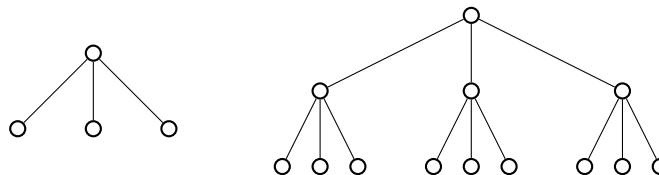


A *rooted tree* is a tree together with a designated vertex r called the *root*. A *perfect k -ary tree* of depth d is defined recursively as follows:

- A perfect k -ary tree of depth 0 is a single vertex.
- For $d \geq 0$, a perfect k -ary tree of depth $d + 1$ is obtained by taking k perfect k -ary trees T_1, \dots, T_k of depth d , a new root vertex r , and adding edges from r to the roots of T_1, \dots, T_k .

Here are diagrams of the perfect ternary (3-ary) trees of depth 1 and 2:



How many vertices $N(d)$ does a perfect k -ary tree of depth d have? When $d \geq 1$, there is one vertex for each of the k subtrees of depth $d - 1$, plus the root vertex. This gives the formula

$$N(d) = k \cdot N(d - 1) + 1$$

for $d \geq 1$, with the “base case” $N(0) = 1$. Plugging in small values of d , this gives

$$N(1) = k \cdot N(0) + 1 = k + 1$$

$$N(2) = k \cdot N(1) + 1 = k(k + 1) + 1 = k^2 + k + 1$$

$$N(3) = k \cdot N(2) + 1 = k(k^2 + k + 1) + 1 = k^3 + k^2 + k + 1$$

and, in general,

$$N(d) = k^d + k^{d-1} + \dots + 1 \quad \text{for every } d \geq 0.$$

You can prove the correctness of this formula by induction on d , but we won't bother. Today we are more interested in “understanding” the value of $N(d)$.

1 Geometric sums

We can evaluate a sum of the form

$$S = x^d + x^{d-1} + \dots + 1$$

for every real number x and positive integer d like this: If we multiply both sides by x , we obtain

$$xS = x^{d+1} + x^d + \dots + x$$

If we now subtract the first expression from the second one, almost all the right hand sides terms cancel out:

$$xS - S = x^{d+1} - 1$$

which simplifies to $(x - 1)S = x^{d+1} - 1$. When $x \neq 1$, we can do a division and obtain the formula

$$x^d + x^{d-1} + \dots + 1 = \frac{x^{d+1} - 1}{x - 1} \quad \text{for every real number } x \neq 1.$$

A sum of this form is called a *geometric sum*.

So the number of vertices in a perfect k -ary tree of depth d is $(k^{d+1} - 1)/(k - 1)$. In particular, for a perfect ternary tree, this number is $(3^{d+1} - 1)/2$. A perfect binary (2-ary) tree of depth d has $2^{d+1} - 1$ vertices.

Annuities You won a prize and you have two options for the prize money. Option A is that you are paid \$5000 per year for the rest of your life. Option B is that you are paid \$80000 today. Which one would you choose?

To answer this question we need to model how the value of money changes over time. If you keep your money in the bank at no interest then option A will pay off for you in twenty years time. If, on the other hand, you want to throw a lavish party right now then option B would make more sense for you. Now suppose that, as a savvy investor, you are quite confident in making a reliable return of $p = 7\%$ per year. How show this affect your choice?

To answer this question we'll calculate how much option A is worth in today's money. The 5K that you will be getting in your zeroth year are worth... well, 5K. For next year's 50K you can reason like this. If you had invested x dollars this year, they would be worth $(1 + p)x$ dollars next year. So today's value of next year's 5K dollars is the amount x for which $(1 + p)x = 5K$, namely $x = 5K/(1 + p)$. By the same reasoning, the 5K you would be getting in two years' time are worth $5K/(1 + p)^2$ today. Continuing this reasoning, you conclude that the value of option A in today's money is

$$5K + \frac{5K}{1 + p} + \frac{5K}{(1 + p)^2} + \dots$$

By the geometric sum formula, the contribution from years zero up to d equals

$$5K \cdot \frac{1/(1 + p)^{d+1} - 1}{1/(1 + p) - 1}$$

In the large n limit, the term $1/(1 + p)^{d+1}$ vanishes and the value converges to $5K \cdot (1 + p)/p$. For $p = 7\%$, the value of option A is about \$76,429. So option B is better.

Here is another way to see the wisdom of option B over option A without evaluating the geometric sum. With a budget of 80K, I can always take out x dollars for spending this year and invest the remaining $80K - x$ dollars aiming to grow them back to 80K by next year. To do this, I need to choose x so that

$$(1 + p) \cdot (80K - x) = 80K,$$

which solves to

$$x = 80K \cdot \left(1 - \frac{1}{1 + p}\right) = 80K \cdot \frac{p}{1 + p}.$$

When p equals 7% we get that x is about 5234 dollars. This improves on the annual 5K in option A.

2 Polynomial sums

In Lecture 3 we proved that

$$1 + 2 + \dots + n = \frac{n(n + 1)}{2}$$

for every integer $n \geq 0$. How did I come up with the expression on the right? Instead of going back to something we already know, let's work out a new one:

$$\text{What is } 1^2 + 2^2 + \dots + n^2?$$

We have to do some guesswork. The sum $1 + 2 + \dots + n$ was a quadratic function in n , perhaps $1^2 + 2^2 + \dots + n^2$ might equal some cubic? Let's make a guess: For all n , there exist real numbers a, b, c, d such that

$$1^2 + 2^2 + \dots + n^2 = an^3 + bn^2 + cn + d.$$

Suppose our guess was correct. Then what are the numbers a, b, c, d ? We can get an idea by evaluating both sides for different values of n :

$$\begin{aligned} 0 &= d && \text{for } n = 0 \\ 1 &= a + b + c + d && \text{for } n = 1 \\ 5 &= 8a + 4b + 2c + d && \text{for } n = 2 \\ 14 &= 27a + 9b + 3c + d && \text{for } n = 3. \end{aligned}$$

I solved this system of equations on the computer and obtained $a = 1/3, b = 1/2, c = 1/6, d = 0$. This suggests the formula

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

for all integers $n \geq 0$. Let us see if we can prove its correctness by induction on n .

We already worked out the base case $n = 0$, so let us do the inductive step. Fix $n \geq 0$ and assume that the equality holds for n . Then

$$1^2 + 2^2 + \dots + (n + 1)^2 = \left(\frac{1}{3}n^3 + \frac{1}{2}n^2 + \frac{1}{6}n\right) + (n + 1)^2 = \frac{1}{3}n^3 + \frac{3}{2}n^2 + \frac{13}{6}n + 1.$$

This indeed equals $\frac{1}{3}(n + 1)^3 + \frac{1}{2}(n + 1)^2 + \frac{1}{6}(n + 1)$. So we have discovered and proved a new theorem:

Theorem 1. For every integer $n \geq 0$, $1^2 + 2^2 + \dots + n^2 = \frac{1}{3}n^3 + \frac{1}{2}n^2 + \frac{1}{6}n$.

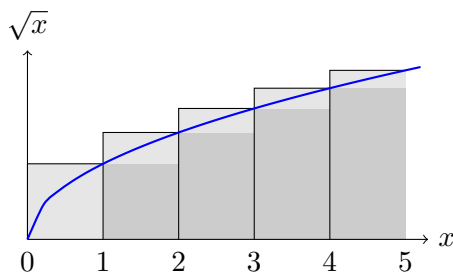
3 Approximating sums

Exact “closed-form” expressions for sums are rather exceptional. Fortunately, we can often obtain very good approximations. One powerful method for approximating sums is the integral method from calculus: It works by comparing the sum with a related integral.

As an example, let us look at the sum

$$S(n) = \sqrt{1} + \sqrt{2} + \dots + \sqrt{n}.$$

The value $S(n)$ can be visualized as the joint area of n bars R_1, \dots, R_n where R_x has base $(x - 1, x)$ and height \sqrt{x} . For example, $S(5)$ equals the area covered by the shaded bars (both light and dark shades) in this plot:



The area under the bars can be lower bounded by the area (i.e., the integral) of the curve $f(x) = \sqrt{x}$ from $x = 0$ to $x = n$:

$$S(5) \geq \int_0^5 \sqrt{x} \, dx.$$

If we remove the area L covered by the lightly shaded bars, the darker shaded area is now dominated by the curve $f(x) = \sqrt{x}$ and so

$$S(5) - L \leq \int_0^5 \sqrt{x} \, dx.$$

The area under L is exactly $\sqrt{5}$: If we stack all of the lightly shaded bars on top of one another, we obtain a column of width 1 and height $\sqrt{5}$. Therefore

$$\int_0^5 \sqrt{x} \, dx \leq S(5) \leq \int_0^5 \sqrt{x} \, dx + \sqrt{5}.$$

By the same reasoning, for every integer $n \geq 1$, we have the inequalities

$$\int_0^n \sqrt{x} \, dx \leq S(n) \leq \int_0^n \sqrt{x} \, dx + \sqrt{n}.$$

We can now use rules from calculus to evaluate the integrals: Recalling that $x^{1/2} = \frac{d}{dx} \frac{2}{3} x^{3/2}$, it follows from the fundamental theorem of calculus that

$$\frac{2}{3} n^{3/2} \leq S(n) \leq \frac{2}{3} n^{3/2} + \sqrt{n}.$$

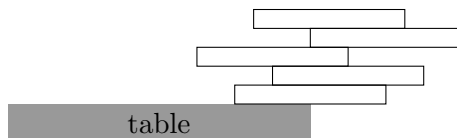
To get a feel for these inequalities, let us plug in a few values of n . (I calculated $S(n)$ by evaluating the sum on the computer.)

n	$\frac{2}{3}n^{3/2}$	$S(n)$	$\frac{2}{3}n^{3/2} + \sqrt{n}$
10	21.082	22.468	24.244
100	666.67	671.46	676.67
1,000	21,081.9	21,097.5	21,113.5
10,000	666,666	666,716	666,766

As n becomes large, the accuracy of these approximations looks quite good.

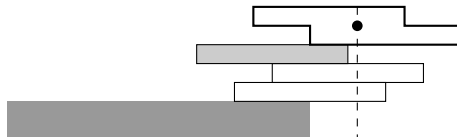
4 Overhang

You have n identical rectangular blocks and you stack them on top of one another at the edge of a table like this:



Is this configuration stable, or will it topple over?

In general, a configuration of n blocks is *stable* if for every i between 1 and n , the center of mass of the top i blocks sits over the $(i + 1)$ st block, where we think of the table as the $(n + 1)$ st block in the stack. For example, the top stack is not stable because the center of mass of the top two blocks does not sit over the third block:



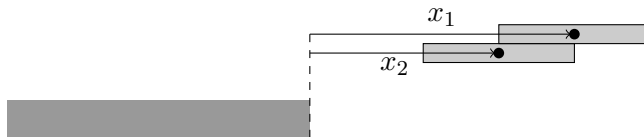
We want to stack our n blocks so that the rightmost block hangs as far over the edge of the table as possible. What should we do? One reasonable strategy is to try to push the top blocks as far as possible away from the table as long as they do not topple over.

We will assume each block has length 2 units and we will use x_i to denote the offset of the center of the i -th block (counting from the top) from the edge of the table:



The offset of a block can be positive, zero, or negative, depending on the position of its center of mass.

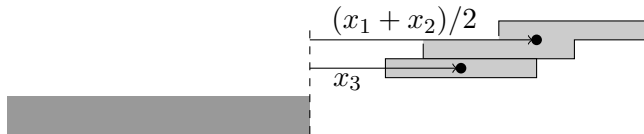
For the top block not to topple over, its center of mass must sit over the second block. To move it as far away from the edge of the table as possible, we should move its center exactly one unit to the right of the center of the second block:



This forces the offsets x_1 and x_2 to satisfy the equation

$$x_1 = x_2 + 1. \tag{1}$$

How about the third block? The center of mass of the first two blocks is at offset $(x_1 + x_2)/2$ from the edge of the table. To push this as far to the right as possible without toppling over the third block



we must set

$$\frac{x_1 + x_2}{2} = x_3 + 1. \tag{2}$$

Continuing our reasoning in this way, for every i between 1 and n , the offset of the center of mass of the top i blocks is $(x_1 + \dots + x_i)/i$. To push this as far to the right without toppling over the $(i + 1)$ st block, we must set

$$\frac{x_1 + x_2 + \dots + x_i}{i} = x_{i+1} + 1 \quad \text{for all } 1 \leq i \leq n. \tag{3}$$

Finally, when $i = n + 1$, we have reached the table whose offset is zero. Since we are thinking of the table as the $(n + 1)$ st block, its centre of mass is one unit left to its edge:

$$x_{n+1} = -1. \tag{4}$$

The overhang of the set of blocks is $x_1 + 1$. To figure out what this number is, we need to solve for x_1 in the system of equations (3-4). Let us develop some intuition first. Equation (1) tells us that $x_2 = x_1 - 1$. Plugging in this formula for x_2 into (2), we get that

$$x_3 = x_1 - \frac{1}{2} - 1.$$

Let's do one more step. Equation (3) tells us that $(x_1 + x_2 + x_3)/3 = x_4 + 1$. Plugging in our formulas for x_2 and x_3 in terms of x_1 we get that

$$\frac{x_1 + (x_1 - 1) + (x_1 - \frac{3}{2})}{3} = x_4 + 1$$

from where

$$x_4 = x_1 - \frac{1 + \frac{3}{2}}{3} - 1 = x_1 - \frac{1}{3} - \frac{1}{2} - 1.$$

At this point it is reasonable to guess that x_{i+1} should equal x_1 minus the sum

$$1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{i}.$$

Let us prove that this guess is correct.

Lemma 2. For all i between 1 and n , $x_i - x_{i+1} = 1/i$.

Proof. If we multiply both sides of the i -th equation (3) by i we obtain

$$x_1 + x_2 + \cdots + x_{i-1} + x_i = i \cdot (x_{i+1} + 1).$$

Under this scaling the $(i - 1)$ st equation is

$$x_1 + x_2 + \cdots + x_{i-1} = (i - 1) \cdot (x_i + 1).$$

Subtracting the two we obtain that

$$x_i = i(x_{i+1} - 1) - (i - 1)(x_i - 1) = ix_{i+1} - (i - 1)x_i + 1$$

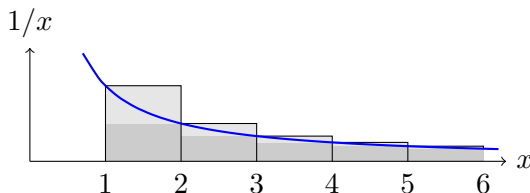
from where, after moving the variables around, we conclude that

$$x_i = x_{i+1} + \frac{1}{i}. \quad \square$$

It follows immediately from this Lemma that

$$x_1 - x_{n+1} = (x_1 - x_2) + (x_2 - x_3) + \cdots + (x_{n+1} - x_n) = 1 + \frac{1}{2} + \cdots + \frac{1}{n}.$$

Since $x_{n+1} = -1$, the overhang $x_1 + 1$ equals exactly this number, which is called the n -th harmonic number and is denoted by $H(n)$. There is no closed-form expression for $H(n)$, but we can obtain an excellent approximation using the integral method. To do this, we compare $H(n)$ with the integral of the function $1/x$:



By similar reasoning as before, the sum $H(n) = 1 + 1/2 + \dots + 1/n$ is given by the area of the first n shaded bars. This area is larger than the integral of $1/x$ from 1 to $n + 1$:

$$H(n) \geq \int_1^{n+1} \frac{1}{x} dx.$$

On the other hand, if we subtract from $H(n)$ the area of the lightly shaded bars, then the integral becomes larger. This area equals $1 - 1/(n + 1)$:

$$H(n) - 1 + \frac{1}{n + 1} \leq \int_1^{n+1} \frac{1}{x} dx.$$

Combining this two inequalities gives the approximation

$$\int_1^{n+1} \frac{1}{x} dx \leq H(n) \leq \int_1^{n+1} \frac{1}{x} dx + 1 - \frac{1}{n + 1}.$$

The antiderivative of $1/x$ is $\ln x$. By the fundamental theorem of calculus it follows that

$$\ln(n + 1) \leq H(n) \leq \ln(n + 1) + 1 - \frac{1}{n + 1}. \quad (5)$$

The left hand side of this inequality tells us that our method of stacking blocks achieves overhang at least $\ln(n + 1)$. The logarithm function is unbounded; given enough blocks, we can grow our stack all the way to New York!

5 Order of growth

In engineering we are often interested in the *asymptotic* behaviour of measures as our problem size becomes large. For example, if you write a program for routing data packets through a network, we might not really care what happens for 2 or 3 packets, but would want to have a good understanding about how fast the program is when we have hundreds or thousands of packets. For this purpose, it is useful to have a quick way of estimating how a function $f(n)$ behaves as n grows large. Usually, we do this by comparing the value of the function f for large n to values of “well-known” functions like n , n^2 , $\log n$, 2^n , or e^n .

The big-oh notation is a handy way to say that a given function does not grow too fast.

Definition 3. For two real-valued functions f and g (defined over the positive reals, or over the positive integers), we say f is $O(g)$ (big-oh of g) if there exists a constant $C > 0$ such that for every sufficiently large input x , $f(x) \leq C \cdot g(x)$.

For example, $13x^4 + 2x^2 + 10x + 1000$ is $O(x^4)$ because when x is large (specifically, at least 1):

$$13x^4 + 2x^2 + 10x + 1000 \leq 13x^4 + 2x^4 + 10x^4 + 1000x^4 = 1025x^4.$$

By the same reasoning, every polynomial is the big-oh of its highest-degree monomial.

Similarly, $\log(16x^5 + 3x + 11)$ is $O(\log x)$ because when x is large,

$$\log(16x^5 + 3x + 11) \leq \log(16x^5 + 3x^5 + 11x^5) \leq \log(30x^5) \leq \log(x^6) \leq 6 \log x.$$

By the same reasoning, the logarithm of any polynomial of x is $O(\log x)$.

The little-oh notation says that asymptotically, one function grows at a significantly slower rate than another one.

Definition 4. For two real-valued functions f and g , we say f is $o(g)$ (little-oh of g) if *for every* constant $c > 0$ and every sufficiently large input x , $f(x) \leq c \cdot g(x)$.

If f is $o(g)$, then f is also $O(g)$, but not necessarily the other way. For example, $\frac{1}{2}x^5 - x^3$ is $O(x^5)$, but it is not $o(x^5)$ because as x gets large, $(\frac{1}{2}x^5 - x^3)/x^5$ converges to $\frac{1}{2}$; when x is sufficiently large, $\frac{1}{2}x^5 - x^3$ will be greater than, say, $\frac{1}{4}x^5$.

For example, $x^{3/2} + 3x^{1/2}$ is $o(x^2)$ because

$$x^{3/2} + 3x^{1/2} \leq 4x^{3/2} \leq \frac{4}{x^{1/2}} \cdot x^2$$

and for every constant $c > 0$, as x becomes large enough, $4/x^{1/2}$ is smaller than c , so $x^{3/2}$ is at most cx^2 . By similar reasoning, every polynomial $p(x)$ of degree d (even one with fractional degrees) is $o(x^e)$ for every constant $e > d$.

Another way to determine order of growth for sufficiently large functions is by taking limits. Assuming that the ratio $f(x)/g(x)$ converges in the limit $x \rightarrow \infty$, we have the relations

$$\begin{aligned} f \text{ is } O(g) & \text{ if } \lim_{x \rightarrow \infty} f(x)/g(x) < \infty, \\ f \text{ is } o(g) & \text{ if } \lim_{x \rightarrow \infty} f(x)/g(x) = 0. \end{aligned}$$

So $17x^4 + 5x^3$ is $O(x^4)$ because $(17x^4 + 5x^3)/x^4$ tends to 17 in the limit $x \rightarrow \infty$, while $x^{3/2} = o(x^2)$ because $x^{3/2}/x^2 = 1/\sqrt{x}$ tends to zero in the limit $x \rightarrow \infty$.

When $0 < B < C$, it is also true that B^x is $o(C^x)$ because we can write $B^x/C^x = (B/C)^x$ tends to zero in the limit.

Theorem 5. For all constants $a, b > 0$, $(\log x)^a$ is $o(x^b)$.

In the statement of this theorem, the base of the logarithm is irrelevant because changing the base from one constant to another only changes the value of the expression $(\log x)^a$ by a constant factor. In the proof we will assume that the logarithm is a natural logarithm.

Proof. When $a = 1$, we can calculate $\lim_{x \rightarrow \infty} (\ln x)/x^b$ using L'Hôpital's rule from calculus. Both numerator and denominator grow to infinity, but this is not true for their derivatives: $\frac{d}{dx} \ln x = 1/x$, while $\frac{d}{dx} x^b = bx^{b-1}$. The ratio of these two numbers is $1/(bx^b)$, which tends to zero as x grows. Therefore

$$\lim_{x \rightarrow \infty} \frac{\ln x}{x^b} = \lim_{x \rightarrow \infty} \frac{1}{bx^b} = 0$$

and so $\log x \leq cx^b$ for every constant $c > 0$ and sufficiently large x .

Now let $a > 0$ be arbitrary and $c > 0$ be an arbitrary constant. By what we just proved,

$$\ln x \leq c^{1/a} x^{b/a}$$

for x sufficiently large, from where

$$(\ln x)^a \leq (c^{1/a} x^{b/a})^a \leq cx^b$$

for x sufficiently large. □

If we set $x = e^y$ and $B = e^b$, we get the following corollary:

Corollary 6. For all constants $a > 0$ and $B > 1$, y^a is $o(B^y)$.

These relations are *transitive*: For example, if f is $o(g)$ and g is $o(h)$ then f is $o(h)$. In fact, if f is $o(g)$ and g is $O(h)$ then f is still $o(h)$. We can therefore use this notation to order sets of functions.

Exercise Suppose we want to order the functions 2^x , 2^{x^2} , x^2 , x^x in terms of their asymptotic growth. By Corollary 6 we know that x^2 is both $o(2^x)$ and $o(2^{x^2})$.

How do 2^x and 2^{x^2} compare to each other? As x is $o(x^2)$ we would expect that 2^x should also be $o(2^{x^2})$. In fact the ratio $2^x/2^{x^2}$ equals 2^{x-x^2} and this goes to zero as the exponent $x - x^2$ tends to $-\infty$. So we have established that

$$x^2 \text{ is } o(2^x) \quad \text{and} \quad 2^x \text{ is } o(2^{x^2}).$$

Where does x^x fit in? To answer this it is usually a sensible strategy to identify whether x^x is a polynomial or an exponential-type function. In this example x^x grows faster than x , x^2 , x^{100} and any polynomial in x , so x^2 is certainly $o(x^x)$. To compare x^x against 2^x and 2^{x^2} it is sensible to rewrite it as a base-2 exponential, namely $x^x = 2^{x \log x}$. Now we see that x is $o(x \log x)$ and $x \log x$ is $o(x^2)$ so x^x should fit between 2^x and 2^{x^2} :

$$x^2 \text{ is } o(2^x), \quad 2^x \text{ is } o(x^x), \quad \text{and} \quad x^x \text{ is } o(2^{x^2}).$$

Big-theta Big-Theta says that two functions have the same order of growth:

Definition 7. We say f is $\Theta(g)$ if f is $O(g)$ and g is $O(f)$.

For example, $x^5 + 7x^3$ is $\Theta(x^5)$ and $(\log x^{1/2})^2 - \log(x^4 + x^2)$ is $\Theta((\log x)^2)$. For sufficiently nice functions,

$$f \text{ is } \Theta(g) \text{ if } 0 < \lim_{x \rightarrow \infty} f(x)/g(x) < \infty.$$

It is customary to abuse the equality sign when talking about order of growth. In books you often see “ $f = O(g)$ ” instead of “ f is $O(g)$ ”. Technically, this is incorrect because f and $O(g)$ are objects of different types: f is a single function while $O(g)$ is not. It is okay to use this notation as long as you are aware of what it means. What you should *not* do is write “equations” like

$$1 + 2 + \dots + n = O(1) + O(1) + \dots + O(n) = (n-1) \cdot O(1) + O(n) = O(n)$$

because it is not clear what they mean and may lead to incorrect conclusions.

6 Order of growth of summations

The order of growth of a summation can be determined by first approximating the sum. For example, from (5) we have that

$$1 \leq \frac{H(n)}{\ln(n+1)} \leq 1 + \frac{1 - 1/(n+1)}{\ln(n+1)}$$

so in the limit $n \rightarrow \infty$, $H(n)/\ln(n+1)$ tends to one. Therefore $H(n)$ is $\Theta(\ln(n+1))$, which is the same as $\Theta(\log n)$.

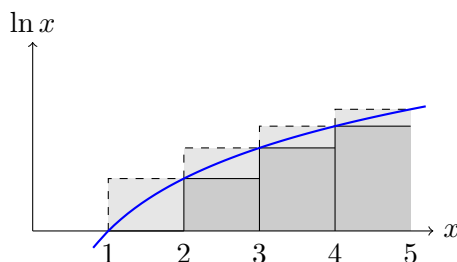
Let's now try to determine the order of growth of

$$n! = 1 \cdot 2 \cdot \dots \cdot n.$$

The expression for the factorial is a product, not a sum. We can turn products into sums by taking logarithms:

$$\ln(n!) = \ln 1 + \ln 2 + \dots + \ln n.$$

This sum is perfect for the integral method.



The sum $\ln(n!) = \ln 1 + \ln 2 + \dots + \ln n$ is the area of the first n dark bars (the very first one has height zero) so it is upper bounded by the integral of $\ln x$ from 1 to n :

$$\ln(n!) \leq \int_1^n \ln x \, dx.$$

To obtain an upper bound, we add the total area of the first n light bars, which is exactly $\ln n$:

$$\ln(n!) - \ln n \leq \int_1^n \ln x \, dx.$$

The antiderivative of $\ln x$ is $x \ln x - x$, so we get that

$$n \ln n - n \leq n! \leq n \ln n - n + \ln n. \tag{6}$$

Both sides of this equation are $\Theta(n \ln n)$ so we have proved that

Theorem 8. $\ln(n!)$ is $\Theta(n \ln n)$.

We cannot conclude from here that $n!$ is $\Theta(e^{n \ln n}) = \Theta(n^n)$: If f is $O(g)$ it does not mean that e^f is $O(e^g)$. What we can do is exponentiate both sides of (6) to obtain

$$e^{n \ln n - n} \leq n! \leq e^{n \ln n - n + \ln n}$$

which we can simplify to

$$\left(\frac{n}{e}\right)^n \leq n! \leq n \cdot \left(\frac{n}{e}\right)^n.$$

The left and right hand side do not match asymptotically. The “naive” integral method cannot determine the order of growth of $n!$. Using an enhancement that you can read about in Section 7 it is possible to prove that the truth is in the middle:

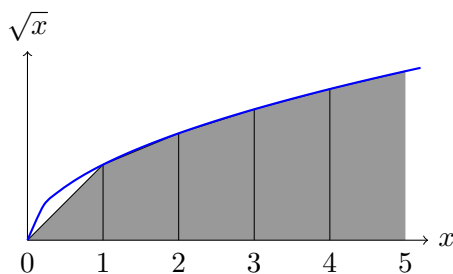
Theorem 9. $n!$ is $\Theta(\sqrt{n} \cdot (n/e)^n)$.

7 Stirling's formula*

The estimate we obtained for $S(n) = 1 + \sqrt{2} + \cdots + \sqrt{n}$ has the undesirable property that the errors $S(n) - \frac{2}{3}n^{3/2}$ and $\frac{2}{3}n^{3/2} + \sqrt{n} - S(n)$ appear to grow to infinity as n becomes larger. You may notice, however, that the average of the upper and lower approximation looks much better:

n	$S(n)$	$\frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n}$
10	22.468	24.663
100	671.463	676.666
1,000	21,097.456	21,096.662
10,000	666,716.459	666,716.666

The estimate $\frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n}$ can be obtained more methodically as an application of the **trapezoidal rule**, which estimates the shaded area in this picture by the integral under the curve in the following plot:



Indeed, the shaded area can be calculated by adding the areas under the first n trapezoids to obtain

$$A = \frac{1}{2}(\sqrt{0} + \sqrt{1}) + \frac{1}{2}(\sqrt{1} + \sqrt{2}) + \cdots + \frac{1}{2}(\sqrt{n-1} + \sqrt{n}) = S(n) - \frac{1}{2}\sqrt{n}.$$

This area is upper bounded by the area under the curve $f(x) = \sqrt{x}$ between $x = 0$ and $x = n$, so

$$A \leq \int_0^n \sqrt{x} dx = \frac{2}{3}n^{3/2}$$

from where

$$S(n) \leq \frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n}.$$

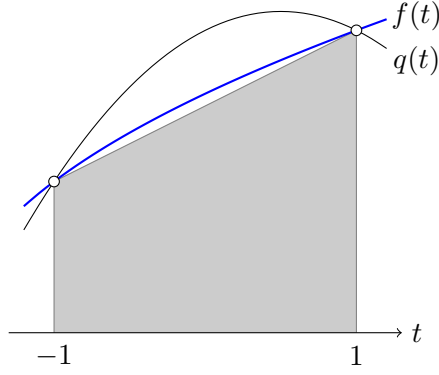
It appears that the approximation error $\frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} - S(n)$ is a small fraction, something like $1/4$, regardless of the value of n . Is this really the case?

To bound the error we need to understand the contribution by a single trapezoid. For convenience we will shift and scale the trapezoid to have x -coordinates -1 and 1 . The following clever lemma bounds the error by the second derivative f'' of f .

Lemma 10. For every function f from the real interval $[-1, 1]$ to the real numbers so that $-B \leq f''(t) \leq 0$ for all $-1 \leq t \leq 1$,

$$f(-1) + f(1) \leq \int_{-1}^1 f(t) dt \leq f(-1) + f(1) + \frac{2B}{3}. \quad (7)$$

Proof. Since $f''(x) \leq 0$, the area under the curve $f(x)$ for $x \in [-1, 1]$ is at most the area of the shaded trapezoid in the graph below, which is exactly $f(-1) + f(1)$.



Let $q(t) = -Bt^2/2 + Ct + D$ be the (unique) quadratic function such that $q(-1) = f(-1)$, $q(1) = f(1)$, and $q''(t) = -B$. We will argue shortly that $f(t)$ must be upper bounded by $q(t)$ for all $t \in [-1, 1]$ as in the graph. Assuming this is the case, we can upper bound the area under f by the area under q to get

$$\int_{-1}^1 f(t)dt \leq \int_{-1}^1 q(t)dt = \left(-\frac{Bt^3}{6} + \frac{Ct^2}{2} + Dt \right) \Big|_{-1}^1 = -\frac{B}{3} + 2D = -B + 2D + \frac{2B}{3}.$$

The value $-B + 2D$ is exactly $q(1) + q(-1)$, which in turn equals $f(1) + f(-1)$, giving the desired upper bound.

It remains to argue that $q(t)$ indeed upper bounds $f(t)$ for $t \in [-1, 1]$. Suppose for contradiction that $q(x) < f(x)$ for some x between -1 and 1 . We consider two cases. If $q'(x) \leq f'(x)$, then $q'(t)$ can be at most $f'(t)$ for all $t \in [x, 1]$ because $q''(t) \leq f''(t)$, so q' drops at a faster rate than f' . Therefore q must drop at a faster rate than f , so if $q(x) < f(x)$ then $q(1)$ must be less than $f(1)$, a contradiction. If, on the other hand, $q'(x) > f'(x)$ then the same reasoning applied to the functions $q(-t)$ and $f(-t)$ also leads to a contradiction. \square

Applying the change of variables $x = i + (t + 1)/2$ in (7) gives that for every i ,

$$\frac{f(i) + f(i+1)}{2} \leq \int_i^{i+1} f(x)dx \leq \frac{f(i) + f(i+1)}{2} + \frac{B_i}{3} \quad (8)$$

assuming that $0 \leq f''(x) \leq B_i$ for every x between i and $i + 1$.

If f is the function $f(x) = \sqrt{x}$, its second derivative is $f''(x) = -\frac{1}{4}x^{-3/2}$, so it takes value between $-\frac{1}{4}i^{-3/2}$ and zero in the interval $[i, i + 1]$. Applying (8) gives that

$$\frac{1}{2}(\sqrt{i} + \sqrt{i+1}) \leq \int_i^{i+1} \sqrt{x}dx \leq \frac{1}{2}(\sqrt{i} + \sqrt{i+1}) + \frac{1}{12i^{3/2}}.$$

The right hand side is meaningless when $i = 0$, so we start at $i = 1$ instead. Adding up these inequalities for i taking integer values from 1 up to $n - 1$ gives

$$S(n) - \frac{1}{2}(\sqrt{1} + \sqrt{n}) \leq \int_1^n \sqrt{x}dx \leq S(n) - \frac{1}{2}(\sqrt{1} + \sqrt{n}) + E(n)$$

where

$$E(n) = \frac{1}{12} + \frac{1}{12 \cdot 2^{3/2}} + \cdots + \frac{1}{12 \cdot (n-1)^{3/2}}.$$

Calculating the integral and simplifying gives that

$$\frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} - \frac{1}{6} - E(n) \leq S(n) \leq \frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} - \frac{1}{6}.$$

This estimate appears even better than the previous one as subtracting $1/6 \approx 0.166$ seems to bring the estimate even closer to the true value of the sum $S(n)$.

It remains to analyze the error term $E(n)$. This can be done by another integral bound:

$$E(n) \leq \int_1^{n+1} \frac{1}{12x^{3/2}} dx = \frac{1}{12} + (-\frac{1}{6}x^{-1/2})|_1^{n+1} \leq \frac{1}{12} + \frac{1}{6} = \frac{1}{4}$$

so the estimate is always within $1/4$ of the true value.

In fact, we can make the approximation error arbitrarily small by starting the summation at a later point. For instance, the same method tells us that

$$\frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} - 16.5 \leq \sqrt{9} + \sqrt{10} + \dots + \sqrt{n} \leq \frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} - 16.5 + E(9, n)$$

where

$$E(9, n) = \frac{1}{12 \cdot 9^{3/2}} + \frac{1}{12 \cdot 10^{3/2}} + \dots + \frac{1}{12 \cdot n^{3/2}} \leq \int_8^\infty \frac{1}{12x^{3/2}} dx = \frac{1}{12 \cdot 8^{2/3}} \leq 0.0834.$$

Therefore the value

$$\frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} - 16.5 + (\sqrt{1} + \dots + \sqrt{8}) = \frac{2}{3}n^{3/2} + \frac{1}{2}\sqrt{n} + 0.194 \pm 0.0001$$

is guaranteed to never exceed $S(n)$ by more than 0.0209. For example, it gives the estimate ≈ 666716.473 for $S(10000)$.

Stirling's formula Let's use the trapezoidal rule to estimate the sum

$$F(n) = \ln 1 + \ln 2 + \dots + \ln n.$$

Since the second derivative of $\ln x$ is $-1/x^2$, its value is between $-1/i^2$ and zero on every interval $[i, i+1]$, where $i \geq 1$. Inequality (8) gives that for every $i \geq 1$

$$\int_i^{i+1} \ln x dx = \frac{\ln i + \ln(i+1)}{2} + \varepsilon_i \quad \text{where} \quad 0 \leq \varepsilon_i \leq \frac{1}{12i^2}.$$

Therefore

$$\int_1^n \ln x dx = F(n) - \frac{1}{2} \ln n + E(n)$$

where

$$E(n) = \varepsilon_1 + \dots + \varepsilon_n \leq \frac{1}{12 \cdot 1^2} + \frac{1}{12 \cdot 2^2} + \dots + \frac{1}{12 \cdot n^2}.$$

The antiderivative of $\ln x$ is $x \ln x - x + 1$, so

$$F(n) = n \ln n - n + \frac{1}{2} \ln n + 1 - E(n).$$

Exponentiating both sides and rearranging terms gives

$$e^{F(n)} = e^{n \ln n} \cdot e^{-n} \cdot e^{(\ln n)/2} \cdot e^{1-E(n)} = e^{1-E(n)} \cdot \sqrt{n} \cdot \left(\frac{n}{e}\right)^n.$$

On the other hand, $e^{F(n)} = e^{\ln 1 + \dots + \ln n} = n!$. Using the integral method to upper bound $E(n)$ by $2/3$ we can conclude that

$$e^{1/3} \cdot \sqrt{n} \cdot \left(\frac{n}{e}\right)^n \leq n! \leq e \cdot \sqrt{n} \cdot \left(\frac{n}{e}\right)^n.$$

so that $n! = \Theta(\sqrt{n} \cdot (n/e)^n)$, which proves Theorem 9.

This is not the last word. Since $E(n)$ is a sum of positive numbers $\varepsilon_1, \varepsilon_2, \dots$ which is at most $2/3$, it must converge to some constant c as n approaches infinity. Therefore it must be that the expression

$$n! / \sqrt{n} \cdot \left(\frac{n}{e}\right)^n$$

approaches the constant $C = e^{1-c}$ as n approaches infinity. Using the **law of large numbers** for the binomial distribution from probability theory and a bit more calculus, one can show that $C = \sqrt{2\pi}$ and obtain this beautiful formula:

Theorem 11 (Stirling's formula). $\lim_{n \rightarrow \infty} n! / \sqrt{2\pi n} \cdot (n/e)^n = 1$.

References

This lecture is based on Chapter 13 of the text *Mathematics for Computer Science* by E. Lehman, T. Leighton, and A. Meyer. The variant of the integral methods described in the textbook is slightly different from the one in these notes.

Surprisingly, if we allow for the blocks to be stacked not only on top of one another but also side by side, the overhang can be much improved. If you are interested, see the amazing work *Overhang* by Mike Paterson and Uri Zwick.