

Street-fighting Mathematics

Sanjoy Mahajan
MIT

Copyright 2008 Sanjoy Mahajan

Contents

1	Dimensions	3
2	Extreme cases	13
3	Discretization	31
4	Picture proofs	45
5	Taking out the big part	57
6	Analogy	80
7	Operators	86
	References	91

1

Dimensions

Dimensions, often called units, are familiar creatures in physics and engineering. They are also helpful in mathematics, as I hope to show you with examples from differentiation, integration, and differential equations.

1.1 Free fall

Dimensions are often neglected in mathematics. Calculus textbooks state many problems in this form:

A ball falls from a height of h **feet**. Neglecting air resistance, estimate its speed when it hits the ground, given a gravitational acceleration of g **feet per second squared**.

The units, highlighted with boldface type, have been separated from g or h , making g and h pure numbers. That artificial purity ties one hand behind your back, and to find the speed you are almost forced to solve this differential equation:

$$\frac{d^2y}{dt^2} = -g, \text{ with } y(0) = h \text{ and } \dot{y}(0) = 0,$$

where $y(t)$ is the ball's height at time t , $\dot{y}(t)$ is its velocity, and g is the strength of gravity (an acceleration). This second-order differential equation has the following solution, as you can check by differentiation:

$$\begin{aligned}\dot{y}(t) &= -gt, \\ y(t) &= -\frac{1}{2}gt^2 + h.\end{aligned}$$

The ball hits the ground when $y(t) = 0$, which happens when $t_0 = \sqrt{2h/g}$. The speed after that time is $\dot{y}(t) = -gt_0 = -\sqrt{2gh}$. This derivation has many

Dimensions

4

spots to make algebra mistakes: for example, not taking the square root when solving for t_0 , or dividing rather than multiplying by g when finding the speed.

Here's the same problem written so that dimensions help you:

A ball falls from a height h . Neglecting air resistance, estimate its speed when it hits the ground, given a gravitational acceleration of g .

In this statement of the problem, the dimensions of h and g belong to the quantities themselves. The reunion helps you guess the final speed, without solving differential equations. The dimensions of h are now length or L for short. The dimensions of g are length per time squared or LT^{-2} ; and the dimensions of speed are LT^{-1} . The only combination of g and h with the dimensions of speed is

$$\sqrt{gh} \times \text{dimensionless constant.}$$

An estimate for the speed is therefore

$$v \sim \sqrt{gh},$$

where the \sim means 'equal except perhaps for a dimensionless constant'. Besides the minus sign (which you can guess) and the dimensionless factor $\sqrt{2}$, the dimensions method gives the same answer as does solving the differential equation – and more quickly, with fewer places to make algebra mistakes. The moral is:

Do not rob a quantity of its intrinsic dimensions.

Its dimensions can guide you to correct answers or can help you check proposed answers.

1.2 Integration

If ignoring known dimensions, as in the first statement of the free-fall problem, hinders you in solving problems, the opposite policy – specifying unknown dimensions – can aid you in solving problems. You may know this Gaussian integral:

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

What is the value of

1.2 Integration

5

$$\int_{-\infty}^{\infty} e^{-\alpha x^2} dx,$$

where α is a constant? The integration variable is x so after you evaluate the integral over the limits, the x disappears; but α remains. The result contains only α and maybe dimensionless numbers, so α is the only quantity in the result that could have dimensions. For dimensional analysis to have a prayer of helping, α needs dimensions. Otherwise you cannot say whether, for example, the result should contain α or contain α^2 ; both choices have identical dimensions. Guessing the answer happens in three steps: (1) specifying the dimensions of α , (2) finding the dimensions of the result, and (3) using α to make a quantity with the dimensions of the result.

In the first step, finding the dimensions of α , it is more intuitive to specify the dimensions of the integration variable x , and let that specification decide the dimensions of α . Pretend that x is a length, as its name suggests. Its dimensions and the exponent $-\alpha x^2$ together determine the dimensions of α . An exponent, such as the 7 in 2^7 , says how many times to multiply a quantity by itself. The notion ‘how many times’ is a pure number; the number might be negative or fractional or both, but it is a pure number:

An exponent must be dimensionless.

Therefore αx^2 is dimensionless, and the dimensions of α are L^{-2} . A convenient shorthand for those words is

$$[\alpha] = L^{-2},$$

where [quantity] stands for the dimensions of the quantity.

The second step is to find the dimensions of the result. The left and right sides of an equality have the same dimensions, so the dimensions of the result are the dimensions of the integral itself:

$$\int_{-\infty}^{\infty} e^{-\alpha x^2} dx.$$

What are the dimensions of an integral? An integral sign is an elongated ‘S’, standing for *Summe*, the German word for sum. The main principle of dimensions is:

Dimensions

6

You cannot add apples to oranges.

Two consequences are that every term in a sum has identical dimensions and that the dimensions of a sum are the dimensions of any term. Similarly, given the kinship of summation and integration, the dimensions of the integral are the dimensions of $e^{-\alpha x^2} dx$. The exponential, despite the fierce-looking exponent of $-\alpha x^2$, is just the pure number e multiplied by itself several times. Since e has no dimensions, e^{anything} has no dimensions. So the exponential factor contributes no dimensions to the integral. However, the dx might contribute dimensions. How do you know the dimensions of dx ? If you read d as ‘a little bit of’, then dx becomes ‘a little bit of x ’. A little bit of length is still a length. More generally:

dx has the same dimensions as x .

The product of the exponential and dx therefore has dimensions of length, as does the integral – because summation and its cousin, integration, cannot change dimensions.

The third step is to use α to construct a quantity with the dimensions of the result, which is a length. The only way to make a length is $\alpha^{-1/2}$, plus perhaps the usual dimensionless constant. So

$$\int_{-\infty}^{\infty} e^{-\alpha x^2} dx \sim \frac{1}{\sqrt{\alpha}}.$$

The twiddle \sim means ‘equal except perhaps for a dimensionless constant’. The missing constant is determined by setting $\alpha = 1$ and reproducing the original integral:

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

Setting $\alpha = 1$ is a cheap trick. Several paragraphs preceding exhorted you not to ignore the dimensions of quantities; other paragraphs were devoted to deducing that α had dimensions of L^{-2} ; and now we pretend that α , like 1, is dimensionless?! But the cheap trick is useful. It tells you that the missing dimensionless constant is $\sqrt{\pi}$, so

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

1.3 Taylor and MacLaurin series

7

The moral of the preceding example is:

Assign dimensions to quantities with unspecified dimensions.

In this example, by assigning dimensions to x and α , we got enough information to guess the integral.

1.3 Taylor and MacLaurin series

The preceding example applied dimensions to integrals. Dimensions also help you remember Taylor series, a result based on derivatives. The idea of Taylor series is that if you know a function and all its derivatives at one point, you can approximate the function at other points. As an example, take $f(x) = \sqrt{x}$. You can use Taylor series to approximate $\sqrt{10}$ by knowing $f(9)$ and all the derivatives $f'(9)$, $f''(9)$, \dots .

The MacLaurin series, a special case of Taylor series when you know $f(0)$, $f'(0)$, \dots , looks like:

$$f(x) = f(0) + \text{stuff}$$

What is the missing stuff? The first principle of dimensions can help, that you cannot add apples to oranges, so all terms in a sum have identical dimensions. The first term is the zeroth derivative $f(0)$. The first term hidden in the ‘stuff’ involves the first derivative $f'(0)$, and this new term must have the same dimensions as $f(0)$. To draw a conclusion from this sameness requires understanding how differentiation affects dimensions.

In the more familiar notation using differentials,

$$f'(x) = \frac{df}{dx}.$$

So the derivative is a quotient of df and dx . You can never – well, with apologies to Gilbert & Sullivan, hardly ever – go astray if you read d as ‘a little bit of’. So df means ‘a little bit of f ’, dx means ‘a little bit of x ’, and

$$f'(x) = \frac{df}{dx} = \frac{\text{a little bit of } f}{\text{a little bit of } x}.$$

Using the [quantity] notation to stand for the dimensions of the quantity, the dimensions of $f'(x)$ are:

Dimensions

8

$$[f'(x)] = \frac{[\text{a little bit of } f]}{[\text{a little bit of } x]}.$$

Since a little bit of a quantity has the same dimensions as the quantity itself,

$$[f'(x)] = \frac{[\text{a little bit of } f]}{[\text{a little bit of } x]} = \frac{[f]}{[x]}.$$

Differentiating with respect to x is, for the purposes of dimensional analysis, equivalent to dividing by x .

So $f'(x)$ has the same dimensions as f/x .

This strange conclusion is worth testing with a familiar example. Take distance x as the function to differentiate, and time as the independent variable. The derivative of $x(t)$ is $\dot{x}(t) = dx/dt$. [Where did the prime go, as in $x'(t)$? When the independent variable is time, a dot instead of a prime is used to indicate differentiation.] Are the dimensions of $\dot{x}(t)$ the same as the dimensions of x/t ? The derivative $\dot{x}(t)$ is velocity, which has dimensions of length per time or LT^{-1} . The quotient x/t also has dimensions of length per time. So this example supports the highlighted conclusion.

The conclusion constrains the missing terms in the MacLaurin series. The first missing term involves $f'(0)$, and the term must have the same dimensions as $f(0)$. It doesn't matter what dimensions you give to $f(x)$; the principle of not adding apples to oranges applies whatever the dimensions of $f(x)$. Since its dimensions do not matter, choose a convenient one, that $f(x)$ is a volume. Do not, however, let x remain unclothed with dimensions. If you leave it bare, dimensions cannot help you guess the form of the MacLaurin series: If x is dimensionless, then differentiating with respect to x does not change the dimensions of the derivatives. Instead, pick convenient dimensions for x ; it does not matter which dimensions, so long as x has some dimensions. Since the symbol x often represents a length, imagine that this x is also a length.

The first derivative $f'(0)$ has dimensions of volume over length, which is length squared. To match $f(0)$, the derivative needs one more power of length. The most natural object to provide the missing length is x itself. As a guess, the first-derivative term should be $xf'(0)$. It could also be $xf'(0)/2$, or $xf'(0)$ multiplied by any dimensionless constant. Dimensional analysis cannot tell you that number, but it turns out to be 1. The series so far is:

$$f(x) = f(0) + xf'(0) + \dots$$

1.4 Cheap differentiation

9

Each successive term in a MacLaurin (or Taylor) series contains a successively higher derivative. The first term used $f(0)$, the zeroth derivative. The second term used $f'(0)$, the first derivative. The third term should use the second derivative $f''(0)$. The dimensions of the second derivative are volume over length squared. because each derivative divides f by one length. Compared to the volume, $f''(0)$ lacks two lengths. The most natural quantity to replace those lengths is x^2 , so the term should be $x^2 f''(0)$. It could be multiplied by a dimensionless constant, which this method cannot find. That number turns out to be $1/2$, and the term is $x^2 f''(0)/2$. The series is now

$$f(x) = f(0) + x f'(0) + \frac{1}{2} x^2 f''(0) + \dots$$

You can guess the pattern. The next term uses $f^{(3)}(0)$, the third derivative. It is multiplied by x^3 to fix the dimensions and by a dimensionless constant that turns out to be $1/6$:

$$f(x) = f(0) + x f'(0) + \frac{1}{2} x^2 f''(0) + \frac{1}{6} x^3 f^{(3)}(0) + \dots$$

The general term is

$$\frac{x^n f^{(n)}(0)}{n!},$$

for reasons that will become clearer in ?? on analogies and operators. This example illustrates how, if you remember a few details about MacLaurin series – for example, that each term has successively higher derivatives – then dimensional analysis can fill in the remainder.

1.4 Cheap differentiation

The relation $[f'(x)] = [f] / [x]$ suggests a way to estimate the *size* of derivatives. Here is the differential equation that describes the oscillations of a mass connected to a spring:

$$m \frac{d^2 x}{dt^2} + kx = 0,$$

where m is the mass, x is its position, t is time, and k is the spring constant. In the first term, the second derivative $d^2 x / dt^2$ is the acceleration a of the mass, so $m(d^2 x / dt^2)$ is ma or the force. And the second term, kx , is the force exerted by the spring. In working out what the terms mean, we have also

Dimensions

10

checked that the terms have the same dimensions (here, dimensions of force). So the equation is at least dimensionally correct.

Here's how to estimate the size of each term. The dimensions of d^2x/dt^2 comes from dividing the dimensions of x by the dimensions of t^2 . The size of d^2x/dt^2 is estimated by dividing the size of x by the size of t^2 . Why not instead divide the dimensions of x^2 by those of t^2 ? The numerator, after all, has a d^2 in it. To answer that question, return to the maxim: d means 'a little bit of'. So dx means 'a little bit of x ', and $d^2x = d(dx)$ means 'a little bit of a little bit of x '. The numerator, therefore does not have anything to do with x^2 . Instead, it has the same dimensions as x . Another way of saying the same idea is that differentiation is a linear operation.

Even if x/t^2 is a rough estimate for the second derivative, x and t are changing: How do you know what x and t to use in the quotient? For x , which is in the numerator, use a **typical value** of x . A typical value is the oscillation amplitude x_0 . For t , which is in the denominator, use the time in which the numerator changes significantly. That time – call it τ – is related to the oscillation period. So

$$\frac{dx}{dt} \sim \frac{\text{typical } x}{\tau} \sim \frac{x_0}{\tau},$$

and

$$\frac{d^2x}{dt^2} = \frac{d}{dt} \left(\frac{dx}{dt} \right) \sim \frac{1}{\tau} \frac{x_0}{\tau} = \frac{x_0}{\tau^2}.$$

Now we can estimate both terms in the differential equation:

$$\begin{aligned} m \frac{d^2x}{dt^2} &\sim m \frac{x_0}{\tau^2}, \\ kx &\sim kx_0, \end{aligned}$$

The differential equation says that the two terms add to zero, so their sizes are comparable:

$$m \frac{x_0}{\tau^2} \sim kx_0.$$

Both sides contain one power of the amplitude x_0 , so it divides out. That cancellation always happens in a linear differential equation. With x_0 gone, it cannot affect the upcoming estimate for τ . So

1.5 Free fall revisited

11

In ideal spring motion – so-called simple harmonic motion – the oscillation period is **independent** of amplitude.

After cancelling the x_0 , the leftover is $k \sim m/\tau^2$ or $\tau \sim \sqrt{m/k}$. A quantity related to the time τ is its reciprocal $\omega = \tau^{-1}$, which has dimensions of inverse time or T^{-1} . Those dimensions are the dimensions of frequency. So

$$\omega = \tau^{-1} \sim \sqrt{\frac{k}{m}}.$$

When you solve the differential equation honestly, this ω is exactly the angular frequency (angle per time) of the oscillations. The missing constant, which dimensional analysis cannot compute, is 1. In this case, dimensional analysis, cheap though it may be, gives the exact frequency.

1.5 Free fall revisited

The ball that fell a height h was released from rest. What if it had an initial velocity v_0 . What is its impact velocity v_{final} ?

1.6 What you have learned

- Preserve dimensions in quantities with dimensions: Do not write ‘ g meters per second squared’; write g .
- Choose dimensions for quantities with arbitrary dimensions, like for x and α in

$$\int_{-\infty}^{\infty} e^{-\alpha x^2} dx.$$

- Exponents are dimensionless.
- You cannot add apples to oranges: Every term in an equation or sum has identical dimensions. Another consequence is that both sides of an equation have identical dimensions.
- The dimensions of an integral are the dimensions of everything inside it, including the dx . This principle helps you guess integrals such as the general Gaussian integral with $-\alpha x^2$ in the exponent.

Dimensions**12**

- The dimensions of a derivative $f'(x)$ are the dimensions of f/x . This principle helps reconstruct formulas based on derivatives, such as Taylor or MacLaurin series.
- The size of df/dx is roughly

$$\frac{\text{typical size of } f}{x \text{ interval over which } f \text{ changes significantly}}$$

See the short and sweet book by Cipra [1] for further practice with dimensions and with rough-and-ready mathematics reasoning.

2

Extreme cases

The next item for your toolbox is the method of **extreme cases**. You can use it to check results and even to guess them, as the following examples illustrate.

2.1 Fencepost errors

Fencepost errors are the most common programming mistake. An index loops over one too many or too few items, or an array is allocated one too few memory locations – leading to a buffer overrun and insecure programs. Since programs are a form of mathematics, fencepost errors occur in mathematics as well. The technique of extreme cases helps you find and fix these errors and deduce correct results instead.

Here is the sum of the first n odd integers:

$$S = \underbrace{1 + 3 + 5 + \cdots}_{n \text{ terms}}$$

Odd numbers are of the form $2k + 1$ or $2k - 1$. Quickly answer this question:

Is the last term $2n + 1$ or $2n - 1$?

For a general n , the answer is not obvious. You can figure it out, but it is easy to make an algebra mistake and be off by one term, which is the difference between $2n - 1$ and $2n + 1$. An extreme case settles the question. Here is the recipe for this technique:

Extreme cases

14

1. Pick an extreme value of n , one where the last term in the sum is easy to determine.
2. For that n , determine the last term.
3. See which prediction, $2n - 1$ or $2n + 1$ (or perhaps neither), is consistent with this last term.

The most extreme value of n is 0. Since n is the number of terms, however, the meaning of $n = 0$ is obscure. The next most extreme case is $n = 1$. With only one term, the final (and also first) term is 1, which is $2n - 1$. So the final term, in general, should be $2n - 1$. Thus the sum is

$$S = 1 + 3 + 5 + \cdots + 2n - 1.$$

Using sigma notation, it is

$$S = \sum_{k=0}^{n-1} (2k + 1).$$

This quick example gives the recipe for extreme-cases reasoning; as a side benefit, it may help you spot bugs in your programs. The sum itself has an elegant picture, which you learn in [Section 4.1](#) in the chapter on pictorial proofs. The rest of this chapter applies the extreme-cases recipe to successively more elaborate problems.

2.2 Integrals

An integral from the [Chapter 1](#), on dimensions, can illustrate extreme cases as well as dimensions. Which of these results is correct:

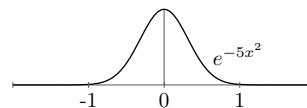
$$\int_{-\infty}^{\infty} e^{-\alpha x^2} dx = \begin{cases} \sqrt{\alpha\pi} \\ \text{or} \\ \sqrt{\frac{\pi}{\alpha}} \end{cases} ?$$

Dimensional analysis answered this question, but forget that knowledge for the moment so that you can practice a new technique.

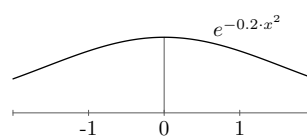
2.2 Integrals

15

You can make the correct choice by looking at the integrand $e^{-\alpha x^2}$ in the two extremes $\alpha \rightarrow \infty$ and $\alpha \rightarrow 0$. As α becomes large, the exponent $-\alpha x^2$ becomes large and negative even when x is only slightly greater than zero. The exponential of a large negative number is nearly zero, so the bell curve narrows, and its area shrinks. As $\alpha \rightarrow \infty$, the area and therefore the integral should shrink to zero. The first option, $\sqrt{\alpha\pi}$, instead goes to infinity. It must be wrong. The second option, $\sqrt{\pi/\alpha}$, goes to infinity and could be correct.



The complementary test is $\alpha \rightarrow 0$. The function flattens to the horizontal line $y = 1$; its integral over an infinite range is infinity. The first choice, $\sqrt{\pi\alpha}$, fails this test because instead it goes to zero as $\alpha \rightarrow 0$. The second option, $\sqrt{\pi/\alpha}$, goes to infinity and passes the test. So the second option passes both tests, and the first option fails both tests. This increases my confidence in $\sqrt{\pi/\alpha}$ while decreasing it, nearly to zero, in $\sqrt{\pi\alpha}$.



If those were the only choices, *and I knew that one choice was correct*, I would choose $\sqrt{\pi/\alpha}$. However, if the joker who wrote the problem included $\sqrt{2/\alpha}$ among the choices, then I need a third test to distinguish between $\sqrt{2/\alpha}$ and $\sqrt{\pi/\alpha}$. For this test, use a third extreme case: $\alpha \rightarrow 1$. Wait, how is 1 an extreme case? Infinity and zero are extreme, but 1 lies between those two so it cannot be an extreme.

Speaking literally, 1 is a special case rather than an extreme case. So extend the meaning of extreme with poetic license and include special cases. The tool, named in full, would be the ‘method of extreme and special cases’. Or, since extreme cases are also special, it could be the ‘method of special cases’. The first option, although correct, is unwieldy. The second option, although also sharing the merit of correctness, is cryptic. It does not help you think of special cases, whereas ‘extreme cases’ does help you: It tells you to look at the extremes. So I prefer to keep the name simple – extreme cases – while reminding myself that extreme cases include special cases like $\alpha \rightarrow 1$.

In the $\alpha \rightarrow 1$ limit the integral becomes

$$I \equiv \int_{-\infty}^{\infty} e^{-x^2} dx,$$

where the \equiv notation means ‘is defined to be’ (rather than the perhaps more common usage in mathematics for modular arithmetic). It is the Gaussian integral and its value is $\sqrt{\pi}$. The usual trick to compute it is to evaluate the square of the integral:

Extreme cases

16

$$I^2 = \left(\int_{-\infty}^{\infty} e^{-x^2} dx \right) \times \left(\int_{-\infty}^{\infty} e^{-x^2} dx \right).$$

In the second factor, change the integration variable to y , making the product

$$I^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\alpha x^2} e^{-\alpha y^2} dx dy.$$

It looks like the integral has become more complicated, but here comes the magic trick. The exponentials multiply to give $e^{-(x^2+y^2)}$, integrated over all x and y – in other words, over the whole plane. And $e^{-(x^2+y^2)} = e^{-r^2}$. So the square of the Gaussian integral is also, in polar coordinates, the integral $\int_{\text{plane}} e^{-r^2} dA$, where dA is the element of area $r dr d\theta$:

$$I^2 = \int_0^{2\pi} \int_0^{\infty} e^{-r^2} \underbrace{r dr d\theta}_{dA}.$$

This integral is doable because the r contributed by the dA is the derivative, except for a factor of 2, of the r^2 in the exponent:

$$\int e^{-r^2} r dr = \frac{1}{2} e^{-r^2} + C,$$

and

$$\int_0^{\infty} e^{-r^2} r dr = \frac{1}{2}.$$

The $d\theta$ integral contributes a factor of 2π so $I^2 = 2\pi/2 = \pi$ and the Gaussian integral is its square root:

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

The only choice consistent with all three extreme cases, even with $\sqrt{2/\alpha}$ among them, is

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

This integral could also be guessed by dimensions, as explained in [Section 1.2](#). Indeed dimensions tell you more than extreme cases do. Dimensions refutes $\sqrt{\pi}/\alpha$ or $\sqrt{\pi}/\alpha^2$, whereas both choices pass the three extreme-case tests:

2.3 Pendulum

17

- $\alpha \rightarrow 0$ Both choices correctly limit to ∞ .
- $\alpha \rightarrow \infty$ Both choices correctly limit to 0.
- $\alpha \rightarrow 1$ Both choices correctly limit to $\sqrt{\pi}$.

Extreme cases, however, has the virtue of being quick. You do not need to find the dimensions for x or α (or invent the dimensions), then find the dimensions of dx and of the result. Extreme cases immediately refutes $\sqrt{\pi\alpha}$. The technique's other virtues become apparent in the next problem: how a pendulum's period varies with amplitude.

2.3 Pendulum

In physics courses, the first problem on oscillations is the ideal spring. Its differential equation is

$$m \frac{d^2x}{dt^2} + kx = 0,$$

where k is the spring constant. Dividing by m gives

$$\frac{d^2x}{dt^2} + \frac{k}{m}x = 0.$$

A consequence of this equation, which we derived in [Section 1.4](#), is that the oscillation period is independent of the amplitude. That property is characteristic of a so-called simple-harmonic system. The oscillation period is:

$$T = 2\pi\sqrt{\frac{m}{k}}.$$

Before moving on to the pendulum, pause to make a sanity check. To make a sanity check, ask yourself: 'Is each portion of the formula reasonable, or does it come out of left field.' [For the non-Americans, left field is one of the distant reaches of a baseball field, and to come out of left fields means an idea come out of nowhere and surprises everyone with how crazy it is.] One species of sanity checking is to check dimensions. Are the dimensions on both sides correct? In this case they are. The dimensions of spring constant are force per length because $F = kx$, so $[k] = \text{MT}^{-2}$. So the dimensions of $\sqrt{m/k}$ are simply time, which is consistent with being an oscillation period

Extreme cases

18

T . [Sorry about the almost-ambiguous notation with T (italic) representing period and T (roman) representing the time dimension.]

Another species of sanity checking is checking extreme cases. Is it reasonable, for example, that m is in the numerator? To decide, check an extreme case of mass. As the mass goes to infinity, the period should go to infinity because the spring has a hard time moving the monstrous mass; and behold, the formula correctly predicts that $T \rightarrow \infty$. Is it reasonable that spring constant k is in the denominator? Check an extreme case of k . As $k \rightarrow 0$, the spring becomes pathetically weak, and the period should go to infinity. Indeed, the formula predicts that $T \rightarrow \infty$. What about the 2π ? To find this constant, either solve the differential equation honestly or use a trick invented by Huygens, which I will explain in lecture if you remind me.

Once the spring has been beaten half to death in physics class, the pendulum is sprung on you. We will study how the period of a pendulum depends on its amplitude – on the maximum angle of the swing, normally called θ_0 . First, let's derive the differential equation for the pendulum, then deduce properties of its solution without solving it. Just as force fights to linearly accelerate an object with mass, torque fights to angularly accelerate an object with moment of inertia. Compare the following formulas:

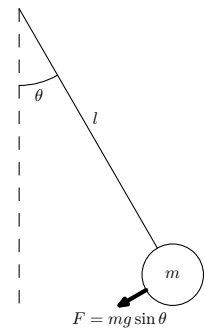
$$\text{force} = \text{mass} \times \text{linear acceleration},$$

$$\text{torque} = \text{moment of inertia} \times \text{angular acceleration}.$$

The first formula is Newton's second law, so you can easily remember it. The second formula follows from the first by analogy, which is the technique of **Chapter 6**. Torque is like force; moment of inertia is like mass; and angular acceleration is like linear acceleration.

The moment of inertia of the bob is $I = ml^2$, and angular acceleration is $\alpha \equiv d^2\theta/dt^2$ (again using \equiv to mean 'is defined to be'). The tangential force trying to restore the pendulum bob to the vertical position is $F = mg \sin \theta$. Or is it $mg \cos \theta$? Decide using extreme cases. As $\theta \rightarrow 0$, the pendulum becomes directly vertical hanging downward, and the tangential force F goes to zero. Since $\sin \theta \rightarrow 0$ as $\theta \rightarrow 0$, the force should contain $\sin \theta$ rather than $\cos \theta$.

The torque, which is the force times the lever arm l , is $F l = mgl \sin \theta$. Putting all three pieces together:



2.3 Pendulum

19

$$\underbrace{-mgl \sin \theta}_{\text{torque}} = \underbrace{ml^2}_I \times \underbrace{\frac{d^2\theta}{dt^2}}_\alpha,$$

where the minus sign in the torque reflects that it is a restoring torque. The mass divides out to produce the pendulum differential equation:

$$\frac{d^2\theta}{dt^2} + \frac{g}{l} \sin \theta = 0.$$

This pendulum equation looks similar to the spring equation

$$\frac{d^2x}{dt^2} + \frac{k}{m}x = 0.$$

Comparing the two equations produces these analogies:

$$\begin{aligned} x &\rightarrow \theta, \\ \frac{k}{m} &\rightarrow \frac{g}{l}, \\ x &\rightarrow \sin \theta. \end{aligned}$$

The first two lines are fine, but the third line contradicts the first one: x cannot map to θ and to $\sin \theta$.

Extreme cases help. Sure, θ and $\sin \theta$ are not identical. However, in the extreme case $\theta \rightarrow 0$, which means that the oscillation angle θ also goes to zero, the two alternatives θ and $\sin \theta$ are identical (a picture proof is given in ??), For small amplitudes, in other words, the pendulum is almost a simple-harmonic system, which would have a constant period. By analogy with the spring equation, the pendulum's period is

$$T = 2\pi \sqrt{\frac{l}{g}},$$

because the pendulum differential equation has g/l where the spring differential equation has k/m . This extreme case is further analyzed in **Chapter 3** using the technique of discretization.

In the Gaussian integral with α , one extreme case was $\alpha \rightarrow 0$ and another was $\alpha \rightarrow \infty$. So try that extreme case here, and see what you can deduce. Not much, since an infinite angle is not informative. However, the idea of a large amplitude is suggestive and helpful. The largest meaningful amplitude – set by the angle of release – is 180° or, in radians, $\theta_0 = \pi$. That angle requires a rod as the pendulum ‘string’, so that the pendulum does not collapse. Such

Extreme cases

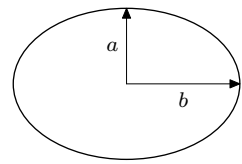
20

a pendulum balanced at $\theta_0 = \pi$ hangs upside down forever. So $T \rightarrow \infty$ when $\theta_0 \rightarrow \pi$. Therefore the period should *increase* as amplitude increases. It could decrease initially, for small θ_0 , then increase as θ_0 gets near π . That behavior would be nasty. The physical world, at least as a first assumption, does not play such tricks on us.

2.4 Ellipse

Now try extreme cases and dimensions on these candidate formulas for the area A of an ellipse:

- ab^2
- $a^2 + b^2$
- a^3/b
- $2ab$
- πab



Let's take them one by one.

- ab^2 . This product has dimensions of length cubed rather than length squared, so it flunks the dimensions test and does not even graduate to the extreme-cases tests. But the other choices have correct dimensions and require more work.
- $a^2 + b^2$. Try an extreme ellipse: a super-thin one with $a = 0$. This case satisfies the first step of the recipe:

Pick an extreme value where the result is easy to determine without solving the full problem.

Now do the second step:

For that extreme case, determine the result.

When $a = 0$ the ellipse has zero area no matter what b is. The third step is:

2.4 Ellipse

21

Determine the prediction in this extreme case, and compare it with the actual value from the second step.

When $a = 0$, the candidate $A = a^2 + b^2$ becomes $A = b^2$. It can be zero, but alas only when $b = 0$. So the candidate fails this extreme-case test except when $a = 0$ and $b = 0$: a boring case of the ellipse shrinking to a point.

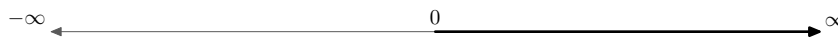
- a^3/b . This candidate passes the thin-ellipse test with $a \rightarrow 0$. When $a \rightarrow 0$, the predicted and actual areas are zero no matter the value of b . Perhaps the candidate is correct. However, it must pass all tests – and even then it may be wrong. If $a \rightarrow 0$ is a reasonable test, then by symmetry $b \rightarrow 0$ should also be worth trying. This test pushes the candidate off the stage. When $b \rightarrow 0$, which produces an infinitely thin vertical ellipse with zero area, the candidate predicts an infinite area whereas the actual area is zero. Although the candidate passed the first test, it fails the second test.
- $2ab$. This candidate is promising. When $a \rightarrow 0$ or $b \rightarrow 0$, the actual and predicted areas are zero. So the candidate passes both extreme-case tests. Both $a \rightarrow 0$ and $b \rightarrow 0$ are literal extreme cases. Speaking figuratively, $a = b$ is also an extreme case. When $a = b$, the candidate predicts that $A = 2a^2$ or, since $a = b$, that $A = 2b^2$. When $a = b$, however, the ellipse is a circle with radius a , and that circle has area πa^2 rather than $2a^2$. So the prediction fails.
- πab . This candidate passes all three tests. Just like $A = 2ab$, it passes $a \rightarrow 0$ and $b \rightarrow 0$. Unlike $A = 2ab$, this candidate also passes the $a = b$ test (making a circle). With every test that a candidate passes, confidence in it increases. So you can be confident in this candidate. And indeed it is correct.

This example introduces extreme cases in a familiar problem, and one where you have choices to evaluate. We next try a three-dimensional problem and guess the answer from scratch. But before moving on, I review the extreme-case tests and discuss how to choose them. Two natural extremes are $a \rightarrow 0$ and $b \rightarrow 0$. However, where did the third test $a \rightarrow b$ originate, and how would one think of it? The answer is symmetry, a useful trick. Actually it's a method: 'a method is a trick I use twice' (George Polya). Symmetry

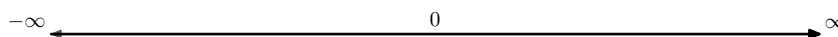
Extreme cases

22

already helped us think of trying $b \rightarrow 0$ after we tried $a \rightarrow 0$. So the following use of it is the second application. Since a and b are lengths, it is natural to compare them by forming their (dimensionless) ratio a/b . The range of a/b is between 0 and ∞ :



The immediately interesting values in this range are its endpoints 0 and ∞ . However, this range is a runt. It is asymmetric, incomplete, and lives on only the right one-half of the real line. To complete the range so that it extends from $-\infty$ to ∞ , take the logarithm of a/b . Here are the possible values of $\ln(a/b)$:

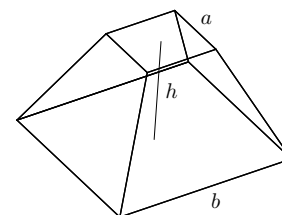


The interesting values on this line are again the endpoints, which are $-\infty$ and ∞ , but also a new one: the middle point, 0. The interesting values of a/b are 0, 1, and ∞ . These points are the three extreme cases for testing the candidate ellipse areas:

$$\begin{aligned} a/b = 0 &\rightarrow b = 0, \\ a/b = \infty &\rightarrow a = 0, \\ a/b = 1 &\rightarrow a = b. \end{aligned}$$

2.5 Truncated pyramid

In the ellipse example, extreme cases helped us evaluate candidates for the area. The next example shows you how to use extreme cases to find a result. Beyond area, the next level of complexity is volume, and the result we look for is the volume of the truncated pyramid formed by slicing off a chunk of the familiar pyramid with a square base. It has therefore a square base and square top that, for simplicity, we assume is parallel to the base. Its height is h , the side length of the base is b , and the side length of the top is a . Guess its volume by finding a formula that meets all the extreme-case tests!



In doing so do not forget the previous technique: dimensions. Any formula must have dimensions of length cubed, so forget about candidate volumes like $V = a^2b^2$ or $V = a^2bh$. But a^2b^2/h would pass the dimensions test.

2.5 Truncated pyramid

23

What are the extreme cases? The simplest is $h \rightarrow 0$, producing a pyramid with zero volume. So a^2b^2/h , although having the correct dimensions, fails because it bogusly produces an infinite volume. Plausible candidates – those producing zero volume – could be ha^2 or h^2a . To choose between those two, think about how the volume must depend on the height. Chop the pyramid into little vertical slivers. When you double the height, you double the height of each sliver, which doubles the volume. So the volume should be proportional to height:

$$V \propto h.$$

A few extreme-cases tests refine this guess. The remaining variables are a and b . The ellipse had only a and b . In the ellipse, a and b are equivalent lengths. Interchanging a and b rotates the ellipse 90° but preserves the same shape and area. For the truncated pyramid, interchanging a and b flips the pyramid 180° but preserves the shape and area. So a and b in the truncated pyramid might have the same interesting extreme cases as do a and b in the ellipse: $a \rightarrow 0$, $b \rightarrow 0$, and $a \rightarrow b$. So let's apply each test in turn, ensuring that the formulas developed in the stepwise process meet all the tests so far investigated.

- $a \rightarrow 0$. This limit shrinks the top surface from a square to a point, making the truncated pyramid an ordinary pyramid with volume $hb^2/3$. This formula also passes the $V \propto h$ test. So $V = hb^2/3$ is a reasonable guess for the truncated volume. Continue testing it.
- $b \rightarrow 0$. This limit shrinks the bottom surface from a square to a point, producing an upside-down-but-otherwise-ordinary pyramid. The previous candidate $V = hb^2/3$ predicts a zero volume, no matter what a is, so $V = hb^2/3$ cannot be correct. The complementary alternative $V = ha^2/3$ passes the $b \rightarrow 0$ test. Great!

Alas, it fails the first test $a \rightarrow 0$. One formula, $V = hb^2/3$, works for $a \rightarrow 0$; the other formula, $V = ha^2/3$, works for $b \rightarrow 0$. Can a candidate pass both tests? Yes! Add the two half-successful candidates:

$$V = \frac{1}{3}ha^2 + \frac{1}{3}hb^2 = \frac{1}{3}h(a^2 + b^2).$$

Two alternatives that also pass both extreme-cases tests, but are not as easy to dream up, are

Extreme cases

24

$$V = \frac{1}{3}h(a+b)^2.$$

and

$$V = \frac{1}{3}h(a-b)^2.$$

- $a \rightarrow b$. In this limit, the pyramid becomes a rectangular prism with height h and base area b^2 (or a^2). So its volume is $V = hb^2$. The hard-won candidate $V = h(a^2 + b^2)/3$, designed to pass the two previous extreme cases, fails this one. Nor do the two alternatives pass. One candidate that does pass is $V = hb^2$. However, it is asymmetric: It treats b specially, which is particularly absurd when $a = b$. What about $V = ha^2$? It treats a specially. What about $V = h(a^2 + b^2)/2$? It is symmetric and passes the $a = b$ test, but it fails the $a \rightarrow 0$ and $b \rightarrow 0$ tests.

We need to expand our horizons. One way to do that is to compare the three candidates that passed $a \rightarrow 0$ and $b \rightarrow 0$:

$$V = \frac{1}{3}h(a^2 + b^2) = \frac{1}{3}h(a^2 + b^2),$$

$$V = \frac{1}{3}h(a + b^2) = \frac{1}{3}h(a^2 + 2ab + b^2),$$

$$V = \frac{1}{3}h(a - b^2) = \frac{1}{3}h(a^2 - 2ab + b^2).$$

The expanded versions share the a^2 and b^2 terms in the parentheses, while differing in the coefficient of the ab term. The freedom to choose that coefficient makes sense. The product ab is 0 in either limit $a \rightarrow 0$ or $b \rightarrow 0$. So adding any amount of ab in the parentheses will not affect the $a \rightarrow 0$ and $b \rightarrow 0$ tests. With just the right coefficient of ab , the candidate might also pass the $a = b$ test. Therefore, find the right coefficient n be in

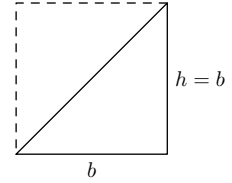
$$V = \frac{1}{3}h(a^2 + nab + b^2).$$

Use the extreme (or special) case $a = b$. Then, the candidate becomes $V = h(2 + n)b^2/3$. To make this volume turn into the correct limit hb^2 , the numerical factor $(2 + n)/3$ should equal 1 meaning that $n = 1$ is the solution:

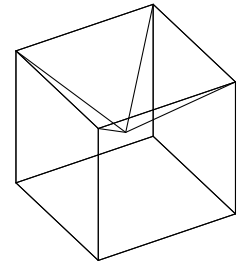
$$V = \frac{1}{3}h(a^2 + ab + b^2).$$

2.6 The magic one-third

You may wonder about the factor of one-third in the volumes of a truncated or regular pyramid. An extreme-case trick explains its origin. First I explain the trick in fewer dimensions: another example of analogy, a technique worthy of its own chapter ([Chapter 6](#)). Instead of immediately explaining the one-third in the volume of a pyramid, which is a difficult three-dimensional problem, first find the corresponding constant in a two-dimensional problem: the area A of a triangle with base b and height h . Its area is $A \sim bh$. What is the constant? Choose a convenient triangle: perhaps a 45-degree right triangle where $h = b$. Two such triangles form a square with area b^2 , so $A = b^2/2$ when $h = b$. The constant in $A \sim bh$ is therefore $1/2$ and $A = bh/2$. Now use the same construction in three dimensions.



What pyramid, when combined with itself perhaps several times, makes a familiar shape? Only the aspect ratio h/b matters in the following discussion. So choose b conveniently, and then choose h to make a pyramid with the clever aspect ratio. The goal shape is suggested by the square pyramid base. Another solid with the same base is a cube. Perhaps several pyramids can combine into a cube of side b . To ease the upcoming arithmetic, I choose $b = 2$. What should h be? To decide, imagine how the cube will be constructed. Each cube has six faces, so six pyramids might make a cube with each pyramid base forming one face of the cube and each pyramid tip facing inwards, meeting in the center of the cube. For the points to meet in the center of the cube, the height must be $h = 1$. So six pyramids with $a = 0$ (meaning that they are not truncated), $b = 2$, and $h = 1$ make a cube with side length 2. The volume of one pyramid is



$$V = \frac{\text{cube volume}}{6} = \frac{8}{6} = \frac{4}{3}.$$

The volume of the pyramid is $V \sim hb^2$, and I choose the missing constant so that the volume is $4/3$. Since $hb^2 = 4$ for these pyramids, the missing constant is $1/3$:

$$V = \frac{1}{3}hb^2 = \frac{4}{3}.$$

So that the general, truncated pyramid agrees with the ordinary pyramid in the limit that $a \rightarrow 0$, the constant for the truncated pyramid is also one-third:

Extreme cases

26

$$V = \frac{1}{3} h(a^2 + ab + b^2).$$

2.7 Drag

The final application of extreme-cases reasoning is to solutions of these nasty nonlinear, coupled, partial-differential equations:

$$\frac{\partial \mathbf{v}}{\partial t} + (\mathbf{v} \cdot \nabla) \mathbf{v} = -\frac{1}{\rho} \nabla p + \nu \nabla^2 \mathbf{v}, \quad (3 \text{ eqns})$$

$$\nabla \cdot \mathbf{v} = 0. \quad (1 \text{ eqn})$$

The top three equations are the Navier–Stokes equations of fluid mechanics, and the bottom equation is the continuity equation. In the four equations is the answer to the following question:

When you drop a paper cone (like a coffee filter) and a smaller cone with the same shape, which falls faster?

Solving those equations is a miserable task, which is why we will instead use our two techniques: dimensions and then extreme cases. For the moment, assume that each cone instantly reaches terminal velocity; that approximation is reasonable but we will check it in ?? using the technique of discretization. So we need to find the terminal velocity. It depends on the weight of the cone and on the drag force F resisting the motion.

To find the force, we use dimensions and add a twist to handle problems like this one that have an infinity of dimensionally correct answers. The drag force depends on the object's speed v ; on the fluid's density ρ ; on its kinematic viscosity ν ; and on the object's size r . Now find the dimensions of these quantities and find all dimensionally correct statements that are possible to make about F . Size r has dimensions of L . Terminal velocity v has dimensions of LT^{-1} . Drag force F has dimensions of mass times acceleration, or MLT^{-2} . Density ρ has dimensions of ML^{-3} . The dimensions of viscosity ν are harder. In the problem set, you show that it has dimensions of L^2T^{-1} . If you look for combinations of ν , ρ , and r , and v that produce dimensions of force, an infinite number of solutions appear, whereas in previous examples using dimensions, only one possibility had the correct dimensions.

Hence the need for a more advanced method to handle the infinite possibilities here. Return to the first principle of dimensions: *you cannot add*

2.7 Drag

27

apples to oranges. The requirement that the sides of an equation match dimensionally is one consequence of the apples-and-oranges principle. Another consequence is that every term in an equation must have the same dimensions. So imagine any true statement about drag force:

$$A + B = C$$

where A , B , and C might be messy combinations of the variables. Then divide each term by A :

$$\frac{A}{A} + \frac{B}{A} = \frac{C}{A}.$$

Because A , B , and C have the same dimensions, each ratio is dimensionless. So you can take any (true) statement about drag force and rewrite it in dimensionless form. No step in this argument depended on the details of drag. It required only that apples must be added to apples. So:

You can write any true statement about the world in dimensionless form.

Furthermore, you can construct any dimensionless expression using dimensionless groups: products of the variables where the product has no dimensions. Since you can write any true statement in dimensionless form, and can write any dimensionless form using dimensionless groups:

You can write any true statement about the world using dimensionless groups.

In the problem of free fall, with variables v , g , and h , the dimensionless group is v/\sqrt{gh} , perhaps raised to a power. With only one group, the only dimensionless statement has the form:

the one group = dimensionless constant,

which results in $v \sim \sqrt{gh}$.

For the drag, what are some dimensionless groups? One group is $F/\rho v^2 r^2$, as you can check by working out its dimensions. A second group is rv/ν . Any other group, it turns out, can be formed from these two groups. With two groups, the most general dimensionless statement is

Extreme cases

28

one group = $f(\text{other group})$,

where f is a dimensionless function. It has a dimensionless argument and must return a dimensionless value because the left side of the equation is dimensionless. Using $F/\rho v^2 r^2$ as the first group:

$$\frac{F}{\rho v^2 r^2} = f\left(\frac{rv}{\nu}\right).$$

The second group, which is the quantity in the parentheses, is the **Reynolds number** and is often written Re . It measures how turbulent the fluid flow is. To find the drag force F , we have to find the function f . It is too hard to determine fully – it would require solving the Navier–Stokes equations – but it might be possible in extreme cases. The extreme cases here are $Re \rightarrow 0$ and $Re \rightarrow \infty$.

Let's hope that the falling cones are in one of those limits! To decide, evaluate Re for the falling cone. From experience, even before you drop the cones to decide which falls faster, either cone falls at roughly $v \sim 1 \text{ m s}^{-1}$. Its size is roughly $r \sim 0.1 \text{ m}$. And the viscosity of the fluid (air) in which it falls is $\nu \sim 10^{-5} \text{ m}^2 \text{ s}^{-1}$, which you can find by looking it up in a table by an online search, or by applying these approximation methods to physics and engineering problems (the theme of another course and book on approximation). So

$$Re \sim \frac{\overbrace{0.1 \text{ m}}^r \times \overbrace{1 \text{ m s}^{-1}}^v}{\underbrace{10^{-5} \text{ m}^2 \text{ s}^{-1}}_\nu} \sim 10^4.$$

So $Re \gg 1$, and we are safe in looking just at that extreme case. Even if the estimate for the speed and size are inaccurate by, say, a factor of 3 each, the Reynolds number is at least 1000, still much larger than 1.

To decide what factors are important in the high-Reynolds-number limit, look at the form of the Reynolds number: rv/ν . One way to send it to infinity is the limit $\nu \rightarrow 0$. Viscosity, therefore, becomes irrelevant as $Re \rightarrow \infty$, and in that limit the drag force F should not depend on viscosity. Although the conclusion is mostly correct, there are subtle lies in the argument. To clarify these subtleties required two hundred years of mathematical and physical development in both theory and experiment. So I will skip the truth, and hope that you are content at least for the moment with almost-truth, especially since it gives the same answer as the truth.

2.7 Drag

29

Let's look at how the requirement of independence from ν constrains the general dimensionless form:

$$\frac{F}{\rho v^2 r^2} = f(Re)$$

The left side does not contain viscosity ν . The right side might because Re contains ν . So if any Reynolds number shows up on the right side, then viscosity will appear on the right side, with no viscosity on the left side with which to cancel it. And that situation would violate the extreme-case result that, in the $Re \rightarrow \infty$ limit, the drag force is independent of viscosity. So the right side must be independent of Re . Since f depended only on the Reynolds number, which has just been stricken off the list of allowed dependencies, the right side $f(Re)$ is a dimensionless constant. Therefore,

$$\frac{F}{\rho v^2 r^2} = \text{dimensionless constant},$$

or

$$F \sim \rho v^2 r^2.$$

And now we have the result that we need to find the relative terminal velocity of the large and small cones. The cones reach terminal speed when the drag force balances the weight. The weight is proportional to the area of the paper, so it is proportional to r^2 . The drag force is also proportional to r^2 , as we just found. To summarize:

$$\underbrace{\rho v^2 r^2}_F \propto \underbrace{r^2}_{\text{weight}}.$$

The factor of r^2 on each side divides out, so

$$v^2 \propto \frac{1}{\rho},$$

showing that

The cones' terminal velocity is independent of its size.

That result is indeed what we found in class by doing the experiment. So, without having to solve the Navier–Stokes differential equations, experiment and cheap theory agree!

2.8 What you have learned

The main theme of this chapter is the recipe for extreme-cases reasoning for checking and guessing the answers to complicated problems:

1. Pick an extreme value where the result is easy to determine without solving the full problem; for example, for the ellipse, its area is easy when $a = 0$ or $b = 0$.
2. For that extreme case, determine the result. For the ellipse, the area is zero when either $a = 0$ or $b = 0$.
3. Determine the prediction in this extreme case, and compare it with the actual value from the second step. So, for the ellipse, any candidate for the area had better go to zero when $a = 0$ or $b = 0$.

Extreme cases also complements the technique of dimensions, once the problems become too complicated for the naive methods of the previous chapter. That symbiosis was illustrated in computing the relative terminal velocities of the falling cones. The general recipe is based on the maxim that **You can write any true statement about the world using dimensionless groups**. It leads to the following problem-solving plan for finding, say, drag force F :

1. Find the quantities on which F depends, and find the dimensions of F and of those quantities.
2. Make dimensionless groups from those quantities.
3. Write the result in general dimensionless form:

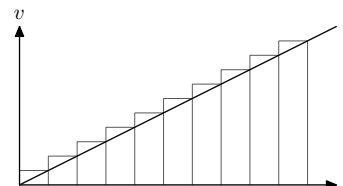
$$\text{group containing } F = f(\text{other groups}).$$

4. Use extreme-cases reasoning to guess the form of the dimensionless function f .

3

Discretization

Discretization takes the fundamental idea of calculus and pushes it to the opposite extreme from what calculus uses. Calculus was invented to analyze changing processes such as orbits of planets or, as a one-dimensional illustration, how far a ball drops in time t . The usual computation



$$\text{distance} = \text{velocity} \times \text{time}$$

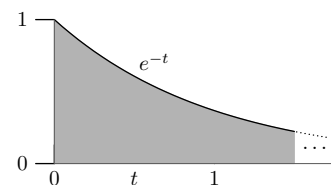
fails because the velocity changes (it increases linearly with time). However – and this next step is the fundamental idea of calculus – over a short time interval, its velocity is almost constant and the usual distance formula works for each short interval. Each short distance is the area of one rectangle, and the total distance fallen is approximately the combined area of the rectangles. To eliminate the error, calculus uses the extreme case of infinite rectangles, ever thinner (shorter intervals) until each shrinks to zero width. Then the approximation of constant speed becomes exact. Discretization uses the opposite extreme: one maybe two fat rectangles. This limitation means the error may not be zero, but it drastically simplifies any computations.

3.1 Exponential decay

The first example is this integral:

$$\int_0^{\infty} e^{-t} dt.$$

Since the derivative of e^t is e^t , the indefinite integral is easy to find exactly, and the limits make the computation even simpler. In an example where the exact answer is known, we can



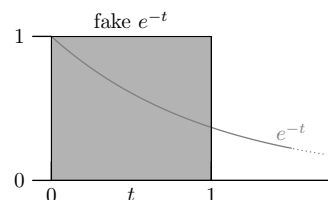
Discretization

32

adjust the free parameters in the method of discretization until the method produces accurate values. So, replace the complicated, continuous, smooth exponential decay e^{-t} by a rectangle, and do the integral by finding the area of the rectangle.

With one rectangle, the approximate function remains constant until it abruptly falls to and remains zero. Finding the area of the rectangle requires choosing its height and width. A natural height is the maximum of e^{-t} , which is 1. A natural width is the time interval until $f(t) = e^{-t}$ changes significantly. A similar idea appeared in [Section 1.4](#) to approximate a derivative df/dx . Its numerator df was estimate as a typical value of $f(x)$. Its denominator dx became the x interval over which $f(x)$ changes significantly. For an exponential, a natural definition for significant change is to changes by a factor of e . When $f(t) = e^{-t}$, this change happens when t goes from t to $t + 1$. So the approximating rectangle, whose height we've chosen to be 1, also has unit width. It is a unit square with unit area. And this rectangle exactly estimates the integral since

$$\int_0^{\infty} e^{-t} dt = 1.$$

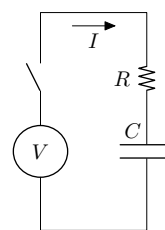


3.2 Circuit with exponential decay

In [Chapter 1](#) on dimensions, I insisted that declaring quantities prematurely dimensionless ties one hand behind your back. In the previous example I committed that sin by making the exponent be $-t$. Since an exponent is dimensionless, my choice made t dimensionless as well.

A more natural interpretation of t is as a time. So here is a similar example but where t has dimensions, which are useful in making and checking the approximations. Let's first investigate the initial conditions, just before the switch closes. No current is flowing since the circuit is not yet a closed loop. Furthermore, because the circuit has been waiting forever, the capacitor has had completely discharged. So capacitor has no charge on it. The charge determines the voltage across the capacitor by

$$Q = CV_C,$$



3.2 Circuit with exponential decay

33

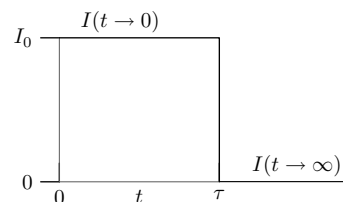
where Q is the charge on the capacitor, C is its capacitance, and V_C is the resulting voltage. [See the classics on circuits [2] and electromagnetism [3] for more on capacitors.] So just before the switch closes, the capacitor has zero voltage on it ($V_C = 0$).

At time $t = 0$, I close the switch, which connects the resistor and capacitor to the source voltage V (which is constant). Since V_C starts at zero, the voltage drops in the resistor is the whole source voltage V :

$$V_R = V \quad (\text{initially}),$$

where V_R is the voltage across the resistor. This voltage drop is caused by a current I flowing through the resistor (which then flows through the capacitor). Ohm's law says that $V_R = IR$. Initially $V_R = V$ so the initial current is $I_0 = V/R$. This current charges the capacitor and increases V_C . As V_C increases, V_R decreases – which decreases the current I , which decreases how fast V_C increases, which . . . Finding the current is a problem for calculus, in particular for a differential equation.

Instead, let's guess the current using dimensions, extreme cases, and the new technique of discretization. First apply extreme cases. At the $t = 0$ extreme, the current is $I_0 = V/R$. At the $t = \infty$ extreme, no current flows: The capacitor accumulates enough charge so that $V_C = V$, whereupon no voltage drops across the resistor. From Ohm's law again, a zero voltage drop is possible only if no current flows.

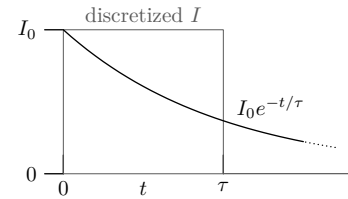


Between those extremes, we guess I using discretization. Pretend that I stays at its $t = 0$ value of I_0 for a time τ , then drops to its $t = \infty$ value of $I = 0$. So τ is the time for the current to change significantly. To determine τ , use dimensions. It can depend only on R and C . [It could depend on V , but because the system is linear, the time constants do not depend on amplitude.] The only way to combine R and C into a time is the product RC . A reasonable guess for τ is therefore $\tau = RC$. In this picture, the discretized current stays at V/R until $t = \tau$, then falls to 0 and remains zero forever.

Discretization

34

No physical current changes so abruptly. To guess the true current, use discretization in reverse. The exponential decay of [Section 3.1](#) produced the same rectangular shape after discretizing. So perhaps the true current here is also an exponential decay. In the other example, the function was e^{-t} , and the changeover from early- to late-time behavior happened at $t = 1$ (in that example, t had no dimensions). By $t = 1$, the exponential decay e^{-t} had changed significantly (by a factor of e). For this circuit, the corresponding changeover time is τ . To change by a factor of e in time τ , the current should contain $e^{-t/\tau}$. The initial current is $I = I_0$, so the current should be



$$I = I_0 e^{-t/\tau} = \frac{V}{R} e^{-t/\tau}.$$

Having a solution, even a guess, turns the hard work of solving differential equations into the easier work of verifying a solution.

To test the guess for I , I derive the differential equation for the current. The source voltage V drops only in the resistor and capacitor, so their voltage drops must add to V :

$$V = V_R + V_C.$$

The capacitor voltage is $V_C = Q/C$. The resistor voltage is $V_R = IR$, so

$$V = IR + \frac{Q}{C}.$$

It seems that there are too many variables: V and C are constants, but I and Q are unknown. Fortunately current I and charge Q are connected: charge is the time integral of current and $I = dQ/dt$. Differentiating each term with respect to time simplifies the equation:

$$0 = R \frac{dI}{dt} + \frac{1}{C} \underbrace{\left(\frac{dQ}{dt} \right)}_I = R \frac{dI}{dt} + \frac{I}{C}.$$

Move the R to be near its companion C (divide by R):

$$0 = \frac{dI}{dt} + \frac{I}{\underbrace{RC}_\tau} = \frac{dI}{dt} + \frac{I}{\tau}.$$

Dimensions, extreme cases, and reverse discretization produced this current:

3.3 Population

35

$$I = I_0 e^{-t/\tau}.$$

Amazing! It satisfies the differential equation:

$$\frac{d}{dt} \left(I_0 e^{-t/\tau} \right) + \frac{I_0 e^{-t/\tau}}{\tau} = 0$$

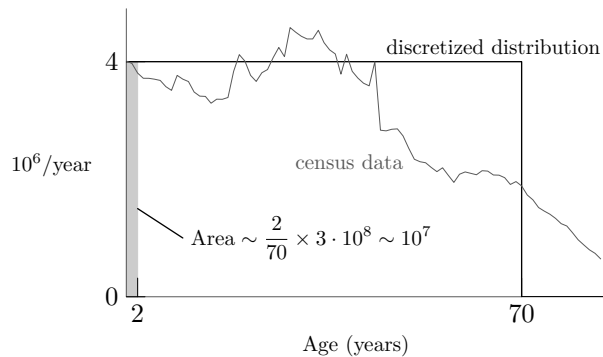
because the time derivative brings down a factor of $-1/\tau$, making the first and second terms equal except for a minus sign.

3.3 Population

Not all problems are exponential decays. In the next example, the true functions are unknown and exact answers are not available. The problem is to estimate the number of babies in the United States. To specify the problem, define babies as children less than two years old. One estimate comes from census data, which is accurate within the limits of statistical sampling. You integrate the population curve over the range $t = 0$ to 2 years. But that method relies on the massive statistical efforts of the US census bureau and would not work on a desert island. If only the population were constant (didn't depend on age), then the integrals are easy! The desert-island, back-of-the-envelope method is to replace the complicated population curve by a single rectangle.

How high is the rectangle and how wide is it? The width τ , which is a time, has a reasonable estimate as the average life expectancy. So $\tau \sim 70$ years. How high is the rectangle? The height does not have such an obvious direct answer as the width. In the exponential-decay examples, the height was the the initial value, from which we found the area. Here, the procedure reverses. You know the area – the population of the United States – from which you find the height. So, with the area being $3 \cdot 10^8$, the height is

$$\text{height} \sim \frac{\text{area}}{\text{width}} \sim \frac{3 \cdot 10^8}{75 \text{ years}},$$



Discretization

36

since the width is the life expectancy, for which we used 70 years. How did it become 75 years? The answer is by a useful fudge. The new number 75 divides into 3 (or 300) more easily than 70 does. So change the life expectancy to ease the mental calculations. The inaccuracies caused by that fudge are no worse than in replacing the complex population curve by a rectangle. So

$$\text{height} \sim 4 \cdot 10^6 \text{ year}^{-1}.$$

Integrating a rectangle of that height over the infancy duration of 2 years gives

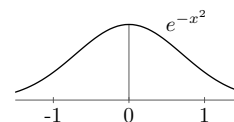
$$N_{\text{babies}} \sim \underbrace{4 \cdot 10^6 \text{ years}^{-1}}_{\text{height}} \times \underbrace{2 \text{ years}}_{\text{infancy}} = 8 \cdot 10^6.$$

Thus roughly 8 million babies live in the United States. From this figure, you can estimate the landfill volume used each year by disposable diapers (nappies).

3.4 Full width at half maximum

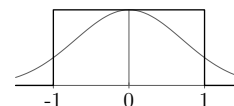
The Gaussian integral

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



has appeared in several examples, and you've seen the trick (in [Section 2.2](#)) of squaring it to show that its value is $\sqrt{\pi}$. The exponential in the integrand is a difficult, continuous function. Except over the infinite integration range, the integral has no closed form, which is why statistics tables enumerate the related error function numerically. I introduce that evidence to show you how difficult the integral is without infinite limits, and it is not easy even with infinite limits.

Pretend therefore that you forget the trick. You can approximate the integral using discretization by replacing the integrand with a rectangle. How high and how wide should the rectangle be? The recipe is to take the height as the maximum height of the function and the width as the distance until the function falls significantly. In the exponential-decay examples, significant meant changing by a factor of e . The maximum of e^{-x^2} is at $x = 0$ when $e^{-x^2} = 1$, so the approximating rectangle has unit height. It falls to $1/e$ when $x = \pm 1$, so the approximating rectangle has width 2 and therefore area 2. This estimate is



3.4 Full width at half maximum

37

half decent. The true value is $\sqrt{\pi} = 1.77\dots$, so the error is about 13%: a reasonable trade for one line of work.

Another recipe, also worth knowing because it is sometimes more accurate, arose in the bygone days of spectroscopy. Spectroscopes measure the wavelengths (or frequencies) where a molecule absorbed radiation and the corresponding absorption strengths. These data provided an early probe into the structure of atoms and molecules, and was essential to the development of quantum theory [4]. An analogous investigation occurs in today's particle accelerators – colloquially, atom smashers – such as SLAC in California and CERN and in Geneva: particles, perhaps protons and neutrons, collide at high energies, showering fragments that carry information about the structure of the original particles. Or, to understand how a finely engineered wristwatch works, hammer it and see what the flying shards and springs reveal.

The spectroscope was a milder tool. A chart recorder plotted the absorption as the spectroscope swept through the wavelength range. The area of the peaks was an important datum, and whole books like [5] are filled with these measurements. Over half a century before digital chart recorders and computerized numerical integration, how did one compute these areas? The recipe was the FWHM.

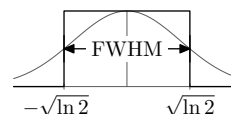
FWHM = full width at half maximum

Unpack the acronym in slow motion:

1. **M.** Find the maximum value (the peak value).
2. **HM.** Find one-half of the maximum value, which is the half maximum.
3. **FWHM.** Find the two wavelengths – above and below the peak – where the function has fallen to one-half of the maximum value. The full width is the difference between the above and below wavelengths.

The FWHM approximation recipe replaces the peak by a rectangle with height equal to the peak height and width equal to the the width estimated by the preceding three-step procedure.

Try this recipe on the Gaussian integral and compare the estimate with the estimate from the old recipe of finding where the function changed by a factor of e . The Gaussian has maximum height 1 at $x = 0$. The half maximum is $1/2$, which



Discretization

38

happens when $x = \pm\sqrt{\ln 2}$. The full width is then $2\sqrt{\ln 2}$, and the area of the rectangle – which estimates the original integral – is $2\sqrt{\ln 2}$. Here, side by side, are the estimate and the exact integral:

$$\int_{-\infty}^{\infty} e^{-x^2} dx = \begin{cases} \sqrt{\pi} = 1.7724\dots & \text{(exact),} \\ 2\sqrt{\ln 2} = 1.6651\dots & \text{(estimate).} \end{cases}$$

The FWHM estimate is accurate to 6%, twice as accurate as the previous recipe. It's far better than one has a right to expect for doing only two lines of algebra!

3.5 Stirling's formula

The FWHM result accurately estimates one of the most useful quantities in applied mathematics:

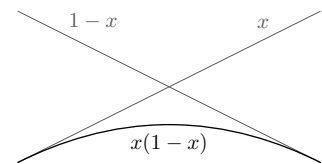
$$n! \equiv n \times (n-1) \times (n-2) \times \dots \times 2 \times 1.$$

We meet this quantity again as a picture proof in [Section 4.6](#). Here we estimate $n!$ by discretizing an integral representing $n!$:

$$\int_0^{\infty} t^n e^{-t} dt = n!$$

You may not yet know that this integral is $n!$; you can show it either with integration by parts or see ?? on generalization to learn differentiation under the integral sign. For now accept the integral representation on faith, with a promise to redeem the trust in a later chapter.

To approximate the integral, imagine what the integrand $t^n e^{-t}$ looks like. It is the product of the increasing function t^n and the decreasing function e^{-t} . Such a product usually peaks. A familiar example of this principle is the product of the increasing function x and the decreasing function $1-x$, over the range $x \in [0, 1]$ where both functions are positive. The product rises from and then falls back to zero, with a peak at $x = 1/2$.



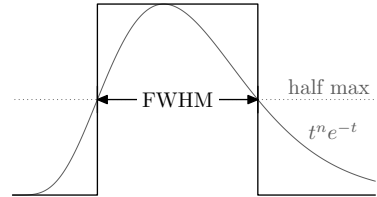
You can check that the product $t^n e^{-t}$ has a peak by looking at its behavior in two extreme cases: in the short run $t \rightarrow 0$ and in the long run $t \rightarrow \infty$. When $t \rightarrow 0$, the exponential is 1, but the polynomial factor t^n wipes it out by multiplying by zero. When $t \rightarrow \infty$, the polynomial factor t^n pushes the product to infinity while the exponential factor e^{-t} pushes it to zero.

3.5 Stirling's formula

39

An exponential beats any polynomial. To see why and avoid the negative exponent $-t$ muddying this issue, compare instead e^t with t^n as $t \rightarrow \infty$. The Taylor series for e^t contains all powers of t , so it is like an infinite-degree polynomial. So e^t/t^n goes to infinity once t gets large enough. Similarly, its reciprocal $t^n e^{-t}$ goes to zero as $t \rightarrow \infty$. Being zero at also $t = 0$, the product is zero at both extremes and positive elsewhere. Therefore it peaks in between. Maybe it has more than one peak, but it should have at least one peak. Furthermore, as n increases, the t^n polynomial factor strengthens, so the e^{-t} requires a larger t to beat down the t^n . Therefore, as n increases the peak moves right.

With $t^n e^{-t}$ having a peak, the FWHM recipe can approximate its area. The recipe requires finding the height (the maximum of the function) and the width (the FWHM) of the approximating rectangle. To find these parameters, slurp the t^n into the exponent:



$$t^n e^{-t} = e^{n \ln t} e^{-t} = e^{n \ln t - t}.$$

The exponent $f(t) \equiv n \ln t - t$ is interesting. As $t \rightarrow 0$, the $\ln t$ takes $f(t)$ to $-\infty$. As $t \rightarrow \infty$, the $-t$ takes $f(t)$ again to $-\infty$. Between these limits, it peaks. To find the maximum, set $f'(t) = 0$:

$$f'(t) = \frac{n}{t} - 1 = 0,$$

or $t_{\text{peak}} = n$. As we predicted, the peak moves right as n increases. The height of the peak is one item needed to estimate the rectangle's area. At the peak, $f(t)$ is $f(n) = n \ln n - n$, so the original integrand, which is $e^{f(t)}$, is

$$e^{f(t_{\text{peak}})} = e^{f(n)} = e^{n \ln n - n} = \frac{n^n}{e^n} = \left(\frac{n}{e}\right)^n.$$

To find the width, look closely at how $f(t)$ behaves near the peak $t = n$ by writing it as a Taylor series around the peak:

$$f(t) = f(n) + f'(n)(t - n) + \frac{1}{2}f''(n)(t - n)^2 + \dots$$

The first derivative is zero because the expansion point, $t = n$, is a maximum and there $f'(n) = 0$. So the second term in the Taylor series vanishes. To evaluate the third term, compute the second derivative of f at $t = n$:

$$f''(n) = -\frac{n}{t^2} = -\frac{1}{n}.$$

Discretization

40

So

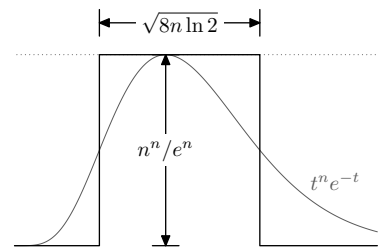
$$f(t) = \underbrace{n \ln n - n}_{f(n)} + \frac{1}{2} \times \underbrace{\left(-\frac{1}{n}\right)}_{f''(n)} (t-n)^2 + \dots$$

The first term gives the height of the peak that we already computed. The second term says how the height falls as t moves away from n . The result is an approximation for the integrand:

$$e^{f(t)} = \left(\frac{n}{e}\right)^n e^{-(t-n)^2/2n}.$$

The first factor is a constant, the peak height. The second factor is the familiar Gaussian. This one is centered at $t = n$ and contains $1/2n$ in the exponent but otherwise it's the usual Gaussian with a quadratic exponent. It falls by a factor of 2 when $(t-n)^2/2n = \ln 2$, which is when

$$t_{\pm} = n \pm \sqrt{2n \ln 2}.$$



The FWHM is $t_+ - t_-$, which is $\sqrt{8n \ln 2}$. The estimated area under $e^{f(t)}$ is then

$$\left(\frac{n}{e}\right)^n \times \sqrt{8n \ln 2}.$$

As an estimate for $n!$, each piece is correct except for the constant factor. The more accurate answer has $\sqrt{2\pi}$ instead of $\sqrt{8 \ln 2}$. However, 2π is roughly $8 \ln 2$ so the approximate is, like the estimate the vanilla Gaussian integral (coincidence?), accurate to 6%.

3.6 Pendulum period

The period of a pendulum is by now a familiar topic in this book. Its differential equation becomes tractable with a bit of discretization. The differential equation that describes pendulum motion is

$$\frac{d^2\theta}{dt^2} + \frac{g}{l} \sin \theta = 0$$

This nonlinear equation has no solution in terms of the usual functions – to put it more precisely, in terms of elementary functions. But you can

3.6 Pendulum period

41

understand a lot about how it behaves by discretizing. If only the equation were

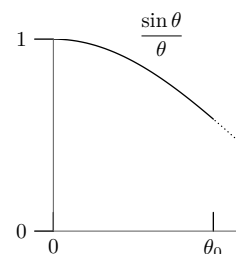
$$\frac{d^2\theta}{dt^2} + \frac{g}{l}\theta = 0.$$

This equation is linear, and therefore possible to solve without too much misery – I hesitate to say that any differential equation is ‘easy’ – and its solution is an oscillation with angular frequency $\omega = \sqrt{g/l}$:

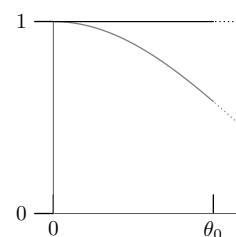
$$\theta(t) = \theta_0 \cos\left(\frac{g}{l}t\right).$$

Its period is $2\pi\sqrt{g/l}$, which is independent of amplitude θ_0 .

The complexity of the unapproximated pendulum equation arises because it has the torque-producing factor $\sin\theta$ instead of its approximation θ . The two functions match perfectly as $\theta \rightarrow 0$. But as θ grows – which happens with large amplitudes – the equivalence becomes less accurate. One way to compare them is to look at their ratio $(\sin\theta)/\theta$. As expected, when $\theta = 0$, the ratio is 1. As θ grows, the ratio falls, making the simple-harmonic approximation less accurate. We can discretize to find a more accurate approximation than the usual simple-harmonic one, yet still produce a linear differential equation. The upcoming figures illustrate making and refining that approximation.



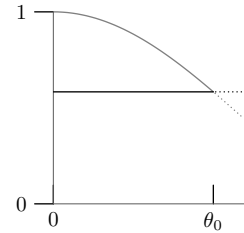
We need a discrete approximation to the difficult function $\sin\theta$ in the range $[0, \theta_0]$. Look at the two functions θ and $\sin\theta$ after dividing by θ ; we are taking out the common big part, the topic of **Chapter 5**. The difficult function becomes $(\sin\theta)/\theta$. The other function, a straight line, is the simple harmonic approximation, and is a useful zeroth approximation. But it does not produce any change in period as a function of amplitude (since the height of the replacement line is independent of θ_0).



Discretization

42

The next approximation does fix that problem. Use a flat line with height $(\sin \theta_0)/\theta_0$. This line is the minimum height of $(\sin \theta)/\theta$. Why is that choice an improvement on the first approximation, using the maximum height of 1? Because in this choice, the height varies with amplitude, so the period varies with amplitude: This choice explains a physical effect that the first choice approximated into oblivion. In this second approximation, the torque term $(g/l) \sin \theta$ becomes



$$\frac{g}{l} \theta \frac{\sin \theta_0}{\theta_0}.$$

Starting from the simple-harmonic approximation, this choice is equivalent to replacing gravity by a slightly weaker gravity:

$$g \rightarrow g \times \frac{\sin \theta_0}{\theta_0}.$$

The Taylor series for sin gives

$$\frac{\sin \theta_0}{\theta_0} \approx 1 - \frac{\theta_0^2}{6}.$$

The fake g is then

$$g_{\text{fake}} = g \left(1 - \frac{\theta_0^2}{6} \right).$$

Using this fake g , the period becomes

$$T \approx 2\pi \sqrt{\frac{l}{g_{\text{fake}}}}.$$

To compute $g_{\text{fake}}^{-1/2}$ requires another Taylor series:

$$(1 + x)^{-1/2} \approx 1 - \frac{x}{2}.$$

Then

$$T \approx 2\pi \sqrt{\frac{l}{g}} \left(1 + \frac{\theta_0^2}{12} \right).$$

3.7 What have you learnt

43

This period is an overestimate because it assumed the weakest torque adjustment factor: scaling the torque by the value of $(\sin \theta)/\theta$ at the endpoints of the swing when $\theta = \pm\theta_0$. The next approximation comes from using an intermediate height for the replacement line. Equivalently, say that the pendulum spends half its flight acting like a spring, where the torque contains just θ ; and half its flight where the torque has the term $\theta(\sin \theta_0)/\theta_0$. Then the period is an average of the simple-harmonic period $T = 2\pi\sqrt{l/g}$ with the preceding underestimate:

$$T = 2\pi\sqrt{\frac{l}{g}} \left(1 + \frac{\theta_0^2}{24}\right).$$

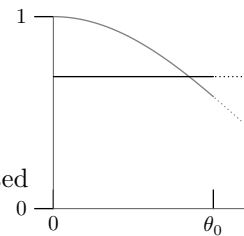
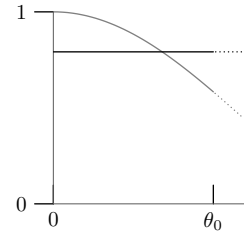
The next step – and here I am pushing this method perhaps farther than is justified – is to notice that the pendulum spends most of its time where it moves the slowest. So it spends most of time near the endpoints of the swings, where the simple-harmonic approximation is the least accurate. So the endpoint-only underestimate estimate for T should be weighted slightly more than the simple-harmonic overestimate. The most recent estimate weighted these pieces equally. To improve it, count the endpoint estimate, say, twice and the center estimate once. This recipe has a further justification in that there are two endpoints and only one center! Then the period becomes

$$T = 2\pi\sqrt{\frac{l}{g}} \left(1 + \frac{\theta_0^2}{18}\right)$$

The true coefficient, which comes from doing a power-series solution, is $1/16$ so this final weighted estimate is very accurate!

3.7 What have you learnt

Discretization makes hard problems simple. The recipe is to replace a complicated function by a rectangle. The art is in choosing the height and width of the rectangle, and you saw two recipes. In both, the height is the maximum of the original function. In the first, easier recipe, the width is the range over which the function changes by a factor of e ; this recipe is useful for linear exponential decays. The second recipe, the FWHM, is useful for messy functions like spectroscopy absorption peaks and Gaussians. In that



Discretization**44**

recipe, the width is the width over which the function goes from one-half the maximum and then returns to that value.

4

Picture proofs

Do you ever walk through a proof, understand each step, yet not believe the theorem, not say ‘Yes, of course it’s true’? The analytic, logical, sequential approach often does not convince one as well as does a carefully crafted picture. This difference is no coincidence. The analytic, sequential portions of our brain evolved with our capacity for language, which is perhaps 10^5 years old. Our pictorial, Gestalt hardware results from millions of years of evolution of the visual system and cortex. In comparison to our visual hardware, our symbolic, sequential hardware is an ill-developed latecomer. Advertisers know that words alone do not convince you to waste money on their clients’ junk, so they spend zillions on images. This principle, which has higher applications, is the theme of this chapter.

4.1 Adding odd numbers

Here again is the sum from [Section 2.1](#) that illustrated using extreme cases to find fencepost errors:

$$S = \underbrace{1 + 3 + 5 + \cdots}_{n \text{ terms}}$$

Before I show the promised picture proof, let’s go through the standard method, proof by induction, to compare it later to the picture proof. An induction proof has three pieces:

1. Verify the *base case* $n = 1$. With $n = 1$ terms, the sum is $S = 1$, which equals n^2 . QED (Latin for ‘quite easily done’).
2. Assume the *induction hypothesis*. Assume that the sum holds for n terms:

Picture proofs

46

$$\sum_1^n (2k - 1) = n^2.$$

This assumption is needed for the next step of verifying the sum for $n + 1$ terms.

3. Do the *induction step* of verifying the sum for $n + 1$ terms, which requires showing that

$$\sum_1^{n+1} (2k - 1) = (n + 1)^2.$$

The sum splits into a new term and the old sum:

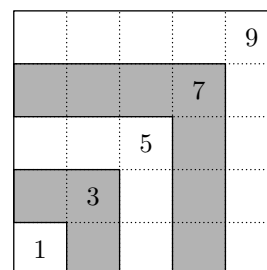
$$\sum_1^{n+1} (2k - 1) = \underbrace{2n + 1}_{\text{new term}} + \sum_1^n (2k - 1).$$

The sum on the right is n^2 courtesy of the induction hypothesis. So

$$\sum_1^{n+1} (2k - 1) = 2n + 1 + n^2 = (n + 1)^2.$$

The three parts of the induction proof are complete, and the theorem is proved. However, the parts may leave you feeling that you follow each step but do not see *why* the theorem is true.

Compare it against the picture proof. Each term in the sum S adds one odd number represented as the area of an L-shaped piece. Each piece extends the square by one unit on each side. Adding n terms means placing n pieces and making an $n \times n$ square. [Or is it an $(n - 1) \times (n - 1)$ square?] The sum is the area of the square, which is n^2 . Once you understand this picture, you never forget why adding the first n odd numbers gives the perfect square n^2 .



4.2 Geometric sums

Here is a familiar series:

$$S = 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \cdots$$

4.3 Arithmetic mean–geometric mean inequality

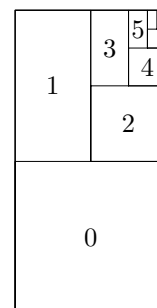
47

The usual symbolic way to evaluate the sum is with the formula for a geometric series. You can derive the formula using a trick. First compute $2S$ by multiplying each term by 2:

$$2S = 2 + 1 + \underbrace{\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \cdots}_S.$$

This sum looks like S , except for the first term 2. So $2S = 2 + S$ and $S = 2$.

The result, though correct, may seem like magic. Here then is a picture proof. A square with unit area represents the first term, which is $1/2^0$ (and is labelled 0). The second term is a $1 \times 1/2$ rectangle representing $1/2^1$ (and is labelled against by the exponent 1). The third term is a $1/2 \times 1/2$ square placed in the nook. The fourth term is, like the second term, a rectangle. With every pair of terms, the empty area between all the rectangles and three-quarters of the 1×2 outlining rectangle fills in. In the limit, the sum fills the 1×2 rectangle, showing that $S = 2$.



4.3 Arithmetic mean–geometric mean inequality

A classic inequality is the arithmetic mean–geometric mean inequality. Here are a few numerical examples before the formal statement. Take two numbers, say, 1 and 2. Their arithmetic mean is 1.5. Their geometric mean is $\sqrt{1 \times 2} = 1.414 \dots$. Now try the same operations with 2 and 3. Their arithmetic mean is 2.5, and their geometric mean is $\sqrt{2 \times 3} = 2.449 \dots$. In both cases, the geometric mean is smaller than the arithmetic mean. This pattern is the theorem of the arithmetic mean and geometric mean. It says that when $a, b \geq 0$, then

$$\underbrace{\frac{a+b}{2}}_{\text{AM}} \geq \underbrace{\sqrt{ab}}_{\text{GM}}$$

where AM means arithmetic mean and GM means geometric mean.

It has at least two proofs: symbolic and pictorial. A picture proof is hinted at by the designation of \sqrt{ab} as the geometric mean. First, however, I prove it symbolically. Look at $(a - b)^2$. Since it is a square,

Picture proofs

48

$$(a - b)^2 \geq 0.$$

Expanding the left side gives $a^2 - 2ab + b^2 \geq 0$. Now do the magic step of adding $4ab$ to both sides to get

$$a^2 + 2ab + b^2 \geq 4ab.$$

The left side is again a perfect square, whose perfection suggests taking the square root of both sides to get

$$a + b \geq 2\sqrt{ab}.$$

Dividing both sides by 2 gives the theorem:

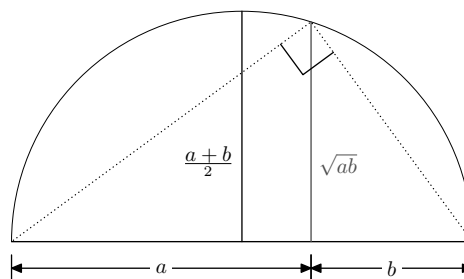
$$\underbrace{\frac{a+b}{2}}_{\text{AM}} \geq \underbrace{\sqrt{ab}}_{\text{GM}}$$

Maybe you agree that, although each step is believable (and correct), the sequence of all of them seems like magic. The little steps do not reveal the structure of the argument, and the *why* is still elusive. For example, if the algebra steps had ended with

$$\frac{a+b}{4} \geq \sqrt{ab},$$

it would not have seemed obviously wrong. We would like a proof whose result *could not have been otherwise*.

Here then is a picture proof. Split the diameter of the circle into the lengths a and b . The radius is $(a+b)/2$, which is the arithmetic mean. Now we need to find the geometric mean, whose name is auspicious. Look at the second half chord rising from the diameter where a and b meet. It is also the height of the dotted triangle, and that triangle is a right triangle. With right triangles everywhere, similar triangles must



come in handy. Let the so-far-unknown length be x . By similar triangles,

$$\frac{x}{a} = \frac{b}{x},$$

4.3 Arithmetic mean–geometric mean inequality

49

so $x = \sqrt{ab}$, showing that the half chord is the geometric mean. That half chord can never be greater than the radius, so the geometric mean is never greater than the arithmetic mean. For the two means to be equal, the geometric-mean half chord must slide left to become the radius, which happens only when $a = b$. So the arithmetic mean equals the geometric mean when $a = b$.

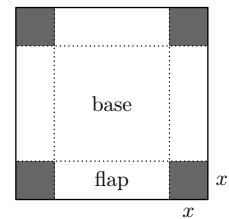
Compare this picture proof with the symbolic proof. The structure of the picture proof is there to see, so to speak. The only non-obvious step is showing that the half chord is the geometric mean \sqrt{ab} , the geometric mean. Furthermore, the picture shows why equality between the two means results only when $a = b$: Only then does the half chord become the radius.

Here are two applications of the AM–GM inequality to problems from introductory calculus that one would normally solve with derivatives. In the first problem, you get $l = 40$ m of fencing to mark off a rectangular garden. What dimensions does the garden have in order to have the largest area? If a is the length and b is the width, then $l = 2(a + b)$, which is $4 \times$ AM. The area is ab , which is $(\text{GM})^2$. Since $\text{AM} \geq \text{GM}$, the consequence in terms of this problem's parameters is

$$\text{AM} = \frac{l}{4} \geq \sqrt{\text{area}} = \text{GM}.$$

Since the geometric mean cannot be larger than $l/4$, which is constant, the geometric mean is maximized when $a = b$. For maximum area, therefore choose $a = b = 10$ m and get $A = 100$ m².

The next example in this genre is a more difficult three-dimensional problem. Start with a unit square and cut out four identical corners, folding in the four edges to make an open-topped box. What size should the corners be to maximize the box volume? Call x the side length of the corner cutout. Each side of the box has length $1 - 2x$ and it has height x , so the volume is



$$V = x(1 - 2x)^2.$$

For lack of imagination, let's try the same trick as in the previous problem. Two great mathematicians, George Polya and Gabor Szego, commented that, 'An idea which can be used once is a trick. If it can be used more than once it becomes a method.' So AM–GM, if it helps solve the next problem, gets promoted from a mere trick to the more exalted method.

Picture proofs

50

In the previous problem, the factors in the area were a and b , and their sum $a + b$ was constant because it was fixed by the perimeter. Then we could use AM–GM to find the maximum area. Here, the factors of the volume are x , $1 - 2x$, and $1 - 2x$. Their sum is $2 - 3x$, which is not a constant; instead it varies as x changes. This variation means that we cannot apply the AM–GM theorem directly. The theorem is still valid but it does not tell us what we want to know. We want to know the largest possible volume. And, directly applied, the theorem says that the volume is never less than the cube of the arithmetic mean. Making the volume equal to this value does not guarantee that the maximum volume has been found, because the arithmetic mean is changing as one changes x to maximize the geometric mean. The largest volume may result where the GM is not equal to the changing AM. In the two-dimensional problem, this issue did not arise because the AM was already constant (it was a fixed fraction of the perimeter).

If only the factor of x were a $4x$, then the $3x$ would disappear when computing the AM:

$$4x + (1 - 2x) + (1 - 2x) = 2.$$

As Captain Jean-luc Picard of *The Next Generation* says, ‘Make it so.’ You can produce a $4x$ instead of an x by studying $4V$ instead of V :

$$4V = 4x \times 1 - 2x \times 1 - 2x.$$

The sum of the factors is 2 and their arithmetic mean is $2/3$ – which is constant. The geometric mean of the three factors is

$$(4x(1 - 2x)(1 - 2x))^{1/3} = (4V)^{1/3}.$$

So by the AM–GM theorem:

$$\text{AM} = \frac{2}{3} \geq (4V)^{1/3} = \text{GM},$$

so

$$V \leq \frac{1}{4} \left(\frac{2}{3} \right)^3 = \frac{2}{27}.$$

The volume equals this constant maximum value when the three factors $4x$, $1 - 2x$, and $1 - 2x$ are equal. This equality happens when $x = 1/6$, which is the size of the corner cutouts.

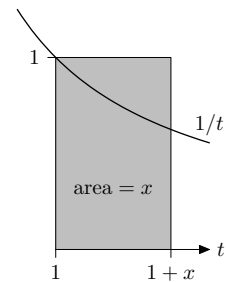
4.4 Logarithms

Pictures explain the early terms in many Taylor-series approximations. As an example, I derive the first two terms for $\ln(1+x)$. The logarithm function is defined as an integral

$$\ln(1+x) = \int_1^{1+x} \frac{dt}{t}.$$

An integral, especially a definite integral, suggests an area as its picture. As a first approximation, the logarithm is the area of the shaded, circumscribed rectangle. The rectangle, although it overestimates the integral, is easy to analyze: Its area is its width (which is x) times its height (which is 1). So the area is x . This area is the first pictorial approximation, and explains the first term in the Taylor series

$$\ln(1+x) = x - \dots$$

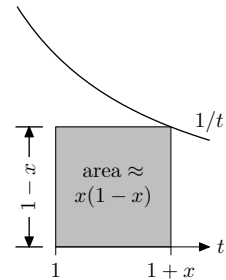


An alternative to overestimating the integral is to underestimate it using the inscribed rectangle. Its width is still x but its height is $1/(1+x)$. For small x ,

$$\frac{1}{1+x} \approx 1-x,$$

as you can check by multiplying both sides by $1+x$:

$$1 \approx 1-x^2.$$

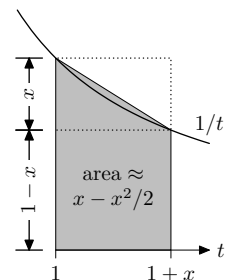


This approximation is valid when x^2 is small, which happens when x is small. Then the rectangle's height is $1-x$ and its area is $x(1-x) = x - x^2$.

For the second approximation, average the over- and underestimate:

$$\ln(1+x) \equiv \text{area} \approx \frac{x + (x - x^2)}{2} = x - \frac{x^2}{2}.$$

These terms are the first two terms in the Taylor series for $\ln(1+x)$. The picture for this symbolic average is a trapezoidal area, so this series of pictures explains the first two terms. Its error lies in making the smooth curve $1/t$ into a straight line, and this error produces the higher-order terms in the series – but they are difficult to compute just using pictures.



Picture proofs

52

Alternatively you can derive all the terms from the binomial theorem and the definition of the logarithm. The logarithm is

$$\ln(1+x) \equiv \int_1^{1+x} \frac{dt}{t} = \int_0^x \frac{1}{1+t} dt.$$

The binomial theorem says that

$$\frac{1}{1+t} = 1 - t + t^2 - t^3 + \dots,$$

so

$$\ln(1+x) = \int_0^x (1 - t + t^2 - t^3 + \dots) dt.$$

Now integrate term by term; although this procedure produces much gnashing of the teeth among mathematicians, it is usually valid. To paraphrase a motto of the Chicago police department, 'Integrate first, ask questions later.' Then

$$\ln(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots.$$

The term-by-term integration shows you the entire series. Understand both methods and you will not only remember the logarithm series but will also understand two useful techniques.

As an application of the logarithm approximation, I estimate $\ln 2$. A quick application of the first two terms of the series gives:

$$\ln(1+x) \approx x - \frac{x^2}{2} \Big|_{x=1} = 1 - \frac{1}{2} = \frac{1}{2}.$$

That approximation is lousy because x is 1, so squaring x does not help produce a small $x^2/2$ term. A trick, however, improves the accuracy. Rewrite $\ln 2$ as

$$\ln 2 = \ln \frac{4/3}{2/3} = \ln \frac{4}{3} - \ln \frac{2}{3}.$$

Then approximate $\ln(4/3)$ as $\ln(1+x)$ with $x = 1/3$ and approximate $\ln(2/3)$ as $\ln(1+x)$ with $x = -1/3$. With $x = \pm 1/3$, squaring x produces a small number, so the error should shrink. Try it:

$$\begin{aligned} \ln \frac{4}{3} &= \ln(1+x) \Big|_{x=1/3} \approx \frac{1}{3} - \frac{1}{2} \cdot \left(\frac{1}{3}\right)^2, \\ \ln \frac{2}{3} &= \ln(1+x) \Big|_{x=-1/3} \approx -\frac{1}{3} - \frac{1}{2} \cdot \left(-\frac{1}{3}\right)^2. \end{aligned}$$

4.5 Geometry

53

When taking the difference, the quadratic terms cancel, so

$$\ln 2 = \ln \frac{4}{3} - \ln \frac{2}{3} \approx \frac{2}{3} = 0.666\dots$$

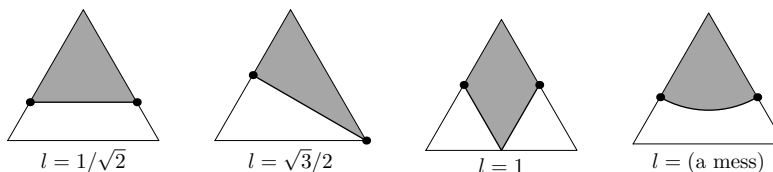
The true value is $0.697\dots$, so this estimate is accurate to 5%!

4.5 Geometry

The following pictorial problem has a natural pictorial solution:

How do you cut an equilateral triangle into two equal halves using the shortest, not-necessarily-straight path?

Here are several candidates among the infinite set of possibilities for the path.



Let's compute the lengths of each bisecting path, with length measured in units of the triangle side. The first candidate encloses an equilateral triangle with one-half the area of the original triangle, so the sides of the smaller, shaded triangle are smaller by a factor of $\sqrt{2}$. Thus the path, being one of those sides, has length $1/\sqrt{2}$. In the second choice, the path is an altitude of the original triangle, which means its length is $\sqrt{3}/2$, so it is longer than the first candidate. The third candidate encloses a diamond made from two small equilateral triangles. Each small triangle has one-fourth the area of the original triangle with side length one, so each small triangle has side length $1/2$. The bisecting path is two sides of a small triangle, so its length is 1. This candidate is longer than the other two.

The fourth candidate is one-sixth of a circle. To find its length, find the radius r of the circle. One-sixth of the circle has one-half the area of the triangle, so

$$\underbrace{\pi r^2}_{A_{\text{circle}}} = 6 \times \frac{1}{2} A_{\text{triangle}} = 6 \times \frac{1}{2} \times \underbrace{\frac{1}{2} \times 1 \times \frac{\sqrt{3}}{2}}_{A_{\text{triangle}}}.$$

Picture proofs

54

Multiplying the pieces gives

$$\pi r^2 = \frac{3\sqrt{3}}{4},$$

and

$$r = \sqrt{\frac{3\sqrt{3}}{4\pi}}.$$

The bisection path is one-sixth of a circle, so its length is

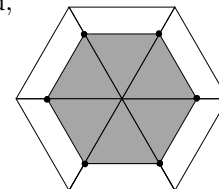
$$l = \frac{2\pi r}{6} = \frac{\pi}{3} \sqrt{\frac{3\sqrt{3}}{4\pi}} = \sqrt{\frac{\pi\sqrt{3}}{12}}.$$

The best previous candidate (the first picture) has length $1/\sqrt{2} = 0.707\dots$. Does the mess of π and square roots produce a shorter path? Roll the drums. . . :

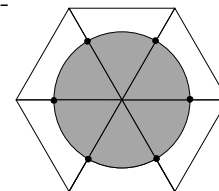
$$l = 0.67338\dots,$$

which is less than $1/\sqrt{2}$. So the circular arc is the best bisection path *so far*. However, is it the best among all possible paths? The arc-length calculation for the circle is messy, and most other paths do not even have a closed form for their arc lengths.

Instead of making elaborate calculations, try a familiar method, symmetry, in combination with a picture. Replicate the triangle six times to make a hexagon, and also replicate the candidate path. Here is the result of replicating the first candidate (the bisection line going straight across). The original triangle becomes the large hexagon, and the enclosed half-triangle becomes a smaller hexagon having one-half the area of the large hexagon.



Compare that picture with the result of replicating the circular-arc bisection. The large hexagon is the same as for the last replication, but now the bisected area replicates into a circle. Which has the shorter perimeter, the shaded hexagon or this circle? The **isoperimetric theorem** says that, of all figures with the same area, the circle has the smallest perimeter. Since the circle and the smaller hexagon enclose the same area – which is three times the area of one triangle – the circle has a smaller perimeter than the hexagon, and has a smaller perimeter than the result of replicating any other path!



4.6 Summing series

Now let's look for a second time at Stirling's approximation to n factorial. In [Section 3.5](#), we found it by approximating the integral

$$\int_0^{\infty} t^n e^{-t} dt = n!.$$

The next method is also indirect, by approximating $\ln n!$:

$$\ln n! = \sum_1^n \ln k.$$

This sum is the area of the rectangles. That area is roughly the area under the smooth curve $\ln k$. This area is

$$\int_1^n \ln k dk = k \ln k - k = n \ln n - n + 1.$$

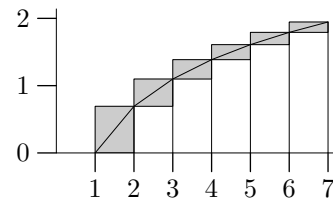
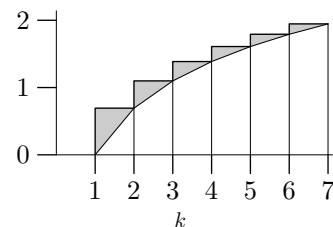
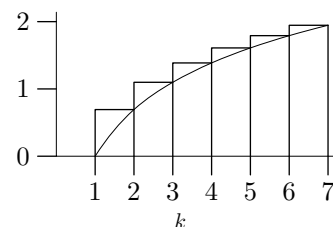
Before making more accurate approximations, let's see how this one is doing by taking the exponential to recover $n!$:

$$n! \approx \frac{n^n}{e^n} \times e.$$

The n^n and the $1/e^n$ factors are already correct. The next pictorial correction make the result even more accurate.

The error in the integral approximation come from the pieces protruding beyond the $\ln k$ curve. To approximate the area of these protrusions, pretend that they are triangles. If $\ln k$ were made of linear segments, there would be no need to pretend; even so the pretense is only a tiny lie. The problem become one of adding up the shaded triangles.

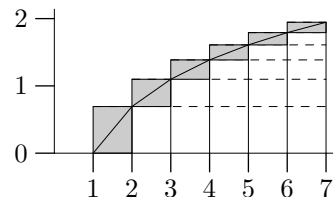
The next step is to double the triangles, turning them into rectangles, and remembering to repay the factor of 2 before the end of the derivation.



Picture proofs

56

The final step is to hold your right hand at the $x = 7$ line to catch the shaded pieces as you shove them rightward with your left hand. They stack to make the $\ln 7$ rectangle. So the total overshoot, after paying back the factor of 2, is $(\ln 7)/2$. For general n , the overshoot is $(\ln n)/2$. The integral $\int_1^n \ln k \, dk$ provides $n \ln n - n$ (from the upper limit) and 1 from the lower limit. So the integral and graph together produce



$$\ln n! \approx n \ln n - n + 1 + \underbrace{\frac{\ln n}{2}}_{\text{protrusions}}$$

or

$$n! \approx e\sqrt{n} \left(\frac{n}{e}\right)^n.$$

Stirling's formula is

$$n! \approx \sqrt{2\pi n} \left(\frac{n}{e}\right)^n.$$

The difference between the pictorial approximation and Stirling's formula is the factor of e that should be $\sqrt{2\pi}$. Except for this change of only 8%, a simple integration and graphical method produce the whole formula.

The protrusion correction turns out to be the first term in an infinite series of corrections. The later corrections are difficult to derive using pictures, just as the later terms in the Taylor series for $\ln(1+x)$ are difficult to derive by pictures (we used integration and the binomial theorem for those terms). But another technique, analogy, produces the higher corrections for $\ln n!$. That analysis is the subject of [Section 7.3](#), where the pictorial, protrusion correction that we just derived turns out to be the zeroth-derivative term in the Euler–MacLaurin summation formula.

5

Taking out the big part

Taking out the big part, the technique of this chapter, is a species of successive approximation. First do the most important part of the analysis: the big part. Then estimate changes relative to this big part. This hygienic approach keeps calculations clean enough to do mentally. Here are a few examples beginning with products, powers, and roots, then moving to exponentials and fierce integrals.

5.1 Multiplication

Suppose you have to estimate 31.5×721 . A first estimate comes from rounding 31.5 to 30 and 721 to 700:

$$31.5 \times 721 \approx 30 \times 700 = 21000.$$

This product is the big part whose estimation is the first step. In the second step, estimate the correction. You could estimate the correction directly by expanding the product:

$$31.5 \times 721 = (30 + 1.5) \times (700 + 21).$$

Expanding produces four terms:

$$30 \times 700 + 1.5 \times 700 + 30 \times 21 + 1.5 \times 21.$$

Taking out the big part

58

What a mess! Using fractional or relative changes cleans up the calculation. The first step is to estimate the fractional change in each factor: 31.5 is 5% more than 30, and 721 is 3% more than 700. So

$$31.5 \times 721 = \underbrace{30 \times (1 + 0.05)}_{31.5} \times \underbrace{700 \times (1 + 0.03)}_{721}.$$

Reorder the pieces to combine the fractional changes:

$$\underbrace{30 \times 700}_{\text{big part}} \times \underbrace{(1 + 0.05) \times (1 + 0.03)}_{\text{correction factor}}.$$

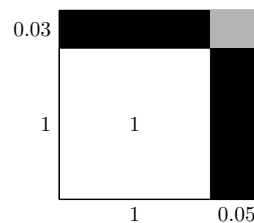
The big part is already evaluated, so the problem reduces to estimating the correction factor. An algebraic method gives

$$(1 + 0.05) \times (1 + 0.03) = 1 \times 1 + 0.05 \times 1 + 1 \times 0.03 + \underbrace{0.05 \times 0.03}_{\text{tiny}}.$$

Because the last term is the product of two corrections, each small, it is smaller than the other terms. Ignoring it gives

$$(1 + 0.05) \times (1 + 0.03) \approx 1 + 0.05 + 0.03 = 1.08.$$

This algebra has an elegant picture. The unit square represents the 1×1 product. Enlarge its width by 0.05 to $1 + 0.05$, and enlarge its height by 0.03 to $1 + 0.03$. The new rectangle has area $(1 + 0.05) \times (1 + 0.03)$, which is the sought-after product. The four pieces of the figure correspond to the four terms in the expansion of $(1 + 0.05) \times (1 + 0.03)$. Relative to the unit square, the new rectangle has a thin rectangle on the right that has area 0.05 and a thin rectangle on top that has area 0.03. There's also an adjustment of 0.05×0.03 for the gray rectangle. It is tiny compared to the long rectangles, so neglect it. Then the area is roughly $1 + 0.05 + 0.03$, which is a geometric proof that the correction factor is roughly



$$1 + 0.05 + 0.03 = 1.08.$$

It represents an 8% increase. The uncorrected product is 21000, and 8% of it is 1680, so

$$31.5 \times 721 = 21000 \times \text{correction factor} \approx 21000 + 1680 = 22680.$$

5.1 Multiplication

59

The true value is 22711.5, so the estimate is low by 0.15%, which is the area of the tiny, gray rectangle.

This numerical example illustrates a general pattern. Suppose that you can easily find the product xy , as in the preceding example with $x = 30$ and $y = 700$, and you want a nearby product $(x + \Delta x)(y + \Delta y)$, where $\Delta x \ll x$ and $\Delta y \ll y$. Call $\Delta(xy)$ the change in the product xy due to the changes in x and y :

$$(x + \Delta x)(y + \Delta y) = xy + \Delta(xy).$$

To find the new product, you could find $\Delta(xy)$ (since xy is easy). But do not expand the product directly:

$$(x + \Delta x)(y + \Delta y) = xy + x\Delta y + y\Delta x + \Delta x \Delta y.$$

Instead, extract the big part of the product and study the correction factor. The big part is xy , so extract xy by extracting x from the first factor and y from the second factor. The correction factor that remains is

$$\left(1 + \frac{\Delta x}{x}\right) \left(1 + \frac{\Delta y}{y}\right) = 1 + \underbrace{\frac{\Delta x}{x} + \frac{\Delta y}{y} + \frac{\Delta x \Delta y}{x y}}_{\text{frac. change in } xy}.$$

The $\Delta x/x$ is the fractional change in x . The $\Delta y/y$ is the fractional change in y . And the $(\Delta x/x)(\Delta y/y)$, the product of two tiny factors, is tiny compared to fractional changes containing one tiny factor. So, for small changes:

$$\begin{aligned} \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } xy \end{array}\right) &\simeq \frac{\Delta x}{x} + \frac{\Delta y}{y} \\ &= \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } x \end{array}\right) + \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } y \end{array}\right). \end{aligned}$$

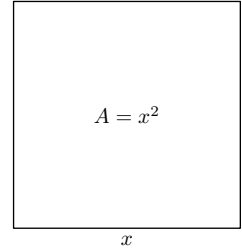
In other words, for small changes:

The fractional change in a product is the sum of fractional changes in its factors.

The simplicity of this rule means that fractional changes simplify computations.

5.2 Squares

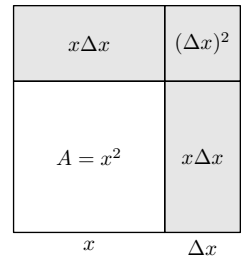
Squares are a particular kind of product, so we could approximate squares using the preceding algebra or pictures. Instead I derive the picture from scratch, to practice with pictures and to introduce the notion of low-entropy expressions. Let A be the area of a square and x be the length of its side, so $A = x^2$. Now imagine increasing x to $x + \Delta x$, producing an area $(x + \Delta x)^2$. This analysis is useful if you can choose x to be a number whose square you know; then Δx is the change to get to the number whose square you want to compute. For example, if you want to compute 9.2^2 , set $x = 9$ and $\Delta x = 0.2$ and find how much the area increases. The algebraic approach is to expand



$$(x + \Delta x)^2 = x^2 + 2x\Delta x + (\Delta x)^2.$$

An alternative approach is to elaborate the picture.

The new area is shaded and has three parts. As long as $\Delta x \ll x$, the tiny corner square is small compared to the two rectangles. So the change in area is



$$\Delta A \approx \underbrace{x\Delta x}_{\text{top rect.}} + \underbrace{x\Delta x}_{\text{right rect.}} = 2x\Delta x.$$

But this result is difficult to remember because it has *high-entropy* [6]. The combination of x and Δx seem arbitrary. If ΔA had turned out to be x^2 or $(\Delta x)^2$, it would also have seemed reasonable. A high-entropy form has variables scattered all over, in a seemingly unconstrained arrangement. A low-entropy form groups together relevant variables to make a form that is easy to understand and therefore to remember.

To turn $\Delta A = 2x\Delta x$ into low-entropy form, divide by $A = x^2$. This choice has two reasons. The first reason is the theme of this chapter: take out the big part. You know how to square x , so A or x^2 is the big part. To take it out, divide the left side ΔA by A and the right side $2x\Delta x$ by x^2 . The second reason comes the method of **Chapter 1**: dimensions. There are many dimensions in the world, so requiring an expression to be dimensionless eliminates this freedom and reduces the entropy:

5.2 Squares

61

Expressions with dimensions have higher entropy than expressions without dimensions.

The high-entropy result has dimensions of area; to make it dimensionless, divide both sides by an area. For the left side ΔA , the natural, related quantity is the area A . For the right side $2x\Delta x$, the natural, related quantity is the area x^2 . So two reasons – taking out the big part and dimensions – suggest dividing by $A = x^2$. A method with two justifications is probably sound, and here is the result:

$$\frac{\Delta A}{A} \approx \frac{2x\Delta x}{x^2} = 2\frac{\Delta x}{x}.$$

Each side has a simple interpretation. The left side, $\Delta A/A$, is the fractional change in area. The right side contains $\Delta x/x$, which is the fractional change in side length. So

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } x^2 \end{array} \right) \approx 2 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } x \end{array} \right).$$

This statement of the result is easier to understand than the high-entropy form. It says that fractional changes produce fractional changes. The only seemingly arbitrary datum to remember is the factor of 2, but it too will make sense after studying cubes and square roots.

Meanwhile you might be tempted into guessing that, because $A = x^2$, the fractional changes follow the same pattern:

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } A \end{array} \right) \approx \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } x \end{array} \right)^2.$$

That reasonable conjecture is wrong! Try a numerical example. Imagine a 10% increase in x , from 1 to 1.1. Then x^2 increases to roughly 1.2, a fractional increase of 0.2. If the candidate formula above were correct, the fractional increase would be only 0.01.

Let's finish the study of squares with 9.2^2 , the numerical example mentioned before. Its big part is $9^2 = 81$. Going from 9 to 9.2 is a fractional increase of $2/90$, so 9.2^2 should increase by $2 \times 2/90 = 4/90$:

$$9.2^2 \approx 81 \times \left(1 + \frac{4}{90} \right) \approx 81 + 3.6 = 84.6.$$

Taking out the big part

62

The exact answer is 84.64, a mere 0.05% higher.

5.3 Fuel efficiency

Section 2.7 used dimensional analysis and an experiment of dropping paper cones to show that drag force is proportional to v^2 , where v is the speed that an object moves through a fluid. This result applied in the limit of high Reynolds number, which is the case for almost all flows in our everyday experience. Highway driving is at a roughly steady speed, so gasoline is burned in fighting drag rather than in lossy, stop-and-go changes of speed. The energy required for a car to travel a distance d at speed v is then

$$E = Fd \propto v^2 d,$$

where F is the drag force. In the 1970's, oil became expensive in Western countries for reasons that were widely misunderstood and often misexplained (maybe intentionally). For a thorough analysis, see [7]. Whatever the causes, the results were hard to avoid. The United States reduced oil consumption by mandating a speed limit of 55 mph on highways. For the sake of this problem, imagine that cars drove at 65 mph before the speed limit was imposed. *By what fraction does the gasoline consumption fall due to the change in speed from 65 to 55 mph?* Pretend that the speed limit does not affect how far people drive. It may be a dubious assumption, since people regulate their commuting by total time rather than distance, but that twist can be the subject of a subsequent analysis (do the big part first).

Fractional changes keep the analysis hygienic. The drag force and the energy consumption are proportional to $v^2 d$, and the distance d is, by assumption, constant. So $E \propto v^2$ and

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } E \end{array} \right) = 2 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } v \end{array} \right).$$

A drop in v from 65 to 55 mph is a drop of roughly 15% so the energy consumption drops by $2 \times 15\% = 30\%$. It is a large reduction in automotive oil consumption. Considering the large fraction of oil consumed by car travel, this 30% drop in highway oil consumption produces a substantial reduction in total oil consumption.

5.4 Third powers

The next example extends the analysis to the volume of a cube with side length x . The usual question recurs: If x increases by Δx , what happens to the volume V ? If you do not use fractional changes, you can try to guess what happens by analogy with the change in area. Perhaps

$$\Delta V \sim x^2 \Delta x$$

or maybe

$$\Delta V \sim x(\Delta x)^2?$$

Both choices have a volume on each side, so their dimensions are correct, and dimensions do not favor either choice. In short, it's a pain to remember how to distribute the three powers of length on the right side. Should the x get all of them, two of them, one of them, or none?

Instead of trying to remember the high-entropy form, work it out from scratch, rewrite it as a fractional change, and see how simple and low-entropy it becomes. The full ΔV is

$$\Delta V = (x + \Delta x)^3 - x^3 = 3x^2 \Delta x + 3x(\Delta x)^2 + (\Delta x)^3.$$

The terms with the higher powers of Δx are the smallest, so ignore them. This approximation leaves

$$\Delta V \approx 3x^2 \Delta x.$$

The fractional change is

$$\frac{\Delta V}{V} \approx \frac{3x^2 \Delta x}{x^3} = 3 \frac{\Delta x}{x}.$$

This result has the same form as the fractional change in area but with a factor of 3. In words:

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } x^3 \end{array} \right) \approx 3 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } x \end{array} \right).$$

The factor of 3 comes from the exponent of x in $V = x^3$, just as the 2 came from the exponent of x in $A = x^2$. Let's look at two examples.

For the first example, estimate 6.3^3 . The big part is $6^3 = 216$. Since 0.3 is 5% larger than 6, its cube is $3 \times 5\% = 15\%$ larger than 6^3 :

Taking out the big part**64**

$$6.3^3 \approx 216 \times (1 + 0.15).$$

To calculate 216×0.15 , first calculate the big part 200×0.15 , which is 30. Then increase the result by 8% of 30, because 216 is 8% larger than 200. Since 8% of 30 is 2.4:

$$216 \times 0.15 = 30 + 2.4 = 32.4$$

Then

$$6.3^3 \approx 216 + 32.4 = 248.4.$$

The true value is 250.047, which is only 0.7% larger.

The second example comes from the physics of wind energy. The power produced by a wind turbine is related to the force exerted by the wind, which is (like the drag force) proportional to v^2 . Since power is force times velocity, it should be proportional to v^3 . Therefore a 10% increase in wind speed increases generated power by 30%! The hunt for fast winds is one reason that wind turbines are placed high in the atmosphere (for example, on cliffs) or at sea, where winds are faster than near land surfaces.

5.5 Reciprocals

The preceding examples used positive exponents. To explore fractional changes in new territory, try a negative exponent. This example is about the simplest one: reciprocals, where $n = -1$. Suppose that you want to estimate $1/13$ mentally. The big part is $1/10$ because 10 is a nearby factor of 10, which means its reciprocal is easy. So $1/13 \approx 0.1$. To get a more accurate approximation, take out the big part $1/10$ and approximate the correction factor:

$$\frac{1}{13} = \frac{1}{10} \times \frac{1}{1 + 0.3}.$$

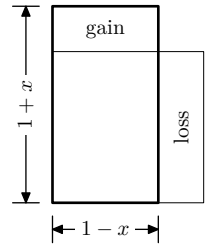
The correction factor is close to 1, reflecting that most of the result is in the big part $1/10$. The correction factor has the form $(1+x)^{-1}$, where $x = 0.3$. It is therefore approximately $1 - x$ as I hope the following example and picture will convince you. If a book is discounted 10% and shipping costs add 10% of the discounted price, the final total is almost exactly the original price. Try an example with a \$20 book. It gets reduced to \$18 but shipping adds \$1.80, for a total of \$19.80. Except for the tiny error of \$0.20, a 10% increase and a 10% decrease cancel each other. In general

5.5 Reciprocals

65

$$\underbrace{(1-x)}_{\text{decrease}} \times \underbrace{(1+x)}_{\text{increase}} \approx 1.$$

The picture confirms the algebra. Relative to the original unit square, the new $(1-x) \times (1+x)$ rectangle loses a rectangle on the right with area x and gains a rectangle on the top, also with area x . So the gain cancels the loss, keeping the area at 1. The error in this tally is the tiny square with area x^2 ; however, as long as x^2 is small, do not worry. That pictorial approximation leads to



$$\frac{1}{1+x} \approx 1-x.$$

In words,

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^{-1} \end{array} \right) = -1 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right).$$

If z increases by 30%, from 1 to 1.3, then z^{-1} decreases by 30%, from 1 to 0.7. So $1/1.3 = 0.7$ and

$$\frac{1}{13} = \frac{1}{10} \times \frac{1}{1.3} \approx 0.1 \times 0.7 = 0.07.$$

The error in the approximation comes from the neglected x^2 term in the reciprocal $(1+x)^{-1}$. To reduce the error, reduce x by making the big part a close approximation. Massage the original fraction to make the denominator close to 1/100:

$$\frac{1}{13} \times \frac{8}{8} = \frac{8}{104} = \frac{8}{100} \times \frac{1}{1.04}.$$

The big part $8/100 = 0.08$ is still easy, and the correction factor 1.04 has a smaller x : only 0.04. A 4% increase in a denominator produces a 4% decrease in the quantity itself, so

$$\frac{1}{13} \approx 0.08 - 4\%,$$

where the -4% means 'subtract 4% of the previous quantity'. To find the 4%, mentally rewrite 0.08 as 0.0800. Since 4% of 800 is 32, reduce the 0.08 by 0.0032:

$$\frac{1}{13} \approx 0.0800 - 0.0032 = 0.0768.$$

Taking out the big part

66

To make an even more accurate value, multiply $1/13$ by $77/77$ to get $77/1001$. The big part is 0.077 and the correction factor is a reduction by 0.1% , which is 0.00077 . The result is 0.076923 . For comparison, the true value is $.0769230769\dots$

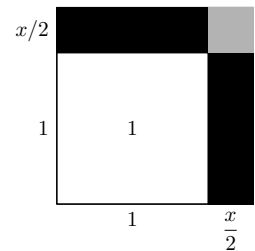
The second application follows up the reduction in gasoline consumption due to a 55-mph speed limit, analyzed in [Section 5.3](#). How much does the reduction in energy consumption increase fuel efficiency? Fuel efficiency is inversely proportional to energy consumption, so the -30% change in energy consumption produces a $+30\%$ change in fuel efficiency. It is often measured in miles per gallon, and a typical value for highway driving may be 35 mph. The 55 mph speed limit would increase it to roughly 45 mph, a larger increase than the legally mandated engineering increases over the last few decades.

5.6 Square roots

After positive and negative integer exponents, the next frontier is fractional exponents. The most common example is square roots, so let's apply these methods to $\sqrt{10}$. First take out the big part from $\sqrt{10}$. The big part is from the number whose square root is easy, which is 9. So factor out $\sqrt{9}$:

$$\sqrt{10} = \sqrt{9} \times \sqrt{1 + \frac{1}{9}}.$$

The problem reduces to estimating $\sqrt{1+x}$ with $x = 1/9$ in this case. Reversing the analysis for squaring in [Section 5.2](#) produces a recipe for square roots. For squaring, the problem was to find the area given the side length. Here the problem is to find the side length $\sqrt{1+x}$ given that the area is $1+x$. Relative to the unit square, the three shaded areas that make an L contribute the extra area x . The width of the vertical rectangle, or the height of the horizontal rectangle, is the change in side length. To find those dimensions, study the areas. Most of the contribution comes from the two dark rectangles, so ignore the tiny gray square. In that approximation, each rectangle contributes an area $x/2$. The rectangles measure $1 \times \Delta x$ or $\Delta x \times 1$, so their small dimension is roughly $\Delta x = x/2$. Thus the side length of the enclosing square is $1 + x/2$. This result produces the first square-root approximation:



$$\sqrt{1+x} \approx 1 + \frac{x}{2}.$$

5.6 Square roots

67

The right side represents a fractional increase of $x/2$, so

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } \sqrt{z} \end{array} \right) \simeq \frac{1}{2} \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right),$$

or in words

A fractional change in z produces one-half the fractional change in \sqrt{z} .

This result is the missing piece in estimating $\sqrt{10}$. The missing step was $\sqrt{1+x}$ with $x = 1/9$. Using the approximation,

$$\sqrt{1 + \frac{1}{9}} \approx 1 + \frac{1}{18}.$$

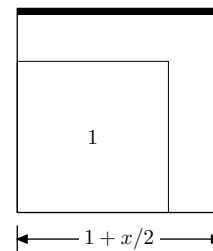
So increase the big part (which is 3) by $1/18$:

$$\sqrt{10} \approx 3 \times \left(1 + \frac{1}{18} \right) = 3 \frac{1}{6} = 3.166\dots$$

The true value is $3.1622\dots$; the estimate is accurate to 0.14%, a reasonable trade for three lines of work.

A few more lines and a refined picture increase the accuracy. The previous analysis ignored the tiny gray square. But now we know enough about the diagram to account for it, or at least to account for most of it. Neglecting the tiny square produced a square of side $1 + x/2$, which has area $1 + x$ plus the area of the tiny square. The tiny square is $x/2$ on each side so its area is $x^2/4$. The error in the first approximation $\sqrt{1+x} = 1 + x/2$ arises from this extra area.

To fix the approximation, shrink the big square slightly, just enough to remove an L-shaped shaded piece with area $x^2/4$. The dimensions of the L cannot be determined exactly – or else we could take square roots exactly – but it is solvable almost exactly using the knowledge from the earlier approximations. The analysis is by successive approximations. The L has two arms, each almost a thin rectangle that is as long or tall as the whole square, which means a length of $1 + x/2$. The ‘almost’ comes from ignoring the miniscule corner square where the two arms overlap. In this approximation, each arm has area $x^2/8$ in order that the L have area $x^2/4$. Since each sliver has length $1 + x/2$, the widths are



Taking out the big part

68

$$\text{width} = \frac{\text{area}}{\text{length}} = \frac{x^2/8}{1+x/2}.$$

The $1+x/2$ in the denominator is a fractional increase in the denominator of $x/2$, so it is a fractional decrease of $x/2$ in the numerator:

$$\frac{x^2/8}{1+x/2} \approx \frac{x^2}{8} \left(1 - \frac{x}{2}\right) = \frac{x^2}{8} - \frac{x^3}{16}.$$

This result is the thin width of the either rectangle arm. So shrink each side of the old square by $x^2/8 - x^3/16$, giving the next approximation to $\sqrt{1+x}$:

$$\sqrt{1+x} = 1 + \frac{x}{2} - \frac{x^2}{8} + \frac{x^3}{16}.$$

The cubic term $x^3/16$ is a bonus. We tried to compute the approximation after $1+x/2$, which presumably would give the coefficient of the $x^2/8$ term, yet we get the x^3 coefficient for free!

For mental calculation, I often neglect the cubic term. And, consistent with taking out the big part, I represent the $x^2/8$ as an adjustment on the next biggest part, which is the $x/2$ term:

$$\sqrt{1+x} = 1 + \frac{x}{2} \left(1 - \frac{x}{4}\right).$$

This formula gives the next approximation for $\sqrt{10}$. The zeroth approximation is $\sqrt{10} = 3$, which is the big part. The next approximation includes the $x/2$ term to give

$$\sqrt{10} = 3 + \frac{1}{6}.$$

The correction is $1/6$. With $x = 1/9$, the correction needs reducing by $x/4 = 1/36$. Because $1/36$ of $1/6$ is $1/216$, the next approximation is

$$\sqrt{10} = 3 + \frac{1}{6} - \frac{1}{216}.$$

For $1/216$ use fractional changes to approximate it: 216 is 8% larger than 200, so

$$\frac{1}{216} \approx \underbrace{\frac{1}{200}}_{0.0050} - 8\%.$$

The percentage is not hard: $8\% \times 50 = 4$, so

$$\frac{1}{216} \approx \underbrace{0.0050}_{-0.0004} = 0.0046.$$

5.7 In general

69

Thus

$$\sqrt{10} \approx 3 + 0.1666 - 0.0046 \approx 3.1626.$$

The true value is $3.162277\dots$, so the estimate is accurate to 0.01%.

Estimating square roots often benefits from a trick to speed convergence of the series. To see the need for the trick, try to estimate $\sqrt{2}$ using the preceding approximations. The big part is $\sqrt{1}$, which is no help. What remains is the whole problem: $\sqrt{1+x}$ with $x = 1$. Its first approximation is

$$\sqrt{2} \approx 1 + \frac{x}{2} = \frac{3}{2}.$$

Compared to the true value $1.414\dots$ this approximation is large by 6%. The next approximation includes the $x^2/8$ term:

$$\sqrt{2} \approx 1 + \frac{x}{2} - \frac{x^2}{8} = \frac{11}{8} = 1.375,$$

which is small by roughly 3%. The convergence is slow because $x = 1$, so successive terms do not shrink much despite the growing powers of x . If only I could shrink x ! The following trick serves this purpose:

$$\sqrt{2} = \frac{\sqrt{4/3}}{\sqrt{2/3}}.$$

Each square root has the form $\sqrt{1+x}$ where $x = \pm 1/3$. Retain up to the $x/2$ term:

$$\sqrt{2} = \frac{\sqrt{4/3}}{\sqrt{2/3}} \approx \frac{1 + 1/6}{1 - 1/6} = \frac{7}{5} = 1.4.$$

This quick approximation is low by only 1%! With the $x^2/8$ correction for each square root, the approximation becomes $\sqrt{2} \approx 83/59 = 1.406\dots$, which is low by 0.5%. The extra effort to include the quadratic term is hardly worth only a factor of 2 in accuracy.

5.7 In general

Look at the patterns for fractional changes. Here they are, in the order that we studied them:

Taking out the big part

70

$$\begin{aligned} \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^2 \end{array} \right) &\simeq 2 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right), \\ \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^3 \end{array} \right) &\simeq 3 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right), \\ \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^{-1} \end{array} \right) &\simeq -1 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right), \\ \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^{1/2} \end{array} \right) &\simeq 1/2 \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right). \end{aligned}$$

The general pattern is

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^n \end{array} \right) \simeq n \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right).$$

Before trying to prove it, check an easy case that was not part of the data used to make the generalization: $n = 1$. The fractional changes in z and z^1 are identical, so the pattern works. You can also check it when n is a nonnegative integer. In that case, z^n is a product of n factors of z . The product principle from [Section 5.1](#) is that the fractional change in a product is the sum of fractional changes in its factors. With n identical factors, the sum is indeed n times the fractional change in each factor.

The shortest proof for general n is by logarithmic differentiation. As the name says: First take the logarithm and then differentiate. The logarithm of $f = z^n$ is $n \log z$. Differentiating, or rather taking the differential, gives

$$\frac{df}{f} = n \frac{dz}{z}.$$

That result is exact for infinitesimal changes ($dz = 0$). For finite changes, use Δz instead of dz and turn the equals sign into an \approx :

$$\frac{\Delta f}{f} \approx n \frac{\Delta z}{z},$$

which is the symbolic expression of the general pattern:

The fractional change in z^n is n times the fractional change in z .

5.8 Seasons

An application of these results is to evaluate a common explanation for seasons. It is often said that, because the earth is closer to the sun in the summer than in the winter, summers are warmer than winters. The earth–sun distance does vary throughout the year because the earth orbits in an ellipse rather than a circle. As the distance varies, so does the solar flux, which is the amount of solar energy per unit area hitting the surface. The flux radiates back to space as blackbody radiation, the subject of numerous physics textbooks. The blackbody flux is related to the surface temperature. So the changing the earth–sun distance changes the earth’s surface temperature. How large is the effect and is it enough to account for the seasons?

The cleanest analysis is, not surprisingly, via fractional changes starting with the fractional change in earth–sun distance. In polar coordinates, the equation of an ellipse is

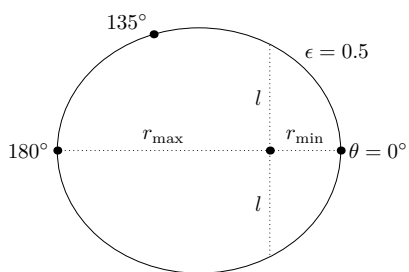
$$r = \frac{l}{1 + \epsilon \cos \theta},$$

where ϵ is the eccentricity, θ is the polar angle, and l is the semi-latus rectum (proportional to the angular momentum of the orbit). The diagram shows an orbit with eccentricity of 0.5, much exaggerated compared to the earth’s orbit in order to show the elliptical nature of the orbit. The distance varies from $r_{\min} = l/(1 + \epsilon)$ to $r_{\max} = l/(1 - \epsilon)$. Going from $r_{\min} = l/(1 + \epsilon)$ to l is a fractional increase of roughly ϵ . Going from l to $r_{\max} = l/(1 - \epsilon)$ is another fractional increase of ϵ , so the earth–sun distance varies by roughly 2ϵ . The earth’s orbit has $\epsilon = 0.016$ or 1.6%, meaning that the distance varies by 3.2%. As a check on that number, here is the relevant orbital data:

$$\begin{aligned} r_{\min} &= 1.471 \cdot 10^8 \text{ km}, \\ r_{\max} &= 1.521 \cdot 10^8 \text{ km}. \end{aligned}$$

These distances differ by roughly 3.2%.

The second step is to estimate the fractional change in flux produced by this fractional change in distance. The total solar power P spreads over a giant sphere with surface area $A = 4\pi d^2$. The power per area, which is flux, is $P/A \propto d^{-2}$. Because of the -2 exponent, a distance increase of 3.2% produces a flux decrease of 6.4%.



Taking out the big part

72

The third step is to estimate the fractional change in temperature produced by this fractional change in incoming flux. The outgoing flux is blackbody radiation, and it equals the incoming flux. So the outgoing flux also changes by 6.4%. Statistical mechanics – the Stefan–Boltzmann law – says that blackbody flux F is proportional to T^4 , where T is the surface temperature:

$$F = \sigma T^4.$$

The σ is the Stefan–Boltzmann constant, a ghastly combination of the quantum of action \hbar , the speed of light c , Boltzmann’s constant k_B , and $\pi^2/60$. But its composition is not relevant, because we are interested only in the fractional change in T . The freedom comes from using fractional changes, and is one of the most important reasons to use them. Since $T \propto F^{1/4}$, if flux changes by 6.4%, then T changes by 6.4%/4 or 1.6%. To find the actual change in temperature, multiply this percentage by the surface temperature T . Do not fall into the trap of thinking that, in winter anyway, the temperature is often 0 °C, so the change ΔT is also 0 °C! The blackbody flux $F \propto T^4$ depends on T being an absolute temperature: measured relative to absolute zero. On one such scale, the Kelvin scale, $T = 300$ K so a 1.6% variation is about 5 K. The reference points of the Celsius and Kelvin scales are different, but their degrees are the same size, so a 5 K difference is also a 5 °C difference. This change is too small to account for the difference between summer and winter, making the proposed explanation for seasons implausible. The explanation has other flaws, such as not explaining how Australia and Europe have opposite seasons despite being almost exactly equidistant from the sun. If orbital distance changes do not produce seasons, what does?

5.9 Exponentials

The preceding examples investigated the approximation

$$(1 + x)^n \simeq 1 + nx$$

where the exponent n was a positive integer, negative integer, and even a fraction. The examples used moderate exponents: 1/2 for the square roots, -1 for reciprocals, and -2 and $1/4$ for the seasons. Now push n to an extreme, but skillfully. If you simply make n huge, then you end up evaluating quantities like 1.1^{800} , which is not instructive. Instead, let n grow but shrink

5.10 Extreme cases

73

x in parallel to keep nx fixed. An intuitive value for nx is 1, and these examples keep $nx = 1$ while increasing n :

$$\begin{aligned} 1.1^{10} &= 2.59374\dots, \\ 1.01^{100} &= 2.70481\dots, \\ 1.001^{1000} &= 2.71692\dots \end{aligned}$$

In each case, $nx = 1$ so the usual approximation is

$$(1+x)^n = 2 \approx 1 + nx = 2,$$

which is significantly wrong. The problem lies in nx growing too large. In the examples with moderate n , the product nx was much smaller than 1. So new mathematics happens when nx grows beyond that limited range.

To explain what happens, guess features of the solution and then find an explanation related to those features. The sequence starting with 1.1^{10} seems to approach $e = 2.718\dots$, the base of the natural logarithms. That limit suggests that we study not $(1+x)^n$ but rather its logarithm:

$$\ln(1+x)^n = n \ln(1+x).$$

As long as x itself is not large (nx can still be large), then $\ln(1+x) \approx x$. So $n \ln(1+x) \approx nx$ and

$$(1+x)^n \approx e^{nx}.$$

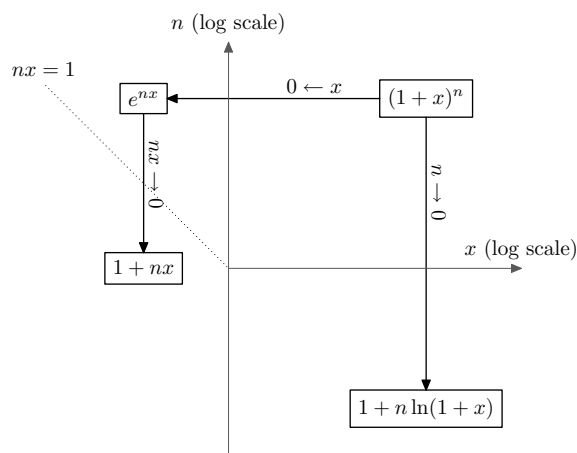
When $nx \ll 1$, then e^{nx} approximates to $1+nx$, which reproduces the familiar approximation $1+nx$. When nx grows large, the approximation $e^{nx} = 1+nx$ fails, and you have to use e^{nx} itself.

5.10 Extreme cases

The general n^{th} power $(1+x)^n$ has several extreme cases depending on n , x , and nx . One limit is taking $n \rightarrow 0$. Then $(1+x)^n$ turns into $1 + n \ln x$, whose proof is left as an exercise for you. The other two limits have been the subject of the preceding analyses. When $x \rightarrow 0$, the limit is e^{nx} . If $nx \rightarrow 0$ in addition $x \rightarrow 0$, then e^{nx} limits to $1 + nx$, which is the result from the first examples in this chapter. Here is a pictorial summary:

Taking out the big part

74



Here are a few numerical examples of these limits:

limit	x	n	$(1+x)^n \approx$
$n \rightarrow 0$	1	0.1	$1 + 0.1 \ln 2$
$x \rightarrow 0$	0.1	30	e^3
$x, nx \rightarrow 0$	0.1	3	1.3

These limits come in handy in the next problem.

5.11 Daunting integral

As a physics undergraduate, I spent many late nights in the department library eating pizza while doing problem sets. The graduate students, in the same boat for their courses, would share their favorite mathematics and physics problems, which included the following from the former USSR. The Landau institute for theoretical physics required an entrance exam of ‘mathematical preliminaries’. One preliminary was to evaluate

$$\int_{-\pi/2}^{\pi/2} \cos^{100} t \, dt$$

to within 5% in less than 5 minutes, without a calculator or computer! That $\cos^{100} t$ looks frightening. Normal techniques for trigonometric functions do not help. For example, this identity is useful when integrating $\cos^2 t$:

5.11 Daunting integral

75

$$\cos^2 t = \frac{1}{2}(\cos 2t + 1).$$

Here it would produce

$$\cos^{100} t = \left(\frac{\cos 2t + 1}{2} \right)^{50},$$

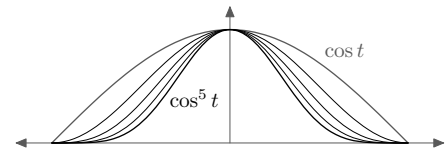
which becomes a trigonometric monster after expanding the 50th power. The answer is to approximate; after all, we need an answer accurate only to 5%. An approximation for $\cos t$ is $\cos t = 1 - t^2/2$. So

$$\cos^{100} t \simeq \left(1 - \frac{t^2}{2} \right)^{100},$$

which looks like $(1 + x)^n$ with $x = -t^2/2$ and $n = 100$. In the range $t \approx 0$ where the approximation for cosine is valid, it is the extreme case $x \rightarrow 0$ of $(1 + x)^n$, which is e^{nx} . So

$$\cos^{100} t = \left(1 - \frac{t^2}{2} \right)^{100} = e^{-50t^2}.$$

The integrand has the general form $e^{-\alpha t^2}$, which is the Gaussian analyzed in [Section 2.2](#) and [Section 3.4](#). This simple conclusion, that a high power of a cosine becomes a Gaussian, seems hard to believe, but the computer-generated plots of $\cos^n t$ for $n = 1 \dots 5$ show the cosine curve turning into the Gaussian bell shape as n increases. A plot is not a proof, but it increases confidence in a surprising result.



The argument has a few flaws but do not concern yourself with them now. Follow Bob Marley: Don't worry, be happy. In other words, approximate first and (maybe) ask questions later *after* getting an answer. To promote this *sang froid* or courage, I practice what I preach and defer the analysis of the flaws. If the limits were infinite, the integral would be

$$\int_{-\infty}^{\infty} e^{-\alpha t^2} dt,$$

which is doable. Alas, our limits are $-\pi/2$ to $\pi/2$ rather than from $-\infty$ to ∞ . Do not worry; just extend the limits and justify it at the end. The infinite-range integral of the Gaussian is

Taking out the big part

76

$$\int_{-\infty}^{\infty} e^{-\alpha t^2} dt = \sqrt{\frac{\pi}{\alpha}}.$$

For $\cos^{100}t$, the parameter is $\alpha = 50$ so the original integral becomes

$$\int_{-\pi/2}^{\pi/2} \cos^{100}t dt \approx \int_{-\infty}^{\infty} e^{-50t^2} dt = \sqrt{\frac{\pi}{50}}.$$

Since $50 \approx 16\pi$, the integral is $\sqrt{1/16} = 0.25$. The exact answer is

$$\int_{-\pi/2}^{\pi/2} \cos^n t dt = 2^{-n} \binom{n}{n/2} \pi,$$

whose proof I leave as a fun exercise for you. For $n = 100$, the result is

$$\frac{12611418068195524166851562157\pi}{158456325028528675187087900672} = 0.25003696348037\dots$$

The `maxima` program, which computed this exact rational-fraction multiple of π , is free software originally written at MIT as the `Macsyma` project. Using a recent laptop (circa 2006) with an Intel 1.83 GHz Core Duo CPU, `maxima` required roughly 20 milliseconds to compute the exact result. Our estimate of $1/4$ used a method that requires less than, say, thirty seconds of human time (with practice), and it is accurate to almost 0.01%. Not a bad showing for wetware.

In order to estimate accurately the computation times for such integrals, I tried a higher exponent:

$$\int_{-\pi/2}^{\pi/2} \cos^{10000}t dt$$

In 0.26 seconds, `maxima` returned a gigantic rational-fractional multiple of π . Converting it to a floating-point number gave $0.025065\dots$, which is almost exactly one-tenth of the previous answer. That rescaling makes sense: Increasing the exponent by a factor of 100 increases the denominator in the integral by $\sqrt{100} = 10$.

Now look at the promised flaws in the argument. Here are the steps in slow motion, along with their defects:

1. Approximate $\cos t$ by $1 - t^2/2$. This approximation is valid as long as $t \approx 0$. However, the integral ranges from $t = -\pi/2$ to $t = \pi/2$, taking t beyond the requirement $t \approx 0$.

5.11 Daunting integral

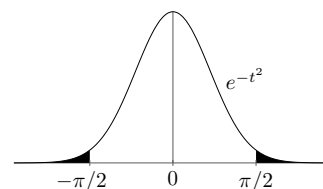
77

2. Approximate $(1 - t^2/2)^n$ as $e^{-nt^2/2}$. This approximation is valid when $t^2/2 \approx 0$. Again, however, t ranges beyond that limited domain.
3. Replace the difficult limits $-\pi/2 \dots \pi/2$ with the easier ones $-\infty \dots \infty$. The infinite limits permit the polar-coordinates trick of [Section 2.2](#) – which I call a trick because I’ve never seen a different problem that uses it. However, what justifies extending the limits?

All three flaws have an justification in the high exponent (100 in this case). Raising $\cos t$ to a high power means that the result is close to zero when $\cos t$ drops even slightly below 1. For example, when $t = 0.5$, its cosine is $0.877 \dots$ and $\cos^{100}t \approx 2 \cdot 10^{-6}$. The exponential approximation e^{-50t^2} is roughly $3.7 \cdot 10^{-6}$, which seems inaccurate: The error is almost 100%! But that error is a relative error or fractional error. The absolute error is roughly $2 \cdot 10^{-6}$. It is fine to make large relative errors where the integrand is tiny. In the region where the integrand contributes most of the area, which is $t \approx 0$, steps 1 and 2 of the approximation are valid. In the other regions, who cares?!

The same argument justifies the third step: extending the limits to infinity. It would be foolhardy to extend the limits in the original integral to give

$$\int_{-\infty}^{\infty} \cos^{100}t \, dt.$$



because each hump of $\cos^{100}t$ contributes equal area and the extended limits enclose an infinity of humps. But this objection disappears if you extend the limits after making the first two approximations. Those approximations give

$$\int_{-\pi/2}^{\pi/2} e^{-50t^2} \, dt.$$

Because the Gaussian e^{-50t^2} is miniscule at and beyond $t = \pm\pi/2$, it is safe to extend the limits to $-\infty \dots \infty$. The figure shows the tails of $e^{-t^2/2}$, and they are already small. In the faster-decaying function e^{-50t^2} , the tails are so miniscule that they would be invisible at any feasible printing resolution.

I do not want to finish the example with a verification. So try a small additional investigation. It arose because of the high accuracy of the approximation when 100 or 10000 is the exponent of the cosine. I wondered how well the approximation does in the other extreme case, when the exponent is small. To study the accuracy, define

Taking out the big part

78

$$f(n) \equiv \int_{-\pi/2}^{\pi/2} \cos^n t \, dt.$$

The preceding approximations produce the approximation

$$f_0(n) = \sqrt{\frac{2\pi}{n}},$$

as you can check by trying the exponents $n = 100$ and $n = 10000$. The fractional error is

$$\frac{f_0(n)}{f(n)} - 1.$$

Here are a few values computed by `maxima`:

n	$f_0(n)/f(n) - 1$
1	0.2533141373155
2	0.1283791670955
3	0.0854018818374
4	0.0638460810704
5	0.0509358530746
6	0.0423520253928
7	0.0362367256182
8	0.0316609527730
9	0.0281092532666
10	0.0252728978367
100	0.0025030858398
1000	0.0002500312109
10000	0.0000250003124

Particularly interesting is the small fractional error when $n = 1$, a case where you can confirm `maxima`'s calculation by hand. The exact integral is

$$f(1) = \int_{-\pi/2}^{\pi/2} \cos^1 t \, dt.$$

So $f(1) = 2$, which compares to the approximation $f_0(1) = \sqrt{2\pi} \approx 2.5$. Even with an exponent as small as $n = 1$, which invalidates each step in the approximation, the error is only 25%. With $n = 2$, the error is only 13% and from there it is, so to speak, all downhill.

5.12 What you have learned

Take out the big part, and use fractional changes to adjust the answer. Using this procedure keeps calculations hygienic. The fundamental formula is

$$(1 + x)^n \simeq 1 + nx,$$

or

$$\left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z^n \end{array} \right) \simeq n \times \left(\begin{array}{c} \text{fractional} \\ \text{change} \\ \text{in } z \end{array} \right).$$

When the exponent n times the fractional change x grows too large (becomes comparable with 1), you need a more accurate approximation:

$$(1 + x)^n \simeq e^{nx}.$$

6

Analogy

When the going gets tough, the tough lower their standards. It is the creed of the sloppy, the lazy, and any who want results. This chapter introduces a technique, reasoning by **analogy**, that embodies this maxim. It works well with extreme-case reasoning.

6.1 Tetrahedral bond angle

What is the bond angle in methane, CH_4 ? The carbon sits at the centroid of a regular tetrahedron, and the hydrogens sit at the vertices. Trigonometry and analytic geometry solve this problem, but let's try analogy. Three dimensions is hard to visualize and figure out, so lower your standards: Look for a simpler problem that preserves its essentials. What is essential is not always obvious, and you might solve several simpler variants before discovering those features.

Let's try the simplest change, going from three dimensions to two dimensions. The two-dimensional version of the problem is to find the bond angle in a planar molecule, for example NH_3 smashed into a plane. The bond angle is one-third of a full circle or 120° . The center of the bond angle is the centroid is the centroid of the object, so its location might be relevant in solving the problem. Who knows where a tetrahedron's centroid is; but the triangle has a centroid one-third of the way from one edge to the opposite vertex.

Here is a table with this data, where d is the number of dimensions. It's hard to generalize from such sparse data, reflected by the question marks in the tetrahedron row. Here is where extreme-cases

	d	centroid	θ
triangle	2	1/3	120°
tetrahedron	3	?	?

reasoning helps. You can get **free data** by extending the analogy to a yet

6.1 Tetrahedral bond angle

81

more extreme problem. If two dimensions are easier than three, then one dimension should be easier than two.

In one dimension, the object is a line. The centroid is one-half of the way between the endpoints. The bond angle is 180° . And the table now is more complete. The bond angle has several generalizations to $d = 3$, depending on what

<i>shape</i>	<i>d</i>	centroid	θ
line	1	1/2	180°
triangle	2	1/3	120°
tetrahedron	3	?	?

pattern underlies it. For example, if the pattern is $\theta = (240 - 60d)^\circ$, then $\theta(d = 3) = 60^\circ$. Having made a conjecture, it is important to test your conjecture. To conjecture and not to test – the great mathematician and mathematics teacher George Polya [8] says that to do so is the mark of a savage! So: Is that conjecture reasonable? It's dubious because, first, the angle is less than 90° . If the molecule were CH_6 , with the carbon at the center of a cube and the hydrogens at the faces of a cube, then the bond angle would be exactly 90° . With only four hydrogens, rather than six, the bond angle should be larger than 90° . So 60° seems to be a dubious conjecture. For a second reason that it is dubious, the try a more extreme case: four dimensions. Then, according to the $(240 - 60d)^\circ$ conjecture, the bond angle would be zero, which is nonsense. So the conjecture is dubious on several grounds.

Let's make another conjecture. What about $\theta = 360^\circ/(d+1)$? That conjecture fits $d = 1$ and $d = 2$. For $d = 3$ it predicts $\theta = 90^\circ$. By the reasoning that rejected the previous conjecture, this angle is too small. Furthermore, it means that for $d = 4$, the angle drops below 90° . That's also not reasonable.

To help find another conjecture, it's time for a new idea. Instead of guessing the bond angle directly, guess a function of it that makes it easier to guess.

<i>shape</i>	<i>d</i>	centroid	θ	$\cos \theta$
line	1	1/2	180°	-1
triangle	2	1/3	120°	-1/2
tetrahedron	3	?	?	?

The bond angle, if we solve it honestly, will come from the dot product of two vectors (the vectors from a vertex to the centroid of the opposite face). Dot products produce cosines, so perhaps $\cos \theta$ is easier to guess than θ itself. This idea adds a column to the table.

One possible pattern for $\cos \theta$ is -2^{1-d} , which fits the $d = 1$ and $d = 2$ data. For $d = 3$ it predicts $\cos \theta = -1/4$, which means $\theta > 90^\circ$, an excellent result. In the extreme case of $d \rightarrow \infty$ it predicts that $\theta \rightarrow 90^\circ$. Let's check

Analogy

82

that result. The d -dimensional problem has a carbon at the center and $d + 1$ hydrogens at the vertices of the object. That bond angle should be more than 90° : The problem with 90° bonds has $2d$ hydrogens, each at center of the $2d$ faces of a d -dimensional cube. And $d + 1$ hydrogens should be more spread out than $2d$ hydrogens. So the -2^{1-d} is not reasonable, although it got off to a good start.

To find another conjecture, look at the pattern in the centroid column. It is $1/(d + 1)$. So $1/(d + 1)$ or $1/d$ might be a reasonable fit for $\cos \theta$. Perhaps $\cos \theta = -1/d$? That fits the $d = 1$ and $d = 2$ data, and predicts $\cos \theta = -1/3$ and $\theta \approx 109.47^\circ$. The only problem is that this conjecture also predicts that $\theta \rightarrow 90^\circ$ as $d \rightarrow \infty$. So maybe that's okay?

Anyway, the more likely conjecture, because it respects the pattern in the centroid column, is that $\cos \theta = -1/d$. Let's see if we can check that. Yes! But first see if we can check the centroid conjecture, since the $\cos \theta$ one depends on it. And we can check that too. It says that the height is $1/(d + 1)$ of the way from the base. Hmm, $d + 1$ – that's how many hydrogens there are. And 1, the numerator, is how many hydrogens are not on the base. Indeed, the average height of the $d + 1$ vertices is $1/(d + 1)$ – which explains the centroid location.

Now, knowing where the centroid is, look at a cross-section of the tetrahedron. The cosine of the complement of θ is

$$\cos(180^\circ - \theta) = \frac{1/(d + 1)}{d/(d + 1)} = \frac{1}{d}.$$

Since $\cos \theta = -\cos(180^\circ - \theta)$, the result is

$$\cos \theta = -\frac{1}{d}.$$

The final verifications are elegant arguments, ones that we might not have thought of on first try. That's okay. Here's what friends who went to the US Math Olympiad training session told me they were taught: Find the answer by any cheap method that you can find; once you know, or are reasonably sure of the answer, you often can then find a more elegant method and never mention the original cheap methods.

I agree with that philosophy, except for one point. It is worthwhile mentioning the cheap methods, because, just as they were useful in this problem, they will be useful in other problems.

6.2 Steiner's plane problem

A famous problem is Steiner's plane problem: Into how many regions do five planes divide space? There are lots of answers to this question, some boring. If the planes are parallel, for example, they make six regions. If the planes are not parallel, the number grows. But the number of regions depends on how 'unparallel' the planes are. So assume that the planes are in a random orientation, to remove the chance of a potential region being wiped out by a silly coincidence.

Five planes are hard to imagine and hard to build. An analogous problem is the same question with four planes. That's still hard, however. So try three planes. That's also hard, so try two planes. That's easy: four regions. Don't forget the more extreme case of one plane: two regions. And more free data comes from the most extreme case of zero planes: one region. So, starting with $n = 0$ planes, the number of regions is: 1, 2, 4, ... Are those powers of two, and is the next number in the sequence 8? Start with two planes making four regions. Place the third plane to cut the other two, so that it splits each region into two pieces – making eight regions total. So 8 is indeed the next number. Is 16 and then 32 next? That is represented in the following table:

n	0	1	2	3	4	5
r	1	2	4	8	16?	32?

So we have a conjecture, an educated guess, for $n = 5$. Its conjectural nature is reflected in the question marks. But how to test it? We still cannot easily visualize four planes, let alone five planes.

Analogy comes to the rescue again. If fewer planes were easier to solve than more planes, fewer dimensions might also help. So let's study the same problem in two dimensions. What is the analogous problem that preserves the essentials? It cannot be placing n planes in a plane. Rather, we should also reduce the dimensionality of the placed object: Place n lines in a plane, in random orientations and positions. How many planar regions does that make? Having learnt the lesson of free data, start with $n = 0$ lines giving 1 region. One line makes two regions; two lines makes four regions. It looks like powers of two again.

Let's test it with three lines. Here's a picture. They make *seven* regions, not eight. So the conjecture fails. Let's do four lines and count carefully. That's 11 regions, remote from the next power of two, which would have been 16. Here are the results for the two-dimensional region:

Analogy

84

n	0	1	2	3	4	5
r	1	2	4	7	11	?

Let's combine the two- and three-dimensional data:

n	0	1	2	3	4	5
r_2	1	2	4	7	11	?
r_3	1	2	4	8	16?	32?

Now once again, use extreme cases and get free data. With data for two and three dimensions, why not include data for one dimension?! In one dimension the problem is, after putting n points on a line, how many regions (line segments) do they make? That's a fencepost problem, so be careful not to be off by one. When $n = 0$, there's only one segment – the whole infinite line. Each dot divides one segment into two, so it increases r by one. So there will be $r = n + 1$ regions.

n	0	1	2	3	4	5	n
r_1	1	2	3	4	5	6	$n + 1$
r_2	1	2	4	7	11	?	
r_3	1	2	4	8	16?	32?	

Now we have lots of data! Can you spot a pattern? Look at the connected entries, where $4 + 7 = 11$:

n	0	1	2	3	4	5	n
r_1	1	2	3	4	5	6	$n + 1$
r_2	1	2	4	7	11	?	
r_3	1	2	4	8	16?	32?	

That pattern holds wherever there is data to check it against. For example, $3 + 4 = 7$. Or $4 + 4 = 8$. If that's true, then in two dimensions when $n = 5$, then $r = 16$. In three dimensions, when $n = 4$, there are $r = 15$ regions (one less than the prediction of $r = 2^n$). And with five planes, there will be 26

6.2 Steiner's plane problem

85

regions. So, that's our conjecture, which now has lots going for it. Let's now be even more extreme and get one more row of free data: 0 dimensions. In 0 dimensions, the object is a point, and there's only one point no matter how many -1-dimensional objects subdivide it. So $r = 1$ always. Then:

n	0	1	2	3	4	5	n
r_0	1	1	1	1	1	1	1
r_1	1	2	3	4	5	6	$n + 1$
r_2	1	2	4	7	11	?	
r_3	1	2	4	8	16?	32?	

And the new row, for 0 dimensions, continues the pattern.

For fun let's fit polynomials to the data we have – before making the conjectured leap. The zeroth row is fit by $r = 1$, a zeroth-degree polynomial. The first row is fit by $r = n + 1$, a first-degree polynomial. A natural generalization of this pattern is that the second row should be fit by a second-degree polynomial: a quadratic. A quadratic requires three data points, so use $n = 0 \dots 2$. The polynomial that fits r_2 for these points is

$$r_2(n) = \frac{1}{2}n^2 + \frac{1}{2}n + 1.$$

Does this quadratic fit the other, certain data points? For $n = 3$, it predicts $r = 7$, which is right. For $n = 4$ it predicts $r = 11$, which is also right. So we can probably trust its prediction for $n = 5$, which is $r = 16$ – in agreement with the prediction from adding numbers.

Carrying this system farther, the third row should be fit by a cubic, which needs four points for its fit. The cubic, as you can check, that fits the first four points is

$$r_3(n) = \frac{1}{6}n^3 + \dots + 1$$

It predicts $r(4) = 15$ and $r(5) = 26$, so once again the previous conjectures for $r(5)$ get new support. And therefore so does the theory that predicted them.

But why is it true? That problem is left as an exercise for the reader.

7

Operators

This chapter is an extended example of an analogy. In the last chapter, the analogy was often between higher- and lower-dimensional versions of a problem. Here it is between operators and numbers.

7.1 Derivative operator

Here is a differential equation for the motion of a damped spring, in a suitable system of units:

$$\frac{d^2x}{dt^2} + 3\frac{dx}{dt} + x = 0,$$

where x is dimensionless position, and t is dimensionless time. Imagine x as the amplitude divided by the initial amplitude; and t as the time multiplied by the frequency (so it is radians of oscillation). The dx/dt term represents the friction, and its plus sign indicates that friction dissipates the system's energy. A useful shorthand for the d/dt is the operator D . It is an operator because it operates on an object – here a function – and returns another object. Using D , the spring's equation becomes

$$D^2x(t) + 3Dx(t) + x(t) = 0.$$

The tricky step is replacing d^2x/dt^2 by D^2x , as follows:

$$D^2x = D(Dx) = D\left(\frac{dx}{dt}\right) = \frac{d^2x}{dt^2}.$$

The analogy comes in solving the equation. Pretend that D is a number, and do to it what you would do with numbers. For example, factor the equation. First, factor out the $x(t)$ or x , then factor the polynomial in D :

7.2 Fun with derivatives

87

$$(D^2 + 3D + 1)x = (D + 2)(D + 1)x = 0.$$

This equation is satisfied if either $(D + 1)x = 0$ or $(D + 2)x = 0$. The first equation written in normal form, becomes

$$(D + 1)x = \frac{dx}{dt} + x = 0,$$

or $x = e^{-t}$ (give or take a constant). The second equation becomes

$$(D + 2)x = \frac{dx}{dt} + 2x = 0,$$

or $x = e^{-2t}$. So the equation has two solutions: $x = e^{-t}$ or e^{-2t} .

7.2 Fun with derivatives

The example above introduced D and its square, D^2 , the second derivative. You can do more with the operator D . You can cube it, take its logarithm, its reciprocal, and even its exponential. Let's look at the exponential e^D . It has a power series:

$$e^D = 1 + D + \frac{1}{2}D^2 + \frac{1}{6}D^3 + \dots$$

That's a new operator. Let's see what it does by letting it operating on a few functions. For example, apply it to $x = t$:

$$(1 + D + D^2/2 + \dots)t = t + 1 + 0 = t + 1.$$

And to $x = t^2$:

$$(1 + D + D^2/2 + D^3/6 + \dots)t^2 = t^2 + 2t + 1 + 0 = (t + 1)^2.$$

And to $x = t^3$:

$$(1 + D + D^2/2 + D^3/6 + D^4/24 + \dots)t^3 = t^3 + 3t^2 + 3t + 1 + 0 = (t + 1)^3.$$

It seems like, from these simple functions (extreme cases again), that $e^D x(t) = x(t + 1)$. You can show that for any power $x = t^n$, that

$$e^D t^n = (t + 1)^n.$$

Since any function can, pretty much, be written as a power series, and e^D is a linear operator, it acts the same on any function, not just on the powers.

Operators

88

So e^D is the successor operator: It turns the function $x(t)$ into the function $x(t+1)$.

Now that we know how to represent the successor operator in terms of derivatives, let's give it a name, S , and use that abstraction in finding sums.

7.3 Summation

Suppose you have a function $f(n)$ and you want to find the sum $\sum f(k)$. Never mind the limits for now. It's a new function of n , so summation, like integration, takes a function and produces another function. It is an operator, σ . Let's figure out how to represent it in terms of familiar operators. To keep it all straight, let's get the limits right. Let's define it this way:

$$F(n) = (\sum f)(n) = \sum_{-\infty}^n f(k).$$

So $f(n)$ goes into the maw of the summation operator and comes out as $F(n)$. Look at $SF(n)$. On the one hand, it is $F(n+1)$, since that's what S does. On the other hand, S is, by analogy, just a number, so let's swap it inside the definition of $F(n)$:

$$SF(n) = (\sum Sf)(n) = \sum_{-\infty}^n f(k+1).$$

The sum on the right is $F(n) + f(n+1)$, so

$$SF(n) - F(n) = f(n+1).$$

Now factor the $F(n)$ out, and replace it by σf :

$$((S-1)\sigma f)(n) = f(n+1).$$

So $(S-1)\sigma = S$, which is an implicit equation for the operator σ in terms of S . Now let's solve it:

$$\sigma = \frac{S}{S-1} = \frac{1}{1-S^{-1}}.$$

Since $S = e^D$, this becomes

$$\sigma = \frac{1}{1-e^{-D}}.$$

7.4 Euler sum

89

Again, remember that for our purposes D is just a number, so find the power series of the function on the right:

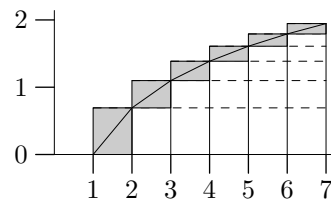
$$\sigma = D^{-1} + \frac{1}{2} + \frac{1}{12}D - \frac{1}{720}D^3 + \dots$$

The coefficients do not have an obvious pattern. But they are the Bernoulli numbers. Let's look at the terms one by one to see what the mean. First is D^{-1} , which is the inverse of D . Since D is the derivative operator, its inverse is the integral operator. So the first approximation to the sum is the integral – what we know from first-year calculus.

The first correction is $1/2$. Are we supposed to add $1/2$ to the integral, no matter what function we are summing? That interpretation cannot be right. And it isn't. The $1/2$ is one piece of an operator sum that is applied to a function. Take it in slow motion:

$$\sigma f(n) = \int_1^n f(k) dk + \frac{1}{2}f(n) + \dots$$

So the first correction is one-half of the final term $f(n)$. That is the result we got with this picture from [Section 4.6](#). That picture required approximating the excess as a bunch of triangles, whereas they have a curved edge. The terms after that correct for the curvature.



7.4 Euler sum

As an example, let's use this result to improve the estimate for Euler's famous sum

$$\sum_1^{\infty} n^{-2}.$$

The first term in the the operator sum is 1, the result of integrating n^{-2} from 1 to ∞ . The second term is $1/2$, the result of $f(1)/2$. The third term is $1/6$, the result of $D/12$ applied to n^{-2} . So:

$$\sum_1^{\infty} n^{-2} \approx 1 + \frac{1}{2} + \frac{1}{6} = 1.666\dots$$

Operators**90**

The true value is $1.644\dots$, so our approximation is in error by about 1%. The fourth term gives a correction of $-1/30$. So the four-term approximation is $1.633\dots$, an excellent approximation using only four terms!

7.5 Conclusion

I hope that you've enjoyed this extended application of analogy, and more generally, this rough-and-ready approach to mathematics.

References

- [1] Barry Cipra. *Mistakes: And How to Find Them Before the Teacher Does*. AK Peters, 3rd edition, 2000.
- [2] P. Horowitz and W. Hill. *The Art of Electronics*. Cambridge University Press, 2nd edition, 1989. ISBN 0521377099.
- [3] Edward M. Purcell. *Electricity and magnetism*, volume 2 of *Berkeley physics course*.. McGraw–Hill, New York, 1985. A classic..
- [4] George Gamow. *Thirty Years that Shook Physics: The Story of Quantum Theory*. Dover, New York, reprint edition, 1985.
- [5] Dwight E. Gray, editor. *AIP Handbook*. McGraw–Hill, New York, 3rd edition, 1972.
- [6] R. D. Middlebrook. Low-entropy expressions: The key to design-oriented analysis. In *IEEE Frontiers in Education*, pages 399–403, Purdue University, 1991. 21st Annual Conference.
- [7] John Malcolm Blair. *The control of oil*. Vintage, 1978.
- [8] George Polya. Let us teach guessing. 1966? MAA.