, n) = 2 \det B$  and the  $(i, j)$  entry of  $B$  equals

$$\begin{aligned} & \binom{m-j+1}{2n-2i} + \binom{m+j-2n}{2n-2i} - \binom{m-j}{2n-2i} - \binom{m+j+1-2n}{2n-2i} \\ &= \binom{m-j}{2n-2i-1} - \binom{m+j-2n}{2n-2i-1}. \end{aligned}$$

For  $j = 1, 2, \dots, n-2$ , subtract column  $j+1$  of  $B$  from column  $j$  of  $B$ . Let  $B'$  be this new matrix. Then  $\det B = \det B'$  and the  $(i, j)$  entry of  $B'$  equals

$$\begin{aligned} & \binom{m-j}{2n-2i-1} - \binom{m+j-2n}{2n-2i-1} - \binom{m-j-1}{2n-2i-1} + \binom{m+j+1-2n}{2n-2i-1} \\ &= \binom{m-j-1}{2n-2i-2} + \binom{m+j-2n}{2n-2i-2} \end{aligned}$$

if  $j < n-1$ . For each  $i$ , the  $(i, n-1)$  entry of  $B'$  is

$$\begin{aligned} & \binom{m-n+1}{2n-2i-1} - \binom{m-n-1}{2n-2i-1} \\ &= \binom{m-n+1}{2n-2i-1} - \binom{m-n}{2n-2i-1} + \binom{m-n}{2n-2i-1} - \binom{m-n-1}{2n-2i-1} \\ &= \binom{m-n}{2n-2i-2} + \binom{m-n-1}{2n-2i-2}. \end{aligned}$$

Thus,  $B' = A(m-2, n-1)$  and  $\det(B') = 2^{n-1}$  by induction. This completes the proof.

Also solved by John Andraos (Canada), Manjul Bhargava (student), J. C. Binz (Switzerland), Carleton College Problem Solvers, Con Amore Problem Group (Denmark), Bill Correll, Jr. (student), Robert L. Doucette, Trinity University Problem Group, and the proposer.

## Two bounding conditions on a series

February 1993

**1416.** Proposed by Stephen Wainger, University of Wisconsin, Madison, Wisconsin, and Jim Wright, Texas Christian University, Fort Worth, Texas.

Let  $(a_k)_{k=1}^{\infty}$  be a sequence of positive numbers, and consider the following two conditions.

(I) There is a constant  $C_1$  such that

$$\sum_{k=1}^n a_k \leq C_1 a_n \quad \text{for all } n \geq 1.$$

(II) There is a constant  $C_2$  such that

$$\sum_{k=n}^{\infty} \frac{1}{a_k} \leq C_2 \frac{1}{a_n} \quad \text{for all } n \geq 1.$$

Which of these conditions (if either) implies the other?

*Solution by the Trinity University Problem Group, Trinity University, San Antonio, Texas.*

Statements (I) and (II) are equivalent.

(I)  $\Rightarrow$  (II). Let  $n \in \mathbb{N}$  be arbitrary and set  $S_n = \sum_{k=1}^n a_k$ . By (I) we have  $S_n + a_{n+1} \leq C_1 a_{n+1}$  from which it follows that

$$a_{n+1} \geq \frac{S_n}{C_1 - 1}.$$

Now combine this inequality with a further appeal to (I) to get  $C_1 a_{n+2} \geq S_n + a_{n+1} + a_{n+2} \geq S_n + S_n/(C_1 - 1) + a_{n+2}$  from which the inequality

$$a_{n+2} \geq \frac{S_n}{C_1 - 1} \left( \frac{C_1}{C_1 - 1} \right)$$

obtains. An easy induction now shows that

$$a_{n+k} \geq \frac{S_n}{C_1 - 1} \left( \frac{C_1}{C_1 - 1} \right)^{k-1}, \quad k \geq 1.$$

Thus,

$$\begin{aligned} \sum_{k=n}^{\infty} \frac{1}{a_k} &= \sum_{k=0}^{\infty} \frac{1}{a_{n+k}} = \frac{1}{a_n} + \sum_{k=1}^{\infty} \frac{1}{a_{n+k}} \\ &\leq \frac{1}{a_n} + \frac{C_1 - 1}{S_n} \sum_{k=1}^{\infty} \left( \frac{C_1 - 1}{C_1} \right)^{k-1} = \frac{1}{a_n} + \frac{C_1 - 1}{S_n} C_1. \end{aligned}$$

But clearly  $S_n \geq a_n$  so

$$\sum_{k=n}^{\infty} \frac{1}{a_k} \leq (C_1^2 - C_1 + 1) \frac{1}{a_n}.$$

Since  $n$  is arbitrary, (II) holds.

(II)  $\Rightarrow$  (I). Again, let  $n \in \mathbb{N}$  be arbitrary and set  $S_n = \sum_{k=n}^{\infty} 1/a_k$ . Condition (II) gives  $S_n + 1/a_{n-1} \leq C_2/a_{n-1}$ , from which

$$\frac{1}{a_{n-1}} \geq \frac{S_n}{C_2 - 1}$$

follows. Further use of (II), in conjunction with the last inequality, gives  $C_2/a_{n-2} \geq S_n + 1/a_{n-1} + 1/a_{n-2} \geq S_n + S_n/(C_2 - 1) + 1/a_{n-2}$  and, after solving for  $1/a_{n-2}$ , the inequality

$$\frac{1}{a_{n-2}} \geq \frac{S_n}{C_2 - 1} \left( \frac{C_2}{C_2 - 1} \right)$$

obtains. Continuing in a similar manner, we get that, for  $1 \leq k \leq n - 1$ ,

$$\frac{1}{a_{n-k}} \geq \frac{S_n}{C_2 - 1} \left( \frac{C_2}{C_2 - 1} \right)^{k-1}.$$

Taking reciprocals, adding, and using the fact that  $S_n \geq 1/a_n$  gives

$$\begin{aligned} a_1 + \cdots + a_n &\leq \frac{C_2 - 1}{S_n} \left[ 1 + \left( \frac{C_2 - 1}{C_2} \right) + \cdots + \left( \frac{C_2 - 1}{C_2} \right)^{n-2} \right] + a_n \\ &\leq \frac{C_2 - 1}{S_n} \sum_{k=0}^{\infty} \left( \frac{C_2 - 1}{C_2} \right)^k + a_n \leq (C_2^2 - C_2 + 1) a_n. \end{aligned}$$

As  $n$  is arbitrary, (I) holds and the proof is complete.

*Also solved by Ed Adams, David Callan, Con Amore Problem Group (Denmark), Hans Georg Killingbergtrø and Ivar Skau (Norway), Kee-Wai Lau (Hong Kong), Howard Morris, Harvey J. Schmidt, Jr., Vu Ha Van (student, Hungary), A. N't Woord (The Netherlands), and the proposers.*

## Billiard balls in collision

February 1993

**1417.** *Proposed by C. Kenneth Fan, Massachusetts Institute of Technology, Cambridge, Massachusetts.*

Let  $N$  billiard balls (of various positive radii and masses) roll about a frictionless, rectangular billiard table. Assume the collisions are elastic. Show that there will either be no collisions at all, or infinitely many.

*Solution by the proposer and Bjorn Poonen, University of California at Berkeley, Berkeley, California.*

It suffices to show that if two billiard balls collide, there will ensue another collision.

Let the dimensions of the table be  $L$  by  $W$ . Consider two balls of radii  $r_1$  and  $r_2$ .

To each ball associate a rectangular tiling of the plane by taking an  $L - 2r_k$  by  $W - 2r_k$  rectangle and reflecting it about its sides until the entire plane is filled. This converts billiard paths of each ball into straight lines, assuming there are no collisions.

Let  $E$  be the product of the two planes so formed. Note that on  $E$ , the motion of both balls is modeled by a single straight line path. Let  $D$  be the obvious lattice whose fundamental domain is a  $2(L - 2r_1)$  by  $2(W - 2r_1)$  by  $2(L - 2r_2)$  by  $2(W - 2r_2)$  hyperblock in  $E$ . The factors of 2 are designed so that points in  $E$  modulo  $D$  refer to the same point on the billiard table.

Let  $C$  denote the set of points in  $E$  which represent impossible states; that is, states where the two balls strictly overlap. Observe that  $C$  is an open subset of  $E$ .

Now suppose our two balls have collided. Let  $R(t)$  represent the ensuing straight line motion of the two balls in  $E$ , extended linearly in both directions, with  $t = 0$  the time of collision. If this straight line is not adhered to for positive  $t$ , it means a collision must take place. Observe that  $R(0) \in \partial C$ , and that  $R(-\delta) \in C$  for some small  $\delta > 0$ . The first condition says that at the moment of collision, the two balls touch, and the second says that there was actually a collision.



A ray starting at one point of a lattice must pass within  $\varepsilon$  of infinitely many other lattice points, so we can pick a time  $T > \delta$  such that  $R(T) - R(0)$  is within  $\varepsilon$  of a vector in  $D$ . Then  $R(T - \delta) - R(-\delta) = R(T) - R(0)$  is within  $\varepsilon$  of a vector in  $D$ , so  $R(T - \delta)$  is in  $C$ , if  $\varepsilon$  is small enough. This implies that a collision takes place.

*Also solved by Hans Georg Killingbergtrø and Ivar Skau (Norway), Andreas Müller (Germany), F. C. Rembis, Western Maryland College Problems Group, and A. N't Woord (The Netherlands).*

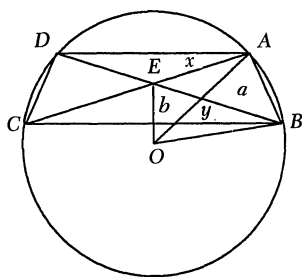
## Answers

*Solutions to the Quickies on page 67.*

**A814.** There is a simple proof by induction. If the assertion is true for all positive integers up to  $n$ , we apply the inductive assumption to  $(n+1)/2$  if  $n+1$  is even and to  $(n+1-3^\beta)/2$  if  $n+1$  is odd, where  $3^\beta \leq n+1 < 3^{\beta+1}$ .

**A815.** Equivalently, we want to show that for a cyclic spherical quadrilateral, the sum of one pair of opposite angles equals the sum of the other pair of opposite angles. Let  $O$  be the pole of the circumcircle of  $ABCD$ . Since triangles  $AOB$ ,  $BOC$ ,  $COD$ , and  $DOA$  are isosceles, the result follows. Note that it does not matter if  $O$  lies in the interior of  $ABCD$  or not.

**A816.** Consider the following figure, and note that  $\angle AEB = \angle AOB$ . Thus,  $A$ ,  $B$ ,  $O$  and  $E$  are concyclic. By Ptolemy's theorem for a cyclic quadrilateral, the sum of the products of the opposite sides equals the product of the diagonals:  $x \cdot 1 + a \cdot b = y \cdot 1$ . Hence,  $y - x = ab$ .



---

# REVIEWS

---

PAUL J. CAMPBELL, *editor*  
Beloit College

*Assistant Editor:* Eric S. Rosenthal, West Orange, NJ. Articles and books are selected for this section to call attention to interesting mathematical exposition that occurs outside the mainstream of the mathematics literature. Readers are invited to suggest items for review to the editors.

vos Savant, Marilyn, Ask Marilyn: The world's most famous math problem has finally been solved ... or has it? *Parade Magazine* (21 November 1993) 16. *The World's Most Famous Math Problem (The Proof of Fermat's Last Theorem and Other Mathematical Mysteries)*, St. Martin's Press, 1993; xii + 80 pp, \$7.95, C\$10.99 (P). ISBN 0-312-10657-2

You were probably asked at holiday parties for your opinion about the Fermat-Wiles Theorem and Marilyn vos Savant's claim that it has not been proved. You may remember that two years ago she revived the "Monty Hall problem" (involving cars and a goat), shook off vilification by Ph.D.'s who defended an incorrect naïve solution, and offered her own somewhat muddled explanations of the conditional probabilities involved.

This time, she has taken on more than self-righteous dunces: "I don't think that the current work [by Wiles] succeeds in proving 'Fermat's last theorem'...—even if no mathematical errors are discovered in it." She casts doubt on non-Euclidean geometry and asserts that Wiles's proof is non-Euclidean—hence dubious. She claims that Bolyai "managed to 'square the circle'—but only by using his own hyperbolic geometry... Has Fermat's last theorem been proved, or has it not? I would say it has not; if we reject a hyperbolic method of squaring the circle, we should also reject a hyperbolic proof of Fermat's last theorem... [I]t is logically inconsistent to *reject* a hyperbolic method of squaring the circle and *accept* a hyperbolic method of proving F.L.T.!"

Where, oh where, to begin? First, let me kick one crutch out from under her argument: Since there are *no* squares in hyperbolic geometry, it's pretty easy to reject a "hyperbolic method" of squaring the circle. What Bolyai did in hyperbolic geometry was "*rhombus*" a circle, showing how to construct (with straightedge and compasses) an equilateral quadrilateral with the same area as the circle of maximum area (see pp. 106–110 of *Non-Euclidean Geometry*, by Robert Bonola; New York: Dover, 1955).

Mathematicians will write her to object to lots of the rest, and Marilyn will no doubt backpedal to dig herself out of this new hole, all the while adding to her notoriety and readership. She is likely, however, to continue to insist that Wiles has not proved FLT (pp. 60–64). She objects to the "very short time" (30 years) of building the new "edifice of mathematics" upon which his proof is based, and to its logic being "widely accepted when almost no one outside of a few specialists understands it." She rejects proofs based on non-Euclidean geometries; she casts doubt on proof by contradiction and proof by induction; she confuses the reader about the difference between the provably impossible and the seemingly impossible; and she claims that Wiles's proof depends on many developments in mathematics that are "recent and poorly understood," naming "Galois theory, modular forms, deformation theory,  $P$ -adic numbers, and theories of  $L$ -functions."

Worse yet, she urges the reader to consider trying to "demolish Einstein's theories of relativity" by proving the parallel postulate (p. 68). She insults Wiles by suggesting that his sense of humor would lead him to "bedevil future generations of mathematicians" by leaving a note that he had such a proof "but the margin is too small to contain it" (p. 70).

Marilyn wrote the book in three weeks, and it appeared less than five months after Wiles's announcement. What effort, what a great shame, a lost opportunity! It communicates some of the excitement of the subject (hence, perhaps, the surprising endorsement from Martin Gardner)—but at the great expense of belittling and betraying mathematics.

Davis, Philip J., *Spirals from Theodorus to Chaos*, A K Peters, 1993; x + 237 pp, \$29.95. ISBN 1-56881-010-5

This book is a greatly amplified version of the author's Hedrick Lectures, delivered at the 75th anniversary meeting of the MAA. It includes a mathematical analysis of spirals, a study of their mathematical history, technical and historical supplements, and figures of more kinds of spirals than you have ever seen before. The name "spiral of Theodorus" derives from Plato's statement that Theodorus of Cyrene discussed the irrationality of  $\sqrt{2}, \dots, \sqrt{17}$ ; Davis applies the name to the spiral generated by a certain difference equation, then generalizes ("Theodorus Goes Wild" is the title of the third lecture). The mathematics involved is not for the faint-hearted, nor is the historical question of what authors from the past understood as they wrote.

Dewdney, A.K., *200% of Nothing: An Eye-Opening Tour through the Twists and Turns of Math Abuse and Innumeracy*, Wiley, 1993; ix + 182 pp, \$19.95. ISBN 0-471-57776-6

Dewdney, in one of his "Computer Recreations" columns for *Scientific American*, once wrote about innumeracy and math abuse. The column brought an avalanche of examples sent in by readers, whom Dewdney dubs his "math detectives." Designed for a popular audience, this book is a taxonomy of math abuse, with an appeal to readers to send in further examples for future editions.

Hills, W. Daniel, and Bruce M. Boghosian, Parallel scientific computation, *Science* 261 (13 August 1993) 856-863.

"This article gives a range of examples of parallel computations that exploit parallelism, summarizes our current understanding of which types of scientific applications are suitable for parallel machines, and discusses some of the issues involved in parallel programming." The fast Fourier transform is examined as an example of a parallel algorithm that requires "nonlocal" communication among the processors.

Beardsley, Tim, Never give a sucker an even break, *Scientific American* (October 1993) 22.

The "tit-for-tat" strategy was the surprise winner in 1980s "tournaments" of strategies for Prisoner's Dilemma. But if the player makes small changes or occasional mistakes in the strategy, then tit-for-tat is usually inferior to the "Pavlov" strategy. Pavlov cooperates after a contest in which either both parties cooperated or both defected, and defects after a contest in which only one player cooperated. Martin Nowak (University of Oxford) and Karl Sigmund (Vienna University) found that in simulations against a wide range of strategies, Pavlov tends to dominate other strategies. They say that the weakness of tit-for-tat is that "mutation allows populations to become more and more cooperative," eventually leading to an invasion by exploiters.

Cipra, Barry A., Mathematics that counts: The FFT: Making technology fly, *SIAM News* (May 1993) 1, 23.

Called by Gil Strang "the most valuable numerical algorithm in our lifetime," the fast Fourier transform is almost 30 years old. The original paper by James W. Cooley (now at the University of Rhode Island) and John W. Tukey (Princeton University and AT&T Bell Laboratories) has been cited over a thousand times. This account gives not details but history and context, including photos of the two discoverers.

Gilbert, George T., Mark I. Krusemeyer, and Loren C. Larson, *The Wohascum County Problem Book*, MAA, 1993; ix + 233 pp, \$26 (P). ISBN 0-88385-316-7

This is a delightful collection of 130 original problems for undergraduates (with solutions). Some problems are set in imaginary Wohascum County in Minnesota, and most require no mathematics beyond calculus. Sample: "If  $n$  is a positive integer, how many real solutions are there, as a function of  $n$ , to  $e^x = x^n$ ?"

Straffin, Philip D., *Game Theory and Strategy*, MAA, 1993; x + 244 pp, \$ (P). ISBN 0-88385-637-9

Based on a course in game theory and strategy taught over the past 15 years by an award-winning teacher, this book arrives in time to mark the fiftieth anniversary of the founding of mathematical game theory. "It is applicable whenever two individuals—or companies, or political parties, or nations—confront situations where the outcome for each depends on the behavior of all. What are the best strategies in such situations? If there are chances of cooperation, with whom should you cooperate, and how should you share the proceeds of cooperation?" Thirty-three chapters, almost all with exercises and answers, treat two-person zero-sum games, two-person non-zero sum games, and  $N$ -person games. [Truth in reviewing disclosure: The author is a colleague in my department.]

Brams, Steven J., Theory of moves, *American Scientist* 81 (November-December 1993) 562-570.

Brams has developed what he calls the "theory of moves" to add a "dynamic dimension" to game theory. In his theory, players start at an initial element in the payoff matrix, then alternately have the opportunity to switch strategies (but a player may decline to do so), until neither player wishes to switch further. The result is a path through the payoff matrix, and the theory allows the players to think ahead and take account the opponents' likely moves—thereby providing a better modeling of how (rational) people tend to behave. Brams applies the theory to "truels" (three-person duels), Prisoners' Dilemma, and the Iran hostage crisis, explaining how his theory can generate better outcomes and better explanations of participants' behavior.

Selvin, Paul, Harrison case: No calm after storm, *Science* 262 (15 October 1993) 324-327.

A seven-year battle was won by mathematician Jenny Harrison last July, when the chancellor of the University of California at Berkeley appointed her a tenured full professor, on recommendation of an independent review panel. Harrison had claimed that she was unfairly denied tenure on the basis of sex discrimination. Despite her reinstatement (with promotion), a backlash against her and the handling of the case by the UC administration has swelled both from within the Berkeley mathematics department and from outside.

Simmons, George F., *Calculus Gems: Brief Lives and Memorable Mathematics*, McGraw-Hill, 1992; xiv + 354 pp, (P). ISBN 0-07-057566-5

"[T]his book has been reconstructed out of two massive appendices in my 1985 calculus book, with many additions, rearrangements and minor adjustments." It provides brief biographies of 33 mathematicians from antiquity through the nineteenth century, plus 26 "nuggets"—miscellaneous topics "for breaking the routine and lifting the spirits." This is an excellent supplement to any calculus book, for humanizing the subject and revealing some of its connections to intellectual and social history.

Moritz, Robert Edouard, *Memorabilia Mathematica, or The Philomath's Quotation Book*, MAA, 1993; xiv + 410 pp, \$29 (P). ISBN 0-88385-321-3. Schmalz, Rosemary (ed.), *Out of the Mouths of Mathematicians": A Quotation Book of Philomaths*, MAA, 1993; x + 294 pp, \$24 (P). ISBN 0-88385-509-7. (Both books are available as a package for \$48.)

The first of these books is a reprint of a 1914 work, with 1140 "anecdotes, aphorisms and passages by famous mathematicians, scientists and writers." The second book is a new compilation, of twentieth-century quotations. Happily, both works are arranged by topic and indexed. Here is a short sample from each: " $\sin^2 \phi$  is odious to me, even though Laplace made use of it . . . let us write  $(\sin \phi)^2$ , but not  $\sin^2 \phi$ , which by analogy should signify  $\sin(\sin \phi)$ ."—Gauss. "The ultimate goal of mathematics is to eliminate all need for intelligent thought."—R.L. Graham, D.E. Knuth, and O. Patashnik.

Cipra, Barry, Nonlinear codes straighten up—and get to work, *Science* 262 (29 October 1993) 658–659.

Linear error-correcting codes—ones in which the sum of two codewords is another codeword—have predominated in coding theory, because they admit fast coding and decoding. Non-linear codes, which can be more efficient but lack the algebraic structure of linear codes, have not been useful because there seemed no way to find fast decoding algorithms. Roger Hammons, Jr., (Hughes Aircraft), Vijay Kumar (University of Southern California), Robert Calderbank and Neil Sloane (AT&T Bell Laboratories), and Patrick Solé (CNRS in France), working in two separate groups, have found unexpected connections between some nonlinear codes and well-known linear codes. In effect, some nonlinear codes are "linear codes in disguise," so decoding them turns out not to be so hard after all.

Sloan, Joseph, Two rectangles are constructible with tangrams: An enumeration proof using Mathematica, *Mathematica in Education* 2 (4) (1993) 7–10.

Tangrams are seven polygonal puzzle pieces cut from a square, and the puzzle consists in discovering how to assemble them (without overlapping) into various shapes. In particular, they can be assembled into the original square and also into a rectangle that is twice as wide as tall. Can they be assembled into other rectangles? Until now, the answer has been that no other rectangle has been found. Sloan converts the geometric problem into an arithmetic one, using Mathematica to enumerate the pieces (in particular, their edge lengths), list all of the possible edge sums, determine the 62 possible sides for rectangles, and reduce the number of possible rectangles; only two rectangles emerge from the search.

Vajda, Steven, *Mathematical Games and How to Play Them*, Ellis Horwood, 1992; x + 128 pp. ISBN 0-13-009275-4

This book offers a collection of solitaire and two-person games that are amenable to mathematical analysis, including an exposition of Nim theory. Also included are two short chapters on game theory and tournament rankings.

Bracewell, R.N., Numerical transforms, *Science* 248 (11 May 1990) 697–704.

Valuable for its perspective on current uses, this article surveys popular and useful transforms, including Laplace, Fourier, Hilbert, Hankel, Mellin, and Abel, as well as fast Fourier, Hartley, fast Hartley, Radon, Walsh, and  $z$ . The wavelet transform does not appear.

# NEWS AND LETTERS

Dear Editor:

I wish to comment on the very nice paper, "On finite abelian groups and parallel edges on polygons" (February 1993, pp.36-39). In this paper S. Szabó gave an elegant group-theoretic proof that *a closed polygon that has vertices coinciding with those of a regular polygon of even order always has at least two parallel edges.*

I found another proof that does not use group theory at all; in fact, it uses only congruence arithmetic. This proof might be of interest since it is short and accessible to beginning students. We begin by labelling the vertices of the regular  $n$ -gon, in order, using the numbers  $0, 1, \dots, n-1$  modulo  $n$ . A line segment connecting vertex  $x$  with vertex  $y$  is labelled by the number  $x+y$  (modulo  $n$ ). Arguing as in the article, we see that two line segments are parallel if and only if their line labels are congruent modulo  $n$ .

A polygon whose vertices coincide with the vertices of the regular  $n$ -gon must have  $n$  edges, even when the polygon has more than one connected component. The sum  $s$  of all the edge labels must be twice the sum of all the vertex labels (since each vertex label is counted twice when calculating the sum of the edge labels) so  $s = 2(0 + 1 + \dots + n-1) \equiv 0 \pmod{n}$ . However, if there are no parallel edges, then all the edge labels must be distinct (modulo  $n$ ), so  $s = 1 + 2 + \dots + n-1$ ; if  $n = 2k$ , this sum is congruent to  $-k \pmod{n}$ , contradicting  $s \equiv 0 \pmod{n}$ . Thus there must be at least a pair of parallel edges.

Iwan Praton  
St. Lawrence University  
Canton, NY 13617

Dear Editor:

Readers interested in generalizing the results "The  $n$ -th power of a  $2 \times 2$  matrix,"

by Kenneth S. Williams (December 1992), might want to refer to chapter 9 in Lancaster and Tismenetsky, *The Theory of Matrices*, Academic Press, 1985. If  $R(x)$  is the remainder after dividing  $x^n$  by the characteristic polynomial of  $A_{2 \times 2}$ , Williams' formulas are just the linear two point interpolation (when  $\alpha \neq \beta$ ) and Taylor's formulas applied to  $R(x)$  (Lagrange and Hermite interpolation), using  $\alpha^n = R(\alpha)$ ,  $\beta^n = R(\beta)$  and if  $\alpha = \beta$ ,  $n\alpha^{n-1} = R'(\alpha)$ , and then interpreting these as matrix polynomial formulas.

R.J. Gregorac  
Iowa State University  
Ames, IA 50010

Dear Editor:

In the recent note, "Extensions of a Sums-of-Squares Problem," by Kelly Jackson, Francis Masat, and Robert Mitchel (February 1993), the following open question was raised: "Given any integer  $K$ , is there an integer  $M$  that is expressible as a sum of  $N$  nonzero squares, where  $1 \leq N \leq K$ ?"

The answer is yes, and the number 169, which played an important part in that note, provides a clue.

In fact, the appearance of Pythagorean Triples satisfying the identity

$$\left(\frac{x^2+1}{2}\right)^2 = \left(\frac{x^2-1}{2}\right)^2 + x^2,$$

with  $x$  odd, in the string  $169 = 13^2 = 12^2 + 5^2 = 12^2 + 4^2 + 3^2$  suggests that we construct recursively such a solution. In particular, the next "step up" using our identity produces  $7225 = 85^2 = 84^2 + 13^2 = 84^2 + 12^2 + 5^2 = 84^2 + 12^2 + 4^2 + 3^2$ .

For the general solution, given  $K$ , we define

$$x_1 = 3, \quad x_i = \frac{x_{i-1}^2 + 1}{2}$$

for  $2 \leq i \leq k$ , and  $M = x_k^2$ . Clearly,  $M$  has the

desired property for  $M = x_k^2 = \left(\frac{x_k^2 - 1}{2}\right)^2$

$$+ x_{k-1}^2 = \left(\frac{x_{k-1}^2 - 1}{2}\right)^2 + \left(\frac{x_{k-2}^2 - 1}{2}\right)^2 + x_{k-2}^2 =$$

$$\dots = \sum_{i=1}^{n-1} \left(\frac{x_{n-i}^2 - 1}{2}\right)^2 + x_1^2.$$

We further note that these squares are all

distinct, by construction, and that there are actually an infinite number of integers with the desired property, since, for example,  $x^2 M$ , for any choice of the integer  $x$ , has the property as well.

Robert B. McNeill  
Northern Michigan University  
Marquette, MI 49855

# EXCURSIONS IN CALCULUS:

## an Interplay of the Continuous and the Discrete

Robert M. Young

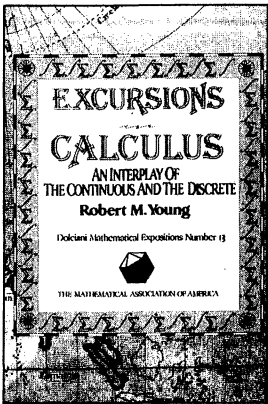
*An excellent source of projects for well motivated students. This list of 463 references is a valuable aid for those who wish to dig deeper.* —CHOICE

*The presentation is clear and the topics very interesting...fully accessible to students for whom the book is intended. The book will be influential in awakening students' awareness for good classical mathematics.* —Paulo Ribenboim

Printed with eight full-color plates.

The purpose of this book is to explore, within the context of elementary calculus, the rich and elegant interplay that exists between the two main currents of mathematics, the continuous and the discrete. Such fundamental notions in discrete mathematics as induction, recursion, combinatorics, number theory, discrete probability, and the algorithmic point of view as a unifying principle are continually explored as they interact with traditional calculus. The interaction enriches both.

The book is addressed primarily to well-trained calculus students and their teachers, but it can serve as a supplement in a traditional calculus course for anyone who wants to see more.



### CONTENTS:

- Infinite Ascent, Infinite Descent: The Principle of Mathematical Induction
- Patterns, Polynomials, and Primes: Three Applications of the Binomial Theorem
- Fibonacci Numbers: Function and Form
- On the Average
- Approximation: from Pi to the Prime Number Theorem
- Infinite Sums: A Potpourri

The problems, taken for the most part from probability, analysis and number theory, are an integral part of the text. Many point the reader toward further excursions. There are over 400 problems presented in this book.

408 pp., 1992, Paperbound  
ISBN 0-88385-317-5  
List: \$39.00 MAA Member: \$31.00  
Catalog Number DOL-13

### ORDER FROM:

The Mathematical Association of America  
1529 Eighteenth Street, NW  
Washington, DC 20036  
1-800-331-1622 Fax (202) 265-2384

Membership Code	Qty.	Catalog Number	Price
-----			
Name _____			
Address _____			
City _____			Total \$ _____
State _____ Zip Code _____			Payment <input type="checkbox"/> Check <input type="checkbox"/> VISA <input type="checkbox"/> MASTERCARD
			Credit Card No. _____
			Signature _____ Exp. Date _____



# CONTENTS

---

## ARTICLES

- 3 Trigonometric Series and Theories of Integration,  
*by Alan D. Gluchoff.*
- 20 Proof without Words: The Cauchy-Schwarz Inequality,  
*by Roger B. Nelsen.*
- 21 Mathematical Certificates, *by John Higgins and  
Douglas Campbell.*
- 28 Math Bite: Pictures, Projections, and Proverbs,  
*by Richard L. Francis.*

## NOTES

- 29 An Advanced Calculus Approach to Finding the Fermat  
Point, *by Mowaffaq Hajja.*
- 34 Proof without Words: Règle des Nombres Moyens,  
*by Roger B. Nelsen.*
- 35 Remarks on Problem B-3 on the 1990 William Lowell  
Putnam Mathematical Competition, *by Jiuqiang Liu  
and Allen J. Schwenk.*
- 40 The Hiding Path, *by Herbert R. Bailey.*
- 45 How Few Transpositions Suffice? ... You Already Know!,  
*by John O. Kiltinen.*
- 47 A Note on the Gaussian Integral, *by Constantine  
Georgakis.*
- 48 Lemniscates and Osculatory Interpolation, *by Donald Teets  
and Patrick Lang.*
- 53 Note on the Evaluation of  $\int_0^x (1/(1+t^2))^n dt$ , *by M. A.  
Gopalan and V. Ravichandran.*
- 55 How Expected is the Unexpected Hanging?,  
*by Dean Clark.*
- 59 Sums of Like Powers of Multivariate Linear Forms,  
*by Ismor Fischer.*
- 62 Smith Numbers, *by Underwood Dudley.*

## PROBLEMS

- 66 Proposals 1438–1442.
- 67 Quickies 814–816.
- 68 Solutions 1413–1417.
- 74 Answers 814–816.

## REVIEWS

- 75 Reviews of recent books and expository articles.

## NEWS AND LETTERS

- 79 Letters to the Editor.

THE MATHEMATICAL ASSOCIATION OF AMERICA  
1529 Eighteenth Street, NW  
Washington, D.C. 20036

