

Computing the Digits in π

Carl D. Offner

September 6, 2009

Contents

1	Why do we care?	2
2	Archimedes	5
3	A digression: means	9
4	Measures of convergence	13
4.1	The rate of convergence of Archimedes' algorithm	13
4.2	Quadratic convergence	14
4.3	The arithmetic-geometric mean	14
5	Leibniz's series	16
6	A digression: acceleration of series	19
6.1	Euler's transformation	20
6.2	Proof that Euler's transformation converges quickly to the original sum	22
7	Quadratically converging approximations to π	27
7.1	Some representative elliptic integrals	27
7.1.1	The elliptic integral of the first kind	27
7.1.2	The elliptic integral of the second kind	30
7.1.3	More representations for $K(k)$ and $E(k)$	31
7.2	Elliptic integrals and the arithmetic-geometric mean	32
7.2.1	The relation of $M(a, b)$ to $K(k)$	32
7.2.2	The arc length of the lemniscate	33
7.3	Landen's transformations	34
7.3.1	Cayley's derivation of Landen's transformations	34
7.3.2	Iteration of Landen's transformation for J	39
7.4	Legendre's relation	40
7.4.1	The functions E' and K'	40
7.4.1.1	The precise behavior of K' near 0	41
7.4.2	Legendre's relation	44
7.5	The algorithm of Brent and Salamin	46

8	Streaming algorithms	48
8.1	Multiplication	48
8.1.1	Proof that we always make progress	51
8.2	A streaming algorithm that computes the digits of π	53
8.2.1	Linear fractional transformations	54
8.2.2	The algorithm	56
8.2.3	Proof that we always make progress	57
9	Algorithmic efficiency	59
9.1	Can the cost of each iteration be bounded?	60
9.2	Can the cost of each iteration be bounded <i>on average</i> ?	60

1 Why do we care?

First, the bottom line: Here is a very abbreviated history, showing the number of correct digits of π computed at various times:

Archimedes	ca. 250 BCE	3	
Liu Hui	263	5	
Van Ceulen	1615	35	
Sharp	1699	71	
Machin	1706	100	
Shanks	1874	707	(527 correct)
Bailey	Jan. 1986	29,360,111	
Kanada	Oct. 2005	more than 1 trillion	

Now it's not hard to understand what it means to compute more and more of the digits of π , faster and faster, but it does come with a question: Who cares? And the answer is perhaps not so obvious.

So here are some possible answers to that question:

The civil engineering answer. I asked a civil engineer how many digits of π he would actually ever need. After thinking about it for a while, he agreed with me that 5 was probably pushing it.

The astrophysics answer. Even physicists don't seem to need all that many digits of π :

... the author is not aware of even a single case of a "practical" scientific computation that requires the value of π to more than about 100 decimal places.

Bailey (1988)

And that estimate may already be absurdly high:

It requires a mere 39 digits of π in order to compute the circumference of a circle of radius 2×10^{25} meters (an upper bound on the distance traveled by a particle moving at the speed of light for 20 billion years, and as such an upper bound on the radius of the universe) with an error of less than 10^{-12} meters (a lower bound for the radius of a hydrogen atom).

Borwein and Borwein (1984)

That's a curious fact, to be sure. But what do astronomers and astrophysicists really need? I spoke to two recently. One had used whatever the π key on her calculator produced—at most, 10 digits. The other thought that 7 or 8 places was the most he'd ever seen the need for. And even in quantum field theory, where there are physical constants that are known to an astounding accuracy of 12 or 13 decimal places, it wouldn't really help to know π to any greater accuracy than that.

And that's probably as precise as any scientist would need.

The computer engineering answer.

In recent years, the computation of the expansion of π has assumed the role as a standard test of computer integrity. If even one error occurs in the computation, then the result will almost certainly be completely in error after an initial correct section. On the other hand, if the result of the computation of π to even 100,000 decimal places is correct, then the computer has performed billions of operations without error. For this reason, programs that compute the decimal expansion of π are frequently used by both manufacturers and purchasers of new computer equipment to certify system reliability.

Bailey (1988)

There is a point to this reasoning—an error of any sort will invalidate all later computations. I think, though, that this is really an argument of convenience rather than of intrinsic strength. That is, as long as we have a good algorithm for computing π , we may as well use it, since it has this nice property. But there's nothing really specific to π in this argument. There are other iterative computations that could exercise the computer's arithmetic capabilities and would serve as well. And since what we're really concerned with here is how to compute π efficiently, any increase in efficiency would actually make the computation less useful for testing a computer, since fewer computations would need to be performed. In addition, it's hard for me to believe that someone would really need to familiarize him or herself with the mathematics of elliptic integrals, as we will do here, solely for the purpose of certifying the correctness of a piece of computer hardware.

The mathematical answer. π is of course important all over the field of mathematics. It comes up in so many contexts that we're naturally curious about its peculiar properties. And one question that has intrigued people for a long time can be loosely worded as follows: "Are the digits in π random?" Well, of course they're not. They're completely determined. There's no chance or uncertainty involved. But really what people mean when they ask this question is something like this: are the digits *statistically* random? That is, as you go farther and farther out in the decimal expansion, does the fraction of 3's approach 1/10? Does the fraction of the pairs 31 approach 1/100? Does the fraction of the triplets 313 approach 1/1000? And so on.

Numbers that have this property are called *normal* numbers. (Perhaps more properly, "normal numbers to base 10", since the question can really be asked with respect to any base.) In 1909 Émile Borel proved that almost all numbers are normal. The term "almost all" has a precise mathematical meaning, but you can think of it in the following colloquial way: if you threw a dart at the number line, your chance of hitting a non-normal number is 0—that is, you might be very lucky and hit a non-normal number, but the fraction of non-normal numbers you'd hit in a long series of such throws would tend to 0.¹ And in fact, almost all numbers are normal to all bases simultaneously.

¹For another model of this, if someone was pulling numbers "randomly" out of a hat, it would never be profitable

That’s really a nice result. But the problem is that practically all the numbers we deal with every day are *not* normal. No rational number, for instance, can be normal, since it has an ultimately repeating decimal expansion. Now π , on the other hand, is not rational, and so perhaps it is normal. No one actually knows whether or not this is true—in fact, very few numbers have ever been proved to be normal, and none of them occur “in the wild”. That is, no number that occurs “naturally”, like e , or $\sqrt{2}$, or π is known to be normal.

Nevertheless, it is widely believed that e , $\sqrt{2}$, and π are all normal, and since so many digits of π have by now been computed, it’s possible to perform statistical analyses on them. And as far as we can tell, the expansion (as far as we know it) *looks* like it’s normal².

Of course that’s not a proof. It might be that starting 100 digits after the last digit that has so far been computed, there are no more 4’s. This would of course be astonishing, and no one seriously even entertains that possibility. But so far it can’t be ruled out.

The computer science answer. The fact that we can efficiently compute many digits of π is due only partly to the wonderful technologies embodied in computer hardware. It is really due much more to some wonderful algorithms that have been discovered that converge unbelievably quickly to π . That is, we are dealing here with algorithms, the core of computer science. And these algorithms are specific to π , of course—if we wanted to compute e , we’d use other algorithms. And algorithms tell us something about the numbers we are computing: a number that can be computed easily and efficiently just “feels” different from one that can only be computed with great difficulty. There is a whole theory of algorithmic complexity based on this idea, and we do feel that we learn something intrinsic about a number when we find out how difficult it is to compute.

So in fact, the computer science answer and the mathematical answer are really closely related. A good algorithm is going to depend on intrinsic mathematical properties of π , and conversely, a good algorithm may well tell us something significant mathematically about π .

The popular cultural answer. Finally—and I only bring this up with some trepidation—it has to be acknowledged that π is just one of those bits of learning that pretty much every one remembers from high school. People who have never heard of e know about π . And it’s part of the popular culture. Everyone knows that π “goes on forever”, and this is widely regarded as at least somewhat mysterious, as indeed it is. You’ve probably heard the jokes about “street math”—mathematical knowledge brought down to a caricature:

“You divide by zero, you die.”

“I knew this kid. He found the exact value of pi. He went nuts.”

All right, so those are the reasons. Now let’s look at what’s been accomplished. We’ll take things pretty much in chronological order. It turns out that the history of this is fascinating in itself. And while I can’t possibly cover this field exhaustively—books have been written about these things, and the topic itself impinges on some rather deep mathematics—I have surveyed the main approaches to the topic³. I have also written out many derivations in much more detail than would normally be done, mainly because I hope that this can be read by advanced undergraduates.

for you to bet that the next number would not be normal, no matter what the odds. OK; that’s enough about gambling.

²Wagon (1985) is a nice short article that goes into these matters at somewhat greater length.

³And I’ve done my best to clarify some tricky points. For instance, section 7.4.1.1 is somewhat similar but more direct and more motivated than some equivalent reasoning in Whittaker and Watson (1927). (For one thing, the actual integral I deal with is simpler and more directly arrived at; Whittaker and Watson perform a seemingly out-of-the-blue transformation, which I believe comes from the deformation of the contour in a contour integral, although

2 Archimedes

Many years ago I read that Archimedes found that π was between $3\frac{10}{71}$ and $3\frac{1}{7}$ by reasoning that the circumference of a circle was greater than that of any inscribed polygon and less than that of any circumscribed polygon—he considered regular polygons of 96 sides and was able to approximate π in this fashion.

This never impressed me much. In the first place, it's not all that great an approximation. In the second place, it just seems like drudgery. And why 96 sides, anyway?

Actually, what Archimedes did was really great. He didn't just figure out what the perimeter of two 96-sided polygons were. What he actually did was come up with an iterative procedure which could in principle be carried out mechanically as far as needed. That is, he produced an iterative algorithm for computing the digits of π . He actually carried out the first 5 steps of this algorithm, corresponding to a polygon of 96 sides, but he could have gone as far as he wanted. And he did all this around 250 BCE.

It's actually not easy to see that this is what he did—the Greeks didn't have algebraic notation, so his reasoning is expressed geometrically and in terms of rational approximations of every quantity. And he certainly didn't have any way to express in formal terms our notion of a recursive or iterative algorithm. This makes his paper extremely hard to read. And to be honest, I haven't gone through it myself, although I did look at it. But enough people have explained what he did that I'm confident that what I'm writing here is correct.

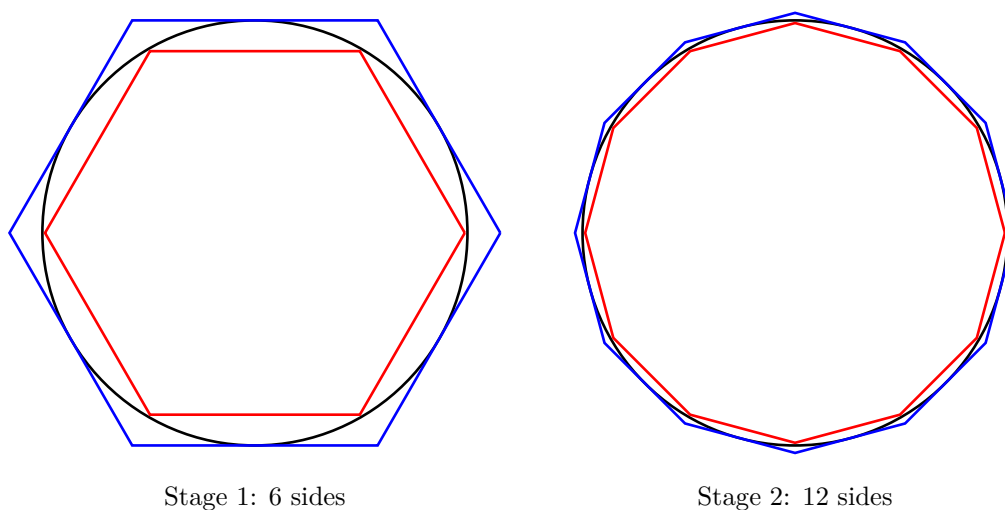


Figure 1: Stages 1 and 2 of the construction used by Archimedes to approximate π .

To set things up conveniently, let us say the circle has radius 1 (so the circumference, which we will

they don't say so, and for another thing, I'm able to motivate naturally the way I arrive at the point of division of the range of the integral). It's also much clearer than an equivalent computation in Cayley (1895) (which is undoubtedly correct but would not be easy to write rigorously today). I did get at least one idea from each of those ways of doing the computation, however.

approximate, will be 2π). It also turns out to be convenient to let the angle at the center of one “wedge” of the polygon at stage n be $2\theta_n$. The edges of the circumscribed and inscribed polygons at stage n will be denoted by α_n and β_n , respectively. Thus in Figure 2 we have

$$\frac{\alpha_n}{2} = \tan \theta_n$$

$$\frac{\beta_n}{2} = \sin \theta_n$$

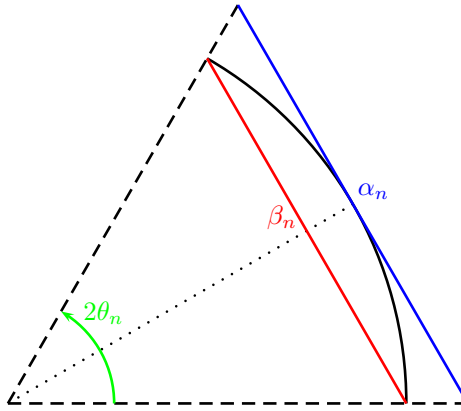


Figure 2: One edge of the **inscribed polygon**, and one edge of the **circumscribed polygon**. If the polygon at stage n has k edges, then $2\theta_n = \frac{2\pi}{k}$. In fact, k only comes in at the very end of the computation, and basically cancels out of the final result.

Let us denote the entire perimeter of the circumscribed polygon at stage n by a_n —this is just α_n multiplied by the number of edges at stage n . Similarly, we will denote the entire perimeter of the inscribed polygon at stage n by b_n , and this is just β_n multiplied by the number of edges at stage n .

We pass from a_n and b_n to a_{n+1} and b_{n+1} by doubling the number of edges in the polygon; that is, by dividing θ_n by 2. Thus to see how a_{n+1} and b_{n+1} are related to a_n and b_n , we can use some of the famous “half-angle formulas” of trigonometry. Identities like these were very familiar to students and practitioners of mathematics a century ago, but have a quaint feel about them today. And since I’m no better at these trigonometric identities than anyone else is, I’ll derive the ones we need. We start with the few that I (and I suppose most people) know:

- (1) $\sin^2 \theta + \cos^2 \theta = 1$
- (2) $\sin 2\theta = 2 \sin \theta \cos \theta$
- (3) $\cos 2\theta = \cos^2 \theta - \sin^2 \theta$

From (3) and (1) we get the identities

$$\cos 2\theta = \cos^2 \theta - \sin^2 \theta = \begin{cases} 1 - 2\sin^2 \theta \\ 2\cos^2 \theta - 1 \end{cases}$$

From those identities we have

$$\tan^2 \theta = \frac{\sin^2 \theta}{\cos^2 \theta} = \frac{1 - \cos 2\theta}{1 + \cos 2\theta} = \frac{1 - \cos^2 2\theta}{(1 + \cos 2\theta)^2} = \frac{\sin^2 2\theta}{(1 + \cos 2\theta)^2}$$

and so

$$(4) \quad \tan \frac{\theta}{2} = \frac{\sin \theta}{1 + \cos \theta} = \frac{\sin \theta \tan \theta}{\sin \theta + \tan \theta}$$

which is one of the two formulas we need.

The other one we'll simply notice to be true by expanding out $\sin \theta = 2 \sin \frac{\theta}{2} \cos \frac{\theta}{2}$ on the right-hand side:

$$(5) \quad \sin^2 \frac{\theta}{2} = \frac{\sin \theta \tan \frac{\theta}{2}}{2}$$

We're now ready to solve for α_{n+1} and β_{n+1} in terms of α_n and β_n : We see that

$$\alpha_{n+1} = 2 \tan \frac{\theta_n}{2} = 2 \frac{\sin \theta_n \tan \theta_n}{\sin \theta_n + \tan \theta_n} = \frac{\alpha_n \beta_n}{\alpha_n + \beta_n}$$

and

$$\beta_{n+1}^2 = 4 \sin^2 \frac{\theta_n}{2} = 2 \sin \theta_n \tan \frac{\theta_n}{2} = \frac{\beta_n \alpha_{n+1}}{2}$$

That is, we have the recursive relations

$$\begin{aligned} \alpha_{n+1} &= \frac{\alpha_n \beta_n}{\alpha_n + \beta_n} \\ \beta_{n+1} &= \sqrt{\frac{\beta_n \alpha_{n+1}}{2}} \end{aligned}$$

We're almost done. What we really need is a similar recursion for the perimeters a_n and b_n of the circumscribed and inscribed polygons. Suppose to make the computation simple we start with a regular hexagon. Then $\beta_0 = 1$ and $\alpha_0 = 2/\sqrt{3}$. So $b_0 = 6$ and $a_0 = 4\sqrt{3}$. Further, $b_n = 6 \cdot 2^n \beta_n$ and $a_n = 6 \cdot 2^n \alpha_n$, so it follows easily that

$$\begin{aligned} a_{n+1} &= \frac{2a_n b_n}{a_n + b_n} \\ b_{n+1} &= \sqrt{b_n a_{n+1}} \end{aligned}$$

and so using this recursion, both a_n and b_n will converge to 2π .

Figure 3 shows how the convergence works out in practice.

n	$a_n/2$	$b_n/2$
1	3.464101615137754587054892	3.000000000000000000000000
2	3.215390309173472477670643	3.105828541230249148186786
3	3.159659942097500483316634	3.132628613281238197161749
4	3.146086215131434971098098	3.139350203046867207135146
5	3.142714599645368298168859	3.141031950890509638111352
6	3.141873049979823871745487	3.141452472285462075450609
7	3.141662747056848526224490	3.141557607911857645516463
8	3.141610176604689538763470	3.141583892148318408668969
9	3.141597034321526151993218	3.141590463228050095738458
10	3.141593748771352027975981	3.141592105999271550544776
11	3.141592927385097033548008	3.141592516692157447592874
12	3.141592722038613818342804	3.141592619365383955189549
13	3.141592670701998047877018	3.141592645033690896672141
14	3.141592657867844419844008	3.141592651450767651704253
15	3.141592654659306032497220	3.141592653055036841691123
16	3.141592653857171436889364	3.141592653456104139264643
17	3.141592653656637788064203	3.141592653556370963662823
18	3.141592653606504375862713	3.141592653581437669762668
19	3.141592653593971022812640	3.141592653587704346287648
20	3.141592653590837684550141	3.141592653589271015418894
21	3.141592653590054349984517	3.141592653589662682701706
22	3.141592653589858516343111	3.141592653589760599522409
23	3.141592653589809557932760	3.141592653589785078727584
24	3.141592653589797318330172	3.141592653589791198528878
25	3.141592653589794258429525	3.141592653589792728479202
26	3.141592653589793493454363	3.141592653589793110966783
27	3.141592653589793302210573	3.141592653589793206588678
28	3.141592653589793254399625	3.141592653589793230494152
29	3.141592653589793242446889	3.141592653589793236470520
30	3.141592653589793239458704	3.141592653589793237964612
31	3.141592653589793238711658	3.141592653589793238338135
32	3.141592653589793238524897	3.141592653589793238431516
33	3.141592653589793238478206	3.141592653589793238454861
34	3.141592653589793238466534	3.141592653589793238460697
35	3.141592653589793238463616	3.141592653589793238462157
36	3.141592653589793238462886	3.141592653589793238462521
37	3.141592653589793238462704	3.141592653589793238462612
38	3.141592653589793238462658	3.141592653589793238462635
39	3.141592653589793238462647	3.141592653589793238462641
40	3.141592653589793238462644	3.141592653589793238462642

Figure 3: The first 40 iterations of the computations in Archimedes' algorithm.

Archimedes stopped with stage 5 (corresponding to regular polygons of 96 sides). His approximations were $3\frac{10}{71} \leq \pi \leq 3\frac{1}{7}$, or equivalently, $3.1408\dots \leq \pi \leq 3.1428\dots$. This is actually only a little less accurate than the computation above—even though Archimedes was approximating every square root (and in fact every arithmetic operation) with fractions, we can see that his rational approximations were very close to what you would get with unlimited precision arithmetic.

But the point, just to make it again, is that the importance of what Archimedes did does not lie in the accuracy of his result, which wasn't that great. The real significance is that he produced for the first time an iterative algorithm capable of producing approximations to π of any desired accuracy.

Archimedes was not the only person to discover such a recursive algorithm, but he seems to have been the first. Around 500 years later Liu Hui in China wrote down a very similar procedure, which played a role in the Chinese mathematical tradition very similar to that which Archimedes' algorithm played in the West.

This algorithm, with only minor changes, formed the basis for all computations of π up through the early 1600's, when 35 digits of π were computed by Ludolph van Ceulen. This feat was considered so remarkable that the value that he got was referred to in parts of Europe for years afterward as the “Ludolphian number”.

For this particular algorithm, it appears that the convergence is “linear” in the sense that the number of correct decimal places is roughly proportional to the number of iterations of the process.

We'll see that this is true. But first we need to take a small side trip.

3 A digression: means

If we have two positive numbers $b \leq a$, there are several different notions of a *mean* we can define. In fact, there is a whole scale of means, parametrized by p :

$$M_p(a, b) = \left(\frac{a^p + b^p}{2} \right)^{1/p}$$

(Actually, this definition can be generalized greatly to talk about means of a finite or even infinite collection of numbers, or of a function, and can accommodate weighted averages as well. But we won't need any of that here, so we'll just stick to the simplest case. Everything we say below, however, generalizes.)

This definition may look strange at first, but it actually encapsulates some well-known “averages”:

- When $p = 1$, we have $M_1(a, b) = \frac{a+b}{2}$, which is the *arithmetic mean* of a and b . It's what people usually mean when they talk about the “average” of two numbers—“you add 'em up and divide by 2”.
- When $p = 2$, we have $M_2(a, b) = \sqrt{\frac{a^2+b^2}{2}}$. This is the *root-mean-square*, which is a mean that is used frequently in measuring physical and statistical errors. If we have a series of measurements $\{a_j : 1 \leq j \leq n\}$ of some quantity whose real value is expected to be m , say, then the “error” of the j^{th} measurement a_j is $|a_j - m|$, and so we could say that the *average*

*error*⁴ is

$$\frac{\sum_{j=1}^n |a_j - m|}{n}$$

However, squares are a lot easier to deal with algebraically than absolute values, and so it turns out to be more convenient to describe the average error by the root-mean-square deviation of $\{a_j\}$ from the expected value:

$$\sqrt{\frac{\sum_{j=1}^n (a_j - m)^2}{n}}$$

- When $p = -1$, we have⁵

$$M_{-1}(a, b) = \frac{1}{\frac{1}{2}\left(\frac{1}{a} + \frac{1}{b}\right)} = \frac{2ab}{a + b}$$

This is the *harmonic mean* of a and b . It comes up in the following kind of problem: Suppose you travel for 100 miles at 50 mph and then another 100 miles at 25 mph. What was your average speed? Well, you might think it was the arithmetic mean of 25 and 50, which is 37.5 mph, but it's not. The total distance you traveled was 200 miles. The total time it took you was $(100/50) + (100/25) = 6$ hours. So your average velocity was $200/6 = 33 \frac{1}{3}$ mph. This in fact is the harmonic mean of 25 and 50, because

$$\begin{aligned} \text{total time} &= \frac{100}{50} + \frac{100}{25} \\ \text{total distance} &= 2 \cdot 100 \end{aligned}$$

and so

$$\begin{aligned} \text{average velocity} &= \frac{2 \cdot 100}{\frac{100}{50} + \frac{100}{25}} \\ &= \frac{1}{\frac{1}{2}\left(\frac{1}{50} + \frac{1}{25}\right)} \\ &= M_{-1}(50, 25) \end{aligned}$$

(By contrast, if you traveled 1 *hour* at 50 mph and then another *hour* at 25 mph, your average velocity would have been the arithmetic mean of 25 and 50.)

Note that the harmonic mean of 25 and 50 is less than the arithmetic mean of 25 and 50. This is true in general, as we'll see shortly.

What is also interesting from our point of view here is that Archimedes' algorithm computes a_{n+1} as the harmonic mean of a_n and b_n .

⁴We use absolute values because we don't care whether the error is positive or negative, and we don't want positive errors to cancel out negative ones.

⁵I know, I didn't specify the values that p could take, but it *can* be negative.

- We know that $M_p(a, b)$ makes sense if $p > 0$ and also if $p < 0$. What happens if we try to set $p = 0$? Of course $M_0(a, b)$ makes no sense as defined above, but we can ask what happens to $M_p(a, b)$ as $p \rightarrow 0$. It turns out (and it is not too hard to prove, but we won't prove it here) that $\lim_{p \rightarrow 0} M_p(a, b) = \sqrt{ab}$. So we define

$$M_0(a, b) = \sqrt{ab}$$

This is the *geometric mean* of a and b . I assume it's called the geometric mean because of the ways it comes up in geometry. For instance, if a square of side s has the same area as a rectangle of sides a and b , then $s^2 = ab$, so s is the geometric mean of a and b . There is also the neat theorem we all learned in high school geometry—and this certainly also goes back to Euclid—that in figure 4 we have $h = \sqrt{ab}$.

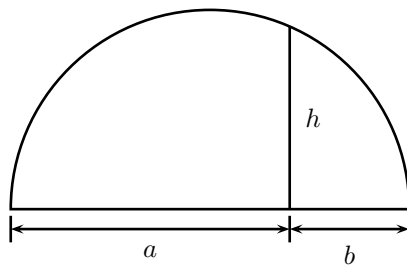


Figure 4: How the geometric mean arises in geometry. $h = \sqrt{ab}$.

- What happens as $p \rightarrow \infty$? Actually, it's pretty easy to compute this: keeping in mind that $0 < b \leq a$, we have

$$M_p(a, b) = \left(\frac{a^p + b^p}{2} \right)^{1/p} = a \left(\frac{1 + \left(\frac{b}{a}\right)^p}{2} \right)^{1/p}$$

Now $\left(\frac{b}{a}\right)^p$ is bounded and so as $p \rightarrow \infty$,

$$\left(\frac{1 + \left(\frac{b}{a}\right)^p}{2} \right)^{1/p} \rightarrow 1$$

Thus it makes sense for us to define

$$M_\infty(a, b) = \max\{a, b\}$$

- Further, we see that for any $p > 0$,

$$(6) \quad M_{-p}(a, b) = \frac{1}{\left(\frac{\left(\frac{1}{a}\right)^p + \left(\frac{1}{b}\right)^p}{2} \right)^{1/p}} = \frac{1}{M_p\left(\frac{1}{a}, \frac{1}{b}\right)}$$

and therefore $\lim_{p \rightarrow -\infty} M_p(a, b) = b$ so we can with a clear conscience define

$$M_{-\infty}(a, b) = \min \{a, b\}$$

Note that all these means M_p satisfy what must be reasonable properties of any mean:

- $M_p(a, b)$ lies between a and b . If $a \neq b$ and $-\infty < p < \infty$, then $M_p(a, b)$ lies strictly between a and b .
- $M_p(a, a) = a$.
- If $a_1 \leq a_2$ then $M_p(a_1, b) \leq M_p(a_2, b)$ and if $b_1 \leq b_2$ then $M_p(a, b_1) \leq M_p(a, b_2)$.
- $M_p(a, b) = M_p(b, a)$.

as well as being homogeneous of degree 1: for any $t > 0$,

$$M_p(ta, tb) = tM_p(a, b)$$

and in particular,

$$M_p(a, b) = aM_p(1, \frac{b}{a})$$

And there is one final inequality which (to my mind, at least) is just wonderful: For $p \leq q$, $M_p(a, b) \leq M_q(a, b)$. We won't prove that here even for these means of two variables, but look what it says for the particular means we have been considering:

$$(7) \quad M_{-\infty}(a, b) \leq M_{-1}(a, b) \leq M_0(a, b) \leq M_1(a, b) \leq M_2(a, b) \leq M_{\infty}(a, b)$$

Some of this we know already. For instance, we know that $M_1(a, b) \leq M_{\infty}(a, b)$. The fact that $M_0(a, b) \leq M_1(a, b)$ is the famous “Theorem of the Arithmetic and Geometric Means”, and in this special case it's quite easy to prove: we just notice that for any real numbers x and y , $(x - y)^2 \geq 0$, so $x^2 - 2xy + y^2 \geq 0$. That is

$$xy \leq \frac{x^2 + y^2}{2}$$

and then setting $x = \sqrt{a}$ and $y = \sqrt{b}$, we are done. So now we know that $M_0(a, b) \leq M_1(a, b) \leq M_{\infty}(a, b)$. The remainder of the inequalities in (7) can be proved similarly or in some cases even more simply by using (6) to “reflect the inequalities in the origin”.

4 Measures of convergence

4.1 The rate of convergence of Archimedes' algorithm

So why were we so interested in these means? Well as we've already noted, in the recurrence relations in Archimedes' algorithm,

$$\begin{aligned} a_{n+1} &= \frac{2a_nb_n}{a_n + b_n} \\ b_{n+1} &= \sqrt{b_na_{n+1}} \end{aligned}$$

we see that $a_{n+1} = M_{-1}(a_n, b_n)$. And b_{n+1} is “almost” the geometric mean of a_n and b_n —in fact, it's the geometric mean of a_{n+1} and b_n . Further, the arithmetic and geometric means will be absolutely crucial later.

In any case, we're now ready to investigate the convergence of Archimedes' algorithm. We will show by induction that

- $b_{j-1} < b_j$ and $a_j < a_{j-1}$ for all $j \geq 1$, and
- $b_j < a_j$ for all $j \geq 0$.

For we know these statements are true when $j = 0$. (Of course the first statement is vacuous in that case.) If then these statements are true for all $j \leq n$ then we have

$$(8) \quad a_{n+1} = \frac{2a_nb_n}{a_n + b_n} \begin{cases} < \frac{2a_nb_n}{2b_n} = a_n \\ > \frac{2a_nb_n}{2a_n} = b_n \end{cases}$$

From the second inequality in (8), it follows that

$$b_{n+1} = \sqrt{b_na_{n+1}} > b_n$$

and also that

$$b_{n+1} = \sqrt{b_na_{n+1}} < a_{n+1}$$

and this completes the inductive proof.

Thus as $n \rightarrow \infty$, b_n increases to a limit b , and a_n decreases to a limit a , with $b \leq a$. Now then we have

$$a_{n+1} - b_{n+1} \leq a_{n+1} - b_n = \frac{2a_nb_n}{a_n + b_n} - b_n = \frac{b_n}{a_n + b_n}(a_n - b_n) = m_n(a_n - b_n)$$

where the multiplier $m_n < 1/2$, which shows two things:

1. $a_n - b_n \rightarrow 0$ as $n \rightarrow \infty$.
2. Each successive iteration of the algorithm increases our accuracy by a factor of m , and in fact, since $m_n \rightarrow 1/2$, ultimately each iteration increases our accuracy by a factor of 2.

Thus each iteration should give us roughly an additional binary digit of accuracy, which is the kind of behavior we saw in practice above. This kind of convergence is called “linear convergence”.

and then repeating this process. Suppose (without loss of generality) that $0 < b < a$. Then by the theorem of the arithmetic and geometric means we have

$$b < b_1 < a_1 < a$$

and so it is clear that $b_n \uparrow$ and $a_n \downarrow$. In fact, they approach the same limit—which we will denote by $M(a, b)$ —as we can see as follows:

$$a_{n+1} - b_{n+1} \leq a_{n+1} - b_n = \frac{a_n + b_n}{2} - b_n = \frac{a_n - b_n}{2}$$

and so by induction we have

$$0 \leq a_n - b_n \leq 2^{-n}(a - b)$$

So we see that the limit $M(a, b)$ well-defined.

The convergence we have just demonstrated is linear convergence. But actually, we gave away too much in the estimation we made. Now that we know that the process converges, we can do better:

$$\begin{aligned} a_{n+1}^2 - b_{n+1}^2 &= \frac{a_n^2 + 2a_nb_n + b_n^2}{4} - a_nb_n \\ (9) \qquad \qquad \qquad &= \frac{a_n^2 - 2a_nb_n + b_n^2}{4} \\ &= \frac{(a_n - b_n)^2}{4} \end{aligned}$$

so since $a_n + b_n \geq 2b_0$ and so is bounded away from 0 for all n , we get

$$a_{n+1} - b_{n+1} = \frac{(a_n - b_n)^2}{4(a_{n+1} + b_{n+1})} = O((a_n - b_n)^2)$$

so we see that the iteration for the arithmetic-geometric mean converges quadratically, just like the square root algorithm.

In particular, suppose that for some N , $|a_N - b_N| \leq 1/2$. Then it follows by induction that

$$(10) \qquad \qquad \qquad a_n - b_n = O\left(2^{-2^{n-N}}\right) \quad \text{for } n \geq N$$

When we come to consider this in more detail, it will be useful to set

$$c_n = \sqrt{a_n^2 - b_n^2}$$

Then from (9), we get

$$c_{n+1}^2 = \frac{(a_n - b_n)^2}{4} = \frac{(a_n - b_n)^2}{4} \cdot \frac{(a_n + b_n)^2}{(a_n + b_n)^2} = \frac{(a_n^2 - b_n^2)^2}{4(a_n + b_n)^2} = \frac{c_n^4}{4(a_n + b_n)^2}$$

and so

$$(11) \qquad \qquad \qquad c_{n+1} = O(c_n^2)$$

Exemplum 3. $a = 1, b = 0,8.$

$^v a = 25,19190\ 722$	$^v b = 0,00000\ 00000\ 0$
$''' a = 12,59595\ 36116\ 78$	$''' b = 0,00000\ 00133\ 367$
$'' a = 6,29797\ 68125\ 07655\ 42373\ 4$	$'' b = 0,00040\ 98644\ 58278\ 08440\ 9$
$' a = 3,14919\ 33384\ 82966\ 75407\ 2$	$' b = 0,05080\ 66615\ 17033\ 24592\ 8$
$a = 1,60000\ 00000\ 00000\ 00000\ 0$	$b = 0,40000\ 00000\ 00000\ 00000\ 0$
$a = 1,00000\ 00000\ 00000\ 00000\ 0$	$b = 0,80000\ 00000\ 00000\ 00000\ 0$
$a' = 0,90000\ 00000\ 00000\ 00000\ 0$	$b' = 0,89442\ 71909\ 99915\ 87856\ 4$
$a'' = 0,89721\ 35954\ 99957\ 93928\ 2$	$b'' = 0,89720\ 92687\ 32734$
$a''' = 0,89721\ 14321\ 16346$	$b''' = 0,89721\ 14321\ 13738$
$a'''' = 0,89721\ 14321\ 15042$	$b'''' = 0,89721\ 14321\ 15042$

Exemplum 4. $a = \sqrt{2}, b = 1.$

$''' a = 19,17024\ 37557\ 69475\ 31905\ 0$	$''' b = 0,00000\ 00009\ 32560\ 02627\ 6$
$'' a = 9,58512\ 18783\ 51017\ 67266\ 3$	$'' b = 0,00013\ 37064\ 06056\ 69181\ 0$
$' a = 4,79262\ 77923\ 78537\ 18223\ 7$	$' b = 0,03579\ 93323\ 67652\ 95745\ 7$
$a = 2,41421\ 35623\ 73095\ 04880\ 2$	$b = 0,41421\ 35623\ 73095\ 04880\ 2$
$a = 1,41421\ 35623\ 73095\ 04880\ 2$	$b = 1,00000\ 00000\ 00000\ 00000\ 0$
$a' = 1,20710\ 67811\ 86547\ 52440\ 1$	$b' = 1,18920\ 71150\ 02721\ 06671\ 7$
$a'' = 1,19815\ 69480\ 94634\ 29555\ 9$	$b'' = 1,19812\ 35214\ 93120\ 12260\ 7$
$a''' = 1,19814\ 02347\ 93877\ 20908\ 3$	$b''' = 1,19814\ 02346\ 77307\ 20579\ 8$
$a'''' = 1,19814\ 02347\ 35592\ 20744\ 1$	$b'''' = 1,19814\ 02347\ 35592\ 20743\ 9$

4.

Habeantur praeter progressionones

$\dots''' a, '' a, ' a, a, a', a'', a'''\dots$
 $\dots''' b, '' b, ' b, b, b', b'', b'''\dots$

duae aliae

$\dots''' c, '' c, ' c, c, c', c'', c'''\dots$
 $\dots''' d, '' d, ' d, d, d', d'', d'''\dots$

simili modo formatae; supponamusque, duos terminos in posterioribus duobus in prioribus *proportionales* esse, e. g. $a:b = c:d$, sive $a:c = b:d = 1:n$. Tunc omnes termini in I ad terminos in III, omnesque in II ad terminos in IV (similiter relative ad a, b, c, d siti ac sitos) in eadem ratione erunt, puta

Figure 7: A page from a manuscript of Gauss (1800). Gauss's computation of $M(\sqrt{2}, 1)$ is the bottom half of the second table on the page.

and integrating term-by-term:

$$\begin{aligned}\arctan z &= \int_0^z \frac{1}{1+t^2} dt \\ &= \int_0^z (1 - t^2 + t^4 - t^6 + t^8 - \dots) dt \\ &= z - \frac{z^3}{3} + \frac{z^5}{5} - \frac{z^7}{7} + \frac{z^9}{9} - \dots\end{aligned}$$

which not only has radius of convergence 1, but actually does converge to $\arctan z$ at 1. This series is known as Gregory's series. Substituting $z = 1$, we get Leibniz's series. (In fact, the series obviously converges at 1, and so by Abel's theorem, it represents $\arctan z$ there.)

Leibniz's series is just beautiful, and it can even be derived from some equally beautiful theorems in number theory—see Hilbert and Cohn-Vossen (1999). For computational purposes, however, it obviously converges abysmally slowly—it would take several hundred terms just to get 2-digit accuracy. But Gregory's series could be used directly:

By the beginning of the eighteenth century Abraham Sharp under the direction of the English astronomer and mathematician E. Halley had obtained 71 correct digits of π using Gregory's series with $x = \sqrt{1/3}$, namely

$$\frac{\pi}{6} = \frac{1}{\sqrt{3}} \left(1 - \frac{1}{3 \cdot 3} + \frac{1}{3^2 \cdot 5} - \frac{1}{3^3 \cdot 7} + \dots \right)$$

Borwein and Borwein (1987), page 339

Even this series, however, did not really converge all that quickly.

A series which converged more quickly was derived by John Machin, who used it to calculate 100 digits of π in 1706. He used some trigonometric identities which we need to recall: From

$$\begin{aligned}\sin(\theta \pm \phi) &= \sin \theta \cos \phi \pm \cos \theta \sin \phi \\ \cos(\theta \pm \phi) &= \cos \theta \cos \phi \mp \sin \theta \sin \phi\end{aligned}$$

we get

$$\begin{aligned}\tan(\theta \pm \phi) &= \frac{\sin \theta \cos \phi \pm \cos \theta \sin \phi}{\cos \theta \cos \phi \mp \sin \theta \sin \phi} \\ &= \frac{\tan \theta \pm \tan \phi}{1 \mp \tan \theta \tan \phi}\end{aligned}$$

and in particular,

$$\tan 2\theta = \frac{2 \tan \theta}{1 - \tan^2 \theta}$$

Machin started by defining an angle θ by

$$\theta = \arctan \frac{1}{5}$$

Then

$$\tan 2\theta = \frac{2 \tan \theta}{1 - \tan^2 \theta} = \frac{5}{12}$$

and thence

$$\tan 4\theta = \frac{2 \tan 2\theta}{1 - \tan^2 2\theta} = 1 + \frac{1}{119}$$

and from that, we get

$$\tan\left(4\theta - \frac{\pi}{4}\right) = \frac{\tan 4\theta - 1}{1 + \tan 4\theta} = \frac{1}{239}$$

and taking the arctan of both sides of this, we get Machin's formula

$$\frac{\pi}{4} = 4 \arctan \frac{1}{5} - \arctan \frac{1}{239}$$

Gregory's series allows this to be evaluated reasonably efficiently: the first term works well with decimal arithmetic, and the second converges rapidly.

Many other series of this form were developed and used. The calculations were all performed by hand (or eventually, by desk calculators, and finally, by early computers). But already by 1976, a better method was available, due independently to Brent and Salamin. We'll look at it next, after another digression.

6 A digression: acceleration of series

Let's go back to Leibniz's series:

$$(12) \quad \frac{\pi}{4} = 1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \dots$$

We have already noted that this series converges so slowly that it is quite useless for any practical purpose. Nevertheless, we might ask if there is some transformation one might apply to this series that would generate an equivalent series that converged faster.

The simplest thing to try, and a very obvious one, is to notice that since this is an alternating series with terms of diminishing absolute value, we can take the terms in pairs and write this series as

$$\begin{aligned} \frac{\pi}{4} &= \left(1 - \frac{1}{3}\right) + \left(\frac{1}{5} - \frac{1}{7}\right) + \left(\frac{1}{9} - \frac{1}{11}\right) \dots \\ &= \frac{2}{3} + \frac{2}{5 \cdot 7} + \frac{2}{9 \cdot 11} + \dots \end{aligned}$$

That is,

$$(13) \quad \frac{\pi}{4} = 2 \sum_{n=1}^{\infty} \frac{1}{(2n-1)(2n+1)}$$

The n^{th} term of this series is $O(\frac{1}{n^2})$, and so it is similar in nature to Euler's wonderful formula

$$\frac{\pi^2}{6} = \sum_{n=1}^{\infty} \frac{1}{n^2}$$

This transformed series (13) is more convenient for computation than Leibniz's series, because it avoids the computational difficulty of adding and subtracting numbers that mainly cancel out—it's hard to do this and maintain a high level of precision. Unfortunately, however, the new series doesn't actually converge any faster than Leibniz's original series does, since each term in (13) is just the sum of two consecutive terms in (12), and after n terms we are still only within a distance of about $1/n$ from the sum.

6.1 Euler's transformation

This didn't stop Euler, however. He found a method of transforming series like this which in many cases dramatically improves the rate of convergence. In effect he recursively applied a transformation similar to the transformation we performed above, and in doing so he came up with a general method that works for a large class of alternating series.

To see what he did it is helpful to have some standard notation: if $(x_0, x_1, x_2, x_3, \dots)$ is any sequence of real (or even complex) numbers, we define the *forward* difference sequence $(\Delta x_0, \Delta x_1, \Delta x_2, \dots)$ by $\Delta x_k = x_{k+1} - x_k$, and we define the *backward* difference sequence $(\nabla x_0, \nabla x_1, \nabla x_2, \dots)$ by $\nabla x_k = x_k - x_{k+1} = -\Delta x_k$. Euler's transformation can actually be written in terms of either Δ or ∇ , but we will find it somewhat more convenient to use ∇ . Note that we can iterate this process: we can take differences of differences, and so on, as in Figure 8. We see immediately that we have

$$\nabla^k a_m = \sum_{j=0}^k (-1)^j \binom{k}{j} x_k$$

With this notation, the Euler transform of an alternating series $\sum_{j=0}^{\infty} (-1)^j a_j$ is produced by the following series of steps:

$$\begin{aligned} S &= a_0 - a_1 + a_2 - a_3 + a_4 - \dots \\ &= \frac{1}{2}a_0 + \frac{1}{2}\left((a_0 - a_1) - (a_1 - a_2) + (a_2 - a_3) - (a_3 - a_4) + \dots\right) \\ &= \frac{1}{2}a_0 + \frac{1}{4}(a_0 - a_1) + \frac{1}{4}\left((a_0 - 2a_1 + a_2) - (a_1 - 2a_2 + a_3) + (a_2 - 2a_3 + a_4) - \dots\right) \end{aligned}$$

That is, Euler replaced the series $\sum_{j=0}^{\infty} a_j$ with the series

$$\sum_{k=0}^{\infty} \frac{\nabla^k a_0}{2^{k+1}}$$

or equivalently,

$$\sum_{k=0}^{\infty} \frac{1}{2^{k+1}} \sum_{j=0}^k (-1)^j \binom{k}{j} a_k$$

$$\begin{array}{cccccc}
x_0 & & & & & \\
& x_0 - x_1 & & & & \\
x_1 & & x_0 - 2x_1 + x_2 & & & \\
& x_1 - x_2 & & x_0 - 3x_1 + 3x_2 - x_3 & & \\
x_2 & & x_1 - 2x_2 + x_3 & & x_0 - 4x_1 + 6x_2 - 4x_3 + x_4 & \\
& x_2 - x_3 & & x_1 - 3x_2 + 3x_3 - x_4 & & \\
x_3 & & x_2 - 2x_3 + x_4 & & x_1 - 4x_2 + 6x_3 - 4x_4 + x_5 & \\
& x_3 - x_4 & & x_2 - 3x_3 + 3x_4 - x_5 & & \\
x_4 & & x_3 - 2x_4 + x_5 & & & \\
& x_4 - x_5 & & & & \\
x_5 & & & & &
\end{array}$$

$$\begin{array}{cccccc}
x_0 = \nabla^0 x_0 & & & & & \\
& \nabla x_0 & & & & \\
x_1 = \nabla^0 x_1 & & \nabla^2 x_0 & & & \\
& \nabla x_1 & & \nabla^3 x_0 & & \\
x_2 = \nabla^0 x_2 & & \nabla^2 x_1 & & \nabla^4 x_0 & \\
& \nabla x_2 & & \nabla^3 x_1 & & \\
x_3 = \nabla^0 x_3 & & \nabla^2 x_2 & & \nabla^4 x_1 & \\
& \nabla x_3 & & \nabla^3 x_2 & & \\
x_4 = \nabla^0 x_4 & & \nabla^2 x_3 & & & \\
& \nabla x_4 & & & & \\
x_5 = \nabla^0 x_5 & & & & &
\end{array}$$

Figure 8: Backward differences. By convention, $\nabla^0 x_k = x_k$ and of course $\nabla = \nabla^1$.

(Just to be clear, note that the series S itself is alternating, but each $a_k \geq 0$.)

The three big questions we need to answer then are:

1. Does the Euler transform converge?
2. If so, does it converge to the same sum as the original series?
3. And finally, assuming it does converge to the same sum, does it converge any faster than the original series?

The utility of Euler's transformation is that in a large class of cases, the answer to all three questions is yes. We'll show this next.

6.2 Proof that Euler's transformation converges quickly to the original sum

Euler didn't actually prove any general theorems about this transformation. He did use it in several specific cases, where he could show that it really did converge to the original sum, and converged much more quickly. We're going to give a general result here showing that the Euler transformation can be applied to a large class of alternating series. The reasoning in the proof of this result is quite elegant and goes ultimately back to Hausdorff. We'll follow Cohen et al. (2000), who have shown how this method can be used to arrive at Euler's transformation (and in fact many others) very quickly.

Suppose—and this is the key assumption—that there is a positive measure μ defined on the interval $[0, 1]$ such that the numbers a_k are the *moments* of μ , by which we simply mean this:

$$(14) \quad a_k = \int_0^1 t^k d\mu(t)$$

This won't always be the case—there are sequences $\{a_k\}$ for which no such μ exists—but in many cases there is such a μ . So assuming that there is such a μ , let us continue⁶: We can use (14) to write the sum of the series as follows:

$$S = \sum_{k=0}^{\infty} (-1)^k a_k = \int_0^1 \left(\sum_{k=0}^{\infty} (-1)^k t^k d\mu(t) \right) = \int_0^1 \frac{1}{1+t} d\mu(t)$$

Now suppose that we have a sequence $\{P_n\}$ of polynomials where P_n has degree n , so we can write

$$P_n(x) = \sum_{j=0}^n p_j^{(n)} x^j$$

and where $P_n(-1) \neq 0$. (We'll see that such sequences of polynomials are not hard to come by.) We know that $P(t) - P(-1)$ must be divisible by $t - (-1)$, i.e., by $t + 1$. Let us set

$$S_n = \frac{1}{P_n(-1)} \int_0^1 \frac{P_n(-1) - P_n(t)}{1+t} d\mu(t)$$

We can split this sum into two parts:

$$S_n = \int_0^1 \frac{1}{1+t} d\mu(t) - \frac{1}{P_n(-1)} \int_0^1 \frac{P_n(t)}{1+t} dt$$

As we have seen, the first term on the right-hand side is just the sum S of the series. Our aim will be to find a sequence of polynomials such that the second integral on the right is very small. For

⁶And by the way, if you are not comfortable with measures, you can assume the situation is somewhat simpler: you can assume that there is a positive "weight function" $w(t)$ such that the numbers a_k are the moments of w :

$$a_k = \int_0^1 t^k w(t) dt$$

and then in what follows simply replace every occurrence of $d\mu(t)$ by $w(t) dt$. This is not a highly restrictive assumption to make, and it is all we need here anyway.

such a sequence of polynomials, then, S_n will be very close to S . If we can do this right, S_n will be easy to express, and will converge to S faster than the original series does.

Now exactly what is S_n ? We have

$$\begin{aligned} S_n &= S - \frac{1}{P_n(-1)} \int_0^1 \frac{P_n(t)}{1+t} dt \\ &= S - \frac{1}{P_n(-1)} \int_0^1 \left(\sum_{j=0}^n p_j^{(n)} t^j \right) \sum_{k=0}^{\infty} (-t)^k d\mu(t) \\ &= S - \frac{1}{P_n(-1)} \int_0^1 \left(\sum_{\kappa=0}^n t^{\kappa} \sum_{j=0}^{\kappa} p_j^{(n)} (-1)^{\kappa-j} + \sum_{\kappa=n+1}^{\infty} t^{\kappa} \sum_{j=0}^n p_j^{(n)} (-1)^{\kappa-j} \right) d\mu(t) \end{aligned}$$

(where κ in the last line is substituted for $k+j$).

Now we're ready to pick the polynomials P_n . In fact, we're going to use $P_n(t) = (1-t)^n$. This works very nicely, because $P_n(-1) = 2^n$, which grows rapidly, and on the other hand $P_n(t)$ is bounded by 1 on the interval $[0, 1]$, so the integral we were concerned about above is easily bounded

$$\frac{1}{P_n(-1)} \int_0^1 \frac{P_n(t)}{t+1} d\mu(t) = O(2^{-n})$$

and we see that it goes to 0 very quickly. So certainly S_n approaches S very quickly. It remains for us to express S_n in a form that is easy to evaluate.

Now with this choice of P_n , we have by the binomial theorem

$$P_n(t) = (1-t)^n = \sum_{j=0}^n \binom{n}{j} (-t)^j$$

That is,

$$p_j^{(n)} = (-1)^j \binom{n}{j}$$

So now we can continue our computation of S_n from where we left off above:

$$S_n = S - \frac{1}{2^n} \int_0^1 \sum_{k=0}^n (-1)^k t^k \sum_{j=0}^k \binom{n}{j} d\mu(t) - \frac{1}{2^n} \int_0^1 \sum_{k=n+1}^{\infty} (-1)^k t^k \sum_{j=0}^n \binom{n}{j} d\mu(t)$$

Now we know that $\sum_{j=0}^n \binom{n}{j} = 2^n$, so the third term on the right is just

$$\frac{1}{2^n} \int_0^1 \sum_{k=n+1}^{\infty} (-1)^k t^k 2^n d\mu(t) = \sum_{k=n+1}^{\infty} (-1)^k a_k$$

and so S minus this term is just $\sum_{k=0}^n (-1)^k a_k$. Thus we have

$$S_n = \sum_{k=0}^n (-1)^k a_k - \frac{1}{2^n} \int_0^1 \sum_{k=0}^n (-1)^k t^k \sum_{j=0}^k \binom{n}{j} d\mu(t)$$

There are a couple of ways we can look at this. The first term on the right-hand side is just the n^{th} partial sum of our original series. Presumably the second term on the right is very small. But actually we want to combine these two terms and get an expression that is convenient to compute. When we do this we get

$$S_n = \sum_{k=0}^n (-1)^k a_k \left(1 - \frac{1}{2^n} \sum_{j=0}^k \binom{n}{j} \right) = \sum_{k=0}^{n-1} (-1)^k a_k \cdot \frac{1}{2^n} \sum_{j=k+1}^n \binom{n}{j}$$

(The outer sum originally has upper bound n , but the inner sum, in either expression, is 0 when $k = n$.)

This may look completely opaque. It turns out, however, that it is equal to the $(n-1)^{\text{th}}$ partial sum of the Euler transform. In fact, that partial sum is

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

and so we want to prove that

$$\frac{1}{2^n} \sum_{j=k+1}^n \binom{n}{j} = \sum_{j=k}^{n-1} \frac{1}{2^{j+1}} \binom{j}{k}$$

or equivalently,

$$(15) \quad \sum_{j=k+1}^n \binom{n}{j} = \sum_{j=k}^{n-1} 2^{n-j-1} \binom{j}{k}$$

This can be proved simply by induction, but as usual, that gives no hint as to why it is really true. There is an elegant combinatorial proof of this, however, using the method of “counting something two different ways”:

Let $S^{(n)}$ be a set with n elements. The left-hand side of (15)—call it $N_k^{(n)}$ —is just the number of subsets of $S^{(n)}$ containing more than k elements. Now let us suppose that the elements of $S^{(n)}$ are numbered from 1 to n . Given any subset T of $S^{(n)}$ containing more than k elements, we can look at which element of $S^{(n)}$ is the $(k+1)^{\text{th}}$ element of T . That element must be in one of the positions $k+1, k+2, \dots, n$ of $S^{(n)}$. Thus, we can write

$$N_k^{(n)} = \sum_{j=k}^{n-1} A_{k,j}^{(n)}$$

where $A_{k,j}^{(n)}$ is the number of subsets of $S^{(n)}$ each of which has the property that its $(k+1)^{\text{th}}$ element occurs in position $j+1$. And we can evaluate $A_{k,j}^{(n)}$ easily by looking at Figure 9.

We can see that any subset T whose $(k+1)^{\text{th}}$ element occurs in position $j+1$ is composed of three pieces:

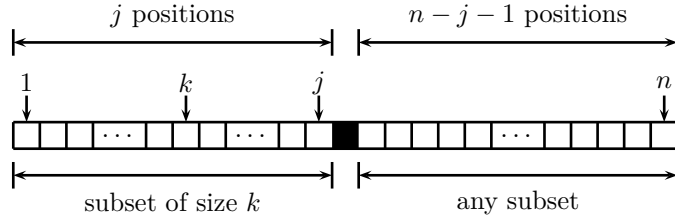


Figure 9: A set with n elements, numbered from 1 to n , containing a subset T whose $(k + 1)^{\text{th}}$ element occurs in position $j + 1$.

- There are k elements in positions 1 through j . These elements can be picked in $\binom{j}{k}$ ways.
- There is 1 element in position $j + 1$.
- There are any number of elements in the remaining positions. That is, we can pick any subset of the remaining positions, and there are 2^{n-j-1} such subsets.

So $A_{k,j}^{(n)}$, which is the number of such subsets T , must be $2^{n-j-1} \binom{j}{k}$. And this is just the term of index j on the right-hand side of (15), so we have derived that identity, and with it, the correctness of the Euler transform for the the class of alternating series the absolute values of whose terms are representable as moments of a positive measure.

Now let us apply the Euler transform to Leibniz's series. We have $a_k = \frac{1}{2k+1}$. We have to make sure that there is a weight function corresponding to this series. And there is: we have

$$\frac{1}{2k+1} = \int_0^1 t^k \frac{1}{2\sqrt{t}} dt$$

(Note that $\frac{1}{2\sqrt{t}}$, even though it is unbounded on $[0, 1]$, still has a finite integral over that interval.)

Therefore the derivation we have just gone through applies, and we know that the Euler transform of this series converges fairly rapidly to its sum. Now what is the Euler transform of this series?

We have

$$\begin{aligned} \nabla^0 a_0 &= 1 \\ \nabla^1 a_0 &= 1 - \frac{1}{3} = \frac{2}{3} \\ \nabla^2 a_0 &= 1 - 2 \cdot \frac{1}{3} + \frac{1}{5} = \frac{8}{15} = \frac{2 \cdot 4}{3 \cdot 5} \\ \nabla^3 a_0 &= 1 - 3 \cdot \frac{1}{3} + 3 \cdot \frac{1}{5} - \frac{1}{7} = \frac{16}{35} = \frac{2 \cdot 4 \cdot 6}{3 \cdot 5 \cdot 7} \end{aligned}$$

and it appears that⁷

$$\nabla^k a_0 = \frac{2 \cdot 4 \cdot \dots \cdot 2k}{3 \cdot 5 \cdot \dots \cdot (2k+1)}$$

We can prove this by induction as follows: by definition we have

$$\nabla^k a_0 = \sum_{j=0}^k \frac{(-1)^j}{2j+1} \binom{k}{j}$$

Let us turn this into a polynomial (this may seem like magic, but it's actually a standard technique); we write

$$f_k(x) = \sum_{j=0}^k (-1)^j \frac{x^{2j+1}}{2j+1} \binom{k}{j}$$

so $f_k(1) = \nabla^k a_0$. Differentiating, we get $f'_k(x) = (1-x^2)^k$. Since $f_k(0) = 0$, then, we must have $f_k(x) = \int_0^x (1-t^2)^k dt$. Starting with the integral for f_k , integrating by parts, simplifying, and setting $x = 1$, we get $(2k+1)f_k(1) = 2kf_{k-1}(1)$, which is exactly the inductive step.

Thus, we have

$$\frac{\pi}{4} = \sum_{k=0}^{\infty} \frac{1}{2^{k+1}} \frac{2^k k!}{3 \cdot 5 \cdot \dots \cdot (2k+1)} = \frac{1}{2} \sum_{k=0}^{\infty} \frac{k!}{3 \cdot 5 \cdot \dots \cdot (2k+1)}$$

and we know that the n^{th} partial sum differs from $\pi/4$ by $O(2^{-n})$.

This is a real improvement over Leibniz's original series (not necessarily in terms of aesthetics, but at least in terms of rate of convergence), but it is yet another linearly convergent algorithm—the number of correct binary digits is proportional to n . In the next section we will see how much faster algorithms were finally arrived at in 1976 by Brent and Salamin. The techniques they used actually go back to Gauss.

Before we leave this topic, however, we note⁸ that another way to write the computations in the series we have just derived is this:

$$\pi = 2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 + \frac{3}{7} \left(\dots \left(2 + \frac{i}{2i+1} \left(\dots \right) \right) \right) \right) \right)$$

We will come back to this in Section 8.2.

⁷This is true provided we make the convention that $\nabla^0 a_0 = 1$. Actually, if we write the formula as

$$\nabla^k a_0 = \frac{2^k k!}{3 \cdot 5 \cdot \dots \cdot (2k+1)}$$

then this works out automatically, because $0! = 1$.

⁸This is obvious if we write out each term, like this:

$$\pi = 2 \left(1 + \frac{1}{3} + \frac{1 \cdot 2}{3 \cdot 5} + \frac{1 \cdot 2 \cdot 3}{3 \cdot 5 \cdot 7} + \frac{1 \cdot 2 \cdot 3 \cdot 4}{3 \cdot 5 \cdot 7 \cdot 9} + \dots \right)$$

7 Quadratically converging approximations to π

We already saw above that the arithmetic-geometric mean converges quadratically. It is a remarkable fact that the arithmetic-geometric mean is related to π in such a way that it can be used as the basis of a quadratically converging algorithm that computes π . The mathematics behind this result was all known by the early 1800's, but was not actually exhibited as an algorithm until 1976, when Brent and Salamin independently and simultaneously derived it, no doubt motivated by the ability to perform these computations on computers.

The connection between the arithmetic-geometric mean and π goes by way of elliptic integrals. These are integrals involving the square root of a cubic or quartic polynomial, which cannot be evaluated in terms of the usual elementary functions such as polynomials, the exponential and logarithmic functions, and the trigonometric functions. They come up in many applications and were widely studied by many mathematicians, well back into the 1700's. There are many forms of such integrals, and there are many transformations that take one form into another. The study of these integrals necessarily involves trains of clever substitutions, which often seem unmotivated to us nowadays, but which were the common practice of mathematicians until the 20th century.

In fact, Legendre showed that by means of these substitutions, all of these elliptic integrals could be expressed in terms of three forms which he called elliptic integrals of the first, second, and third kind. Two of these three kinds will be important to us, and here we show, by way of motivation, how these two forms arise in applications.

7.1 Some representative elliptic integrals

7.1.1 The elliptic integral of the first kind

In first-year college physics courses, it is shown that the period of a simple pendulum is practically independent of the amplitude of its swing, provided the angular amplitude α is small. This depends on the approximation $\sin \alpha \cong \alpha$, which is quite accurate when α is small. In general, however, the solution involves an elliptic integral, and does in fact depend on α .

So let us assume that our pendulum is a point mass m suspended on a massless string of length l . At any time t , let θ denote the angle the string makes with the vertical rest position, and let v denote the velocity of the mass. And as usual, let g denote the gravitational acceleration.

At time t , the energy of the pendulum is the sum of its kinetic and potential energies, so it is

$$\frac{1}{2}mv^2 - mgl \cos \theta$$

and this must be constant. We want to express this in terms of θ . Since $v = l\dot{\theta}$, this becomes⁹

$$\frac{1}{2}ml^2\dot{\theta}^2 - mgl \cos \theta$$

Dividing by the constant $\frac{1}{2}ml^2$, we get the necessarily constant expression

$$\dot{\theta}^2 - 2\frac{g}{l} \cos \theta$$

⁹Ever since Newton, physicists have used dots to represent derivatives with respect to time, and since this is physics, we'll briefly follow that convention.

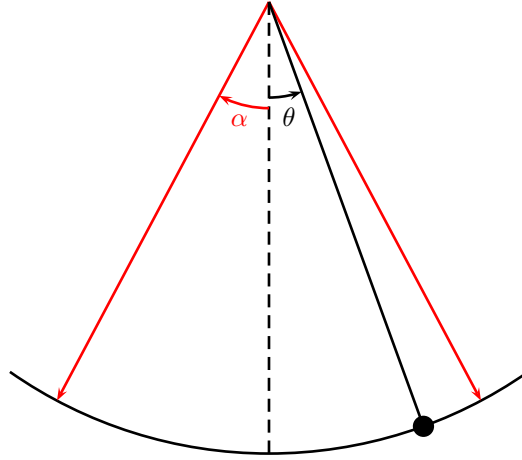


Figure 10: A pendulum. θ is the displacement of the pendulum from its vertical rest position at time t . α represents the amplitude; that is, the greatest possible value of θ ; equivalently, the greatest angular displacement at any time. Note that by convention $\alpha > 0$, even though for clarity α is represented as being negative in this figure. In fact, in the formulas we derive, it really doesn't matter whether α is positive or negative.

Now $\dot{\theta} = 0$ when $\theta = \pm\alpha$ (i.e., when the pendulum reaches its highest points), so the value of the constant expression must be $2\frac{g}{l}\cos\alpha$, and therefore we have

$$\dot{\theta}^2 = 2\frac{g}{l}(\cos\theta - \cos\alpha)$$

Now it is convenient¹⁰ to use the identity $\cos 2\phi = \cos^2\phi - \sin^2\phi = 1 - 2\sin^2\phi$, or equivalently, $\cos\phi = 1 - 2\sin^2\frac{1}{2}\phi$, to rewrite this as

$$\dot{\theta}^2 = 4\frac{g}{l}(\sin^2\frac{1}{2}\alpha - \sin^2\frac{1}{2}\theta)$$

So

$$\frac{d\theta}{dt} = 2\sqrt{\frac{g}{l}}\sqrt{\sin^2\frac{1}{2}\alpha - \sin^2\frac{1}{2}\theta}$$

and so

$$\frac{dt}{d\theta} = \frac{1}{2}\sqrt{\frac{l}{g}} \cdot \frac{1}{\sqrt{\sin^2\frac{1}{2}\alpha - \sin^2\frac{1}{2}\theta}}$$

¹⁰Who would think of making this substitution? Well, in those days, this was pretty standard stuff.

Therefore the time taken to go from the bottom point $\theta = 0$ to a point at amplitude $\theta = \phi$ must be

$$t(\phi) = \frac{1}{2} \sqrt{\frac{l}{g}} \int_0^\phi \frac{d\theta}{\sqrt{\sin^2 \frac{1}{2}\alpha - \sin^2 \frac{1}{2}\theta}}$$

Now it is again convenient to simplify the integrand by introducing a variable ψ given by

$$\sin \frac{1}{2}\theta = \sin \frac{1}{2}\alpha \sin \psi$$

With this substitution, we have

$$d\theta = \frac{2 \sin \frac{1}{2}\alpha \cos \psi}{\cos \frac{1}{2}\theta} d\psi = \frac{2 \sin \frac{1}{2}\alpha \cos \psi}{\sqrt{1 - \sin^2 \frac{1}{2}\alpha \sin^2 \psi}} d\psi$$

and so we get

$$t(\phi) = \sqrt{\frac{l}{g}} \int_{\theta=0}^{\theta=\phi} \frac{d\psi}{\sqrt{1 - \sin^2 \frac{1}{2}\alpha \sin^2 \psi}}$$

The period T of the pendulum is then 4 times the time taken in going from $\theta = 0$ all the way to $\theta = \alpha$, i.e., from $\psi = 0$ to $\psi = \frac{\pi}{2}$. That is

$$T = 4t(\alpha) = 4 \sqrt{\frac{l}{g}} \int_0^{\pi/2} \frac{d\psi}{\sqrt{1 - \sin^2 \frac{1}{2}\alpha \sin^2 \psi}}$$

And so finally,

$$(16) \quad T = 4 \sqrt{\frac{l}{g}} K(\sin \frac{1}{2}\alpha)$$

where

$$(17) \quad K(k) = \int_0^{\pi/2} \frac{d\psi}{\sqrt{1 - k^2 \sin^2 \psi}}$$

is the *complete elliptic integral of the first kind*. (Note that $|k| = |\sin \frac{1}{2}\alpha| \leq 1$.)

By making the substitution $t = \sin \theta$, we see also that

$$(18) \quad K(k) = \int_0^1 \frac{dt}{\sqrt{(1-t^2)(1-k^2 t^2)}}$$

which shows, as promised, that $K(k)$ is the integral of a rational function of a quartic polynomial. (More generally, elliptic integrals are integrals with respect to t of functions of the form $R(t, w)$, where R is a rational function of its two arguments and where w is the square root of a polynomial in t of degree 3 or 4.)

7.1.2 The elliptic integral of the second kind

The term “elliptic integral” comes from the attempt to find the perimeter of an ellipse in terms of its major and minor axes. If the ellipse (see Figure 11) is given by the parametric equations

$$x = a \cos \theta \quad y = b \sin \theta$$

where $a > b > 0$, then the differential of the arc length s is

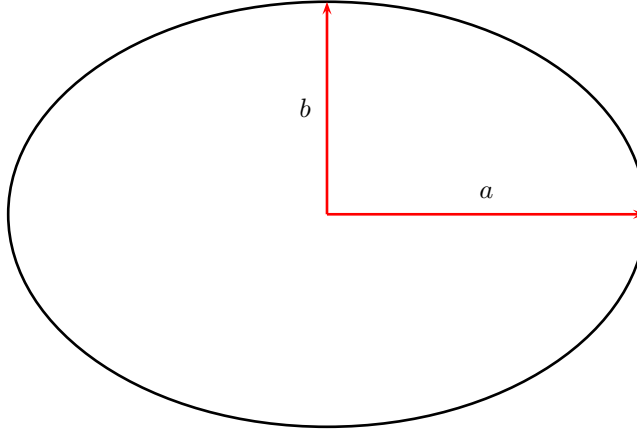


Figure 11: An ellipse.

$$ds = \sqrt{dx^2 + dy^2} = \sqrt{a^2 \sin^2 \theta + b^2 \cos^2 \theta} d\theta$$

and so the perimeter of the ellipse is 4 times the perimeter of a quadrant, and so is

$$s = 4 \int_0^{\pi/2} \sqrt{a^2 \sin^2 \theta + b^2 \cos^2 \theta} d\theta = 4 \int_0^{\pi/2} \sqrt{a^2 \cos^2 \theta + b^2 \sin^2 \theta} d\theta$$

where the last equality is arrived at by replacing θ by $\frac{\pi}{2} - \theta$. This is customarily done, probably to bring this integral into a form parallel with that for K in the next section.

Now letting $e = \sqrt{1 - b^2/a^2}$ denote the eccentricity of the ellipse, we then have the arc length of the ellipse given by

$$s = 4a \int_0^{\pi/2} \sqrt{1 - e^2 \sin^2 \theta} d\theta = 4aE(e)$$

where

$$(19) \quad E(k) = \int_0^{\pi/2} \sqrt{1 - k^2 \sin^2 \theta} d\theta$$

is the *complete elliptic integral of the second kind*. (Note here again that $k = e$, which is between 0 and 1.)

As above, making the substitution $t = \sin \theta$ yields the representation

$$(20) \quad E(k) = \int_0^1 \frac{\sqrt{1-k^2t^2}}{\sqrt{1-t^2}} dt = \int_0^1 \frac{1-k^2t^2}{\sqrt{(1-t^2)(1-k^2t^2)}} dt$$

which shows that $E(k)$ is also an integral with respect to t of a function $R(t, w)$ where R is a rational function of its arguments and w is the square root of a polynomial of order 4 in t .

7.1.3 More representations for $K(k)$ and $E(k)$

First, note that with a , b and e as above, and writing k for e as is conventional here, we have (similarly to the computation above)

$$\int_0^{\pi/2} \frac{d\theta}{\sqrt{a^2 \cos^2 \theta + b^2 \sin^2 \theta}} = \frac{1}{a} \int_0^{\pi/2} \frac{d\theta}{\sqrt{1-k^2 \sin^2 \theta}}$$

That is, defining

$$(21) \quad I(a, b) = \int_0^{\pi/2} \frac{d\theta}{\sqrt{a^2 \cos^2 \theta + b^2 \sin^2 \theta}}$$

we have

$$(22) \quad I(a, b) = \frac{1}{a} K(k)$$

As we'll see, it is customary to set $k' = \sqrt{1-k^2}$ (so that $k^2 + k'^2 = 1$). Then if $a = 1$, then $b = k'$, and so

$$(23) \quad I(1, k') = K(k)$$

It will be useful to us later to have a similar representation for E : the computations in section 7.1.2 show that for $0 < b < a$,

$$(24) \quad J(a, b) = \int_0^{\pi/2} \sqrt{a^2 \cos^2 \theta + b^2 \sin^2 \theta} d\theta = aE(k)$$

where again

$$k = \sqrt{1 - \frac{b^2}{a^2}}$$

And similarly to (23), we have

$$(25) \quad J(1, k') = E(k)$$

Now returning to (21) and making the clever substitution $t = b \tan \theta$, we have

$$dt = b \sec^2 \theta d\theta = b(1 + \tan^2 \theta) d\theta = b\sqrt{1 + \tan^2 \theta} \sqrt{1 + \tan^2 \theta} d\theta = \sqrt{b^2 + t^2} \frac{d\theta}{\cos \theta}$$

and so

$$(26) \quad I(a, b) = \int_0^\infty \frac{1}{\cos \theta \sqrt{a^2 + b^2 \tan^2 \theta}} \frac{\cos \theta}{\sqrt{b^2 + t^2}} dt = \frac{1}{2} \int_{-\infty}^\infty \frac{dt}{\sqrt{(a^2 + t^2)(b^2 + t^2)}}$$

and this is the form we will find useful.

This representation makes it perfectly clear that $I(a, b) = I(b, a)$, which could also have been proved directly by making the substitution $\phi = \frac{\pi}{2} - \theta$ in (21).

7.2 Elliptic integrals and the arithmetic-geometric mean

7.2.1 The relation of $M(a, b)$ to $K(k)$

Gauss discovered that the arithmetic-geometric mean $M(a, b)$ is closely related to the elliptic integral $K(k)$. We can now derive this fact in a relatively short manner.

Following Newman (1985), we start with equation (26) and we make the change of variables

$$t = \frac{x - \frac{ab}{x}}{2}$$

which entails

$$\begin{aligned} dt &= \frac{x^2 + ab}{2x^2} dx \\ t^2 + \left(\frac{a+b}{2}\right)^2 &= \frac{(x^2 + a^2)(x^2 + b^2)}{4x^2} \\ t^2 + ab &= \frac{(x^2 + ab)^2}{4x^2} \end{aligned}$$

so that

$$\begin{aligned} \int_{-\infty}^\infty \frac{dx}{\sqrt{(x^2 + a^2)(x^2 + b^2)}} &= \int_0^\infty \frac{2 dx}{\sqrt{(x^2 + a^2)(x^2 + b^2)}} \\ &= \int_{-\infty}^\infty \frac{dt}{\sqrt{\left(t^2 + \left(\frac{a+b}{2}\right)^2\right)(t^2 + ab)}} \end{aligned}$$

That is to say, $I(a, b) = I(a_1, b_1)$, and by induction and the fact that $I(a, b)$ is a continuous function of (a, b) , we get $I(a, b) = \lim_n I(a_n, b_n) = I(m, m)$, where $m = M(a, b)$. Now

$$I(m, m) = \frac{1}{2} \int_{-\infty}^\infty \frac{dx}{x^2 + m^2} = \frac{1}{2m} \int_{-\infty}^\infty \frac{dy}{y^2 + 1} = \frac{\pi}{2m}$$

and so finally, we see that

$$M(a, b) = \frac{\pi}{2I(a, b)}$$

and in particular,

$$(27) \quad M(1, k') = \frac{\pi}{2K(k)}$$

7.2.2 The arc length of the lemniscate

What seems to have been historically crucial in the development of elliptic integrals and the arithmetic-geometric mean is a topic which might seem much less significant to us today—the arc length of the lemniscate. The lemniscate is a curve in the plane whose equation in polar coordinates is $r^2 = \cos 2\theta$ (see Figure 12). The total arc length of this curve is

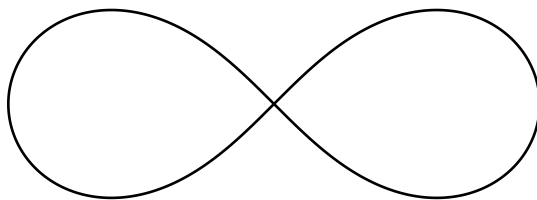


Figure 12: The lemniscate $r^2 = \cos 2\theta$. The equation only makes sense for $-\frac{\pi}{4} \leq \theta \leq \frac{\pi}{4}$ and $\frac{3\pi}{4} \leq \theta \leq \frac{5\pi}{4}$, but that is sufficient to parametrize the entire curve.

$$4 \int_0^{\pi/4} (r^2 + (dr/d\theta)^2)^{1/2} d\theta = 4 \int_0^{\pi/4} (\cos 2\theta)^{-1/2} d\theta$$

Now we make one of the wonderful trigonometric substitutions $\cos 2\theta = \cos^2 \phi$, and this becomes

$$4 \int_0^{\pi/2} (1 + \cos^2 \phi)^{-1/2} d\phi = 4 \int_0^{\pi/2} (2 \cos^2 \phi + \sin^2 \phi)^{-1/2} d\phi = 4I(\sqrt{2}, 1) = 4I(1, \sqrt{2}) = \frac{2\pi}{M(1, \sqrt{2})}$$

Now the lemniscate was a curve that people were widely interested in, and Gauss in fact had been looking at the integral for the arc length. (He had been preceded in this by some very acute observations of Euler and Lagrange.) He knew that it was related to π in certain ways. Gauss was no slouch at computation, and he actually computed the ratio

$$\frac{2}{\pi} \int_0^{\pi/2} \frac{d\phi}{\sqrt{2 \cos^2 \phi + \sin^2 \phi}}$$

by hand to 15 decimal places¹¹. Given what people were interested in at the time, this is perhaps not so amazing. But what came next was: he independently was investigating the arithmetic-geometric mean, and on May 30, 1799 he computed $M(\sqrt{2}, 1)$ and found that its reciprocal equaled the ratio he had computed above. He then proved all this; but the computation actually came first. Now all this didn't exactly come out of the blue, and in fact, Lagrange seems to have derived essentially the same result several years earlier. But he doesn't seem to have realized its significance; and Gauss did.

7.3 Landen's transformations

7.3.1 Cayley's derivation of Landen's transformations

The transformations by means of which we demonstrated the identity $M(1, k') = \pi/2K(k)$ are modern polished versions of some transformations that are originally due to John Landen (1719–1790), an obscure English mathematician who nevertheless contributed one of the foundational results in the field of elliptic integrals. Landen's transformations actually become fully meaningful only in the larger realm of what is called elliptic functions (which are the inverses of elliptic integrals and have an extensive and elegant theory which dominated much of 19th century mathematics). We can't even begin to go there in this survey; nor do we want to reproduce Landen's original work, which would be quite hard to follow, even though elementary.

But we do need his results, and to get them, we will follow Cayley (1895), Chapter XIII. Although Cayley's treatment of all this is fairly direct, it is nonetheless real 19th century mathematics, full of seemingly unmotivated substitutions which miraculously come out just right. They often can be better motivated and understood in the context of elliptic functions, but that would make the story even longer, and in any case, it's not how these results were originally discovered. So my suggestion to the reader who is encountering this intricate series of geometric and algebraic computations for the first time is simply to relax and enjoy the ride.

We will use the notations

$$I(a, b, \theta) = \int_0^\theta \frac{d\phi}{\sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi}}$$

$$J(a, b, \theta) = \int_0^\theta \sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi} d\phi$$

for the *incomplete* elliptic integrals. The *complete* elliptic integrals are those for which $\theta = \pi/2$:

$$I(a, b, \pi/2) = I(a, b)$$

$$J(a, b, \pi/2) = J(a, b)$$

In Figure 13, say the radius of the circle (i.e., the length $OA = OP = OB$) is r . Setting

$$a_1 = \frac{a+b}{2} \quad b_1 = \sqrt{ab}$$

¹¹How did he evaluate the integral? Well, Euler and others had already worked on this, and Gauss was able to evaluate it in terms of some rapidly converging infinite series. This is a topic in itself, and is surveyed in Cox (1984).

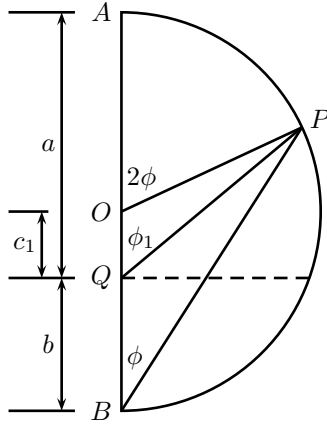


Figure 13: The geometric construction on which Landen's transformation is based. O is the center of the circle. Landen's transformations relate the values of I and J at (a, b, ϕ) with their values at (a_1, b_1, ϕ_1) .

We then have

$$r = \frac{a+b}{2} = a_1 \quad c_1 = \frac{a-b}{2}$$

Note that $a_1^2 - b_1^2 = c_1^2$, and for consistency we will also set $c^2 = a^2 - b^2$.

There are two facts from elementary geometry we need to rely on:

- The angle at O marked 2ϕ really is twice the angle at B marked ϕ .
- A horizontal line (the dashed line in the figure) from Q to the circumference of the circle has length $\sqrt{ab} = b_1$.

Now we have

$$\begin{aligned} QP^2 &= c_1^2 + 2c_1a_1 \cos 2\phi + a_1^2 \\ &= \frac{1}{2}(a^2 + b^2) + \frac{1}{2}(a^2 - b^2) \cos 2\phi \\ &= \frac{1}{2}(a^2 + b^2)(\cos^2 \phi + \sin^2 \phi) + \frac{1}{2}(a^2 - b^2)(\cos^2 \phi - \sin^2 \phi) \\ &= a^2 \cos^2 \phi + b^2 \sin^2 \phi \end{aligned}$$

Therefore

$$\begin{aligned} \sin \phi_1 &= \frac{a_1 \sin 2\phi}{\sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi}} \\ \cos \phi_1 &= \frac{c_1 + a_1 \cos 2\phi}{\sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi}} \end{aligned}$$

Now this yields

$$\begin{aligned}
a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1 &= a_1^2(1 - \sin^2 \phi_1) + b_1^2 \sin^2 \phi_1 \\
&= a_1^2 - \frac{(a_1^2 - b_1^2) \sin^2 2\phi}{a^2 \cos^2 \phi + b^2 \sin^2 \phi} \\
&= \frac{a_1^2(a^2 \cos^2 \phi + b^2 \sin^2 \phi - (a_1^2 - b_1^2) \sin^2 2\phi)}{a^2 \cos^2 \phi + b^2 \sin^2 \phi} \\
&= \frac{a_1^2(a^2 \cos^2 \phi + b^2 \sin^2 \phi - (a^2 - 2ab + b^2) \sin^2 \phi \cos^2 \phi)}{a^2 \cos^2 \phi + b^2 \sin^2 \phi} \\
&= \frac{a_1^2(a^2 \cos^2 \phi + b^2 \sin^2 \phi - a^2 \sin^2 \phi \cos^2 \phi + 2ab \sin^2 \phi \cos^2 \phi - b^2 \sin^2 \phi \cos^2 \phi)}{a^2 \cos^2 \phi + b^2 \sin^2 \phi} \\
&= \frac{a_1^2(a^2 \cos^4 \phi + 2ab \sin^2 \phi \cos^2 \phi + b^2 \cos^4 \phi)}{a^2 \cos^2 \phi + b^2 \sin^2 \phi} \\
&= \frac{a_1^2(a \cos^2 \phi + b \sin^2 \phi)^2}{a^2 \cos^2 \phi + b^2 \sin^2 \phi}
\end{aligned}$$

and also

$$\begin{aligned}
\sin(2\phi - \phi_1) &= \frac{\frac{1}{2}(a - b) \sin 2\phi}{\sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi}} && \text{(this is easy)} \\
\cos(2\phi - \phi_1) &= \frac{a \cos^2 \phi + b \sin^2 \phi}{\sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi}} && \text{(just work it out)}
\end{aligned}$$

and thence

$$\cos(2\phi - \phi_1) = \frac{1}{a_1} \sqrt{a_1 \cos^2 \phi_1 + b_1 \sin^2 \phi_1}$$

Now (see Figure 14) move P through a small angle on the circle to a point P' . Let P'' be the point on the circle centered at Q of radius QP such that P' is on the radius QP'' .

First, as $P' \rightarrow P$ we see that

- The line through P and P' approaches the tangent to the original circle at P .
- The line through P and P'' approaches the tangent to the dashed circle at P .

and therefore as $P' \rightarrow P$,

$$(28) \quad \angle P''PP' \rightarrow \angle QPO = 2\phi - \phi_1$$

Next, QP'' meets the dashed circle at a right angle, and so as $P' \rightarrow P$,

$$(29) \quad \angle P'P''P \rightarrow \text{a right angle}$$

Putting (28) together with (29), we see that as $P' \rightarrow P$,

$$\frac{PP''}{PP'} \rightarrow \cos(2\phi - \phi_1) = \frac{1}{a_1} \sqrt{a_1 \cos^2 \phi_1 + b_1 \sin^2 \phi_1}$$

Note that if $\theta = \pi/2$, then $\theta_1 = \pi$, and this identity reduces to

$$\begin{aligned} I(a, b) &= I(a, b, \pi/2) = \frac{1}{2} I(a_1, b_1, \pi) = \frac{1}{2} \int_0^\pi \frac{d\phi}{\sqrt{a_1^2 \cos^2 \phi + b_1^2 \sin^2 \phi}} \\ &= \int_0^{\pi/2} \frac{d\phi}{\sqrt{a_1^2 \cos^2 \phi + b_1^2 \sin^2 \phi}} = I(a_1, b_1) \end{aligned}$$

which we have already seen to be true.

What we are really interested in, however, is a similar formula involving the function J . Now on the one hand we already know that

$$(32) \quad \sin^2 \phi_1 = \frac{4a_1^2 \sin^2 \phi \cos^2 \phi}{a^2 \cos^2 \phi + b^2 \sin^2 \phi}$$

and if for convenience we write

$$X = a^2 \cos^2 \phi + b^2 \sin^2 \phi$$

then

$$X - a^2 = (b^2 - a^2) \sin^2 \phi$$

$$X - b^2 = (a^2 - b^2) \cos^2 \phi$$

and so

$$(X - a^2)(X - b^2) = -4(a - b)^2 a_1^2 \sin^2 \phi \cos^2 \phi$$

Putting this together with (32), we get

$$(X - a^2)(X - b^2) + X(a - b)^2 \sin^2 \phi_1 = 0$$

That is,

$$X^2 + X(-(a^2 + b^2)(\sin^2 \phi_1 + \cos^2 \phi_1) + (a - b)^2 \sin^2 \phi_1) + a^2 b^2 = 0$$

which, after simplifying the second term, is

$$X^2 - X((a^2 + b^2) \cos^2 \phi_1 + 2ab \sin^2 \phi_1) + a^2 b^2 = 0$$

which can also be written as

$$\left(X - \left(\frac{1}{2}(a^2 + b^2) \cos^2 \phi_1 + ab \sin^2 \phi_1 \right) \right)^2 - \left(\frac{1}{2}(a^2 + b^2) \cos^2 \phi_1 + ab \sin^2 \phi_1 \right)^2 + a^2 b^2 (\cos^2 \phi_1 + \sin^2 \phi_1)^2 = 0$$

and so

$$\begin{aligned} &\left(X - \left(\frac{1}{2}(a^2 + b^2) \cos^2 \phi_1 + ab \sin^2 \phi_1 \right) \right)^2 \\ &= \frac{1}{4}(a^2 + b^2)^2 \cos^4 \phi_1 + (a^2 + b^2)ab \cos^2 \phi_1 \sin^2 \phi_1 + a^2 b^2 \sin^4 \phi_1 \\ &\quad - a^2 b^2 \cos^4 \phi_1 - 2a^2 b^2 \cos^2 \phi_1 \sin^2 \phi_1 - a^2 b^2 \sin^4 \phi_1 \\ &= \frac{1}{4}(a^2 - b^2)^2 \cos^4 \phi_1 + ab(a - b)^2 \cos^2 \phi_1 \sin^2 \phi_1 \\ &= 4c_1^2 \cos^2 \phi_1 (a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1) \end{aligned}$$

Now restoring the value of X , we get

$$\begin{aligned} a^2 \cos^2 \phi + b^2 \sin^2 \phi &= \frac{1}{2}(a^2 + b^2) \cos^2 \phi_1 + ab \sin^2 \phi_1 + 2c_1 \cos \phi_1 \sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1} \\ &= 2(a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1) - b_1^2 + 2c_1 \cos \phi_1 \sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1} \end{aligned}$$

We rewrite this as

$$\frac{a^2 \cos^2 \phi + b^2 \sin^2 \phi}{\sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1}} = 2 \left(\sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1} - \frac{\frac{1}{2}b_1^2}{\sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1}} + c_1 \cos \phi_1 \right)$$

Putting this together with (30), we get

$$\sqrt{a^2 \cos^2 \phi + b^2 \sin^2 \phi} d\phi = \left(\sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1} - \frac{\frac{1}{2}b_1^2}{\sqrt{a_1^2 \cos^2 \phi_1 + b_1^2 \sin^2 \phi_1}} + c_1 \cos \phi_1 \right) d\phi_1$$

and therefore

$$(33) \quad J(a, b, \theta) = J(a_1, b_1, \theta_1) - \frac{1}{2}b_1^2 I(a_1, b_1, \theta_1) + c_1 \sin \theta_1$$

In particular, when $\theta = \pi/2$, $\theta_1 = \pi$, and we have

$$(34) \quad J(a, b) = 2J(a_1, b_1) - b_1^2 I(a_1, b_1)$$

7.3.2 Iteration of Landen's transformation for J

Now let us set $a_0 = a$, $b_0 = b$, and $c_0 = c$. Let us inductively define

$$\begin{aligned} a_{n+1} &= \frac{a_n + b_n}{2} \\ b_{n+1} &= \sqrt{a_n b_n} \\ c_{n+1} &= \frac{a_n - b_n}{2} \end{aligned}$$

so as before, $c_n^2 = a_n^2 - b_n^2$ and $a_{n+1}c_{n+1} = \frac{1}{4}(a_n^2 - b_n^2)$.

We know that a_n and b_n converge extremely rapidly to the arithmetic-geometric mean $M(a, b)$.

We can write (34) as

$$(J(a, b) - a^2 I(a, b)) = 2(J(a_1, b_1) - a_1^2 I(a_1, b_1)) + I(a_1, b_1)(-b_1^2 + 2a_1^2 - a^2)$$

and since

$$2a_1^2 - a^2 - b_1^2 = -\frac{1}{2}(a^2 - b^2) = -\frac{1}{2}c^2$$

this can be written as

$$(35) \quad (J(a, b) - a^2 I(a, b)) = 2(J(a_1, b_1) - a_1^2 I(a_1, b_1)) - \frac{1}{2}c^2 I(a_1, b_1)$$

By induction from (35), and using the fact that $I(a_n, b_n) = I(a, b)$ for all n , this yields

$$(J(a, b) - a^2 I(a, b)) = 2^n (J(a_n, b_n) - a_n^2 I(a_n, b_n)) - \frac{1}{2} I(a, b) \sum_{j=0}^{n-1} 2^j c_j^2$$

Further, we have (using (10) on page 15 for the inequality at the end)

$$\begin{aligned} J(a_n, b_n) - a_n^2 I(a_n, b_n) &= \int_0^{\pi/2} \sqrt{a_n^2 \cos^2 \phi + b_n^2 \sin^2 \phi} d\phi + \int_0^{\pi/2} \frac{a_n^2}{\sqrt{a_n^2 \cos^2 \phi + b_n^2 \sin^2 \phi}} d\phi \\ &= \int_0^{\pi/2} \frac{(b_n^2 - a_n^2) \sin^2 \phi}{\sqrt{a_n^2 \cos^2 \phi + b_n^2 \sin^2 \phi}} d\phi \\ &= O(2^{-2^{n-N}}) \end{aligned}$$

for some N and all $n > N$. Therefore

$$2^n (J(a_n, b_n) - a_n^2 I(a_n, b_n)) = O(2^n 2^{-2^{n-N}})$$

which tends to 0 as $n \rightarrow \infty$, and so

$$(J(a, b) - a^2 I(a, b)) = -\frac{1}{2} I(a, b) \sum_{j=0}^{\infty} 2^j c_j^2$$

or equivalently,

$$(36) \quad J(a, b) = I(a, b) \left(a^2 - \sum_{j=0}^{\infty} 2^{j-1} c_j^2 \right)$$

This will be of great use to us below. And even by itself it tells us something interesting: we know that $I(a, b)$ can be computed by the arithmetic-geometric mean, which converges quadratically. And the numbers $\{c_j\}$ are just numbers that occur in that computation, and they converge to 0 quadratically, as we saw in (11). Therefore, we can compute $J(a, b)$ quadratically as well.

7.4 Legendre's relation

7.4.1 The functions E' and K'

It is convenient to define $k' = \sqrt{1 - k^2}$ so that $k^2 + k'^2 = 1$. (This is a standard notational convention in this subject, and is used in many contexts.)

Now let us define functions E' and K' by $E'(k) = E(k')$ and $K'(k) = K(k')$. (Actually, these definitions of “primed” objects are not just for convenience. But their real significance is only apparent in the larger theory of elliptic functions.)

We will need to know the behavior of these functions as $k \downarrow 0$. It is easy to see that $K(0) = E(0) =$

$\pi/2$. But we actually need a little more: by the binomial theorem,

$$(1 - k^2 \sin^2 \phi)^{1/2} = 1 - \frac{1}{2}k^2 \sin^2 \phi + O(k^4)$$

$$(1 - k^2 \sin^2 \phi)^{-1/2} = 1 + \frac{1}{2}k^2 \sin^2 \phi + O(k^4)$$

and so integrating these two quantities over $[0, \pi/2]$, we get

$$(37) \quad K(k) = \frac{\pi}{2} - \frac{\pi}{8}k^2 + O(k^4)$$

$$(38) \quad E(k) = \frac{\pi}{2} + \frac{\pi}{8}k^2 + O(k^4)$$

Further, we see easily¹² that

$$(39) \quad E'(0) = E(1) = \int_0^{\pi/2} \cos \phi \, d\phi = 1$$

The behavior of K' near 0 is more difficult to handle because it blows up. We need to get a bound on how fast it blows up at 0. In fact, for our needs below, a crude bound will suffice: we have

$$(40) \quad K'(k) = K(k') = \int_0^{\pi/2} \frac{1}{\sqrt{1 - k'^2 \sin^2 \phi}} \, d\phi \leq \int_0^{\pi/2} \frac{1}{\sqrt{1 - k'^2}} \, d\phi = O(1/k)$$

7.4.1.1 The precise behavior of K' near 0

In fact, a lot more is known about the behavior of K' near 0, and although we don't need this result, we'll derive it anyway, for completeness. We'll start with the representation of K' which we get from Section 7.1.3, and then divide the range of integration into two parts:

$$\begin{aligned} K'(k) = K(k') = I(1, k) &= \int_0^\infty \frac{dt}{\sqrt{(1+t^2)(k^2+t^2)}} \\ &= \int_0^\alpha \frac{dt}{\sqrt{(1+t^2)(k^2+t^2)}} + \int_\alpha^\infty \frac{dt}{\sqrt{(1+t^2)(k^2+t^2)}} \end{aligned}$$

We have divided up the range of integration so that we can approximate these two integrals separately. We have not yet decided what α will be, but we can allow it to depend on k . Now for the first integral, we have

$$\frac{1}{\sqrt{1+\alpha^2}} \int_0^\alpha \frac{dt}{\sqrt{k^2+t^2}} \leq \int_0^\alpha \frac{dt}{\sqrt{(1+t^2)(k^2+t^2)}} \leq \int_0^\alpha \frac{dt}{\sqrt{k^2+t^2}}$$

¹²And we also see easily that if we substitute (37) into (16) (page 29), we find that the period of a pendulum with small angular displacement α is

$$T = 2\pi \sqrt{\frac{l}{g}} \left(1 + \frac{1}{8}\alpha^2 + O(\alpha^4) \right)$$

the first term of which is the ordinary approximation.

And for the second integral, let us assume at any rate that $k < \alpha$. Then since in fact $k < \alpha \leq t$, we have $k = (k/t)t \leq (k/\alpha)t$, so

$$\frac{1}{\sqrt{1 + \frac{k^2}{\alpha^2}}} \int_{\alpha}^{\infty} \frac{dt}{t\sqrt{1+t^2}} \leq \int_{\alpha}^{\infty} \frac{dt}{\sqrt{(1+t^2)(k^2+t^2)}} \leq \int_{\alpha}^{\infty} \frac{dt}{t\sqrt{1+t^2}}$$

For both of these approximations to work, then, we are certainly going to need both of the factors on the left hand sides to approach 1 as $k \rightarrow 0$. Thus, we need $\alpha \rightarrow 0$ and $k/\alpha \rightarrow 0$. That serves to narrow things down a bit. Now let's evaluate the approximating integrals. The first one is

$$\int_0^{\alpha} \frac{dt}{\sqrt{k^2+t^2}} = \int_0^{\alpha/k} \frac{ds}{\sqrt{1+s^2}} = \log(\sqrt{1+s^2}+s) \Big|_0^{\alpha/k} = \log\left(\sqrt{1+\frac{\alpha^2}{k^2}} + \frac{\alpha}{k}\right)$$

and the second is

$$\int_{\alpha}^{\infty} \frac{dt}{t\sqrt{1+t^2}} = \log \frac{\sqrt{1+s^2}-1}{s} \Big|_{\alpha}^{\infty} = -\log \frac{\sqrt{1+\alpha^2}-1}{\alpha} = -\log(\sqrt{1+\alpha^2}-1) + \log \alpha$$

where the integrated term at the upper bound of integration is 0, because

$$\frac{\sqrt{1+s^2}-1}{s} = \sqrt{\left(\frac{1}{s^2}+1\right)} - \frac{1}{s} \rightarrow 1 \quad \text{as } s \rightarrow \infty$$

Putting these two results together, we have

$$\begin{aligned} \frac{1}{\sqrt{1+\alpha^2}} \log\left(\sqrt{1+\frac{\alpha^2}{k^2}} + \frac{\alpha}{k}\right) + \frac{1}{\sqrt{1+\frac{k^2}{\alpha^2}}} \left(-\log(\sqrt{1+\alpha^2}-1) + \log \alpha\right) \\ (41) \qquad \qquad \qquad \leq K'(k) \\ \leq \log\left(\sqrt{1+\frac{\alpha^2}{k^2}} + \frac{\alpha}{k}\right) + \left(-\log(\sqrt{1+\alpha^2}-1) + \log \alpha\right) \end{aligned}$$

We could at this point cut a long story short and notice that setting $\alpha = \sqrt{k}$ not only satisfies the two constraints we noted above but also makes the two additional factors on the left-hand side of this inequality the same, which makes the expression a lot simpler to deal with. However, there is another reason we might want to set $\alpha = \sqrt{k}$ as well:

Let us denote the right-hand side of this by U (for “upper bound”), and let us hold k fixed and regard U as a function of α :

$$U(\alpha) = \log\left(\sqrt{1+\frac{\alpha^2}{k^2}} + \frac{\alpha}{k}\right) + \left(-\log(\sqrt{1+\alpha^2}-1) + \log \alpha\right)$$

To get the most effective approximation (at least from above), we certainly want to make $U(\alpha)$ as small as possible. Calculus can help us here: The derivative of $U(\alpha)$ with respect to α can be worked out easily enough; it turns out to be

$$(42) \qquad \qquad \qquad U'(\alpha) = \frac{1}{\sqrt{k^2+\alpha^2}} - \frac{1}{\alpha\sqrt{1+\alpha^2}}$$

Differentiating once again with respect to α yields the second derivative:

$$U''(\alpha) = \frac{-\alpha}{(k^2 + \alpha^2)^{3/2}} + \frac{1 + 2\alpha^2}{\alpha^2(1 + \alpha^2)^{3/2}}$$

and assuming that $0 < k < \alpha \leq 1$, it is straightforward¹³ to see that $U''(\alpha) > 0$. Thus, $U(\alpha)$ is a convex function of α and so if $U'(\alpha) = 0$ for some α in the open interval $(k, 1)$, $U(\alpha)$ will have a minimum in that interval and it will be at that point. So we take the expression for $U'(\alpha)$ in (42) and set it equal to 0, and we get $k = \alpha^2$. (And I imagine that this also gives the tightest bound from below, but it's too painful for me to check that.) And with this choice of α , as already noted above, we do in fact have $k < \alpha < 1$, and also $\alpha \rightarrow 0$ and $k/\alpha \rightarrow 0$ as $k \rightarrow 0$.

Making the substitution $\alpha = \sqrt{k}$ in (41) now yields

$$(43) \quad \frac{1}{\sqrt{1+k}} \log \frac{\sqrt{1+k} + 1}{\sqrt{1+k} - 1} \leq K'(k) \leq \log \frac{\sqrt{1+k} + 1}{\sqrt{1+k} - 1}$$

and so we see that now regarding the right-hand side U as a function of k , we have

$$U(k) = \log \frac{\sqrt{1+k} + 1}{\sqrt{1+k} - 1}$$

and

$$\frac{K'(k)}{U(k)} \rightarrow 1 \quad \text{as } k \rightarrow 0$$

Further, since

$$U(k) = \log \frac{\sqrt{1+k} + 1}{\sqrt{1+k} - 1} = \log \frac{(\sqrt{1+k} + 1)^2}{k} = \log \frac{4 + O(k)}{k} = \log \left(\frac{4}{k} + O(1) \right) = \log \frac{4}{k} + O(k)$$

it follows that

$$\frac{U(k)}{\log \frac{4}{k}} \rightarrow 1 \quad \text{as } k \rightarrow 0$$

so

$$\frac{K'(k)}{\log \frac{4}{k}} \rightarrow 1 \quad \text{as } k \rightarrow 0$$

In this case, however, we can get a stronger result: Since $\frac{1}{\sqrt{1+k}} = 1 - O(k)$, we have from (43) that

$$(1 - O(k)) \left(\log \frac{4}{k} - O(k) \right) \leq K'(k) \leq \log \frac{4}{k} + O(k)$$

and so for k close to 0,

$$\left| K'(k) - \log \frac{4}{k} \right| = O \left(k \log \frac{1}{k} \right)$$

¹³First use $k > 0$. Then simplify and finally appeal to the fact that $\alpha \leq 1$.

and in particular,

$$(44) \quad K'(k) - \log \frac{4}{k} \rightarrow 0 \quad \text{as } k \rightarrow 0$$

which shows much more precisely how $K'(k)$ behaves near 0.

7.4.2 Legendre's relation

First let us derive formulas for the derivatives of E and K with respect to k .

Recall that we have

$$E = \int_0^{\pi/2} \sqrt{1 - k^2 \sin^2 \phi} d\phi \quad K = \int_0^{\pi/2} \frac{d\phi}{\sqrt{1 - k^2 \sin^2 \phi}}$$

and so we immediately get

$$\begin{aligned} \frac{dE}{dk} &= - \int_0^{\pi/2} \frac{k \sin^2 \phi}{\sqrt{1 - k^2 \sin^2 \phi}} d\phi = \frac{E - K}{k} \\ \frac{dK}{dk} &= \int_0^{\pi/2} \frac{k \sin^2 \phi}{(1 - k^2 \sin^2 \phi)^{3/2}} d\phi \end{aligned}$$

We need to go further with dK/dk .

Remember that by convention we have $k'^2 = 1 - k^2$. It is convenient to define a variable t by

$$(45) \quad \frac{k'^2}{1 - k^2 \sin^2 \phi} = 1 - k^2 \sin^2 t$$

so that as ϕ increases from 0 to $\pi/2$, t decreases from $\pi/2$ to 0. Solving for $\sin^2 t$, we get

$$\sin^2 t = \frac{1}{k^2} \left(1 - \frac{k'^2}{1 - k^2 \sin^2 \phi} \right) = \frac{\cos^2 \phi}{1 - k^2 \sin^2 \phi}$$

and so in turn we get

$$\cos t = \sqrt{1 - \sin^2 t} = \frac{k' \sin \phi}{\sqrt{1 - k^2 \sin^2 \phi}}$$

Since (45) is symmetric in t and ϕ , the last two identities hold if we switch t and ϕ . And from that in turn we get

$$\sin \phi \cos \phi = \frac{k' \sin t \cos t}{1 - k^2 \sin^2 t}$$

Further, differentiating and simplifying, using the last identity, we find that

$$d\phi = - \frac{k'}{1 - k^2 \sin^2 t} dt$$

Putting these results together we can derive

$$\begin{aligned}
k \frac{dK}{dk} &= \int_0^{\pi/2} \frac{k^2 \sin^2 \phi}{(1 - k^2 \sin^2 \phi)^{3/2}} d\phi \\
&= \frac{1}{k'^2} \int_0^{\pi/2} \frac{k'^2}{1 - k^2 \sin^2 \phi} \cdot \frac{k^2 \sin^2 \phi}{\sqrt{1 - k^2 \sin^2 \phi}} \cdot \frac{k'}{1 - k^2 \sin^2 t} dt \\
&= \frac{1}{k'^2} \int_0^{\pi/2} (1 - k^2 \sin^2 t) \cdot \frac{k^2 \sin^2 \phi}{\sqrt{1 - k^2 \sin^2 \phi}} \cdot \frac{k'}{1 - k^2 \sin^2 t} dt \\
&= \frac{k^2}{k'} \int_0^{\pi/2} \sin \phi \cdot \frac{\cos t}{k'} dt \\
&= \frac{k^2}{k'^2} \int_0^{\pi/2} \frac{\cos^2 t}{\sqrt{1 - k^2 \sin^2 t}} dt
\end{aligned}$$

On the other hand, we have

$$\begin{aligned}
\frac{E}{k'^2} - K &= \frac{1}{k'^2} \int_0^{\pi/2} \left(\sqrt{1 - k^2 \sin^2 \phi} - \frac{k'^2}{\sqrt{1 - k^2 \sin^2 \phi}} \right) d\phi \\
&= \frac{1}{k'^2} \int_0^{\pi/2} \frac{1 - k^2 \sin^2 \phi - 1 + k^2}{\sqrt{1 - k^2 \sin^2 \phi}} d\phi \\
&= \frac{1}{k'^2} \int_0^{\pi/2} \frac{k^2 \cos^2 \phi}{\sqrt{1 - k^2 \sin^2 \phi}} d\phi
\end{aligned}$$

so that finally we have

$$k \frac{dK}{dk} = \frac{E}{k'^2} - K$$

So we have

$$\begin{aligned}
\frac{dE}{dk} &= \frac{1}{k}(E - K) \\
\frac{dK}{dk} &= \frac{1}{kk'^2}(E - k'^2 K)
\end{aligned}$$

These same equations also hold with every variable “primed”. Further, since $dk'/dk = -k/k'$, we can use the chain rule to get

$$\begin{aligned}
\frac{dE'}{dk} &= \frac{dE'}{dk'} \frac{dk'}{dk} = -\frac{k}{k'^2}(E' - K') \\
\frac{dK'}{dk} &= \frac{dK'}{dk'} \frac{dk'}{dk} = -\frac{1}{kk'^2}(E' - k'^2 K')
\end{aligned}$$

Now we calculate

$$\begin{aligned}
\frac{d}{dk}(EK' + E'K - KK') &= \frac{1}{k}(E - K)K' + E\left(-\frac{1}{kk'^2}(E' - k^2K')\right) \\
&\quad + \left(-\frac{k}{k'^2}\right)(E' - K')K + E'\left(\frac{1}{kk'^2}\right)(E - k'^2K) \\
&\quad - \frac{1}{kk'^2}(E - k'^2K)K' - K\left(-\frac{1}{kk'^2}\right)(E' - k^2K') \\
&= \frac{1}{k}EK' - \frac{1}{k}KK' - \frac{1}{kk'}EE' + \frac{k}{k'^2}EK' \\
&\quad - \frac{k}{k'^2}E'K + \frac{k}{k'^2}K'K + \frac{1}{kk'^2}E'E - \frac{1}{k}E'K \\
&\quad - \frac{1}{kk'^2}EK' + \frac{1}{k}KK' + \frac{1}{kk'^2}E'K - \frac{k}{k'^2}KK' \\
&= E'K\left(-\frac{k}{k'^2} + \frac{1}{kk'^2} - \frac{1}{k}\right) \\
&\quad + EK'\left(\frac{k}{k'^2} - \frac{1}{kk'^2} + \frac{1}{k}\right) \\
&\quad + KK'\left(-\frac{1}{k} + \frac{k}{k'^2} + \frac{1}{k} - \frac{k}{k'^2}\right) \\
&= 0
\end{aligned}$$

Therefore the expression $EK' + E'K - KK'$ must be constant. To see what its value is, we can just let $k \rightarrow 0$. Using (37), (38), (39) and (40), we get (as $k \downarrow 0$)

$$\begin{aligned}
E(k)K'(k) + E'(k)K(k) - K(k)K'(k) &= (E(k) - K(k))K'(k) + E'(k)K(k) \\
&= O(k^2) \cdot O(1/k) + 1 \cdot \frac{\pi}{2} \\
&\rightarrow \frac{\pi}{2}
\end{aligned}$$

and so we have *Legendre's relation*

$$(46) \quad EK' + E'K - KK' = \frac{\pi}{2}$$

7.5 The algorithm of Brent and Salamin

We can now derive the algorithm of Brent and Salamin for computing π . Starting with (23) and (25), we have

$$\begin{aligned}
I(1, k') &= K(k) & J(1, k') &= E(k) \\
I(1, k) &= K'(k) & J(1, k) &= E'(k)
\end{aligned}$$

Now if we start two arithmetic-geometric iterations with

$$\begin{aligned}
a_0 &= 1 & b_0 &= k' \\
a'_0 &= 1 & b'_0 &= k
\end{aligned}$$

and if c_n and c'_n are defined in the usual manner:

$$c_n = \frac{a_{n-1} - b_{n-1}}{2} \quad c'_n = \frac{a'_{n-1} - b'_{n-1}}{2}$$

then from (36) on page 40 we have

$$E(k) = \left(1 - \sum_{j=0}^{\infty} 2^{j-1} c_j^2\right) K(k)$$

$$E'(k) = \left(1 - \sum_{j=0}^{\infty} 2^{j-1} c_j'^2\right) K'(k)$$

Substituting in Legendre's relation for $E(k)$ and $E'(k)$, we get

$$\left(1 - \sum_{j=0}^{\infty} 2^{j-1} c_j^2\right) K(k) K'(k) + \left(1 - \sum_{j=0}^{\infty} 2^{j-1} c_j'^2\right) K(k) K'(k) - K(k) K'(k) = \frac{\pi}{2}$$

But since

$$K'(k) = \frac{\pi}{2M(1, k)}$$

$$K(k) = \frac{\pi}{2M(1, k')}$$

we then get

$$\frac{\pi}{2} = \left(\left(1 - \sum_{j=0}^{\infty} 2^{j-1} c_j^2\right) + \left(1 - \sum_{j=0}^{\infty} 2^{j-1} c_j'^2\right) - 1 \right) \frac{\pi^2}{4M(1, k)M(1, k')}$$

Now the $j = 0$ terms of the two summations sum to $\frac{1}{2}c_0^2 + \frac{1}{2}c_0'^2 = \frac{1}{2}((a_0^2 - b_0^2) + (a_0'^2 - b_0'^2)) = \frac{1}{2}(k^2 + k'^2) = \frac{1}{2}$. Therefore we can write this as

$$\frac{\pi}{2} = \frac{1}{2} \left(1 - \sum_{j=1}^{\infty} 2^j c_j^2 - \sum_{j=1}^{\infty} 2^j c_j'^2 \right) \frac{\pi^2}{4M(1, k)M(1, k')}$$

That is,

$$\pi = \frac{4M(1, k)M(1, k')}{1 - \sum_{j=1}^{\infty} 2^j (c_j^2 + c_j'^2)}$$

Finally, suppose we start with $b_0 = k = \frac{1}{\sqrt{2}}$. Then also $b'_0 = k' = \frac{1}{\sqrt{2}}$, and $c_j = c'_j$ for all j , and so we have

$$\pi = \frac{4M\left(1, \frac{1}{\sqrt{2}}\right)^2}{1 - \sum_{j=1}^{\infty} 2^{j+1} c_j^2}$$

which is the basis for Brent and Salamin’s algorithm for approximating π . In fact to turn it into an algorithm, simply approximate the numerator by $4a_{n+1}^2$ and truncate the sum in the denominator¹⁴. That is, we denote the n^{th} approximant by

$$(47) \quad \pi_n = \frac{4a_{n+1}^2}{1 - \sum_{j=1}^n 2^{j+1} c_j^2}$$

This converges quadratically, as we know it must (since both the numerator and denominator converge quadratically), and we can see this quite dramatically in Figure 15.

n	π_n
1	3.140579250522168248311331268975823311773440237512948335643486693345582758034902907827
2	3.141592646213542282149344431982695774314437223345602794559539484821434767220795264694
3	3.141592653589793238279512774801863974381225504835446935787330702026382137838927399031
4	3.141592653589793238462643383279502884197114678283648921556617106976026764500643061711
5	3.141592653589793238462643383279502884197169399375105820974944592307816406286208998625

Figure 15: Convergence of the approximations to π , using the algorithm of Brent and Salamin.

8 Streaming algorithms

All the algorithms we have seen so far have the property that you have to start by deciding how accurate you want your final result to be. Then you have to insure that each intermediate step is computed to enough accuracy to make the final result come out right.

We will now consider a different class of algorithm. In the algorithms we will consider here, you don’t have to decide how accurate you want your result to be. As it runs, the algorithm spits out consecutive digits of the result (in our case, the consecutive digits of π), and you can just let it run as long as you want. And since all the arithmetic is integer arithmetic, you don’t have to worry about approximating intermediate results.

The discussion here is adapted from Gibbons (2006).

8.1 Multiplication

We start with a simple but instructive algorithm—multiplication. Of course we know how to multiply, and it’s easy to write an algorithm to multiply two numbers given as finite strings in their decimal representation. But suppose we want to multiply a number—say, 87—by some other number which is given to us as a non-terminating decimal. For instance,

$$\frac{12}{35} = .3428571428571428571428571 \dots$$

¹⁴Why a_{n+1} rather than simply a_n in the numerator? Well, to compute c_n in the denominator, we already need to compute a_n and b_n , and at that point, we have a_{n+1} almost for free. It turns out that using a_{n+1} , the algorithm yields a considerably better approximation, although simply using a_n is quite good also—either way yields a quadratically converging sequence of approximations.

Now this is really too simple, because we know that we can represent this number simply by the fraction $12/35$, and so if we wanted to multiply it by 87, say, we would just multiply 87 by $(12/35)$, and then turn the resulting fraction into a decimal if we wanted to.

But suppose we had a more complicated situation: suppose, for instance, we wanted to multiply 87 by $\sqrt{2}$. Now $\sqrt{2}$, being irrational, has a non-repeating decimal expansion:

$$\sqrt{2} = 1.414213562373009504880\dots$$

If we wanted to multiply this by 87, the process would be more complicated.

We’re going to come up below with an algorithm that will work for any such infinite decimal string provided only that the product does not end in an infinite sequence of 9’s. Such a product of course can be represented as a terminating decimal, and in particular is rational. So if we know a priori that the product is irrational, our algorithm will be guaranteed to work. And actually, our algorithm works just fine for terminating decimals as well—it’s only when ordinary multiplication leads to an infinite sequence of 9’s that there is a problem. We’ll explain all this at the end of this section.

To make things simple, let us first see what we need to do to multiply a *terminating* decimal by 87. Let us multiply .987436917 by 87. Normally, one does this “right-to-left”, and this makes sense, since carries propagate to the left. However, we want to perform the multiplication left-to-right because eventually we want to let the decimal be non-terminating—it will be a stream—and so we can only access it from left to right. Figure 16 shows how the multiplication works.

Now the crucial point is this: *we don’t have to wait until the very end to start producing numbers in the output stream*. In fact, in Figure 16, we have indicated in **red** the digits that we know have been “produced for sure” at each stage.

For instance, once we have added the first two products to get 8526, we actually know that the first digit of the output is 8. The reason we know this is that the rest of all the numbers that could possibly be added up are at most $.9999999 \times 87$, which is at most 87, and $87 + 8526 < 9000$, so no matter what comes next, the initial 8 will never get changed to a 9. In our pseudo-code below, this gets implemented as

$$87 + 526 < 1000$$

How do we get 526? It is “8526 mod 1000”. And where do we get 1000? It is the smallest power of 10 with the same number of digits as 8526. It is called **pow** in the pseudo-code in Figure 17.

Warning: it’s not *always* true that each time we add a new partial product in, we can take off the next digit of the new partial sum. We always have to perform the check, as above. If the check does not succeed, then we just go on and add in the next partial product.

Furthermore, exactly this same reasoning holds when the input is an infinite stream.

We can write a program for this in which we multiply a stream of digits (regarding that stream as the decimal representation of a number between 0 and 1) by a multiplier **m**. (In our example, **m** was 87.) In writing this program, we can regard this computation as dealing with three pieces of data:

- The input stream. (This was the finite stream (9 8 7 4 3 6 9 1 7) in our example.) We will call this **strm** in the pseudo-code below.
- The intermediate list being processed. We will call this **a-list** in the pseudo-code below. In our example, **a-list** starts out being the empty list, then is the list (7 8 3), then is the list (8 5 2 6), and then is the list (5 2 6) after we discover that the 8 will never change.

$$\begin{array}{r}
 9 \ 8 \ 7 \ 4 \ 3 \ 6 \ 9 \ 1 \ 7 \\
 8 \ 7 \\
 \hline
 87 \times 9 = \textcolor{blue}{7} \ 8 \ 3 \\
 87 \times 8 = 6 \ 9 \ 6 \\
 \textcolor{red}{8} \ \textcolor{blue}{5} \ 2 \ 6 \\
 87 \times 7 = 6 \ 0 \ 9 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{blue}{8} \ 6 \ 9 \\
 87 \times 4 = 3 \ 4 \ 8 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{red}{9} \ \textcolor{blue}{0} \ 3 \ 8 \\
 87 \times 3 = 2 \ 6 \ 1 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{red}{9} \ 0 \ \textcolor{blue}{6} \ 4 \ 1 \\
 87 \times 6 = 5 \ 2 \ 2 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{red}{9} \ 0 \ \textcolor{blue}{6} \ \textcolor{blue}{9} \ 3 \ 2 \\
 87 \times 9 = 7 \ 8 \ 3 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{red}{9} \ 0 \ \textcolor{red}{7} \ 0 \ \textcolor{red}{1} \ 0 \ 3 \\
 87 \times 1 = 8 \ 7 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{red}{9} \ 0 \ \textcolor{red}{7} \ 0 \ \textcolor{red}{1} \ \textcolor{blue}{1} \ 1 \ 7 \\
 87 \times 7 = 6 \ 0 \ 9 \\
 \textcolor{red}{8} \ \textcolor{red}{5} \ \textcolor{red}{9} \ 0 \ \textcolor{red}{7} \ 0 \ \textcolor{red}{1} \ 1 \ 7 \ 7 \ 9
 \end{array}$$

Figure 16: Multiplying 87 by .987436917. Each **consume** step adds in a partial product (87×9 , 87×8 , and so on). The partial sums (783, 8526, and so on) represent the situation after each **consume** step. The red digits are digits that are known to be fixed at that point and can be output. The blue digits are used below in proving that the algorithm can always make progress.

Both here and in the algorithm, we've left out figuring where the decimal point goes. That could be added easily; it's not important for the algorithm.

- The output stream.

There are also some auxiliary variables that we will find useful:

Let us denote by **a** the value of **a-list**. For instance

$$\begin{array}{ll}
 \mathbf{a-list} & \mathbf{a} \\
 (5 \ 0 \ 3) & 503 \\
 (0 \ 5 \ 0 \ 3) & 503
 \end{array}$$

In other words, **a** is insensitive to initial 0's in **a-list**. Both representations are useful.

Next, given a *non-empty* list **a-list**, we will define a number **pow** such that

- **pow** is a power of 10, and

- **pow** is the smallest number with the same number of digits as **a-list**.

For instance, if **a-list** is (0 0 5 0), then **pow** is 1000.

In terms of these data structures, our algorithm can be described as follows. It consists of two basic procedures:

consume This takes the next element off of **strm**, multiplies it by **m**, and adds it to the number represented by **a-list**.

produce This takes the left-most digit off of **a-list** and incorporates it in the output stream.

Figure 17 illustrates some pseudo-code for these mutually recursive procedures.

(If you have studied compilers, you may find these procedures vaguely reminiscent of the **shift** and **reduce** procedures in an LR parser. The intermediate data structure **a-list** here takes the place of the stack, except that our **a-list** is really more like a queue.)

The whole basis of this computation is the test in **action** which enables us to determine when it is safe to call **produce**—that is, when the most significant digit of **a-list** can never be changed by any future **consume**.

This depends on the fact that we know that “**m times strm**” can never be greater than **m**.

Please note that this is very sketchy pseudo-code. The stream and list operations are written using Scheme notation (**cons-stream**, **stream-cdr**, **cdr**). I have used C notation for expressions and flow of control, and **%** for the remainder operator. Also, I have not specified how parameters are passed to these various functions, or how the system is initialized. But the fundamental ideas are all here.

Figure 18 shows the first few steps of the algorithm.

8.1.1 Proof that we always make progress

One thing that always needs proving in an algorithm such as this is that progress can always be made. We know, for instance, that not every **consume** step is followed by a **produce** step—for instance in our example, 6 and 9 are consumed one after the other without any digit being produced in between. And it is possible even for a long string of digits to be consumed before another one is produced. But is it possible (assuming we are multiplying by an unending decimal string) for production of digits to suddenly stop, and all the rest of the steps to be **consume** steps, without anything further being produced?

The answer is yes, it is possible. For instance, if we multiply $1/9$, which is $.11111111\dots$, by 9, we can see that actually no digit will ever be produced. And in fact, the product of these two numbers is the repeating decimal $.99999999\dots$, which is 1. This is really the only kind of thing that can go wrong, as we can now see:

At the beginning of each step (i.e., at the beginning of the code for the **action** procedure), if **m** has at least as many digits as **a-list** does, then either

- **a-list** is empty, or
- $m + (a \% \text{pow}) \geq \text{pow}$.

```

action:

  If strm is empty, just return the finite stream consisting of the digits
  in a-list.

  Otherwise:

    If ((a-list has at least 1 digit) AND
        (m + (a % pow) < pow))
    {

      ;; Adding "m times strm" to a can't possibly affect the leading
      ;; element of a-list, so we are free to move the leading element of
      ;; a-list to the output stream.

      produce
    }
    else
    {
      consume
    }

produce:

  ;; move the first element of a-list to the return stream
  ;; and delete it from a-list.

  Return what you get when you cons-stream the first element of a-list onto what
  you get when you call action recursively after

    replacing a by (a % pow)
    replacing a-list by (cdr a-list)
    ;; could get leading 0's here

consume:

  return the result of calling action recursively after

    replacing a by (10*a) + (m*(the first element of strm))
    replacing a-list by a represented as a string with 0's prepended if needed
    to make its length > the length of the original a-list.
    replacing strm by (stream-cdr strm)

```

Figure 17: Pseudo-code for stream multiplication

and hence the next step will be a **consume** step. This will increase the length of **a-list**, and so eventually **a-list** will have more digits than **m** does.

In Figure 16, after each **consume** step (which corresponds to adding in one of the partial products), there may be one or more **produce** steps, which append to the output stream of red digits. After those **produce** steps are finished, we have colored blue the digits remaining in **a** that are “to the left” of **m**—that is that are to the left of the greatest digit in **m**. In one case in this example, there is no blue digit. In most of the cases, there is 1 blue digit, and in one case there are two.

A **consume** step always increases the number of blue digits, while a **produce** step turns the left-most blue digit red.

	Output Stream	a	a-list	pow	Input Stream
	()	0	()	—	(9 8 7 4 3 6 9 1 7)
a-list is empty; consume					
	()	783	(7 8 3)	100	(8 7 4 3 6 9 1 7)
$87 + 83 \geq 100$; consume					
	()	8526	(8 5 2 6)	1000	(7 4 3 6 9 1 7)
$87 + 526 < 1000$; produce					
	(8)	526	(5 2 6)	100	(7 4 3 6 9 1 7)
$87 + 26 \geq 100$; consume					
	(8)	5869	(5 8 6 9)	1000	(4 3 6 9 1 7)
$87 + 869 < 1000$; produce					
	(85)	869	(8 6 9)	100	(4 3 6 9 1 7)

Figure 18: The first few steps of the multiplication algorithm.

The crucial point then is this: each **consume** step adds a number to **a**, but the digits of that number do not lie below any of the blue digits in **a**. (That’s because the number that is added is at most **m**, or equivalently, at most $10 * \mathbf{m}$ is added to $10 * \mathbf{a}$.)

Therefore, each blue digit can only change if a carry is propagated into it. If there is more than 1 blue digit, all except the left-most one must be 9. This is because when we add two decimal numbers, no matter how long, the “carry” at each digit is always at most 1. Therefore if any blue digit (other than the left-most one) was less than 9, it would still be ≤ 9 after anything was carried into it, and therefore no carries would propagate to its left, and all the blue digits to its left would be ready to be produced. That is, in such a case we are guaranteed that we can produce something.

Thus, the only way we could stop making progress is if after some point, all the new blue digits created were 9. This would stop any production of red digits. And the only way this could happen is if the product that we are computing ended in a repeating chain of 9’s.

So as long as we are not computing a product which ends in an infinite sequence of 9’s, we are guaranteed that there will be an infinite number of **produce** steps, and with that stipulation, that concludes the proof that the algorithm always makes progress.

8.2 A streaming algorithm that computes the digits of π

We will base our algorithm on the series

$$\pi = 2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 + \frac{3}{7} \left(\cdots \left(2 + \frac{i}{2i+1} \left(\cdots \right) \right) \right) \right) \right)$$

The way you understand this computation is like this: looking at footnote 8 on page 26, it is really a sequence of approximations to π :

$$\begin{aligned}
2 + \frac{1}{3} \left(2 \right) &= 2.6666666666 \dots \\
2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 \right) \right) &= 2.9333333333 \dots \\
2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 + \frac{3}{7} \left(2 \right) \right) \right) &= 3.04761904761 \dots \\
2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 + \frac{3}{7} \left(2 + \frac{4}{9} \left(2 \right) \right) \right) \right) &= 3.09841269841 \dots \\
2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 + \frac{3}{7} \left(2 + \frac{4}{9} \left(2 + \frac{5}{11} \left(2 \right) \right) \right) \right) \right) &= 3.12150072150 \dots \\
2 + \frac{1}{3} \left(2 + \frac{2}{5} \left(2 + \frac{3}{7} \left(2 + \frac{4}{9} \left(2 + \frac{5}{11} \left(2 + \frac{6}{13} \left(2 \right) \right) \right) \right) \right) \right) &= 3.13215673215 \dots \\
&\dots
\end{aligned}$$

Actually (and we'll come back to this later), we didn't have to apply all these computations to the inner value 2—we could have started with *any* value, and the sequence of approximations would still converge to π .

8.2.1 Linear fractional transformations

There is a clever way we can rewrite this formula which will enable us to perform these computations efficiently: We can think of the first computation above as $f_1(2)$, where f_1 is the function given by

$$f_1(x) = 2 + \frac{1}{3}x = \frac{1 \cdot x + 6}{3}$$

The second computation can be regarded as $f_1(f_2(2))$, where f_1 is as just defined, and f_2 is the function given by

$$f_2(x) = \frac{2 \cdot x + 10}{5}$$

and we see that we can write the n^{th} computation as

$$f_1 \circ f_2 \circ \dots \circ f_n(2)$$

where for any k ,

$$f_k(x) = \frac{k \cdot x + (4k + 2)}{2k + 1}$$

This kind of function is what is known as a *linear fractional transformation*. The most general linear fractional transformation is a function of the form

$$f(x) = \frac{ax + b}{cx + d}$$

One of the really nice properties of linear fractional transformations is that it is easy to form their composition, because they correspond naturally to matrices: If we associate the linear fractional transformation

$$f(x) = \frac{a_{11}x + a_{12}}{a_{21}x + a_{22}}$$

with the matrix

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

and if we similarly associate the linear fractional transformation

$$g(x) = \frac{b_{11}x + b_{12}}{b_{21}x + b_{22}}$$

with the matrix

$$B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}$$

then it turns out (and it's easy to prove; if you've never done it before, you should do it now) that $f \circ g$ is the linear fractional transformation that is associated in the same way with the matrix AB (using ordinary matrix multiplication). In particular, the set of linear fractional transformations is closed under composition.

Note that a number of common functions are really linear fractional transformations. For instance:

- The function $f(x) = c$ (i.e., a constant) is just the linear fractional transformation that corresponds to the matrix

$$\begin{pmatrix} 0 & c \\ 0 & 1 \end{pmatrix}$$

(This is generally considered a “degenerate” transformation because it is not invertible.)

- The function $f(x) = cx$ (where c is a constant) is just the linear fractional transformation that corresponds to the matrix

$$\begin{pmatrix} c & 0 \\ 0 & 1 \end{pmatrix}$$

- The function $f(x) = ax + b$ is just the linear fractional transformation that corresponds to the matrix

$$\begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}$$

8.2.2 The algorithm

So now we are ready to construct an algorithm to compute π based on this sequence of products of linear fractional transformations. Remember that we are computing

$$\lim_{n \rightarrow \infty} f_1 \circ f_2 \circ f_3 \circ \cdots \circ f_n(2)$$

where

$$f_k(x) = \frac{k \cdot x + (4k + 2)}{2k + 1}$$

which is the linear fractional transformation corresponding to the matrix

$$\begin{pmatrix} k & 4k + 2 \\ 0 & 2k + 1 \end{pmatrix}$$

We will structure this computation in a manner which is similar to that of the previous problem. That is, we will have

- An input stream **strm**. This is just the stream f_1, f_2, f_3, \dots of linear fractional transformations, each one represented as a matrix.
- An intermediate data structure called **a**. **a** will itself be a linear fractional transformation.
- An output stream consisting of the digits of π .

Each step of the computation is a (recursive) call to the procedure **action**, which in turn calls one of two procedures:

consume: Multiply **a** on the right by the first element of **strm** and make that the new **a**. Replace **strm** by the stream you get when you remove its first element. (Of course, I really mean by this that you make a recursive call to **action**, passing new values of **a** and **strm** as defined here.)

produce: Extract the next decimal digit of π from **a** (and incorporate that digit in the output stream), and multiply **a** on the left by a linear fractional transformation that reflects the fact that it has “lost that value”.

(Actually, these procedures are so simple that they can just as well be inlined. But it’s convenient to think of them this way.)

There are two open questions here, both having to do with **produce**:

- How do we know if we are allowed to take a digit from **a**—and how do we get that digit in any case?
- How do we modify **a** when we do extract a digit from it?

Fortunately, we’re in good luck here because these transformations are really pretty simple. Here are some important properties of this sequence of transformations—the first three are all easy to prove:

- Each fractional linear transformation in **strm** is a monotonic increasing map, and the same is true for any finite product of these fractional linear transformations.
- Each fractional linear transformation in **strm** maps the interval $[3, 4]$ into a sub-interval of $[3, 4]$. And again, the same is true for any finite product of those fractional linear transformations.
- Each fractional linear transformation in **strm** maps any interval into an interval of less than half its length.
- As I mentioned before, in the computations above, we didn't have to apply all the functions to the number 2—we could have used any fixed number. (This is not obvious, but it's true because of the particular properties of the linear fractional transformations in **strm**. We'll prove it in the next section.) It will be convenient for us to use 3 instead of 2.

Now we know that π lies in the interval $[3, 4]$. So we start out by computing products of the first elements of **strm**. We do this until we find a product transformation which maps 3 and 4 into numbers that have the same integral part (that is, into numbers that have the same **floor**).

Actually, this happens immediately, by what we have said. So now we can take that value (which must be 3) and make it the first element of the output stream. Then we have to adjust **a** to reflect the fact that we have “taken away a 3”. We do this by multiplying **a** on the left by the matrix

$$\begin{pmatrix} 10 & -10 \cdot 3 \\ 0 & 1 \end{pmatrix}$$

This just corresponds to multiplying by 10 and subtracting 30—in other words, we shift the decimal expansion to the left 1 digit and delete the left-most digit, which is 3.

Now we proceed. We **consume** successive elements of **strm** until again **a** applied to 3 and **a** applied to 4 yield numbers having the same floor. Note that now the floor is not guaranteed to be 3, because **a** has been modified by multiplying it on the left, as above. Say the floor is n . We then **produce**, passing n on to the output stream, and multiplying **a** on the left by the matrix

$$\begin{pmatrix} 10 & -10 \cdot n \\ 0 & 1 \end{pmatrix}$$

And that's pretty much it. By the way, when I say “apply **a** to 3”, what I mean is “regard **a** as a linear fractional transformation and apply it to 3”.

8.2.3 Proof that we always make progress

How do we know that in the algorithm **a** applied to 3 and **a** applied to 4 will always eventually have the same floor? We'll prove that now, and along with it, we'll prove that numbers we get from this algorithm converge to π no matter what number we start with.

The argument depends on two things: 1) the “contraction” property of these linear fractional transformations in **strm**, and 2) the fact that π is irrational, or even more simply that it does not have a terminating decimal expansion.

First, it is convenient to establish some notation:

- Let $T_n = f_1 \circ f_2 \circ \dots \circ f_n$. So $T_{n+1} = T_n \circ f_{n+1}$.
- Let s_n denote the n^{th} result which we get from the computation starting from 3, and t_n be the n^{th} result starting from 4. That is,

$$s_n = T_n(3)$$

$$t_n = T_n(4)$$

Now it is easy to see that if $x \leq 3$, then $f_1(x) \geq x$ and for $k > 1$, $f_k(x) > x$. And similarly, if $y \geq 4$, then for all $k \geq 1$, $f_k(y) < y$. And since each transformation f_k as well as each product of these transformations, is a strictly increasing function, we have

$$s_{n+1} = T_{n+1}(3) = T_n(f_{n+1}(3)) \geq T_n(3) = s_n$$

$$t_{n+1} = T_{n+1}(4) = T_n(f_{n+1}(4)) \leq T_n(4) = t_n$$

In fact, let's be a little more general: let s_0 be any number ≤ 3 and t_0 be any number ≥ 4 , and set

$$s_n = T_n(s_0)$$

$$t_n = T_n(t_0)$$

Then just as above, we see that s_n increases with n , and t_n decreases with n . And further, as we have seen, $|t_{n+1} - s_{n+1}| \leq \frac{1}{2} |t_n - s_n|$. So certainly s_n and t_n converge to the same number. Further, if we start with $s_0 \leq s'_0 \leq t'_0 \leq t_0$ and produce sequences $\{s_n\}$, $\{s'_n\}$, $\{t'_n\}$, and $\{t_n\}$, then s'_n and t'_n get squeezed in between s_n and t_n and therefore must have the same limit. And this is even true if s_0 and/or t_0 is between 3 and 4. This verifies the claim made earlier: it doesn't matter what number we start with—the algorithm will converge to a unique value.

And actually, we know the value: we have already seen that starting with $s_0 = 2$ we get a limit of π (this was after all, just Euler's transformation applied to Leibniz's formula)

So our algorithm is guaranteed to produce sequences s_n and t_n approaching π from below and above, respectively.

That's almost the whole story. We just have to ensure that the decimal digits in s_n and t_n keep agreeing on longer and longer intervals.

Now the only way this might not happen is if from some point on, s_n and t_n looked something like this:

$$s_n = x_1.x_2x_3x_4x_5 \dots x_p799999993524$$

$$t_n = x_1.x_2x_3x_4x_5 \dots x_p800010325796$$

(Of course, we know that $x_1 = 3$, $x_2 = 1$, and so on. But that's not important for this argument.)

In other words, all the digits up to x_p agree, the next digit differs by 1, and s_n continues with a series of 9's, while t_n continues with some 0's. And further, since s_n and t_n have to get closer and closer together, the sequence of 9's and the sequence of 0's would have to get longer and longer as n became greater.

But if this really was true, then both s_n and t_n would approach the number

$$x_1.x_2x_3x_4x_5\dots x_p8$$

which is certainly a rational number. And we know that π is irrational. So this situation cannot occur, and while we may have to wait awhile (in case something like this does happen for a number of iterations), we will always make progress, just as we did in the multiplication algorithm. The algorithm will continue to generate successive digits of π as long as it is kept running.

And further, assuming that the decimal expansion of π is not overwhelmed by long sequences of 9's¹⁵, we can be assured also that our algorithm will spit out a digit roughly every 3 or 4 iterations, since we are doubling the accuracy with every iteration.

9 Algorithmic efficiency

The Brent-Salamin algorithm, as we have seen, converges quadratically—each successive iteration roughly doubles the number of significant digits. All the other algorithms we have considered converge linearly—the number of significant digits is proportional to the number of iterations.

Of course, these are somewhat crude measures. For instance, we really should include in our estimates the cost of an individual iteration. If the cost goes up rapidly enough, then a quadratically converging algorithm might be swamped by the expense of the later iterations.

In our case, this does not happen—the arithmetic-geometric mean based algorithms are in fact extremely efficient.

Really however, the best way to look at the efficiency of these algorithms is not in terms of the “operational cost”, (i.e., the number of operations needed to compute each successive digit), since the cost of each operation may vary, and in fact typically increases as the algorithm progresses. The real measure of cost is the “cost per bit”, or equivalently, the number of bit operations needed to compute each successive digit. And there is a trivial lower bound here: it costs a small fixed amount simply to output each digit, so the cost to output n digits cannot be smaller than $O(n)$.

Now this may seem like a great disappointment when compared with the quadratic convergence that we have seen in the Brent-Salamin algorithm. Nevertheless, it is really not bad at all. And in fact, in terms of cost per bit, the Brent-Salamin algorithm is an $n \log n$ algorithm, and this is the kind of algorithm we computer scientists have learned to love.

On the other hand, we can look more closely at the other algorithms and see if we have really got the best estimate for their efficiency. In particular, let us look at the streaming algorithm of the last section, which although it is nowhere near as efficient as the Brent-Salamin algorithm, does have some inherent interest of its own.

There is no doubt that just in terms of operational cost, it is a linear algorithm—we know that the interval getting squeezed contracts by a factor of about 1/2 each iteration, so between 3 and 4 iterations are needed for each additional decimal digit to be emitted.

But it's actually more expensive than this would imply, because each iteration deals with larger and larger numbers, which cost more and more to multiply, and so the cost of each iteration actually increases. In fact, at each iteration we multiply by a 2×2 matrix whose greatest element is

¹⁵If π is normal, there *will* be arbitrarily long sequences of 9's in its decimal expansion. But they will not “overwhelm” the sequence of digits—they will be quite rare.

$O((Cn)!)$, for some C . The number of digits in such a number is $O(\log(Cn)!) = O(n \log n)$. So the cost of the n^{th} iteration is some constant times the cost of multiplying a number with n digits by a number with $n \log n$ digits, and using the fastest way of multiplying integers currently known, this is $O(n \log n \log \log n)$. So the cost certainly increases with n .

It then follows that the cost of computing the first n digits is $O(n^2 \log n \log \log n)$. And as we've noted, this compares pretty unfavorably with the cost of computing the first n digits using the Brent-Salamin algorithm—i.e., $O(n \log n)$, which is strikingly better. On the other hand, this is a streaming algorithm, which is really very nice. So we might ask the next question: How good can a streaming algorithm be *in principle*? There are two versions of this question that come naturally to mind, and we consider them in the next two sections.

We may as well assume that we have a streaming algorithm in which each iteration consists of a fixed (or at any rate, bounded) number of computations and which takes a bounded number of iterations to produce each successive digit (at least on average).¹⁶ Then the cost of each digit is pretty much the same as the cost of each iteration.

9.1 Can the cost of each iteration be bounded?

No, it can't. The reason for this is quite simple: if the cost of each iteration was bounded, then in effect we would have a finite state machine: each iteration consists of a finite sequence of operations, and if the cost of each iteration was bounded, the numbers being operated on would have to have a bounded number of digits, so there could be only a finite number of states that the algorithm could be in. It follows that the stream of digits that would be generated would have to be ultimately periodic, which could only be true if π was rational, which it is not.

So certainly the cost per iteration cannot be bounded. On the other hand, we really don't know how slowly the cost has to increase. Perhaps, for instance, there is a streaming algorithm for the digits of π in which the cost of each iteration is $O(n)$ or even $O(\log n)$. I don't know if this is possible or not. And the argument I gave here is quite crude: I only used the fact that π was irrational. So this argument applies to any irrational number, and π is of course a very special one.

There is another interesting question that comes up in this connection: if π really is a normal number, then its digits could be used as a good pseudo-random number sequence. Now the most commonly used algorithms for pseudo-random number generation provide each successive digit at a fixed incremental cost. As we've just seen, this could only be the case if the generated sequence was ultimately periodic, and in fact, these pseudo-random number sequences *are* periodic—usually with extremely long periods. Since the modern algorithms for computing π are so extraordinarily efficient, however, there has been some interest in seeing if the sequence of digits of π could be used instead; we know a priori that this sequence, at least, is not periodic.

9.2 Can the cost of each iteration be bounded *on average*?

I don't know the situation for π , but I do know that for some irrational numbers, the average cost can indeed be bounded. Here are some examples:

¹⁶For what we're doing here, it's not worth reasoning about this in any greater generality.

Liouville numbers Consider the number

$$L = \sum_{n=1}^{\infty} 2^{-n!}$$

The cost of computing the first $n!$ binary digits of this number is a (small) fixed cost for each digit plus the cost of computing “where all the ones go”, i.e., the cost of computing all the factorials up to $n!$, which can be very easily but crudely bounded by $n^2 \log^2 n$. So the cost is $O(n! + n^2 \log^2 n) = O(n!)$, and so the average cost of computing each of the first $n!$ digits is $O(n!)/n! = O(1)$.

We called this number L because it is a *Liouville* number, and perhaps you are thinking that this is a really exceptional case. For instance, it’s nowhere near a normal number—in fact, the fraction of 1’s in the first n digits tends to 0 as n becomes large.

The Thue-Morse number Now of course any number for which we could get a result like this must have *some* special property. But still, we can find numbers which look a lot more “random” and still have an average fixed cost for computing each digit.

One such number is the Thue-Morse number (let us call it T), whose binary expansion is

$$T = .0110100110010110 \dots$$

The simplest definition of this number is that its n^{th} digit is the parity (i.e., the number of ones mod 2) in the binary expansion of n . (Actually, this number isn’t normal either¹⁷, but it certainly appears to be a lot more regular than L .)

Now to compute the n^{th} digit of T , we need to compute the parity of n . To compute this from scratch, we need to know all the digits of n , and there are $\log n$ of them. So just to do this naively would require a cost of

$$\sum_{j=1}^n \log j = \log n! \approx n \log n$$

for the first n digits.

However, we can do better. Let us use two data structures, which hold at the end of iteration $j - 1$

- The last parity computed, i.e., the parity of the binary representation of $j - 1$.
- The binary representation of $j - 1$.

At the j^{th} step we compute the binary representation of j . As we do this, we keep track of whether or not the parity of j changed from that of $j - 1$, and we use that to decide whether or not to toggle the first bit of data.

Now how expensive is this? It’s less expensive than you might think. For instance, suppose that the binary representation of $j - 1$ is

$$j - 1 = 1101000110111$$

¹⁷For instance, as Thue already showed around 1906, it doesn’t contain the strings 000 or 111.

Then to compute the binary representation of j , we just add 1 in the usual way, starting from the right. And it's easy to see that to do this, we just have to replace the rightmost string of the form 01^* by a string 10^* of equal length¹⁸. And in doing this, we also know whether the parity changed or not. In the example here, the parity does not change, because 1000 has the same parity as 0111.

The reason that this is cheap on the average is that most of the time we don't have to inspect every digit in $j - 1$. In fact, half the numbers end in 0 to begin with, so we only have to inspect 1 digit. Half of the rest end in 01, so we only have to inspect 2 digits. And we see that in general, the average number of digits we need to inspect is no more than

$$\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \cdots = \sum_{j=0}^{\infty} \frac{1}{2^j} \cdot j$$

and it's not hard to see that this entire series converges. That is, the average cost per iteration is bounded, even though the individual cost of each iteration is not—there will be some values of j for which we need to inspect all the digits in their binary representation. But there won't be many of them.

Champernowne's number Finally, let us consider Champernowne's number

$$C = .012345678910111213141516 \dots$$

where the decimal expansion is produced by concatenating the decimal representations of the integers starting from 0. This number is known to be normal (at least in base 10).

How easy is it to compute? Well, if we have used all the integers up to n , then we have generated roughly $\sum_{j=1}^n \log j = \log n! \sim n \log n$ decimal places. Further, the average incremental cost of computing each successive number is bounded (by the same argument as above). Hence (roughly speaking) the cost of computing $n \log n$ digits is $O(n \log n + n) = O(n \log n)$, so again the average cost of computing each successive digit is bounded.

Of course, none of this tells us what the cost is for π . No one at this point knows enough about any special properties of π that would help us answer this kind of question.

References

- Akhiezer, N. I. (Naum Il'ich). 1990. *Elements of the Theory of Elliptic Functions*, Translation of Mathematical Monographs, Vol. 79, American Mathematical Society, Providence, RI. [QA343.A3813].
- Almkvist, Gert and Bruce Berndt. 1988. *Gauss, Landen, Ramanujan, the Arithmetic-Geometric Mean, Ellipses, pi, and the Ladies Diary*, American Mathematical Monthly **95**, no. 7, 585–608.
- Bailey, David H. 1988. *The Computation of π to 29,360,000 Decimal Digits Using Borweins' Quartically Convergent Algorithm*, Mathematics of Computation **50**, no. 181, 283–296.
- Berggren, Lennart, Peter B. Borwein, and Jonathan M. Borwein. 2000. *Pi: a Source Book*, 2nd edition, Springer-Verlag, New York. [QA484.P5 2000].
- Borwein, Jonathan M., Peter B. Borwein, and David H. Bailey. 1989. *Ramanujan, Modular Equations, and Approximations to Pi or How to Compute One Billion Digits of Pi*, American Mathematical Monthly **96**, no. 3, 201–219. Reprinted in Berggren et al., *Pi, a Source Book*.

¹⁸I'm using regular expression notation here, of course.

- Borwein, Jonathan M. and Peter B. Borwein. 1984. *The Arithmetic-Geometric Mean and Fast Computation of Elementary Functions*, SIAM Review **26**, no. 3, 351–366. Reprinted in Breggren et al., *Pi, a Source Book*.
- . 1987. *Pi and the AGM: a study in analytic number theory and computational complexity*, Wiley, New York. [QA241.B774 1987].
- Brent, Richard P. 1976. *Fast Multiple-Precision Evaluation of Elementary Functions*, Journal of the Association of Computing Machinery **23**, no. 2, 242–251. Reprinted in Breggren et al., *Pi, a Source Book*.
- Cayley, Arthur. 1895. *An Elementary Treatise on Elliptic Functions*, 2nd edition, Cambridge University Press, Cambridge. [QA343.C38].
- Cohen, Henri, Fernando Rodrigues Villegas, and Don Zagier. 2000. *Convergence Acceleration of Alternating Series*, Experimental Mathematics **9**, no. 1, 3–12.
- Cox, David A. 1984. *The Arithmetic-Geometric Mean of Gauss*, L'Enseignement Mathématique **30**, 275–330. Reprinted in Breggren et al., *Pi, a Source Book*.
- Gauss, Carl Friedrich. 1800. *De Origine Proprietatibusque Generalibus Numerorum Mediorum Arithm. Geometricorum*. Manuscript, unpublished during Gauss's lifetime. Werke, Vol. 3, 361–371.
- . 1863. *Carl Friedrich Gauss Werke*, Königlichen Gesellschaft der Wissenschaften zu Göttingen, Göttingen. [QA3.G3].
- Gibbons, Jeremy. 2006. *Unbounded Spigot Algorithms for the Digits of Pi*, American Mathematical Monthly **113**, no. 4, 318–328.
- Hilbert, David and S. Cohn-Vossen. 1999. *Geometry and the Imagination*, 2nd edition, American Mathematical Society, Providence, Rhode Island. Translated from the German *Anschauliche Geometrie* by P. Nemenyi and originally published in English by Chelsea Publishing. [QA685.H515 1999].
- Lawden, Derek F. 1989. *Elliptic Functions and Applications*, Springer-Verlag, New York-Berlin-Heidelberg. [QA343.L39].
- Newman, Donald J. 1985. *A Simplified Version of the Fast Algorithms of Brent and Salamin*, Mathematics of Computation **44**, no. 169, 207–210. Reprinted in Breggren et al., *Pi, a Source Book*.
- Roy, Ranjan. 1990. *The Discovery of the Series Formula for π by Leibniz, Gregory and Nilakantha*, Mathematics Magazine **63**, 291–306. Reprinted in Breggren et al., *Pi, a Source Book*.
- Salamin, Eugene. 1976. *Computation of π Using Arithmetic-Geometric Mean*, Mathematics of Computation **30**, no. 135, 565–570. Reprinted in Breggren et al., *Pi, a Source Book*.
- Wagon, Stan. 1985. *Is π Normal?*, The Mathematical Intelligencer **7**, no. 3, 65–67. Reprinted in Breggren et al., *Pi, a Source Book*.
- Whittaker, E. T. and G. N. Watson. 1927. *A Course of Modern Analysis*, 4th ed., Cambridge University Press, Cambridge (England). [QA401.W624].