

MATH 647 and MATH 648 Course Notes
Version 1.0

Ty Ghaswala and Paul McGrath
© Faculty of Mathematics, University of Waterloo

April 20, 2026

Contents

I	MATH 647	5
1	Limits	6
1.1	Why Do We Need Limits?	6
1.2	Limit of a Function at a Point	9
1.2.1	Definition of a Limit	11
1.3	Limit Properties	21
1.4	Limits at Infinity	23
1.5	Infinite Limits	27
1.6	The Squeeze Theorem	31
2	Continuity	35
2.1	Introduction to Continuity	35
2.2	Types of Discontinuities	37
2.3	Continuity Properties of Elementary Functions	40
2.4	Rules for Continuous Functions	42
2.5	Intermediate Value Theorem	45
2.6	Extreme Value Theorem	48
3	Differentiation	51
3.1	Instantaneous Velocity Revisited	51
3.2	Definition of the Derivative	51
3.3	Differentiability	53
3.4	The Derivative Function	55
3.5	Differentiation Rules	57
3.6	Differentiating Elementary Functions	60
3.6.1	Trigonometric Functions	60
3.6.2	Exponential Functions	63
3.7	Chain Rule	64
3.8	Implicit Differentiation	68
3.8.1	Logarithmic Functions	71
3.8.2	Inverse Trigonometric Functions	72
3.9	Tangent Lines and Linear Approximations	73
3.10	Newton's Method	78
3.11	Local Extreme Values	83
3.11.1	Locating Local Extreme Values	87
3.12	Related Rates	92
4	The Mean Value Theorem	96
4.1	The Mean Value Theorem	96
4.1.1	Fermat's Theorem	96

4.1.2	Rolle's Theorem	97
4.1.3	Mean Value Theorem	99
4.1.4	Extreme Values Revisited	102
4.2	Indeterminate Forms	106
4.2.1	L'Hôpital's rule	106
4.2.2	Indeterminate Products	109
4.2.3	Indeterminate Differences	110
4.2.4	Indeterminate Powers	111
5	Extra Topics	113
5.1	Optimization	113
5.2	Antidifferentiation	117
5.2.1	Higher-Order Antiderivatives	121
5.3	Modelling	123
5.3.1	Exponential Processes	124
5.3.2	Oscillatory Motion	128
II	MATH 648	131
6	Sequences and Series	132
6.1	Sequences	132
6.1.1	Introduction to Sequences	133
6.1.2	Limit of a Sequence	134
6.1.3	Monotonicity and Boundedness	138
6.2	Series	140
6.2.1	Partial Sums	142
6.2.2	Geometric Series	145
6.2.3	Series Properties	147
6.3	Comparison Tests	148
6.3.1	Comparison Test	148
6.3.2	Limit Comparison Test	150
6.4	Alternating Series	152
6.4.1	Alternating Series Test	153
6.4.2	Alternating Series Approximations	156
6.5	Ratio Test	158
6.5.1	Absolute Convergence	158
6.5.2	The Ratio Test	160
6.6	Power Series	163
6.6.1	Introduction to Power Series	163
6.6.2	Power Series Convergence	164
6.6.3	Manipulating Power Series Representations of Functions	167
6.6.4	Working with Divergent Power Series	170
7	Taylor Polynomials	172
7.1	Introduction to Taylor Polynomials	172
7.2	Taylor Series	174
7.3	Taylor Polynomials	179
7.4	Taylor's Inequality	182
7.5	Taylor Polynomial Approximations in Science	185
7.6	Binomial Series	187

7.6.1	The Binomial Approximation	188
8	Integration	190
8.1	Riemann Sums	190
8.2	Definite Integrals	193
8.2.1	Properties of Definite Integrals	195
8.2.2	Integrability	197
8.3	The Fundamental Theorem of Calculus	201
8.4	Indefinite Integrals	204
8.4.1	Net Change	206
8.5	Integration Techniques	207
8.5.1	Integration by Substitution	207
8.5.2	Integration by Parts	211
8.5.3	Trigonometric Integrals	214
8.5.4	Trigonometric Substitutions	216
8.5.5	Integration by Partial Fractions	217
8.6	Improper Integrals	220
8.6.1	Improper Integrals: Type 1	221
8.6.2	Convergence of Power Functions: Part 1	223
8.6.3	Improper Integral Comparison Test	224
8.6.4	Improper Integrals: Type 2	226
8.6.5	Convergence of Power Functions: Part 2	227
8.7	Probability	228
8.7.1	Probability Density	229
8.7.2	Average Value	231
8.7.3	Expected Value	233
8.8	Series and Taylor Polynomials Revisited	235
8.8.1	Integral Test	235
8.8.2	Taylor's Inequality Revisited	237
8.8.3	Integral Approximations	241
9	Geometrical Applications of Integration	245
9.1	Areas Between Curves	245
9.2	Volumes	249
9.2.1	Volumes by Disks	251
9.2.2	Volumes by Washers	253
9.2.3	Volumes by Shells	254
9.3	Arc Length	257
9.3.1	Surfaces of Revolution	260

Part I

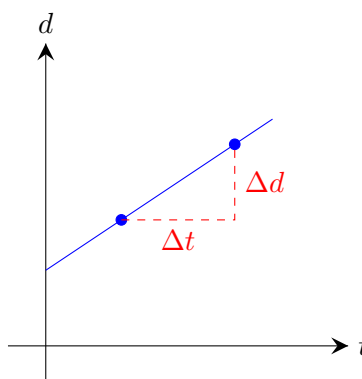
MATH 647

Chapter 1

Limits

1.1 Why Do We Need Limits?

Consider an airplane moving at constant speed in a straight line. If we denote the position of the airplane by d , and time by t , we can depict this graphically as in the plot below. Notice that the graph is linear since the speed of the airplane is constant.



The slope of this graph gives the rate of change of the position of the airplane with respect to time - a quantity better known as the speed of the airplane. We can convince ourselves of this by looking at some interval during which the airplane travels a distance Δd in a time of Δt . The quantity Δd gives the ‘rise’ while Δt gives the ‘run’ over this interval. Their ratio quantifies the slope.

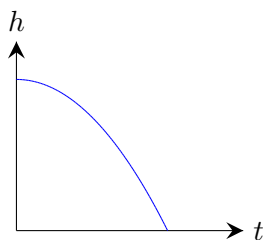
$$\text{slope} = \frac{\Delta d}{\Delta t} = \text{speed}$$

Notice, we could pick any two points on the graph to compute this slope since the airplane is moving at constant speed.

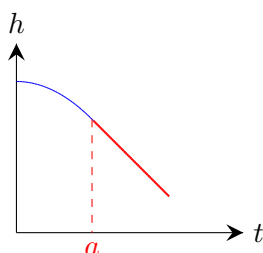
Now, suppose a skydiver jumps out of the airplane. How do we analyze the motion of the skydiver in the vertical (call it the h direction)? Notice we cannot do what we did above with the airplane because the skydiver will *not* be moving at constant speed! First, we note that we need to treat the skydiver’s speed as a function of time since it is constantly changing. This is crucial because it also means that we need to rephrase our inquiry as how

do we find the speed of the skydiver at a specific instant in time. That is, how do we find the **instantaneous speed**?

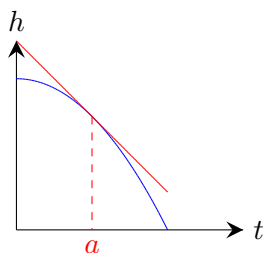
Consider again a graph of the position (denoted h) as a function of time (denoted t).



Let's pause for a moment and think about what exactly we are trying to find when we are looking for the *instantaneous speed* at some time $t = a$. One way to interpret this is to assume that at time $t = a$, all acceleration (due to gravity or air resistance) ceases, and the skydiver continues as whatever speed they are travelling at that moment in time. In this case, after $t = a$ the graph becomes linear like this:



Since the red portion of the graph to the right of $t = a$ is linear, we can compute the slope for that part of the graph. This slope would be the instantaneous speed. Let's look at both the linear function (extended to the left of a) and the original function on the same axes:



So, to find the instantaneous speed of the skydiver at time $t = a$, we need to find the slope of the line tangent to the curve at $t = a$. How can we do this?

At the moment, we can't do this precisely, but we can find some pretty good approximations! Over a finite interval of time, we can calculate the **average speed** of the skydiver. This average speed is given by the ratio of the change in position, Δh , to the time elapsed, Δt . Observe that this average speed is not the same as the instantaneous speed at every instant during the interval but it does provide an approximation for the instantaneous speed at any instant during the interval.

Let's say we use this average speed calculated to approximate the speed at the beginning of the interval which we'll label $t = a$. Let's also label the right endpoint as $t = b$ so that $\Delta t = b - a$. Observe that we can improve the approximation by making the time interval smaller which is equivalent to shifting b closer to a .

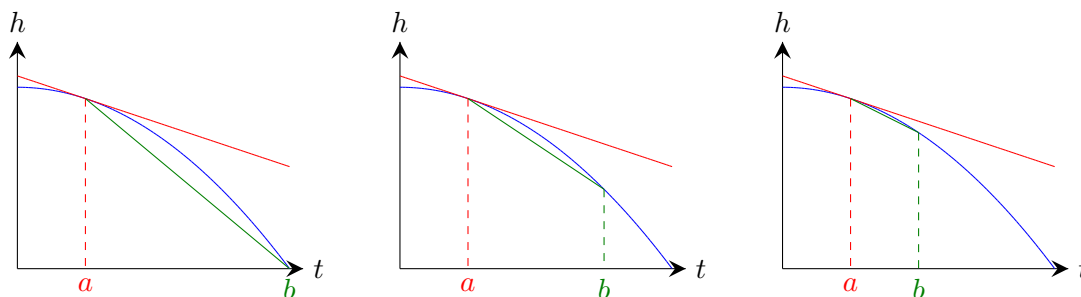


Figure 1.1: Approximating the slope of the tangent line by slopes of smaller and smaller segments

Our intuition tells us that if we could take the interval to be vanishingly small, then we might actually get an *exact* value for the speed at that singular moment in time. That is, the slope of the green line would be the same as the slope of the red tangent line at $t = a$. Indeed, we will see that we can make this idea work and this amounts to computing a derivative. However, we need to be careful. We can't simply set $b = a$ because then $\Delta t = 0$ and the average speed expression becomes undefined.

To avoid this pitfall, we instead need to look at the behaviour of the ratio $\frac{\Delta h}{\Delta t}$ as a whole as the denominator gets arbitrarily small. Inspecting the behaviour of a function in this way is called taking a **limit**.

Let's do this process of approximating the slope of a tangent line at a point with a specific example.

Example 1

We wish to find the slope of the tangent line to the point $(1, 1)$ on the curve $y = x^2$. Let's try to approximate it as we did above. Let $a = 1$, and we will let b get closer and closer to 1.

Let's start with $b = 2$, and notice that $(2, 4)$ is on the curve. Then the slope of the line joining $(1, 1)$ to $(2, 4)$ is

$$\frac{4 - 1}{2 - 1} = 3.$$

What about when $b = 1.5$? Then the slope is given by

$$\frac{(1.5)^2 - 1}{1.5 - 1} = 2.5.$$

When $b = 1.2$ we get the slope as

$$\frac{(1.2)^2 - 1}{1.2 - 1} = 2.2.$$

Continuing like this, we get the following values of the slope joining $(1, 1)$ to (b, b^2) as b gets closer and closer to 1:

b	$\frac{b^2-1}{b-1}$
1.1	2.1
1.05	2.05
1.001	2.001

We might guess that the slope of the tangent line to $y = x^2$ at the point $(1, 1)$ is 2 (and it turns out we would be right!). However, we cannot simply substitute in $b = 1$ to the expression $\frac{b^2-1}{b-1}$, because the fraction would then be undefined.

We will deal with these issues in the next section when we define limits.

EXERCISE

Assuming no wind resistance, write down the equation for the height of a skydiver as a function of time (you may need to revisit some physics to remind yourself of the relevant equations of motion). Sketch this function on the graph. Approximate the instantaneous speed at the beginning, end, and somewhere in the middle of the dive, and make guesses as to what the instantaneous speed is at those points.

1.2 Limit of a Function at a Point

Let's begin with a heuristic definition of a limit. We say L is the limit of a function f as x approaches a if as x gets arbitrarily close to a , f gets arbitrarily close to L .

Consider the graph of the function $f(x)$ below. We see that as we approach the x -value of a from either side of a , the value of the function approaches L . So we say the limit as x tends towards a of $f(x)$ is L .

"That was easy," you say, "since $f(a) = L$, so of course the limit should be L !". It turns out that this isn't always the case. Consider the functions $g(x)$ and $h(x)$ below. Both these functions agree with $f(x)$ for all values except for a . The value $g(a)$ is undefined (so a is not in the domain of g), and $h(a)$ is defined but $h(a) \neq L$. In both these cases we still have that the value of the function approaches L as x approaches a , and in both cases the limit as x approaches a is indeed L .

With this heuristic approach, we require that the same limit is approached, regardless of whether we approach it from the left or the right. This means that, in particular, in order for us to even talk about the limit of a function at a point $x = a$, the function needs to be defined on intervals on both sides of a .

Although we don't have a formal definition at the moment, let's investigate some limits with some specific examples.

Example 2

Let $f(x) = \frac{x^2-4}{x-2}$. This function is not defined at $x = 2$, but it is defined everywhere else (so its domain is $\mathbb{R} \setminus \{2\}$). Let's investigate the behaviour of $f(x)$ around $x = 2$.

Here is a table of values of $f(x)$ as we approach x from both above and below 2.

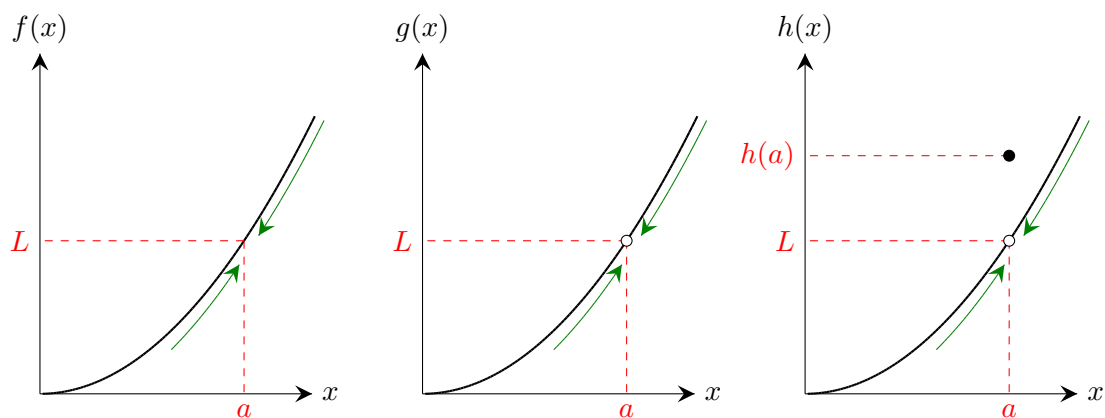


Figure 1.2: Graphs of three different functions, all with a limit of L as x approaches a .

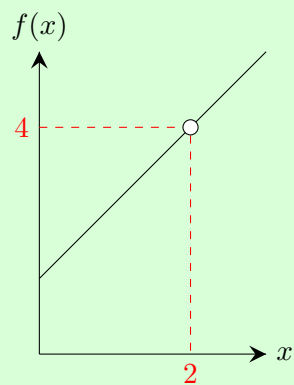
x	$f(x)$	x	$f(x)$
1.9	3.9	2.1	4.1
1.99	3.99	2.01	4.01
1.999	3.999	2.001	4.001
1.9999	3.9999	2.0001	4.0001
1.99999	3.99999	2.00001	4.00001
1.999999	3.999999	2.000001	4.000001

So we see that $f(x)$ gets closer and closer to 4 as x gets closer and closer to 2. With this data, our guess is that the limit of $f(x)$ as x approaches 2 is 4.

Here's another approach. If $x \neq 2$, then we can perform the following manipulation:

$$\frac{x^2 - 4}{x - 2} = \frac{(x - 2)(x + 2)}{x - 2} = x + 2.$$

So the value of the function at any point $x \neq 2$ is simply $x + 2$. Now, as we get closer and closer to $x = 2$ (from above or below), the value of $f(x)$ approaches 4 since the value of $x + 2$ approaches 4. We can also see this looking at the graph of $f(x)$, which looks like the graph of $x + 2$, except at $x = 2$.



Again, our guess is that the limit as x approaches 2 of $f(x)$ is 4.

It is worth reiterating that $f(x) \neq x + 2$ since $f(2)$ is not defined. However, the limit of $f(x)$ at $x = 2$ is defined (or at least we will see it is defined once we have the formal definition!).

EXERCISE

Consider the function

$$f(x) = \begin{cases} x - 1 & \text{if } x \neq 2 \\ 3 & \text{if } x = 2. \end{cases}$$

Take a guess as to what the limit is of $f(x)$ as x approaches 2. What about as x approaches 3?

EXERCISE

The *Heaviside function* is the function

$$H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0. \end{cases}$$

Graph the Heaviside function. What is the limit of $H(x)$ as x approaches 0? What about as x approaches 1?

Remark: The Heaviside function is named after English mathematician and physicist Oliver Heaviside. It also happens to be a function that has a heavy side.

1.2.1 Definition of a Limit

In the previous few examples, we have tried computing the limit of a function at some x value by computing values close to that x value. While this seems like a reasonable way to proceed, it's fraught with danger! Let's see an example where things go wrong, highlighting the need for a better definition.

Example 3

Consider the function

$$f(x) = \begin{cases} 1 & \text{if } x = \frac{1}{n} \text{ where } n \text{ is a non-zero even integer} \\ -1 & \text{if } x = \frac{1}{n} \text{ where } n \text{ is an odd integer} \\ 0 & \text{otherwise.} \end{cases}$$

Let's investigate the behaviour of the function around $x = 0$. First note that $f(0) = 0$.

Now let's see what happens to the function as we get closer and closer to $x = 0$.

x	$f(x)$
$\pm \frac{1}{100}$	1
$\pm \frac{1}{1000}$	1
$\pm \frac{1}{10000}$	1
$\pm \frac{1}{100000}$	1
$\pm \frac{1}{1000000}$	1
$\pm \frac{1}{10000000}$	1

It seems reasonable to guess that the limit as x approaches 0 is 1.

However, little did we know, that as we were doing the computations above, our best friend from high school was also looking at the same problem. Here's what she computed:

x	$f(x)$
$\pm \frac{1}{99}$	-1
$\pm \frac{1}{999}$	-1
$\pm \frac{1}{9999}$	-1
$\pm \frac{1}{99999}$	-1
$\pm \frac{1}{999999}$	-1
$\pm \frac{1}{9999999}$	-1.

Well, now it looks like the limit should be -1 . This is a problem. To make matters worse, our old high school math teacher also did some investigating:

x	$f(x)$
$\pm \frac{3}{64}$	0
$\pm \frac{3}{128}$	0
$\pm \frac{3}{256}$	0
$\pm \frac{3}{512}$	0
$\pm \frac{3}{1024}$	0
$\pm \frac{3}{2048}$	0.

Perhaps you and your friend were wrong, and the limit is really 0. Just when you thought things couldn't get any worse, our old high school best friend's boyfriend also did some computations:

x	$f(x)$
$\pm \frac{1}{100001}$	-1
$\pm \frac{1}{200002}$	1
$\pm \frac{1}{300003}$	-1
$\pm \frac{1}{400004}$	1
$\pm \frac{1}{500005}$	-1
$\pm \frac{1}{600006}$	1.

He says "the function doesn't tend to any particular value as x approaches 0".

So who is right? In order to prevent a long argument, we need a formal definition!

Let's attempt to come up with a more formal definition. Consider the function $f(x) = 2x$,

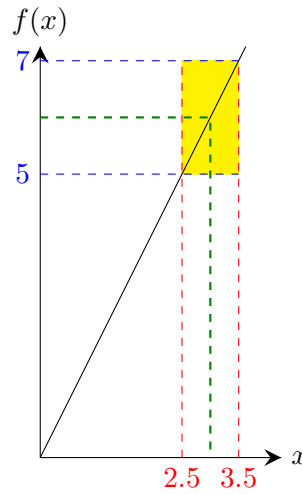
and imagine Michael and Janet are trying to figure out what the limit of the function is as x approaches 3. They both agree it should be 6 (and they will be correct!).

Janet has an idea of how they can convince themselves that the limit is indeed 6. She challenges Michael to give her a small positive number ϵ . She will try to find an interval around $x = 3$, inside which the function stays within ϵ of the proposed limit, which is 6.

“Good idea!” says Michael, “I give you $\epsilon = 1$.”

Janet thinks for a moment, and looks at a graph of the function (below).

“Easy,” she responds, “as long as you’re within 0.5 of 3, then $f(x)$ is within 1 of 6.”



More formally, if $3 - 0.5 < x < 3 + 0.5$, then $6 - 1 < f(x) < 6 + 1$.

“Fine,” Michael responds, “I’ll make it harder by choosing a really small number. What about $\epsilon = 0.5$?”

“Easy,” she says, “as long as x is within 0.25 of 3, then $f(x)$ is within 0.5 of 6.”

“What about $\epsilon = 0.01$?” Michael responds, “try that!”

“Got it,” replies Janet, “if the x value stays within 0.005 of 3, then $f(x)$ stays within 0.01 of 6. In fact, I can always win! If you give me any $\epsilon > 0$, as long as the x value stays within $\frac{\epsilon}{2}$, $f(x)$ is within ϵ of 6.”

Formally, what Janet realised is that for any number $\epsilon > 0$

$$\text{if } |x - 3| < \epsilon, \text{ then } |f(x) - 6| < \frac{\epsilon}{2}.$$

Note that the condition $|x - 3| < \epsilon$ is simply a reformulation of the statement “ x is within ϵ of 3.”

So, informally, Janet and Michael have shown that you can get arbitrarily close to $f(x) = 6$ by getting arbitrarily close to $x = 3$. This is essentially the definition of a limit!

Definition 1.2.1
limit of a function
at a point

Let f be a function defined on some open interval containing a except possibly at a itself. We write

$$\lim_{x \rightarrow a} f(x) = L \quad (1.1)$$

and say the limit of f as x tends to a is L if for every $\epsilon > 0$, there exists a $\delta > 0$ such that if

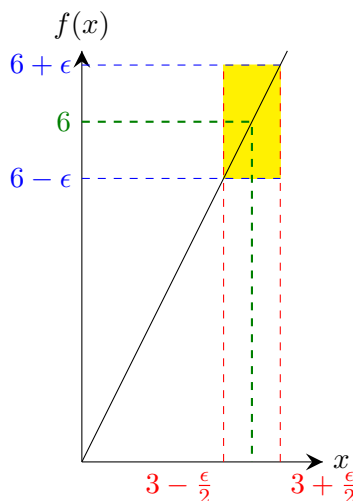
$$0 < |x - a| < \delta, \quad (1.2)$$

then

$$|f(x) - L| < \epsilon. \quad (1.3)$$

Keep in mind that we don't care about the value of $f(a)$ ($x = a$ may not even be in the domain of f). This is why we include the inequality $0 < |x - a|$ in the definition, as it excludes considering $x = a$.

It is also common to write $f(x) \rightarrow L$ as $x \rightarrow a$ in place of $\lim_{x \rightarrow a} f(x) = L$. Let's unpack graphically how this definition works. Consider the graph of $y = 2x$ and focus on the portion of this graph where the y -values are within ϵ of 6 for some arbitrary positive value of ϵ .



Observe that the function takes the value $6 - \epsilon$ at $x = 3 - \frac{\epsilon}{2}$ and the value $6 + \epsilon$ at $x = 3 + \frac{\epsilon}{2}$. Therefore, by taking $\delta = \frac{\epsilon}{2}$, we can ensure that for any x on the interval $(3 - \delta, 3 + \delta)$, the values of $f(x)$ are within ϵ of 6. Crucially, this argument (i.e., that $f(x)$ is within ϵ of 6 as x nears 3) works for arbitrarily small values of ϵ and therefore this must be the limiting value of $f(x)$.

Let's use the definition to recreate the argument Janet came up with above and show $\lim_{x \rightarrow 3} (2x) = 6$.

Example 4 Let $f(x) = 2x$. Use the definition of the limit to show $\lim_{x \rightarrow 3} f(x) = 6$.

Solution: Let ϵ be a positive real number so that $\epsilon > 0$ and take $\delta = \frac{\epsilon}{2}$.

Suppose $|x - 3| < \delta$, then we have

$$\begin{aligned} |f(x) - 6| &= |2x - 6| \\ &= 2|x - 3| \\ &< 2\delta \\ &= \epsilon. \end{aligned}$$

Therefore if $|x - 3| < \delta$, then $|f(x) - 6| < \epsilon$, and we conclude $\lim_{x \rightarrow 3} f(x) = 6$.

You may also be wondering how we decided to take $\delta = \frac{\epsilon}{2}$. One strategy is to write out the proof not knowing what δ should be in advance. Once the end of the proof is reached, it sometimes becomes clear what δ should be. You can then go back and rewrite the proof, with your chosen value of δ (often in terms of ϵ).

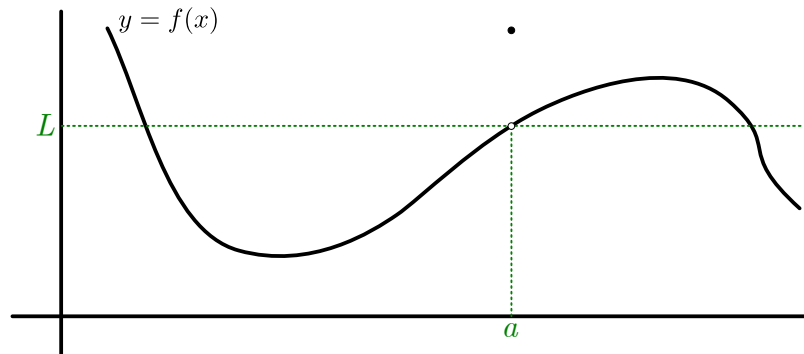
In our example above, we wanted $|f(x) - 6| < \epsilon$ or $|2x - 6| < \epsilon$. If we re write this as $|x - 3| < \frac{\epsilon}{2}$, it gives us a hint as to what δ should be. We can guess that $\delta = \frac{\epsilon}{2}$ could work, and now we have to go back and see if the proof still works. Sure enough, it does!

EXERCISE

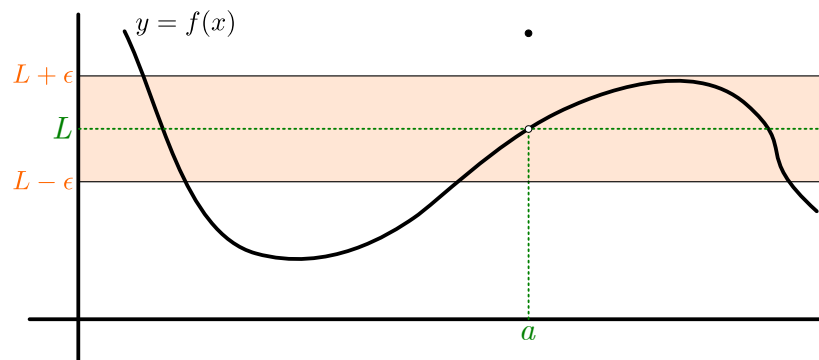
Use the definition of a limit to prove that $\lim_{x \rightarrow 2} \frac{x^2 - 4}{x - 2}$ is 4.

Before moving on, let's explore the definition of a limit graphically in a little more depth.

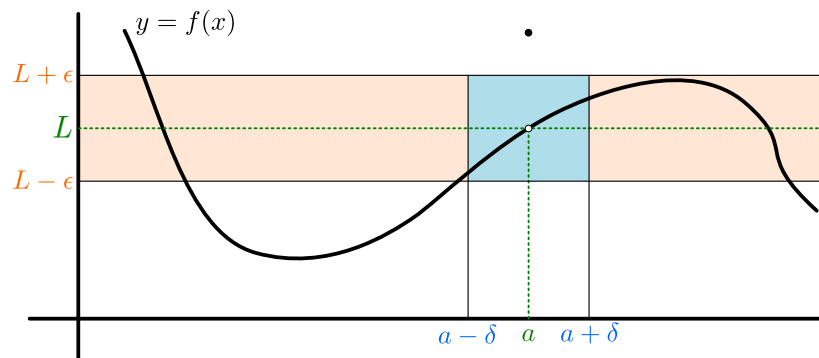
Suppose we are given a function f , and after drawing its graph, we suspect it has a limit of L at $x = a$ (even though it's clear that $f(a) \neq L$).



Now choose some positive number $\epsilon > 0$, which is going to act as a tolerance for how close we need our function to be to L . Our goal is to find a $\delta > 0$ so that if the x value of the function is within δ of a , then $f(x)$ is within ϵ of L . Graphically, we need to find an interval around a so that within that interval, the graph of the function is always between the horizontal lines $y = L + \epsilon$ and $y = L - \epsilon$.

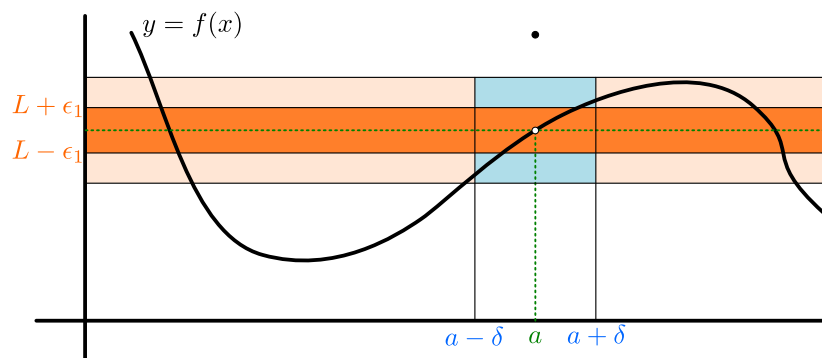


We can choose such a δ as in the following diagram:

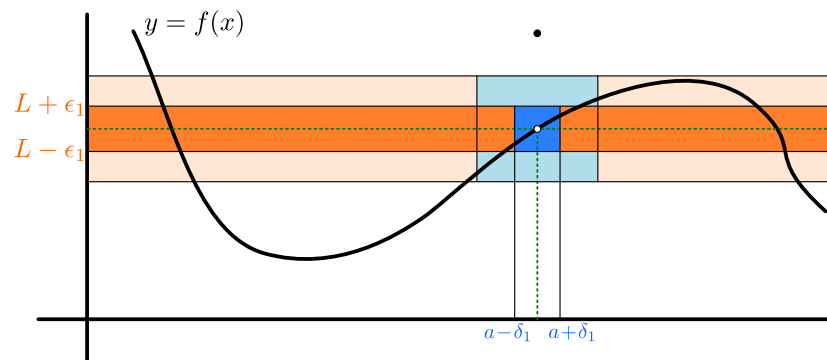


Notice how between the vertical lines $x = a - \delta$ and $x = a + \delta$, the function is between the horizontal lines $y = L - \epsilon$ and $y = L + \epsilon$. This is a graphical interpretation of the statement, if $|x - a| < \delta$, then $|f(x) - L| < \epsilon$. Great! We have shown that for this particular ϵ , there is a δ that does the job.

However, this is not enough to prove that $\lim_{x \rightarrow a} f(x) = L$. The definition requires us to, given *any* $\epsilon > 0$, find a δ . So, someone may now come along with a smaller value of ϵ , call it ϵ_1 , and now our value of δ is now longer good enough. Here is what that may look like:



Notice that there are now parts of the graph of $f(x)$ between $x = a - \delta$ and $x = a + \delta$ that lie outside the bounds of $y = L - \epsilon$ and $y = L + \epsilon$. All is not lost! We can now choose a smaller value, call it δ_1 , which works with ϵ_1 .



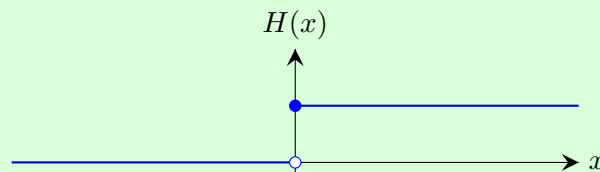
Now the graph of the function between the vertical lines $x = a - \delta_1$ and $a + \delta_1$ lies between the horizontal lines $y = L - \epsilon_1$ and $y = L + \epsilon_1$. This shows that if $|x - a| < \delta_1$, then $|f(x) - L| < \epsilon_1$.

If for *any* choice of ϵ (no matter how small), we can find a δ so that the graph of $f(x)$ between $x = a - \delta$ and $x = a + \delta$ lies between $y = L - \epsilon$ and $y = L + \epsilon$, then we have proved that $\lim_{x \rightarrow a} f(x) = L$.

Example 5

Recall the Heaviside function

$$H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0. \end{cases}$$



We will prove that $\lim_{x \rightarrow 0} H(x)$ does not exist. The key observation here is that the function makes a jump at $x = 0$, and the value of the function changes by 1.

Consider an interval $(-\delta, \delta)$ around $x = 0$, for some $\delta > 0$. Then there are x -values a and b in the interval so that $H(a) = 0$ and $H(b) = 1$ (for example, $a = -\frac{\delta}{2}$, $b = \frac{\delta}{2}$).

Now, suppose we try to show $\lim_{x \rightarrow 0} H(x) = L$. For any $\epsilon > 0$, we would need to find $\delta > 0$ so that for *every* $x \in (-\delta, \delta)$, $L - \epsilon < H(x) < L + \epsilon$. However, $L - \epsilon$ and $L + \epsilon$ are only 2ϵ apart. This spells trouble!

To see why, suppose ϵ is really small, for example $\epsilon = \frac{1}{10}$. Then we need all the values of $H(x)$, where x is in the interval $(-\delta, \delta)$, to lie within a range of $2\epsilon = \frac{1}{5}$. Unfortunately, we know there will always be two values that differ by 1, so this won't be possible! Let's turn this discussion into a proof.

Proof: Let L be a real number. We will show $\lim_{x \rightarrow 0} H(x) \neq L$. Since L is an arbitrary number, this will prove that the limit does not exist.

Let $\epsilon = \frac{1}{10}$. Suppose, towards a contradiction, that there exists $\delta > 0$ so that for all x satisfying $|x| < \delta$, $|H(x) - L| < \frac{1}{10}$. Note that $|\pm \frac{\delta}{2}| < \delta$. Therefore,

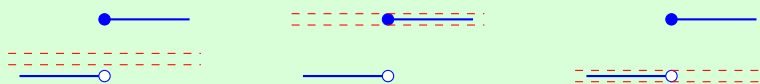
$$\begin{aligned} \left| H\left(\frac{\delta}{2}\right) - H\left(-\frac{\delta}{2}\right) \right| &= \left| H\left(\frac{\delta}{2}\right) - L + L - H\left(-\frac{\delta}{2}\right) \right| \\ &\leq \left| H\left(\frac{\delta}{2}\right) - L \right| + \left| L - H\left(-\frac{\delta}{2}\right) \right| \\ &= \left| H\left(\frac{\delta}{2}\right) - L \right| + \left| H\left(-\frac{\delta}{2}\right) - L \right| \\ &< \frac{1}{10} + \frac{1}{10} \\ &= \frac{1}{5}. \end{aligned}$$

However, $H\left(-\frac{\delta}{2}\right) = 0$, and $H\left(\frac{\delta}{2}\right) = 1$. Therefore, $\left| H\left(\frac{\delta}{2}\right) - H\left(-\frac{\delta}{2}\right) \right| = 1$, a contradiction.

We can now conclude that $\lim_{x \rightarrow 0} H(x)$ does not exist. \square

The contradiction above comes from the fact that $1 \not< 2\epsilon$, so we did not have to choose the particular value $\epsilon = \frac{1}{10}$. Any value of ϵ satisfying $1 \not< 2\epsilon$ would have sufficed.

Here is a graphical interpretation of the proof. No matter what L we choose, when $\epsilon = \frac{1}{10}$, there is no interval around zero so that the graph of the Heaviside function lies within the dotted red lines.



EXERCISE

Revisit the function defined in Example 3 and prove that the limit as $x \rightarrow 0$ of the function does not exist.

LOOKING AHEAD

Later on, we will not only consider limits as x approaches a particular value, but also as x approaches infinity.

What should the limit of $f(x) = \frac{1}{x}$ be as x approaches infinity?

What should the limit of $f(x) = x$ be as x approaches infinity?

What should the limit of $f(x) = \ln(x)$ be as x approaches infinity?

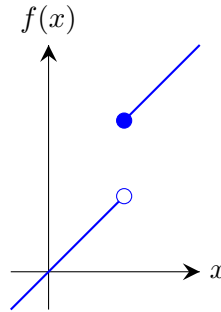
Try to write down a precise definition of what it means for the limit of $f(x)$ to be a as x approaches infinity.

One-Sided Limits

Consider the piecewise-defined function

$$f(x) = \begin{cases} x & \text{if } x < 1 \\ x + 1 & \text{if } x \geq 1. \end{cases}$$

Here is a graph of the function:



The graph of the function suggests that the limit does not exist at $x = 1$.

To show this, we first observe that for every $\delta > 0$, we can find pairs x_1 and x_2 such that $|f(x_1) - f(x_2)| > 1$ while $0 < |x_1 - 1| < \delta$ and $0 < |x_2 - 1| < \delta$ by taking x_1 and x_2 on opposite sides of $x = 1$.

We now use this observation to show the limit does not exist. We'll use proof by contradiction so we begin by assuming that a limit does exist. That is, assume $\lim_{x \rightarrow 1} f(x) = L$ for some real number L .

Next, let $\epsilon = 0.1$ and suppose $\delta > 0$ exists such that if $0 < |x - 1| < \delta$, then $|f(x) - L| < 0.1$.

Now let x_1 and x_2 satisfy the properties above as well as being within a distance of δ of $x = 1$ so that $|f(x_1) - L| < \epsilon$ and $|f(x_2) - L| < \epsilon$. It follows that

$$\begin{aligned} 2\epsilon &> |f(x_1) - L| + |f(x_2) - L| \\ &= |f(x_1) - L| + |L - f(x_2)| \\ &\geq |f(x_1) - L + L - f(x_2)| && \text{Triangle inequality} \\ &= |f(x_1) - f(x_2)| \\ &> 1 \end{aligned}$$

It follows from our assumption that $2\epsilon > 1$ or $\epsilon > 0.5$. This contradicts that ϵ can be any positive real number and, in this case, we took it to be 0.1 which is certainly not greater than 0.5. Since we have reached a contradiction, we conclude that the limit does not exist.

Observe that the key to this argument is that we could pick x_1 and x_2 to be on opposite sides of where we're taking the limit. If we require x_1 and x_2 to be on the same side of $x = 1$, then we would not be able to arrive at a contradiction. This makes sense because it is equivalent to asking if the function is approaching a specific value if we approach $x = 1$ only from one direction. Indeed, if we approach from the left, the function tends towards a value of 1 while if we approach from the right, the function tends towards a value of 2. This motivates us to define one-sided limits.

Definition 1.2.2

**one-sided limit
from the right**

Let f be a function whose domain includes an open interval to the **right** of a . We write

$$\lim_{x \rightarrow a^+} f(x) = L \quad (1.4)$$

and say the limit of f as x tends to a from the right is L if for every $\epsilon > 0$, there exists a $\delta > 0$ such that if

$$0 < x - a < \delta, \quad (1.5)$$

then

$$|f(x) - L| < \epsilon. \quad (1.6)$$

In words, this definition says that the one-sided limit from the right of a function f is L if we can make $f(x)$ as close to L as we like by taking x sufficiently close to a with $x > a$.

EXERCISE

A key difference between the definition above and that of a regular limit is the absence of the absolute value bars in Equation 1.5 compared to Equation 1.2. Convince yourself of why this is and then construct a similar definition a one-sided limit from the left.

REMARK

The one-sided limit from the left is written $\lim_{x \rightarrow a^-} f(x)$ and is defined analogously to the one-sided limit from the right.

Example 6

Recall the Heaviside function defined by

$$H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0. \end{cases}$$

We showed earlier that $\lim_{x \rightarrow 0} H(x)$ does not exist. However, $\lim_{x \rightarrow 0^+} H(x) = 1$ and $\lim_{x \rightarrow 0^-} H(x) = 0$. Let's prove the first of these statements.

Suppose $\epsilon > 0$ is given. Let $\delta = 1$. Then if $0 < x < \delta$, $H(x) = 1$ and so $|H(x) - 1| = 0 < \epsilon$. Therefore, $\lim_{x \rightarrow 0^+} H(x) = 1$.

See if you can prove that $\lim_{x \rightarrow 0^-} H(x) = 0$.

EXERCISE

Prove that $\lim_{x \rightarrow a} f(x) = L$ if and only if the one-sided limits of $f(x)$ both exist and are equal.

1.3 Limit Properties

It's a real pain to always have to use the definition using ϵ and δ to compute limits. Thankfully, we can prove some general results which make our computations much easier.

Example 7

For all $a \in \mathbb{R}$ we have

(i) $\lim_{x \rightarrow a} c = c$, and

(ii) $\lim_{x \rightarrow a} x = a$.

Before proving these, sketch out the graphs of these functions and convince yourself that these are believable statements!

For (i), let $\epsilon > 0$ and let $\delta = 1$. If $0 < |x - a| < \delta$, then

$$|f(x) - c| = |c - c| = 0 < \epsilon.$$

Therefore, $\lim_{x \rightarrow a} c = c$. Note here that we could have chosen any value of $\delta > 0$ and the proof would have worked!

For (ii), let $\epsilon > 0$ and let $\delta = \epsilon$. Suppose x is such that $0 < |x - a| < \delta$. Then

$$|f(x) - a| = |x - a| < \delta = \epsilon.$$

Therefore, $\lim_{x \rightarrow a} x = a$.

With a bit more work, we can show that limits obey a handful of useful arithmetic properties.

Theorem 1 (Limit Properties)

Let f and g be functions and let $a \in \mathbb{R}$. Suppose $\lim_{x \rightarrow a} f(x)$ and $\lim_{x \rightarrow a} g(x)$ both exist, then

(i) $\lim_{x \rightarrow a} c = c$ for $c \in \mathbb{R}$

(ii) $\lim_{x \rightarrow a} x = a$

(iii) $\lim_{x \rightarrow a} (c f(x)) = c \lim_{x \rightarrow a} f(x)$ for $c \in \mathbb{R}$

(iv) $\lim_{x \rightarrow a} (f(x) + g(x)) = \lim_{x \rightarrow a} f(x) + \lim_{x \rightarrow a} g(x)$

(v) $\lim_{x \rightarrow a} (f(x)g(x)) = \left(\lim_{x \rightarrow a} f(x) \right) \left(\lim_{x \rightarrow a} g(x) \right)$

(vi) $\lim_{x \rightarrow a} \left(\frac{f(x)}{g(x)} \right) = \frac{\lim_{x \rightarrow a} f(x)}{\lim_{x \rightarrow a} g(x)}$ provided $\lim_{x \rightarrow a} g(x) \neq 0$

(vii) $\lim_{x \rightarrow a} (f(x))^r = \left(\lim_{x \rightarrow a} f(x) \right)^r$ whenever $\lim_{x \rightarrow a} f(x) > 0$ and $r > 0$

Proof: Properties (i) and (ii) were proved in Example 7. Let's prove (iii).

We are assuming that the limit $\lim_{x \rightarrow a} f(x)$ exists, so let's call this limit L . First, if $c = 0$, then the result follows from (i). So, we can assume $c \neq 0$.

We want to show that $\lim_{x \rightarrow a} (cf(x)) = cL$. The fact that $\lim_{x \rightarrow a} f(x) = L$ means that for every $\epsilon > 0$, there is a $\delta > 0$ with the property that if $0 < |x - a| < \delta$, then $|f(x) - L| < \epsilon$. Let's exploit the existence of such a δ for all ϵ to prove what we desire.

Let $\epsilon > 0$, and let δ be such that if $0 < |x - a| < \delta$, then $|f(x) - L| < \frac{\epsilon}{|c|}$ (take a moment and convince yourself that such a δ exists, this is the key point in the proof!). Now, suppose $x \in \mathbb{R}$ is such that $0 < |x - a| < \delta$. Then

$$|cf(x) - cL| = |c| |f(x) - L| < |c| \frac{\epsilon}{|c|} = \epsilon.$$

Therefore, $\lim_{x \rightarrow a} cf(x) = cL$, completing the proof.

We will leave the remaining properties unproved, but you should try to prove them yourself! \square

Equipped with Theorem 1, we can now compute limits of much more complicated functions.

Example 8

Let's compute $\lim_{x \rightarrow 3} (x^2 + 2x - 4)$. Applying the various limit properties we can manipulate the limit as follows:

$$\begin{aligned} \lim_{x \rightarrow 3} (x^2 + 2x - 4) &= \lim_{x \rightarrow 3} (x^2) + \lim_{x \rightarrow 3} (2x) + \lim_{x \rightarrow 3} (-4) \\ &= \left(\lim_{x \rightarrow 3} (x) \right)^2 + 2 \lim_{x \rightarrow 3} (x) + (-4) \\ &= 3^2 + 2 \cdot 3 - 4 \\ &= 11. \end{aligned}$$

It is not a coincidence that if $f(x) = x^2 + 2x - 4$, then $f(3) = \lim_{x \rightarrow 3} f(x)$.

We can generalize Example 8 to evaluate limits of all polynomials.

Theorem 2 (Limits of Polynomials)

Let $P(x) = c_0 + c_1x + c_2x^2 + \cdots + c_nx^n$, then

$$\lim_{x \rightarrow a} P(x) = P(a)$$

EXERCISE

Prove Theorem 2. (Hint: See Example 8 for inspiration.)

EXERCISE

Suppose $f(x)$ and $g(x)$ are functions so that $\lim_{x \rightarrow 3} f(x) = 4$ and $\lim_{x \rightarrow 3} g(x) = 7$. Evaluate

$$\lim_{x \rightarrow 3} \frac{f(x)g(x)}{f(x) + g(x)}.$$

Example 9

Let $f(x) = \frac{x^2 - 1}{x + 1}$. Evaluate $\lim_{x \rightarrow -1} f(x)$.

It is tempting to just use the limit properties here. However, $\lim_{x \rightarrow -1} (x + 1) = 0$ so we cannot use Property (vi).

Here's a useful trick. The limit only depends on the behaviour of the function around $x = -1$, but not at $x = -1$. When $x \neq -1$ we have

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

So, for any $a \neq -1$, $f(a) = a - 1$. Therefore,

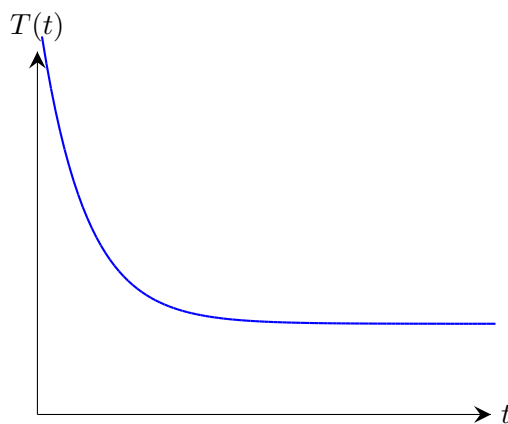
$$\lim_{x \rightarrow -1} f(x) = \lim_{x \rightarrow -1} (x - 1) = -2.$$

1.4 Limits at Infinity

All of our limits so far have looked at the behaviour of a function $f(x)$ as x tends towards some finite value. However, it is very often the case that we want to know the behaviour of a function at arbitrarily large (positive or negative) values of x . In these cases, we write the limit with x approaching ∞ or $-\infty$.

As a concrete example, consider a cup of coffee with an initial temperature of 90°C sitting in a room with ambient temperature 20°C . Our experience tells us that the coffee will cool down by 70° to be at thermal equilibrium with its surroundings. How do we show this mathematically?

The first step is to find a function, T , describing the temperature of the coffee as a function of time, t . This involves setting up and solving a differential equation modelled after Newton's Law of Cooling. This requires a bit more machinery than we have at the moment but, for our example, the result is a function of the form $T(t) = 20 + 70e^{-kt}$ where $k > 0$ is a constant.



To determine the *eventual* temperature of the cup of coffee, we want to evaluate the limit of $T(t)$ as $t \rightarrow \infty$. Using limit properties and rescaling the time variable by letting $x = kt$, we have

$$\lim_{t \rightarrow \infty} T(t) = 20 + 70 \lim_{x \rightarrow \infty} e^{-x}$$

The question then boils down to evaluating $\lim_{x \rightarrow \infty} e^{-x}$. We know decaying exponential functions tend to zero, so we can reason that the limit is zero. With that in mind, we can also make better sense of the temperature function $T(t) = 20 + 70e^{-kt}$. The temperature of the coffee is 20° plus a part which gives the *difference* in temperature between the coffee and its surroundings and eventually decays to zero.

Let's now adapt our definition of a limit to define these limits at infinity. The key difference now is that instead of supposing there will exist some quantity δ near a which will guarantee $f(x)$ is within an amount ϵ of the limit L , we will suppose that x can always be taken large enough such that $f(x)$ is within an amount ϵ of the limit L .

Definition 1.4.1

limit at positive
infinity

We say that the limit of $f(x)$ as x approaches ∞ is L and write

$$\lim_{x \rightarrow \infty} f(x) = L \tag{1.7}$$

if for every $\epsilon > 0$, there exists some real number k such that if

$$x > k, \tag{1.8}$$

then

$$|f(x) - L| < \epsilon. \tag{1.9}$$

A similar definition applies for limits at negative infinity which we denote by $\lim_{x \rightarrow -\infty} f(x)$.

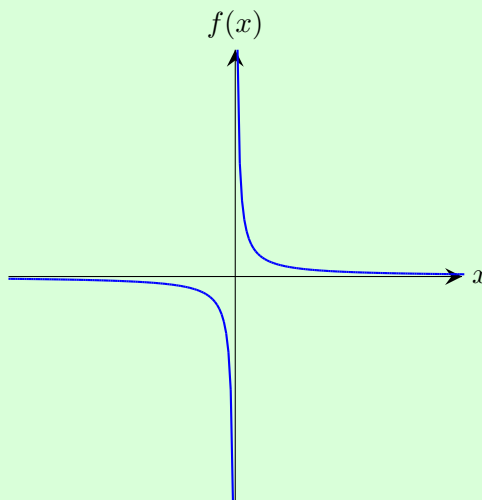
EXERCISE

Write down a formal definition for the meaning of the equation $\lim_{x \rightarrow -\infty} f(x) = L$.

We start with an important example, which will come in handy as we move forward.

Example 10

Consider the function $f(x) = \frac{1}{x}$. Here is its graph:



As x gets really really big, and really really small, $f(x)$ gets closer and closer to 0. Let's prove that

$$\lim_{x \rightarrow \infty} \frac{1}{x} = \lim_{x \rightarrow -\infty} \frac{1}{x} = 0.$$

Let $\epsilon > 0$. We want to find a value k so that $x > k$ implies $|\frac{1}{x}| < \epsilon$. Set $k = \frac{1}{\epsilon}$. Then if $x > k$ we have

$$\left| \frac{1}{x} \right| = \frac{1}{x} < \frac{1}{k} = \epsilon.$$

Therefore, $\lim_{x \rightarrow \infty} f(x) = 0$.

For the limit to $-\infty$, if $\epsilon > 0$ is given, we can let $k = -\frac{1}{\epsilon}$. Then if $x < k$ we have

$$\left| \frac{1}{x} \right| = -\frac{1}{x} < -\frac{1}{k} = \epsilon.$$

Therefore $\lim_{x \rightarrow -\infty} f(x) = 0$.

REMARK

All of the properties from Theorem 1 are valid for limits at infinity.

Example 11

Evaluate $\lim_{x \rightarrow -\infty} \frac{x-1}{x^2-1}$.

Solution: The definition of the limit as $x \rightarrow -\infty$ of a function only depends on the behaviour of the function near $-\infty$. More precisely, if two functions agree for all $x < -1000$, for example, then the two functions have the same limit at $-\infty$.

With that in mind, let's assume that $x < 0$ (so in particular, $x \neq 0$). Then

$$\frac{x-1}{x^2-1} = \frac{\frac{1}{x} - \frac{1}{x^2}}{1 - \frac{1}{x^2}}.$$

Now we can use our properties of limits. We have

$$\begin{aligned} \lim_{x \rightarrow -\infty} \frac{x-1}{x^2-1} &= \lim_{x \rightarrow -\infty} \frac{\frac{1}{x} - \frac{1}{x^2}}{1 - \frac{1}{x^2}} \\ &= \frac{\lim_{x \rightarrow -\infty} \frac{1}{x} - \lim_{x \rightarrow -\infty} \frac{1}{x^2}}{\lim_{x \rightarrow -\infty} 1 - \lim_{x \rightarrow -\infty} \frac{1}{x^2}} \\ &= \frac{0 - \left(\lim_{x \rightarrow -\infty} \frac{1}{x}\right)^2}{1 - \left(\lim_{x \rightarrow -\infty} \frac{1}{x}\right)^2} \\ &= \frac{0}{1} \\ &= 0. \end{aligned}$$

Note that the limit of the denominator is not 0 (it's 1), so we can apply the limit properties without any problems.

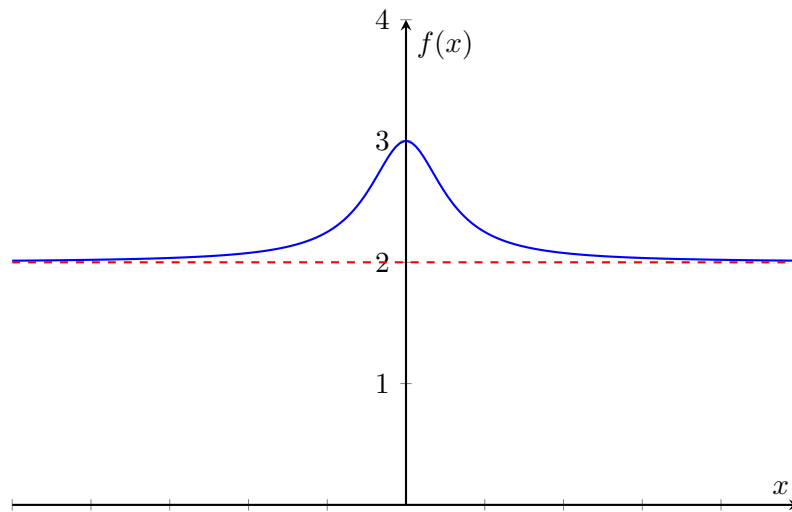
When you have a rational function, that is, a quotient of two polynomials, dividing out by the highest power of x can often be a useful manipulation!

EXERCISE

Evaluate $\lim_{x \rightarrow \infty} \frac{6x^2 + 3}{3x^2 + 1}$.

Horizontal Asymptotes

Observe that the graph of the function $f(x) = \frac{6x^2 + 3}{3x^2 + 1}$ has a horizontal asymptote at $y = 2$.



In other words, the values of $f(x)$ tend to 2 as x gets really big (in this case, positive or negative). It's no coincidence then that the limit in the previous exercise is equal to 2 as this is precisely the behaviour of a function which a limit at infinity is probing.

Definition 1.4.2

horizontal
asymptote

We say that the line $y = L$ is a **horizontal asymptote** of $f(x)$ if $\lim_{x \rightarrow \pm\infty} f(x) = L$.

REMARK

A function may not approach the same value at both ∞ and $-\infty$. This means a function can have zero, one, or two horizontal asymptotes. Moreover, unlike vertical asymptotes which we'll talk about soon, the graph of a function may cross any of its horizontal asymptotes.

EXERCISE

Locate all horizontal asymptotes of $f(x) = \tan^{-1}(x)$.

1.5 Infinite Limits

You might think that we just covered infinite limits in the last section titled limits at infinity but they're not the same thing at all! Consider the expression $\lim_{x \rightarrow a} f(x) = L$. If the value a (i.e., the value x is approaching) is infinite (i.e., $x \rightarrow \pm\infty$), then you have a limit at infinity. In contrast, if the value of the function is getting infinitely large as x tends to a , then we say you have an infinite limit. In this case, we cannot write down a finite value in place of the quantity L . Instead, we replace L with either ∞ or $-\infty$ and say the function diverges as $x \rightarrow a$.

We wish for $\lim_{x \rightarrow a} f(x) = \infty$ to morally mean that $f(x)$ can be made arbitrarily large by taking x to be sufficiently close to a . This is good intuition to have, but let's give a definition with which we can work.

Definition 1.5.1

positive infinite
limit

We say that

$$\lim_{x \rightarrow a} f(x) = \infty$$

if for all real numbers $k > 0$, there exists $\delta > 0$ so that $0 < |x - a| < \delta$ implies $f(x) > k$.

A similar definition holds for negative infinite limits, and for one-sided limits. In particular, if $f(x)$ grows arbitrary large and negative as $x \rightarrow a$, then we write $\lim_{x \rightarrow a} f(x) = -\infty$.

EXERCISE

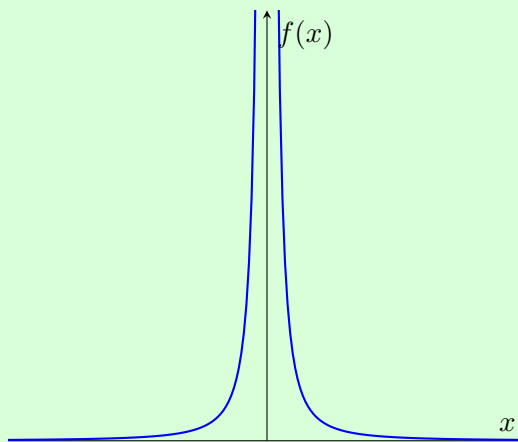
Write down formal definitions for the following limits

1. $\lim_{x \rightarrow a} f(x) = -\infty$.
2. $\lim_{x \rightarrow a^+} f(x) = \infty$.
3. $\lim_{x \rightarrow a^-} f(x) = -\infty$.

Example 12

Determine $\lim_{x \rightarrow 0} \frac{1}{x^2}$.

Before computing this limit, let's graph the function so that we have a guess as to what the limit should be.



As x gets closer to 0, the graph suggests that $f(x)$ goes off to $+\infty$. Let's prove it.

Let $k > 0$ and let $\delta = \frac{1}{\sqrt{k}}$. If $|x| < \delta$, then

$$\left| \frac{1}{x^2} \right| = \frac{1}{x^2} > \frac{1}{\delta^2} = k.$$

Therefore $\lim_{x \rightarrow 0} \frac{1}{x^2} = \infty$.

REMARK

It is important to keep in mind that when we write something like

$$\lim_{x \rightarrow 0} \frac{1}{x^2} = \infty$$

we do **not** mean to suggest that $f(x) = \frac{1}{x^2}$ has a well-defined limit as $x \rightarrow 0$. The reality is precisely the opposite. Since f grows without bound as $x \rightarrow 0$, then the limit does **not** exist. However, it is still useful to write an expression like this to describe in what way the function diverges as $x \rightarrow 0$.

EXERCISE

Consider the function $f(x) = \frac{1}{x-1}$. Determine the one-sided limits of $f(x)$ as $x \rightarrow 1$.

For a function to have an infinite limit as $x \rightarrow a^\pm$, it must be the case that f is defined on some open interval on one or both sides of a . Moreover, the graph of f comes arbitrarily close to the vertical line $x = a$ while growing arbitrarily large, but it must not cross this line otherwise the crossing point would be the (finite-valued) limit. Graphically, this means $y = f(x)$ has a vertical asymptote at $x = a$.

Definition 1.5.2
vertical asymptote

We say that the line $x = a$ is a **vertical asymptote** of $f(x)$ if one or both of the following hold:

- $\lim_{x \rightarrow a^+} f(x) = \pm\infty$, or
- $\lim_{x \rightarrow a^-} f(x) = \pm\infty$.

EXERCISE

Locate all vertical asymptotes of $f(x) = \frac{1}{x^2 - 4}$. Next, inspect the sign of $f(x)$ on either side of each asymptote to help you generate a sketch of $y = f(x)$.

We will also encounter infinite limits when taking limits at infinity. This happens when a function grows arbitrarily large as x gets arbitrarily large. For example, when $f(x) = x$, we have

$$\lim_{x \rightarrow \pm\infty} f(x) = \pm\infty$$

Similarly, for any polynomial of degree $n \geq 1$, the limits at infinity will be infinite. Here is a formal definition.

Definition 1.5.3
Infinite limits at
infinity

We say $\lim_{x \rightarrow \infty} f(x) = \infty$ if for every $k > 0$ there exists $t > 0$ such that $x > t$ implies $f(x) > k$.

EXERCISE

Write down definitions for the following:

- $\lim_{x \rightarrow \infty} f(x) = -\infty$.
- $\lim_{x \rightarrow -\infty} f(x) = \infty$.
- $\lim_{x \rightarrow -\infty} f(x) = -\infty$.

Example 13

Prove that $\lim_{x \rightarrow \infty} x = \infty$.

Proof: Let $f(x) = x$, and let $k > 0$ be given. Let $t = k$. Then if $x > t$, $f(x) = x > t = k$, completing the proof. \square

This proof seems like we didn't do very much at all. It's an important exercise for you to make sure you understand what was proved here!

EXERCISE

Let n be a positive integer. Prove that $\lim_{x \rightarrow \infty} x^n = \infty$.

REMARK

When we want to formally capture a statement like “as x gets close to a ”, we interpret that as $|x - a| < \delta$ for some small $\delta > 0$. The smaller the δ , the closer x is to a .

Similarly, to turn a statement like “as x gets arbitrarily large” into something we can work with mathematically, we say that $x > k$ for some k . The larger k is, the “closer” x is to $+\infty$. Of course, here when we say “closer”, we don't mean it in any precise mathematical sense, but instead in an intuitive sense.

To summarise, a neighbourhood of a point $a \in \mathbb{R}$ is an interval of the form $(a - \delta, a + \delta) \subset \mathbb{R}$. Analogously, a neighbourhood of $+\infty$ is an interval of the form $(k, \infty) \subset \mathbb{R}$.

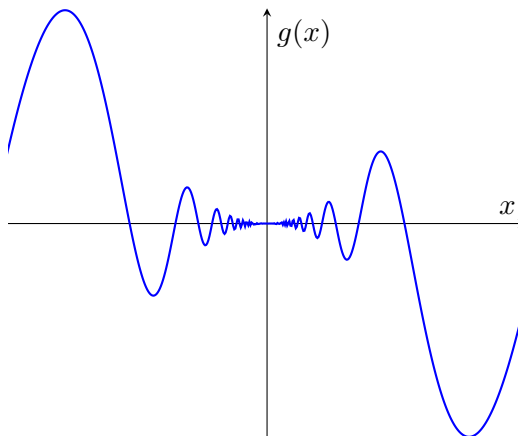
It's worth comparing the definitions of

$$\lim_{x \rightarrow a} f(x) = L, \quad \lim_{x \rightarrow \infty} f(x) = L, \quad \lim_{x \rightarrow a} f(x) = \infty, \quad \text{and} \quad \lim_{x \rightarrow \infty} f(x) = \infty$$

through this lens.

1.6 The Squeeze Theorem

Consider the limit of $g(x) = x^2 \sin\left(\frac{1}{x}\right)$ as $x \rightarrow 0$. A graph of this function is shown below.



It appears that the limit in question is zero and it is, but how do we justify this? It is tempting to use our limit properties (specifically property (v)) and say that $g(x)$ is the product of two functions: x^2 and $\sin\left(\frac{1}{x}\right)$ and then reason that since $x^2 \rightarrow 0$ as $x \rightarrow 0$, then $g(x)$ must go to zero. The problem is that this ignores the possibility that the other function (in this case, $\sin\left(\frac{1}{x}\right)$) behaves in such a way as to make the limit non-zero. Indeed, a more careful review of our limit properties would require that the limit of $\sin\left(\frac{1}{x}\right)$ as $x \rightarrow 0$ exist to take this approach, but this limit does *not* exist. We need to come up with a better argument.

Recall that $-1 \leq \sin(c) \leq 1$ for any c . This means that

$$-1 \leq \sin\left(\frac{1}{x}\right) \leq 1$$

for all x . Let's multiply this inequality by x^2 and note that we don't need to worry about the directions of inequality signs changing since $x^2 \geq 0$ for all x . This gives the expression

$$-x^2 \leq x^2 \sin\left(\frac{1}{x}\right) \leq x^2$$

Importantly, since this inequality holds for all values of x , it must hold in an open interval around $x = 0$ where we want to take the limit of $g(x) = x^2 \sin\left(\frac{1}{x}\right)$. This means, we can take the limit as $x \rightarrow 0$ of each term in the inequality and still have the inequality hold.

$$\lim_{x \rightarrow 0} -x^2 \leq \lim_{x \rightarrow 0} x^2 \sin\left(\frac{1}{x}\right) \leq \lim_{x \rightarrow 0} x^2$$

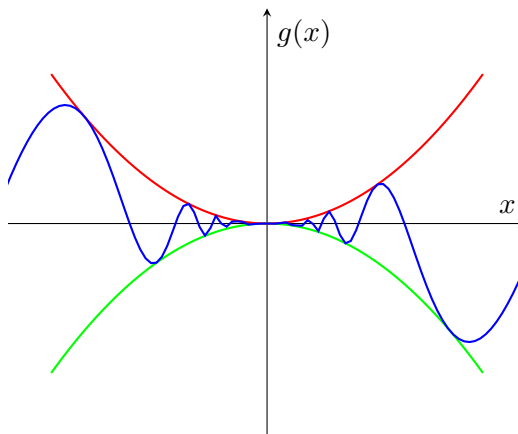
The limits on the left and right sides in this expression are just zero, so we can simplify this to.

$$0 \leq \lim_{x \rightarrow 0} x^2 \sin\left(\frac{1}{x}\right) \leq 0$$

Since the desired limits is not less than zero and not more than zero, it must be exactly zero.

$$\lim_{x \rightarrow 0} x^2 \sin\left(\frac{1}{x}\right) = 0$$

This argument is an application of what is called the Squeeze Theorem (also known as the Sandwich Theorem). The diagram below helps to show why this is a great name for the theorem.



Observe that the functions $f(x) = -x^2$ and $h(x) = x^2$ bound $g(x)$ on either side (i.e., $f(x) \leq g(x) \leq h(x)$) and have equal limits at the point where we're interested in finding the limit of $g(x)$. As such, we can infer that the range of possible values for the limit of $g(x)$ as $x \rightarrow 0$ gets 'squeezed' to a single point. Let's state this theorem more precisely.

Theorem 3 (Squeeze Theorem)

Let f , g , and h be functions defined on some open interval I containing $x = a$ (except possibly at $x = a$). If

$$f(x) \leq g(x) \leq h(x)$$

for all $x \in I$ and

$$\lim_{x \rightarrow a} f(x) = L = \lim_{x \rightarrow a} h(x),$$

then

$$\lim_{x \rightarrow a} g(x) = L.$$

Proof: Let f , g , and h be defined on the open interval I containing $x = a$. Assume that $f(x) \leq g(x) \leq h(x)$ for all $x \in I$ and $\lim_{x \rightarrow a} f(x) = L = \lim_{x \rightarrow a} h(x)$.

By the definition of the limit, for every $\epsilon > 0$, there exists a $\delta_1 > 0$ such that if

$$0 < |x - a| < \delta_1 \quad \implies \quad |f(x) - L| < \epsilon$$

Similarly, there exists a $\delta_2 > 0$ such that if

$$0 < |x - a| < \delta_2 \quad \implies \quad |h(x) - L| < \epsilon$$

Let δ_0 be the lesser of δ_1 and δ_2 and observe that if we replace both δ_1 and δ_2 with δ_0 , then both implications still hold since in one case there is no change and in the other we are actually restricting ourselves to values of x which are *closer* to a .

By the assumption that $f(x) \leq g(x) \leq h(x)$ on I , it follows that

$$f(x) - L \leq g(x) - L \leq h(x) - L$$

Now, since $|f(x) - L| < \epsilon$ for $0 < |x - a| < \delta_0$, it follows that $-\epsilon < f(x) - L$ on this same interval.

Similarly, since $|h(x) - L| < \epsilon$ for $0 < |x - a| < \delta_0$, we have that $h(x) - L < \epsilon$ on this same interval.

This allows us to expand the earlier inequality to read

$$-\epsilon < f(x) - L \leq g(x) - L \leq h(x) - L < \epsilon$$

from which we get

$$|g(x) - L| < \epsilon$$

Therefore, by the definition of a limit, we have

$$\lim_{x \rightarrow a} g(x) = L$$

□

REMARK

The squeeze theorem also holds for one-sided limits and limits at infinity

LOOKING AHEAD

Later on we will deal with derivatives of trigonometric functions. At the heart of a lot of the results relating to trigonometric functions is the important result that

$$\lim_{x \rightarrow \infty} \frac{\sin(x)}{x} = 1.$$

One way to prove this is to use the squeeze theorem. Here we'll walk you through one geometric approach to proving this limit.

- (i) Consider a triangle $\triangle ABC$ right-angled at B with side length $AB = 1$ and $\angle CAB = x$ (where x is in radians of course!). Let D be on the hypotenuse AC so that the length of AD is 1. By comparing various areas, prove that for all $x \in \mathbb{R}$,

$$\frac{1}{2} \tan(x) \geq \frac{1}{2} x \geq \frac{1}{2} \sin(x).$$

- (ii) Prove that for all x , $\cos(x) \leq \frac{\sin(x)}{x} \leq 1$.

- (iii) The cosine function has the property that for all $a \in \mathbb{R}$, $\lim_{x \rightarrow a} \cos(x) = \cos(a)$. Using this fact (and the squeeze theorem), prove that

$$\lim_{x \rightarrow 0} \frac{\sin(x)}{x} = 1.$$

Part (iii) hinged on the fact that to evaluate the limit at a point $x = a$ for the cosine function, you simply evaluate the function at $x = a$! Functions with this property are called continuous.

So far in this course we have proved that some functions are continuous. Which functions do we know are continuous?

Chapter 2

Continuity

2.1 Introduction to Continuity

Our goal in the previous chapter of computing the instantaneous rate of change of a function motivated the need for limits. The same goal, among others, motivates the need to introduce a notion of continuity of functions. Roughly speaking, if a function is continuous over some interval, then you could sketch the graph of that function (over that interval) without lifting your pencil. It seems sensible to expect that a function need to have this property if we want to compute its rate of change. In particular, it only makes sense to define a rate of change if the function exists across some interval of finite width and the value of the function at any point is determined by the value of the function near to that point.

Many familiar phenomenon are continuous. For example, the velocity of a car with respect to time is a continuous function. One believable consequence of the continuity of velocity is that if a car is travelling at 20 km/h at some point in time, and 30 km/h 20 seconds later, then we can conclude that at some point in those first 20 seconds, the car was travelling at 25 km/h. This is an example of something called the *intermediate value theorem*, which we will see in Section 2.5.

This notion of continuity can be defined formally using the limit.

Definition 2.1.1
continuity at a
point

We say a function f is continuous at $x = a$ if

$$\lim_{x \rightarrow a} f(x) = f(a)$$

REMARK

If a function f is continuous at $x = a$, then it is implicit by the definition of continuity that both the $\lim_{x \rightarrow a} f(x)$ exists and $f(a)$ exists.

Recall that, in general, the limit of $f(x)$ as $x \rightarrow a$ doesn't care about the actual value of f at $x = a$. Instead, it tells us what value the function approaches as x tends to a from the left and right. For a continuous function, we have that these two values are the same. In

other words, the value the function takes at a point is what one would expect based on its behaviour approaching that point from either side.

Observe that when a function is known to be continuous at a point, then we can use the definition of continuity to directly evaluate the limit of the function approaching that point using direct substitution.

Example 1

Consider the function $f(x) = 3x^2 + 2x + 1$. Since f is a polynomial, we know that $\lim_{x \rightarrow 2} f(x) = f(2)$. Therefore, f is continuous at $x = 2$. In fact, f is continuous at $x = a$ for all $a \in \mathbb{R}$!

It is natural and very useful to extend the notion of continuity of a function at a point to continuity over an interval of points.

Definition 2.1.2

continuity on an interval

We say a function f is continuous on an interval (a, b) if it is continuous at every point $x \in (a, b)$.

When a function is continuous at every point $x \in \mathbb{R}$, we say the function is continuous on \mathbb{R} or, more simply, just that the function is continuous. So, as we saw above, the polynomial $f(x) = 3x^2 + 2x + 1$ is continuous. We will see below that many familiar functions are continuous on their domains.

EXERCISE

Give an example of a function that is continuous at every real number except for $x = 3$.

The two main theorems in this Chapter, the Intermediate Value Theorem and the Extreme Value Theorem, both need continuity on a *closed* interval. The question then becomes, what does it mean for a function to be continuous on a closed interval $[a, b]$?

Well, the idea is that as you approach a from the right, the function should be continuous. That is, the value of the function at $x = a$ should be the value the function approaches as you get closer and closer to a from the right. The same should be true as you approach b from the left. Thankfully, we already have the notion of a one-sided limit to make this more concrete.

Definition 2.1.3

One-sided continuity

We say a function $f(x)$ is continuous from the right at $x = a$ if $\lim_{x \rightarrow a^+} f(x) = f(a)$. Similarly, $f(x)$ is continuous from the left at $x = a$ if $\lim_{x \rightarrow a^-} f(x) = f(a)$.

With one-sided continuity at our disposal, we can define continuity on a closed interval.

Definition 2.1.4

continuity on a closed interval

A function $f(x)$ is continuous on the closed interval $[a, b]$ if

- $f(x)$ is continuous on (a, b) ,
- $f(x)$ is continuous from the right at $x = a$, and
- $f(x)$ is continuous from the left at $x = b$.

Example 2 Recall again the Heaviside function

$$H(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0. \end{cases}$$

The Heaviside function is continuous on the closed interval $[0, 1]$, but not on the closed interval $[-1, 0]$. Can you see why?

LOOKING AHEAD

Here are four important exercises. If you cannot do any of them, you should try to prove that they are impossible!

1. Find a function $f(x)$ defined on the closed interval $[0, 1]$ with the property that $f(0) = 2$, $f(1) = 4$ but for all $x \in [0, 1]$, $f(x) \neq 3$.
2. Find a continuous function $f(x)$ defined on the closed interval $[0, 1]$ with the property that $f(0) = 2$, $f(1) = 4$ but for all $x \in [0, 1]$, $f(x) \neq 3$.
3. Find a function $f(x)$ defined on $[0, 1]$ that does not have a maximum. That is, there does not exist $t \in [0, 1]$ so that for all $s \in [0, 1]$, $f(t) \geq f(s)$.
4. Find a continuous function $f(x)$ defined on $[0, 1]$ that does not have a maximum. That is, there does not exist $t \in [0, 1]$ so that for all $s \in [0, 1]$, $f(t) \geq f(s)$.

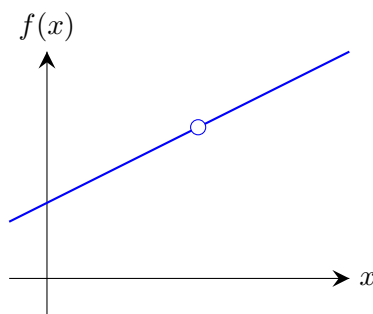
To approach each of the four problems, you may find it helpful to sketch out the graph of a function with the desired properties.

2.2 Types of Discontinuities

We say a function is **discontinuous** at $x = a$ if it is not continuous at $x = a$. Let's look at the different situations which lead to a discontinuity in a function.

Removable Discontinuities

Consider the function $f(x) = \frac{x^2 - 4}{x - 2}$. When $x \neq 2$, $f(x)$ is equivalent to the linear function $g(x) = x + 2$. We can infer from this that the limit of $f(x)$ as $x \rightarrow 2$ is equal to $g(2) = 4$. However, since $x = 2$ is not in the domain of f , $f(2)$ does not exist and therefore f is not continuous at $x = 2$.



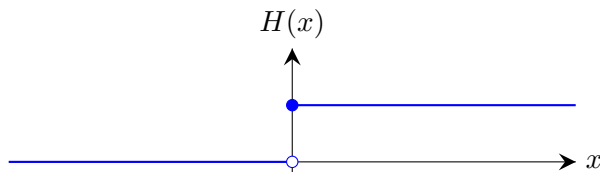
We call this type of discontinuity a **removable** discontinuity because the function can be made continuous at $x = 2$ by defining $f(2)$ to be equal to 4.

In general, suppose $\lim_{x \rightarrow a} f(x)$ exists but f is *not* continuous at $x = a$. Then we can redefine the function at $x = a$ to be $f(a) = \lim_{x \rightarrow a} f(x)$. Redefining the function at this one point then makes f continuous at $x = a$. If we can remove a discontinuity in this way, it's a removable discontinuity.

Jump Discontinuities

Recall the Heaviside step function

$$H(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}$$



Notice that the one-sided limits of $H(x)$ as $x \rightarrow 0$ exist and $H(0)$ exists. However, since the one-sided limits are not equal to one-another, there is an upward “jump” in the value of H at $x = 0$. For this reason, we call this a **jump discontinuity**.

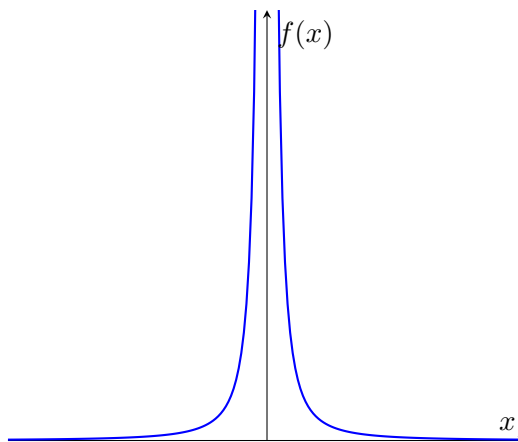
We can quantify the “size” of a jump discontinuity for a function f at $x = a$ by taking the difference in the one-sided limits.

$$\lim_{x \rightarrow a^+} f(x) - \lim_{x \rightarrow a^-} f(x)$$

Observe that for $H(x)$, $\lim_{x \rightarrow 0^+} H(x) = H(0)$. So, even though $H(x)$ is not continuous at $x = 0$, it is continuous from the right at $x = 0$.

Essential Discontinuities

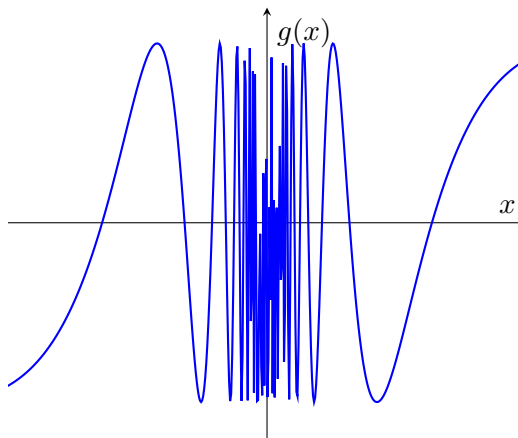
When one or both of the one-sided limits of a function f as $x \rightarrow a$ does not exist, we say the function has an **essential discontinuity**. For example, consider the function $f(x) = \frac{1}{x^2}$.



While we could write $\lim_{x \rightarrow 0} f(x) = \infty$, the limit does **not** actually exist. Similarly, neither of the one-sided limits exist. As such, the function has an essential discontinuity. Note, a discontinuity like this one where one or both one-sided limits are infinite can also be referred to as an **infinite discontinuity**.

In this case, $x = 0$ is *not* in the domain of the function. We could always define $f(0)$ to be some finite value. However, it doesn't matter what we choose as the value of $f(0)$, $f(x)$ would still have an infinite discontinuity at $x = 0$.

Another example of a function with an essential discontinuity is $g(x) = \sin\left(\frac{1}{x}\right)$.



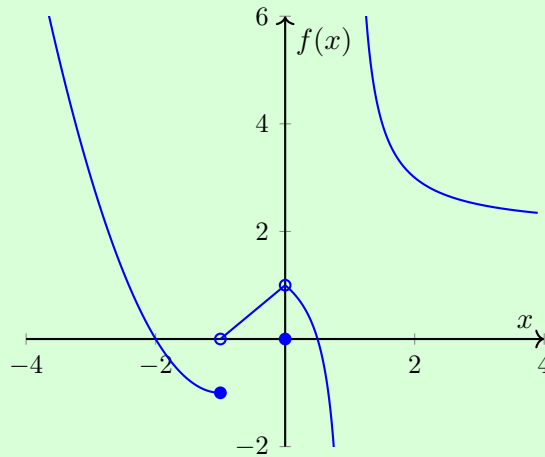
At $x = 0$ the limit of $g(x)$ does not exist because the frequency at which $g(x)$ oscillates between 1 and -1 grows without bound as $x \rightarrow 0$ (try using this observation to formally prove that the limit does not exist).

This type of essential discontinuity may also be referred to as an **oscillatory discontinuity**.

Example 3 Consider the piecewise-defined function

$$f(x) = \begin{cases} (x+1)^2 - 1 & \text{if } x \leq -1 \\ x+1 & \text{if } -1 < x < 0 \\ 0 & \text{if } x = 0 \\ \frac{1}{x-1} + 2 & \text{if } x > 0 \end{cases}$$

Here is a graph of the function:



The function is continuous at all points in \mathbb{R} except for $x = -1$, $x = 0$, and $x = 1$.

At $x = -1$, there is a jump discontinuity since $\lim_{x \rightarrow -1^-} f(x) = -1$ and $\lim_{x \rightarrow -1^+} f(x) = 0$.

At $x = 0$, there is a removable discontinuity since $\lim_{x \rightarrow 0} f(x)$ exists (and is equal to 1).

At $x = 1$, there is a vertical asymptote, and therefore an essential discontinuity.

EXERCISE

Suppose a function f has domain \mathbb{R} , is continuous everywhere except for at $x = a$, and $\lim_{x \rightarrow a^+} f(x)$ and $\lim_{x \rightarrow a^-} f(x)$ exist. What types of discontinuities could f have?

EXERCISE

Classify the discontinuities of the function $f(x) = \frac{x-1}{x^3-x}$.

2.3 Continuity Properties of Elementary Functions

As continuity will play an important role for much of what is to come, it will be useful to explore where our most common elementary functions are continuous.

Polynomial Functions

For any polynomial $P(x) = c_0 + c_1x + c_2x^2 + \cdots + c_nx^n$, we recall that for all $a \in \mathbb{R}$

$$\lim_{x \rightarrow a} P(x) = P(a)$$

Therefore, every polynomial is continuous at each point in \mathbb{R} .

Trigonometric Functions

Intuitively, we expect that $\sin(x)$ and $\cos(x)$ are continuous at each point in \mathbb{R} . In order to show this explicitly, we will rely on the following facts. We will not prove these facts in the notes, but you certainly have the tools to do them yourself!

Fact 1 We have $\lim_{x \rightarrow 0} \sin(x) = 0$ and $\lim_{x \rightarrow 0} \cos(x) = 1$.

Since $\sin(0) = 0$ and $\cos(0) = 1$, this fact shows that the sine and cosine functions are continuous at $x = 0$. Let's now show that the sine function is continuous on all of \mathbb{R} .

For any $a \in \mathbb{R}$, we have

$$\begin{aligned} \lim_{x \rightarrow a} \sin(x) &= \lim_{y \rightarrow 0} \sin(y + a) \\ &= \lim_{y \rightarrow 0} (\sin(y) \cos(a) + \sin(a) \cos(y)) \\ &= \cos(a) \lim_{y \rightarrow 0} \sin(y) + \sin(a) \lim_{y \rightarrow 0} \cos(y) \\ &= \cos(a) \sin(0) + \sin(a) \cos(0) \\ &= \sin(a) \end{aligned}$$

In the calculation above, we made the substitution $y = x - a$ and then proceeded to use the angle summation formula for sine.

EXERCISE

Prove that $\cos(x)$ is continuous at each point in \mathbb{R} .

EXERCISE

The tangent function, $\tan(x)$, is discontinuous at $x = \frac{\pi}{2} + n\pi$ for all $n \in \mathbb{Z}$. What types of discontinuities are these?

Exponential and Logarithmic Functions

Functions of the form $f(x) = b^x$ for positive real numbers b are continuous at every point in \mathbb{R} .

Logarithmic functions $g(x) = \log_b(x)$ where the base b is a positive real number are continuous on their domains which includes all $x > 0$.

EXERCISE

What is the domain of $f(x) = \ln(x^2)$? Classify any discontinuities f has over the real numbers.

Many other familiar functions are continuous on their domains. This includes power functions with non-integer exponents and rational functions. One example of this is the square root function $f(x) = \sqrt{x}$, which is continuous on its domain $[0, \infty)$. What this means is that $f(x)$ is continuous at all real numbers $x > 0$ and continuous from the right at $x = 0$.

Using the fact that the square root function is continuous on its domain, try the following exercise.

EXERCISE

Evaluate the following limits:

$$1. \lim_{x \rightarrow 0} \frac{\sqrt{x+9} - 3}{x}$$

$$2. \lim_{x \rightarrow 1} \arctan \frac{1-x}{2(1-\sqrt{x})}$$

2.4 Rules for Continuous Functions

We have seen that just about every familiar function that's easy to write down is continuous on its domain. Of course, we can take these basic functions and create lots of other functions by composing, adding, multiplying, and taking quotients.

It is often the case that we can determine continuity properties of a function by looking at the continuity properties of related and/or component functions.

Theorem 2 (Continuity of Sums, Differences, and Products)

Let f and g be continuous functions at $x = a$, then each of the following is continuous at $x = a$

(i) $f + g$

(ii) $f - g$

(iii) fg

Before embarking on the proof (well, the proof of (i) anyway), let's first recall precisely what we mean by the function $f + g$.

We are given functions f and g . In order for me to tell you what the function $f + g$ is, I have to tell you what it does to every input. We *define* the function $f + g$ by declaring that

$$(f + g)(x) = f(x) + g(x)$$

for all x such that x is in the domain of both f and g .

Proof of (i): We will prove (i) here and leave the other two as an exercise.

Suppose f and g are continuous at $x = a$. That means

$$\lim_{x \rightarrow a} f(x) = f(a) \quad \text{and} \quad \lim_{x \rightarrow a} g(x) = g(a).$$

Then

$$\lim_{x \rightarrow a} (f + g)(x) = \lim_{x \rightarrow a} (f(x) + g(x)) = \lim_{x \rightarrow a} f(x) + \lim_{x \rightarrow a} g(x) = f(a) + g(a) = (f + g)(a)$$

completing the proof. □

The key step in the proof was the second equals sign, where we used one of the limit properties from Section 1.3.

If there is a rule about limits, then there may be a way to use that rule to prove a corresponding fact about continuity. With that in mind, have a go at proving parts (ii) and (iii) from Theorem 2 and Theorem 3

Theorem 3 (Continuity of Quotients)

Let f and g be continuous functions at $x = a$ with $g(a) \neq 0$, then $\frac{f}{g}$ is continuous at $x = a$.

EXERCISE

Determine for what values of $x \in \mathbb{R}$ the function $f(x) = \frac{x + 3}{x^2 + 5x + 6}$ is continuous.

Our next goal is to prove that continuity plays nicely with composition of functions. We would like to say that the composition of continuous functions is a continuous function.

In order to do this, we first need an intermediary result.

Theorem 4 (Limits of Compositions)

If f is continuous at $x = b$ and $\lim_{x \rightarrow a} g(x) = b$, then $\lim_{x \rightarrow a} f(g(x)) = f(b)$.

We can think of this theorem as allowing us to take the limit inside continuous functions. Another way of writing the conclusion of the theorem is by saying

$$\lim_{x \rightarrow a} f(g(x)) = f\left(\lim_{x \rightarrow a} g(x)\right).$$

Proof: Let $\epsilon > 0$. The end goal is to find $\delta > 0$ so that $0 < |x - a| < \delta$ implies $|f(g(x)) - f(b)| < \epsilon$.

We are given the assumption that $\lim_{x \rightarrow b} f(x) = f(b)$ (that is, f is continuous at $x = b$). Therefore, there exists $\tilde{\delta} > 0$ so that if $0 < |x - b| < \tilde{\delta}$, then $|f(x) - f(b)| < \epsilon$.

The other assumption provided is that $\lim_{x \rightarrow a} g(x) = b$. Therefore, there exists $\delta > 0$ so that $0 < |x - a| < \delta$ implies $|g(x) - b| < \tilde{\delta}$ (here $\tilde{\delta}$ is playing the role ϵ would usually play).

Combining these two implications gives us that if $0 < |x - a| < \delta$, then $|g(x) - b| < \tilde{\delta}$, which in turn implies $|f(g(x)) - f(b)| < \epsilon$. We may now conclude that $\lim_{x \rightarrow a} f(g(x)) = f(b)$, completing the proof. \square

Theorem 5 (Continuity of Compositions)

Let f and g be functions such that g is continuous at $x = a$ and f is continuous at $g(a)$, then $h = f \circ g$ is continuous at $x = a$.

Proof: Since g is continuous at $x = a$, $\lim_{x \rightarrow a} g(x) = g(a)$. Applying Theorem 4 yields

$$\lim_{x \rightarrow a} f(g(x)) = f(g(a))$$

completing the proof. \square

Example 4 Determine $\lim_{x \rightarrow 0} \cos(x + \sin(x))$.

Solution: Let $f(x) = \cos(x)$ and $g(x) = x + \sin(x)$. Since f and g are both continuous on \mathbb{R} , then $f(g(x))$ is continuous on \mathbb{R} . Therefore,

$$\lim_{x \rightarrow 0} \cos(x + \sin(x)) = \cos(0 + \sin(0)) = 1$$

Another way we can create functions out of other functions is by taking inverses.

Recall that the inverse of f is g if $f \circ g(s) = s$ and $g \circ f(t) = t$ for all s in the domain of g and all t in the domain of f .

Not every function has an inverse, but some do! A nice example to have in mind is $f(x) = e^x$ and $g(x) = \ln(x)$. These are inverses of each other. Note that the inverse of e^x is *not* e^{-x} .

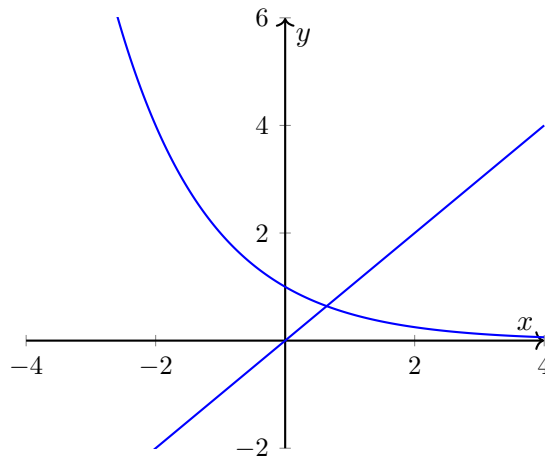
Theorem 6 (Continuity of Inverses)

Let f be a function with inverse g . If f is continuous at $x = a$ and $f(a) = b$, then g is continuous at b .

We shall not give the proof of this here, but you should try to prove it for yourself. The result is believable when you recall that if g and f are inverses of each other, that the graphs of g and f are reflections of one another across the line $y = x$. See if you can use this observation to come up with a proof of Theorem 6.

2.5 Intermediate Value Theorem

Consider the equation $2^{-x} = x$. Does this equation have any real-valued solutions?



It is clear from the plot of $y = 2^{-x}$ and $y = x$ that these two curves intersect somewhere on the interval $(0, 1)$. Therefore, the equation $2^{-x} = x$ does have a solution. How might we argue this though if we could not readily generate the graph above.

Here is another way to argue that the equation $2^{-x} = x$ has a solution on the interval $(0, 1)$. Let $f(x) = 2^{-x} - x$. Observe that $f(0) = 1$ and $f(1) = -\frac{1}{2}$. This means the function f changes from positive to negative on the interval $[0, 1]$. Since the function f is the sum of two continuous functions, it is also a continuous function. Therefore, in order to change from positive to negative, it must be equal to exactly zero at least once on the interval $(0, 1)$. Let x_0 be a value of $x \in (0, 1)$ such that $f(x_0) = 0$. It follows that $2^{-x_0} = x_0$. That is, the given equation has a solution on the interval $(0, 1)$.

EXERCISE

Let $f(x) = \frac{1}{x-2}$. Compute $f(1)$ and $f(3)$. Can we conclude that there exists a value of $x \in (1, 3)$ such that $f(x) = 0$? (Hint: Sketch the graph of $y = f(x)$.)

In the previous exercise, you should discover that it is critical that a function is continuous on a particular interval to infer that it must take values between the values which it takes at the end points of that interval. However, when continuity can be assured, this idea is a very powerful one and is called the **Intermediate Value Theorem**.

Theorem 7 (Intermediate Value Theorem (IVT))

If f is continuous on the interval $[a, b]$ and N is a number strictly between $f(a)$ and $f(b)$, then there exists $c \in (a, b)$ such that $f(c) = N$.

The proof of the Intermediate Value Theorem is beyond the scope of this course.

Example 5 Consider the function

$$f(x) = x^5 - \frac{7}{2}x^4 + \frac{1}{2}x^3 + \frac{25}{4}x^2 - \frac{51}{16}x - \frac{27}{32}.$$

This is one ugly looking polynomial, but at least it's a polynomial, so it's continuous everywhere! Let's see if we can learn anything about its roots.

By evaluating the function we see

$$f(-2) < 0$$

$$f(-1) > 0$$

$$f(0) < 0$$

$$f(1) > 0$$

$$f(2) < 0$$

$$f(3) > 0.$$

So, by the Intermediate Value Theorem, we know there is at least one root in each of the following intervals:

$$(-2, -1), \quad (-1, 0), \quad (0, 1), \quad (1, 2), \quad \text{and} \quad (2, 3).$$

Since a degree-5 polynomial has at most 5 real roots, we can conclude that this polynomial has exactly 5 distinct real roots.

We don't know what they are, but we know they are there!

EXERCISE

Prove that there exists $x \in (1, 3)$ such that $h(x) = \sqrt{x+1} + \sqrt{3x}$ takes the value 4.

The Intermediate Value Theorem lends itself to some neat physical applications under the assumption that things that seem continuous in the real world are indeed mathematically continuous. The next two examples are instances of this.

Example 6

A square table with equal length legs at each corner and perpendicular to the tabletop is placed on an uneven surface. Initially, one of the table legs does not touch the uneven ground making the table wobbly. Argue that by rotating the table, you can find an orientation such that all four legs touch the ground simultaneously thereby removing the wobble.

Solution: Label the legs counterclockwise from above as 1,2,3,4 with legs 1,2,3 in contact with the ground while 4 is the one that is not. This means, the table wobbles with legs 1 and 3 always on the ground (i.e., if you push down on leg 4, leg 2 comes up off the ground). Importantly, in this configuration, for legs 2 and 4 to simultaneously touch the ground, we'd need to push one or both of legs 1 and 3 into the ground.

Now, rotate the table counterclockwise while keeping legs 1,2, and 3 in contact with the ground. Note, this might mean that leg 4 enters the ground which may not be physically sensible but we can allow it mathematically as follows. Denote the height of the fourth leg

by $h(\theta)$ where θ is the angle of rotation. Note, $h(0) > 0$ and a negative h means the fourth table leg extends into the ground.

After rotating 90 degrees, we have effectively swapped legs 1 and 3 with legs 2 and 4 relative to the initial configuration. If legs 1 and 3 are at ground level, then based on the argument above, one of both of legs 2 and 4 must be in the ground. Since leg 2 is forced to be at ground level, then it must be leg 4 which is in the ground now. Mathematically, $h(\frac{\pi}{2}) < 0$. This means the height of the fourth leg which should be varying continuously as we rotate the table has gone from positive to negative. By the Intermediate Value Theorem, there must have been some value of $\theta \in (0, \frac{\pi}{2})$ such that $h(\theta) = 0$. In other words, some orientation such that all four legs were simultaneously at ground level and the table was not wobbly.

There are of course various things implicitly assumed in order to make this argument valid. What are those assumptions?

Verify that the argument is reasonable by taking an actual four-legged piece of furniture on some uneven ground and rotating it! If your family asks what you're doing, you can respond with "calculus".

Example 7

Prove that there exist a pair of diametrically opposite points on the equator with the exact same temperature.

Solution: Let's assume that temperature on the surface of the Earth varies continuously with respect to position. What this means physically is that if you start somewhere on Earth and the temperature is 10° , and you walk to somewhere where the temperature is 20° , then at some point you must have been somewhere where the temperature was 15° , for example.

Label each point on the equator by its longitude θ (in radians). So, for example, the points $\theta = 0$ and $\theta = \pi$ are diametrically opposite from each other, as are the points $\theta = \frac{\pi}{2}$ and $\theta = \frac{3\pi}{2}$.

Let $T(\theta)$ be the temperature at the point on the equator with longitude θ . To make the mathematics a little neater, we will allow θ to take on any real number, and insist that θ and $\theta + 2\pi k$ for any integer k represent the same point on the equator.

It follows that $T(\theta) = T(\theta + 2\pi)$ for all $\theta \in \mathbb{R}$.

Now, let $H(\theta) = T(\theta) - T(\theta + \pi)$. That is, $H(\theta)$ is the difference in temperature between the point with longitude θ , and the point diametrically opposite.

Since we are assuming that T is continuous, it follows that H is continuous. We have

$$H(0) = T(0) - T(\pi) = T(2\pi) - T(\pi) = -H(\pi).$$

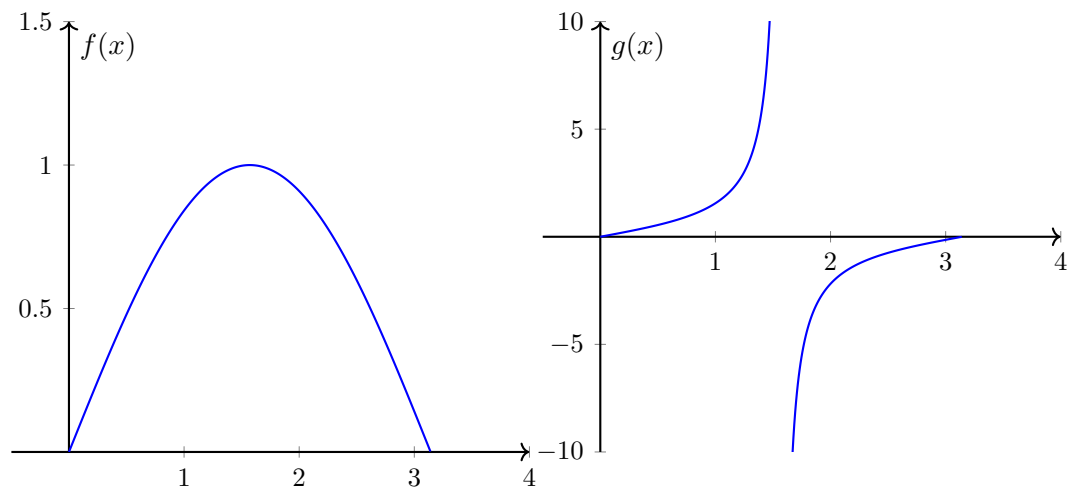
If $H(0) = 0$, we're done, since $T(0) = T(\pi)$ so there are two diametrically opposite points with equal temperature.

If $H(0) \neq 0$, then 0 is strictly between $H(0)$ and $H(\pi) = -H(0)$. Therefore, by the Intermediate Value Theorem, there is some $\mu \in (0, \pi)$ so that $H(\mu) = 0$. That value μ corresponds to a point on the equator whose diametrically opposite point (given by $\mu + \pi$) has the same temperature.

What other examples can you think of?

2.6 Extreme Value Theorem

Consider the functions $f(x) = \sin(x)$ and $g(x) = \tan(x)$ on the domain $D = [0, \pi]$.



What is the maximum value that $f(x)$ takes on this interval? What is the minimum value? What about for $g(x)$?

To answer these questions, we should first clarify what we mean by maximum and minimum in this context. In particular, we are interested in identifying the **absolute** maximum and minimum points. (This is in contrast to **local** maximum and minimum points which we will discuss later in the course.)

Definition 2.6.1 absolute max/min

Let f be a function with domain D and let $c \in D$.

- If $f(c) \geq f(x)$ for all $x \in D$, then we say $f(c)$ is the absolute maximum of f on D .
- If $f(c) \leq f(x)$ for all $x \in D$, then we say $f(c)$ is the absolute minimum of f on D .

According to this definition, if we can identify a value which a function takes on a given domain as being greater than or equal to every other value the function takes in that domain, then we call that value the absolute maximum. Likewise for the absolute minimum. Note, we call the maximum and minimum values of a function the **extreme values**. Sometimes we also refer to absolute maximum and minimum values as **global** maximum and minimum values.

Going back to $f(x) = \sin(x)$ with $D = [0, \pi]$, we see that f takes its absolute maximum value of 1 when $x = \frac{\pi}{2}$. Every other value f takes on D is strictly less than 1.

What about an absolute minimum? Of course, f takes its minimum value on D of 0 at $x = 0$ and $x = \pi$. The definition does not preclude having an extreme value at more than one place, so we would say the absolute minimum value of f is 0.

Let's look now at $g(x) = \tan(x)$ on the same domain. Does g have an absolute maximum on this domain? g grows arbitrary large as $x \rightarrow \frac{\pi}{2}$ from the left, it is not possible to identify a value of g which is greater than or equal to every other value g takes on the same domain.

Therefore, g does *not* have an absolute maximum on D . For similar reasons, g does not have an absolute minimum value on D either.

It is clear that the infinite discontinuity in g at $x = \frac{\pi}{2}$ is the reason that g fails to have an absolute maximum or minimum in this case. In contrast, $f(x) = \sin(x)$ was continuous on D . In fact, a function being continuous on an interval guarantees the existence of extreme values.

Theorem 8 (Extreme Value Theorem)

Let f be a function which is continuous on the closed interval $[a, b]$, then f has an absolute maximum value and an absolute minimum value in $[a, b]$.

The proof of the extreme value theorem is beyond the scope of this course.

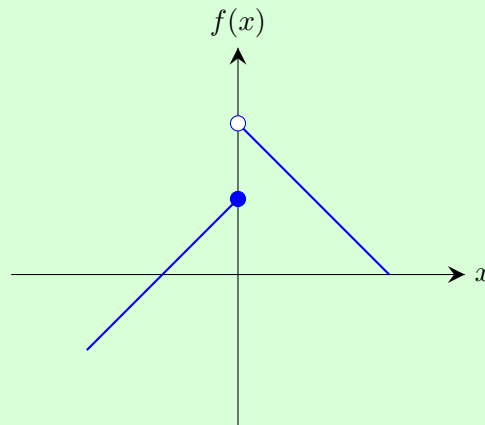
Example 8

Let $f(x)$ be the following piecewise-defined function with domain $[-2, 2]$.

$$f(x) = \begin{cases} x + 1, & -2 \leq x \leq 0 \\ -x + 2 & 0 < x \leq 2 \end{cases}$$

Determine if f has an absolute maximum or minimum on its domain.

Solution: First, to help guide us, here is the graph of the function:



Let's start with the easy one - the absolute minimum. At $x = -2$, the function takes the value -1 and every other value the function takes on its domain is greater than this. Without appealing to the graph, we could argue this as follows.

If $-2 \leq x \leq 0$, then $-1 \leq x + 1 \leq 1$. Thus, the part of the graph where $f(x) = x + 1$ has a minimum value of -1 . Similarly, if $0 < x \leq 2$, then $-2 \leq -x < 0$ and so $0 \leq -x + 2 < 2$. So, for the part of the graph where $f(x) = -x + 2$, the minimum value attained by f is 0 . Therefore, the minimum value attained by f on the interval $[-2, 2]$ is -1 .

In contrast, an absolute maximum does not exist. This is because of the nature of the jump discontinuity of f at $x = 0$. The function f tends to 2 as $x \rightarrow 0$ from the right but never actually takes a value of 2 . Therefore, the absolute maximum, if it existed, would have to be some number strictly less than two. But, for any such number we write down, we can always find a value which is both closer to 2 and a value which f takes for some value of $x > 0$.

EXERCISE

Let $f(x) = \frac{1}{x^2-1}$ on the open interval $(-1, 1)$. Determine if f has an absolute maximum or minimum on its domain. Are your findings consistent with the Extreme Value Theorem?

Chapter 3

Differentiation

3.1 Instantaneous Velocity Revisited

Let's return now to the problem of computing instantaneous velocity. Recall, given a function $d(t)$ describing the displacement, d , of an object as a function of time, t , the *average* velocity over the time interval from $t = t_1$ to $t = t_2$ is given by

$$v_{\text{avg}} = \frac{d(t_2) - d(t_1)}{t_2 - t_1}$$

This is just the slope (i.e., rise over run) of a secant line connecting the points $(t_1, d(t_1))$ and $(t_2, d(t_2))$ on the graph of d vs. t . We can use this expression to approximate the *instantaneous* velocity at some fixed time, say $t = a$, by setting $t_1 = a$. Then,

$$v_{\text{avg}} = \frac{d(t_2) - d(a)}{t_2 - a}$$

gives an approximation of the instantaneous velocity at $t = a$ which tends to improve as t_2 gets closer to a . However, we cannot just set $t_2 = a$ since v_{avg} is discontinuous there. Fortunately, this is a *removable* discontinuity, so we can take the limit of $v(t_2)$ as $t_2 \rightarrow a$ and assign the limiting value to be the instantaneous velocity of the object at $t = a$.

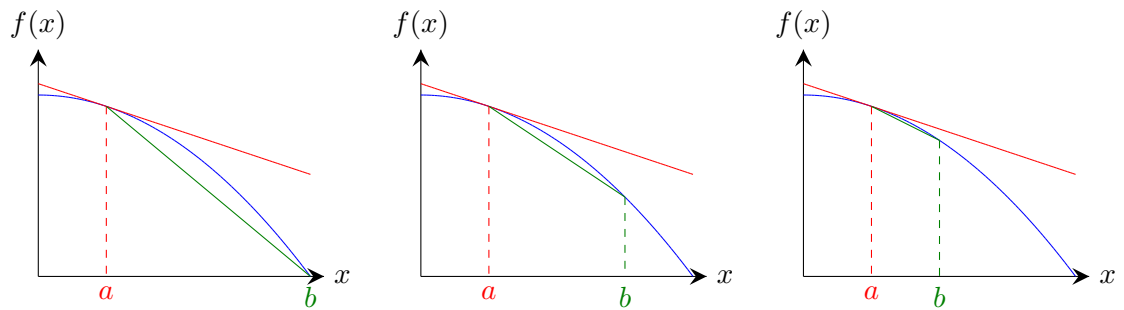
$$v_{\text{inst}}(a) = \lim_{t_2 \rightarrow a} \frac{d(t_2) - d(a)}{t_2 - a}$$

EXERCISE

The vertical position, z , of a particle released at rest is described as a function of time, t , by $z(t) = -\frac{1}{2}gt^2$ where g is the magnitude of the acceleration due to gravity. Find the instantaneous velocity at $t = 5$.

3.2 Definition of the Derivative

The approach above for computing instantaneous velocity can be applied in a more general context to compute the instantaneous rate of change of an arbitrary function $f(x)$ at $x = a$. Essentially, we write down an expression for the average rate of change over some interval (i.e., the slope of a secant line) and then take the limit as the interval goes to zero (i.e., the secant line becomes a tangent line).



This means the instantaneous rate of change of an arbitrary function $f(x)$ at $x = a$ is given by the limit

$$\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$$

The instantaneous rate of change of a function is such an important feature that we call this expression the *derivative* of $f(x)$ at $x = a$.

Definition 3.2.1
derivative at a point

The derivative of a function $f(x)$ at $x = a$ is given by

$$\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$$

provided this limit exists.

REMARKS

- Implicit in this definition is that the quantity $f(a)$ exists. If a is not in the domain of f , then the derivative at $x = a$ does not exist.
- The derivative of $f(x)$ at $x = a$ can be denoted $f'(a)$.
- The quantity $\frac{f(x) - f(a)}{x - a}$ is commonly referred to as the **difference quotient**.
- It is often convenient to make the substitution $x = h + a$ which gives the equivalent definition

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a + h) - f(a)}{h}$$

Example 1

Determine the derivative of $f(x) = x^3$ at $x = 2$.

Solution: By the definition of the derivative, we have

$$f'(2) = \lim_{x \rightarrow 2} \frac{x^3 - 2^3}{x - 2}$$

Observe that the numerator can be factored as $x^3 - 2^3 = (x - 2)(x^2 + 2x + 4)$. This allows us to rewrite

$$\begin{aligned} f'(2) &= \lim_{x \rightarrow 2} \frac{(x - 2)(x^2 + 2x + 4)}{x - 2} \\ &= \lim_{x \rightarrow 2} (x^2 + 2x + 4) \\ &= 2^2 + 2(2) + 4 \\ &= 12 \end{aligned}$$

EXERCISE

Determine the derivative of $g(x) = x^4$ at $x = 1$.

Hint: $x^4 - a^4 = (x - a)(x^3 + ax^2 + a^2x + a^3)$.

EXERCISE

Show that the derivative of $h(x) = ax^2 + bx + c$ exists for an arbitrary $x = x_0$. In other words, show that the derivative of every quadratic function exists everywhere.

3.3 Differentiability

When the derivative of a function $f(x)$ exists at $x = a$, we say that f is **differentiable** at $x = a$.

There is an important connection between continuity and differentiability. Suppose that $f(x)$ is differentiable at $x = a$, then we know that the following limit exists.

$$\lim_{x \rightarrow a} \frac{f(x) - f(a)}{x - a}$$

As $x \rightarrow a$, the denominator of the difference quotient clearly tends to zero. That is,

$$\lim_{x \rightarrow a} x - a = 0$$

Therefore, the only way that this limit can exist is if the numerator of the difference quotient also tends to zero.

$$\lim_{x \rightarrow a} f(x) - f(a) = 0$$

Equivalently, we have

$$\lim_{x \rightarrow a} f(x) = f(a)$$

This is precisely the definition of continuity for $f(x)$ at $x = a$. We have just proved the following theorem.

Theorem 1 (Differentiability Implies Continuity)

If f is differentiable at $x = a$, then f is continuous at $x = a$.

REMARKS

- It is important to note that the previous theorem can only be applied in one direction. In particular, continuity does *not* imply differentiability.
- The contrapositive of the theorem can also be useful. That is, if f is *discontinuous* at $x = a$, then f is *not differentiable* at $x = a$.

The derivative of a function $f(x)$ at a point $x = a$ gives the slope of a tangent line to the graph of that function at the point $(a, f(a))$. Therefore, to determine if a function is differentiable at a point, it can be helpful to visualize its graph and decide if a line can be drawn tangent to the graph at that point unambiguously with a finite slope.

Example 2

Let $f(x) = |x^2 - 4|$. For what values of x is $f(x)$ differentiable?

Solution:

We can rewrite $f(x)$ as a piecewise function:

$$f(x) = \begin{cases} x^2 - 4, & x < -2 \text{ or } x > 2 \\ 4 - x^2, & -2 \leq x \leq 2 \end{cases}$$

Observe that when $x \neq \pm 2$, $f(x)$ is a quadratic function and you should have established in a previous exercise that quadratic functions are differentiable everywhere. What about at $x = 2$ and $x = -2$?

Consider first $x = 2$. For $f(x)$ to be differentiable at $x = 2$, the limit of the difference quotient must exist. We need to check the left-side and right-side limits:

The right-side derivative at $x = 2$ is:

$$\lim_{x \rightarrow 2^+} \frac{f(x) - f(2)}{x - 2} = \lim_{x \rightarrow 2^+} \frac{(x^2 - 4) - 0}{x - 2} = \lim_{x \rightarrow 2^+} x + 2 = 4$$

Note that we were able to replace $f(x) = |x^2 - 4|$ with $x^2 - 4$ in this limit since $x > 2$ when we are taking a limit as x approaches 2 from the right.

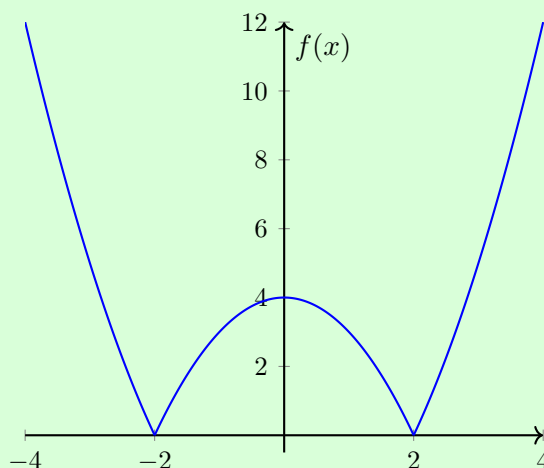
The left-side derivative at $x = 2$ is:

$$\lim_{x \rightarrow 2^-} \frac{f(x) - f(2)}{x - 2} = \lim_{x \rightarrow 2^-} \frac{(4 - x^2) - 0}{x - 2} = - \lim_{x \rightarrow 2^-} x + 2 = -4$$

Since the left-side and right-side limits are not equal, the limit $\lim_{x \rightarrow 2} \frac{f(x) - f(2)}{x - 2}$ does not exist at $x = 2$. Therefore, $f(x)$ is not differentiable at $x = 2$.

A similar calculation shows that $f(x)$ is also not differentiable at $x = -2$.

The graph of $f(x) = |x^2 - 4|$ reveals that this result is not surprising.



Observe that the curve $y = f(x)$ has cusps at $x = 2$ and $x = -2$. At these points, the slope of a tangent to the curve changes abruptly. It is therefore impossible to define a unique tangent line to the graph at those points.

Therefore, $f(x)$ is differentiable for all $x \in \mathbb{R}$ except for $x = 2$ and $x = -2$.

EXERCISE

Let $g(x) = \tan(x)$. For what values of x is $g(x)$ differentiable?

EXERCISE

Let

$$f(x) = \begin{cases} x, & x < 0 \\ x + 1, & x \geq 0 \end{cases}$$

Sketch $y = f(x)$. Determine using the definition of the derivative as needed for what values of x is $f(x)$ differentiable?

3.4 The Derivative Function

The derivative of a function gives us useful information about the function. Derivatives can be computed at specific points as needed or, recognizing that functions are often differentiable over non-trivial intervals, we can alternatively define a new function - the **derivative function** - whose output is the desired derivative. When we need to know the derivative of a function at multiple points, this is almost always more efficient but the function itself is useful too as it describes the rate of change of the original function.

Definition 3.4.1
derivative function

Suppose f is differentiable on the interval I , then the derivative function of f , denoted f' , is given by

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

and we say $f'(x)$ is the derivative of f with respect to x .

Observe that the derivative function is defined using the same (alternative form) difference quotient that is used to define the derivative at a point $x = a$. The key modification is that we do not specify the point at which the limit is being taken but instead leave it arbitrary as x .

REMARKS

- The derivative function $f'(x)$ can be equivalently denoted in *Leibniz notation* as $\frac{df}{dx}$.
- It is often convenient to consider the derivative as an “operation” acting on a function. This may be denoted by $(f(x))'$ or, in Leibniz notation, as $\frac{d}{dx}f$.
- The derivative of f evaluated at the point $x = a$ can also be denoted by $\left. \frac{df}{dx} \right|_{x=a}$.

The derivative operation can be repeatedly applied to a function to produce *higher-order derivatives*. The second derivative of $f(x)$ is denoted by

$$f''(x) \quad \text{or} \quad \frac{d^2f}{dx^2}$$

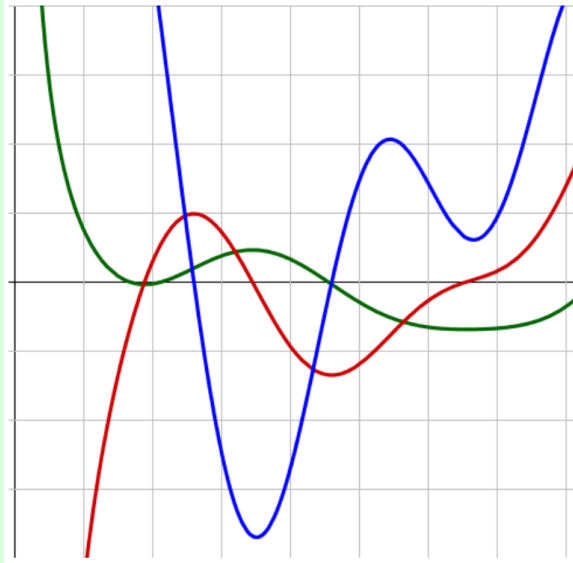
This notation can be generalized for the n -th derivative as

$$f^{(n)}(x) \quad \text{or} \quad \frac{d^n f}{dx^n}$$

It holds true that $f^{(n+1)}(x)$ describes the rate of change of $f^{(n)}$. For example, suppose $x(t)$ gives the position, x , of an object as a function of time, t , then $x'(t)$ gives the velocity, $x''(t)$ gives the acceleration, and so on.

Example 3

The following plot shows the graphs of $f(x)$, $f'(x)$, and $f''(x)$ for some function $f(x)$. Determine which graph corresponds to the plot of each of $f(x)$, $f'(x)$, and $f''(x)$.



Solution: The function $f'(x)$ tells us the rate of change of $f(x)$ at a given value of x . This means that if $f(x)$ is increasing (decreasing), then its derivative should be positive (negative) valued. Moreover, when $f(x)$ is changing from increase to decreasing or decreasing to increasing, its derivative is instantaneously zero.

By inspection, the red curve is positive (negative) whenever the green curve is increasing (decreasing). This means that the red curve represents the derivative of the function represented by the green curve.

Similarly, the blue curve is positive (negative) whenever the red curve is increasing (decreasing). This means that the blue curve represents the derivative of the function represented by the red curve.

Combining these two observations, it must be the case that the green curve is the graph of $f(x)$, the red curve is the graph of $f'(x)$, and the blue curve is the graph of $f''(x)$.

EXERCISE

Let $f(x) = x^n$ for some positive integer n . Prove that $f^{(n+1)}(x) = 0$.

(You may find the binomial expansion formula useful $(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k$.)

3.5 Differentiation Rules

Suppose we are asked to determine the derivative of $f(x) = x^{100}$. Applying the definition of the derivative yields:

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$\begin{aligned}
&= \lim_{h \rightarrow 0} \frac{(x+h)^{100} - x^{100}}{h} \\
&= \lim_{h \rightarrow 0} \frac{(x^{100} + 100x^{99}h + \cdots + 100xh^{99} + h^{100}) - x^{100}}{h} \\
&= \lim_{h \rightarrow 0} \frac{(x^{100} - x^{100}) + h(100x^{99} + \cdots + 100xh^{98} + h^{99})}{h} \\
&= \lim_{h \rightarrow 0} (100x^{99} + \cdots + 100xh^{98} + h^{99}) \\
&= 100x^{99}
\end{aligned}$$

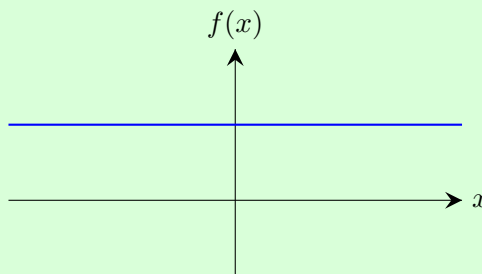
That's great, but it was a lot of work to differentiate just a simple monomial function. In this section, we present some useful rules to make the process of computing derivatives more efficient.

Example 4 Let $f(x) = c$ where c is an arbitrary real constant. Determine $f'(x)$.

Solution: Let us find the derivative function using its definition.

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} \frac{c - c}{h} = 0$$

Therefore, the derivative of any constant function is zero everywhere. Geometrically, this makes sense since the slope of the graph of a constant function is zero.



The following theorem gives a handful of useful differentiation properties.

Theorem 2 (Arithmetic Differentiation Rules)

Let f and g be functions and let c be a real number. Suppose that at $a \in \mathbb{R}$, $f'(a)$ and $g'(a)$ exist. Then the following hold:

- (i) $(c)' = 0$
- (ii) $(cf(a))' = cf'(a)$
- (iii) $(f(a) + g(a))' = f'(a) + g'(a)$ (Sum Rule)
- (iv) $(f(a)g(a))' = f'(a)g(a) + f(a)g'(a)$ (Product Rule)
- (v) $\left(\frac{f(a)}{g(a)}\right)' = \frac{f'(a)g(a) - f(a)g'(a)}{(g(a))^2}$ provided $g(a) \neq 0$ (Quotient Rule)

We prove the sum rule and leave the rest as exercises.

Proof of (iii): Suppose $f'(a)$ and $g'(a)$ exist, meaning the limits

$$\lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h} \quad \text{and} \quad \lim_{h \rightarrow 0} \frac{g(a+h) - g(a)}{h}$$

both exist. Then

$$\begin{aligned} (f(a) + g(a))' &= \lim_{h \rightarrow 0} \frac{(f(a+h) + g(a+h)) - (f(a) + g(a))}{h} \\ &= \lim_{h \rightarrow 0} \frac{(f(a+h) - f(a)) + (g(a+h) - g(a))}{h} \\ &= \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h} + \lim_{h \rightarrow 0} \frac{g(a+h) - g(a)}{h} \\ &= f'(a) + g'(a) \end{aligned}$$

□

EXERCISE

Let $f(x) = mx + b$ where m and b are arbitrary real constants. Without doing any calculations, what do you expect $f'(x)$ to be? Find $f'(x)$ using the definition of the derivative function and properties above as needed to check your hypothesis.

Calculations like the one in the previous exercise can be repeated for higher degree polynomials and made simpler by knowing the following rule for the derivative of a power function.

Theorem 3 (Power Rule)

If $f(x) = x^n$ where n is a non-zero real number, then

$$f'(x) = n x^{n-1}$$

This rule will be proved later in the course once we have implicit differentiation as a tool to wield.

REMARK

Observe that this statement asserts that the power rule works even when the exponent is not an integer.

EXERCISE

Prove the power rule in the case that n is a positive integer. It may be helpful to use the product rule and argue by induction.

Similarly to when we did continuity, the arithmetic differentiation rules allow us to differentiate polynomials and rational functions (that is, a quotient of a polynomial by a polynomial).

EXERCISE

Let

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

Find the derivative function $f'(x)$ and find all the x values satisfying $f'(x) = 0$.

3.6 Differentiating Elementary Functions

The derivative is defined formally through a limiting process. Let's now look at what this process yields for some common elementary functions. In doing so, we will see that, in practice, the limit-based definition of the derivative can often be skipped over. So, moving forward, we will usually only need to appeal to the formal (limit-based) definition of the derivative in special cases.

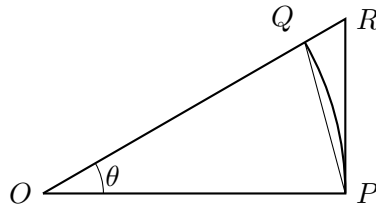
3.6.1 Trigonometric Functions

To determine the derivatives of $\sin(x)$ and $\cos(x)$, we will need the Fundamental Trigonometric Limit

$$\lim_{\theta \rightarrow 0} \frac{\sin(\theta)}{\theta} = 1$$

The proof of this limit involves a fun exercise in geometry.

Proof: Consider a unit circle centred at the origin O . Place points P and Q on the unit circle with P at $(1, 0)$ and Q at $(\cos(\theta), \sin(\theta))$ with $\theta > 0$. Observe that the angle between lines OP and OQ is θ . Next, extend line OQ to meet the vertical line $x = 1$ at $(1, \tan(\theta))$; we call this point R .



We now define three areas:

- The area of $\triangle OPQ$ is $A_{\triangle OPQ} = \frac{1}{2} \sin(\theta)$ since it has base length 1 and height $\sin(\theta)$.
- The area of the sector of the circle subtended by the angle θ is $A_{\theta} = \frac{1}{2}\theta$.
- The area of $\triangle OPR$ is $A_{\triangle OPR} = \frac{1}{2} \tan(\theta)$ since it has base length 1 and height $\tan(\theta)$.

Observe that these areas satisfy $A_{\triangle OPQ} \leq A_\theta \leq A_{\triangle OPR}$. This gives us the inequality

$$\sin(\theta) \leq \theta \leq \tan(\theta)$$

Dividing by $\sin(\theta)$ and taking the reciprocal gives

$$\cos(\theta) \leq \frac{\sin(\theta)}{\theta} \leq 1$$

Taking the right-side limit as $\theta \rightarrow 0^+$ of each term in the inequality and applying the Squeeze theorem gives

$$\lim_{\theta \rightarrow 0^+} \frac{\sin(\theta)}{\theta} = 1$$

A similar argument can be made with $\theta < 0$ to show that the left-side limit is also 1. Putting these together, we have the desired result

$$\lim_{\theta \rightarrow 0} \frac{\sin(\theta)}{\theta} = 1$$

□

Example 5

Evaluate the following limit.

$$\lim_{h \rightarrow 0} \frac{\cos(h) - 1}{h}$$

Solution: We cannot simply substitute $h = 0$ into the limit because this would give the indeterminate form $\frac{0}{0}$. Instead, we will multiply the argument of the limit by

$$\frac{\cos(h) + 1}{\cos(h) + 1} = 1$$

Note that since we are inserting this into a limit as $h \rightarrow 0$, we do not need to worry about the denominator of this expression going to zero.

So, we have

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{\cos(h) - 1}{h} &= \lim_{h \rightarrow 0} \frac{\cos(h) - 1}{h} \frac{\cos(h) + 1}{\cos(h) + 1} \\ &= \lim_{h \rightarrow 0} \frac{\cos^2(h) - 1}{h(\cos(h) + 1)} \\ &= - \lim_{h \rightarrow 0} \frac{\sin^2(h)}{h(\cos(h) + 1)} \\ &= - \left(\lim_{h \rightarrow 0} \frac{\sin(h)}{h} \right) \left(\lim_{h \rightarrow 0} \frac{\sin(h)}{\cos(h) + 1} \right) \\ &= -(1)(0) \\ &= 0 \end{aligned}$$

Observe that in the second last step, we used the Fundamental Trigonometric Limit to evaluate $\lim_{h \rightarrow 0} \frac{\sin(h)}{h}$.

Therefore, we have

$$\lim_{h \rightarrow 0} \frac{\cos(h) - 1}{h} = 0$$

The previous example illustrates a scenario where the Fundamental Trigonometric Limit proves useful. It also turns out that this specific result is needed to determine the derivative of $\sin(x)$. Let's do that now. To begin, we apply the definition of the derivative function

$$\begin{aligned} \frac{d}{dx}(\sin(x)) &= \lim_{h \rightarrow 0} \frac{\sin(x+h) - \sin(x)}{h} \\ &= \lim_{h \rightarrow 0} \frac{\sin(x)\cos(h) + \cos(x)\sin(h) - \sin(x)}{h} \\ &= \sin(x) \left(\lim_{h \rightarrow 0} \frac{\cos(h) - 1}{h} \right) + \cos(x) \left(\lim_{h \rightarrow 0} \frac{\sin(h)}{h} \right) \end{aligned}$$

In the second line of the calculation above, we used the identity $\sin(A+B) = \sin(A)\cos(B) + \cos(A)\sin(B)$. We were then able to separate the expression into two limits and factor out of these two limits $\sin(x)$ and $\cos(x)$, respectively, since the limit is only with respect to h . As for limits themselves, we showed in the previous example that the first one is zero and the second one is the Fundamental Trigonometric Limit which is equal to 1. It follows that,

$$\frac{d}{dx}(\sin(x)) = \cos(x)$$

EXERCISE

Prove using the limit definition of the derivative that $\frac{d}{dx}(\cos(x)) = -\sin(x)$.

Now that we know the derivatives of $\sin(x)$ and $\cos(x)$, we can use differentiation rules to find the derivatives of other trigonometric functions.

Example 6 Determine the derivative of $\tan(x)$.

Solution: We begin by writing $\tan(x)$ in terms of $\sin(x)$ and $\cos(x)$ and then we apply

the Quotient Rule.

$$\begin{aligned}
 \frac{d}{dx} \tan(x) &= \frac{d}{dx} \left(\frac{\sin(x)}{\cos(x)} \right) \\
 &= \frac{\cos(x) \frac{d}{dx}(\sin(x)) - \sin(x) \frac{d}{dx}(\cos(x))}{\cos^2(x)} \\
 &= \frac{\cos^2(x) + \sin^2(x)}{\cos^2(x)} \\
 &= \frac{1}{\cos^2(x)} \\
 &= \sec^2(x)
 \end{aligned}$$

EXERCISE

Determine the derivative of $\csc(x)$.

3.6.2 Exponential Functions

Consider the base- a exponential function $f(x) = a^x$ with $a > 0$. With $a > 1$, this function can be used to describe exponential growth while for $0 < a < 1$, this function can be used to describe exponential decay. If we'd like to know the rate of growth or decay for such a function, we will need its derivative.

Applying the limit definition of the derivative, we have

$$\begin{aligned}
 f'(x) &= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \\
 &= \lim_{h \rightarrow 0} \frac{a^{x+h} - a^x}{h} \\
 &= \lim_{h \rightarrow 0} a^x \frac{a^h - 1}{h} \\
 &= a^x \lim_{h \rightarrow 0} \frac{a^h - 1}{h}
 \end{aligned}$$

Observe that we can extract the a^x term from the limit since the limit is being taken with respect to h . However, it is not clear how to evaluate this limit. We can prescribe geometric meaning to it though.

It is always true that $f(0) = a^0 = 1$. So, we can rewrite the limit above as

$$\begin{aligned}
 \lim_{h \rightarrow 0} \frac{a^h - 1}{h} &= \lim_{h \rightarrow 0} \frac{a^h - a^0}{h} \\
 &= \lim_{h \rightarrow 0} \frac{f(h) - f(0)}{h - 0}
 \end{aligned}$$

This is the limit definition of $f'(0)$. Therefore, we have

$$f'(x) = a^x \cdot f'(0)$$

Base- e Exponential Function

We define the irrational number e ($= 2.718\dots$) to be the base of the exponential function whose derivative at $x = 0$ is equal to 1 - that is, $f'(0) = 1$ for this particular exponential function. The derivative is then given by

$$(e^x)' = e^x$$

or, put more simply, the function is its own derivative.

This also allows us to define the following limit.

$$\lim_{h \rightarrow 0} \frac{e^h - 1}{h} = 1$$

We've run into a bit of a roadblock with determining the derivative of the exponential function a^x when the base is *not* equal to e . However, once we have the Chain Rule at our disposal, this derivative along with many others will become readily available. So, let's look at that now.

3.7 Chain Rule

Consider a drop of water on a kitchen counter top. Let's model this drop as a hemisphere of water with radius r . Since it is half of a sphere, its volume as a function of radius is $V(r) = \frac{2}{3}\pi r^3$. After a while, you'll notice the drop getting smaller due to evaporation. Suppose we want to model the rate of change of the volume with respect to time, $\frac{dV}{dt}$. This is hard to measure directly, but we might more easily take some radius measurements and model $\frac{dr}{dt}$. How does this rate of change relate to the one we want to find?

Well, we have a formula relating the volume and the radius. If we know how quickly the radius is changing at a particular time, then we just need to scale this by another quantity which tells us how quickly the volume would change with respect to the radius. In other words, we need to multiply $\frac{dr}{dt}$ by $\frac{dV}{dr}$ to get $\frac{dV}{dt}$.

To put this in more concrete terms, suppose that at a particular instant, say $t = 10$ min after we start observing, the radius of the drop is 5 mm and the radius is changing at a rate of 0.1 mm/min. We can compute $\frac{dV}{dr} = 2\pi r^2$ so we also get that at that instant in time $\frac{dV}{dr}\Big|_{r(10)=5} = 2\pi(5)^2 = 50\pi$ mm². Then, the rate of change of the volume of the droplet with respect to time is

$$\begin{aligned} \frac{dV}{dt}\Big|_{t=10} &= \frac{dV}{dr}\Big|_{r(10)=5} \cdot \frac{dr}{dt}\Big|_{t=10} \\ &= (50\pi \text{ mm}^2) (0.1 \text{ mm/s}) \\ &= 5\pi \text{ mm}^3/\text{min} \\ &\approx 15.7 \text{ mm}^3/\text{min} \end{aligned}$$

What we were essentially dealing with in the example above was finding the derivative of a composite function. We know the function $V(r)$ and we also know there also exists

some function $r(t)$. It is sensible to be interested in the composition of these two functions $V(r(t))$. It turns out that, in general, these functions will satisfy

$$\frac{dV}{dt} = \frac{dV}{dr} \cdot \frac{dr}{dt}$$

and we can generalize this result to other composite functions. We call the generalization the Chain Rule.

Theorem 4 (Chain Rule)

Let $g(x)$ be differentiable at $x = a$ and let $f(y)$ be differentiable at $y = g(a)$. Then $h(x) = f \circ g(x) = f(g(x))$ is differentiable at $x = a$ and

$$h'(a) = f'(g(a)) \cdot g'(a)$$

REMARKS

1. When the functions f and g are composed in the order $f \circ g$, the function f is sometimes called the outer function and the function g is something called the inner function.
2. We can write the chain rule in Leibniz notation as

$$\left. \frac{df}{dx} \right|_{x=a} = \left. \frac{df}{dg} \right|_{u=g(a)} \cdot \left. \frac{dg}{dx} \right|_{x=a}$$

3. The term $f'(g(a))$ is computed by first taking the derivative of f with respect to g and then evaluating this derivative at $g(a)$.
4. If the differentiability requirements for f and g are satisfied on some interval in x , then the derivative function of h can be written as

$$h'(x) = f'(g(x)) \cdot g'(x)$$

We will use this more general statement frequently to differentiate composite functions.

A rigorous proof of the chain rule is beyond the scope of this course. However, the basic idea underlying the proof can be understood by generalizing the example at the beginning of this section. If a change Δx in x gives rise to a change Δg in g , then $\Delta g \approx \frac{dg}{dx} \Delta x$. Similarly, if a change Δg in g gives rise to a change Δf in f , then $\Delta f \approx \frac{df}{dg} \Delta g$. It follows that

$$\Delta f \approx \frac{df}{dg} \Delta g \approx \frac{df}{dg} \cdot \frac{dg}{dx} \Delta x \quad \implies \quad \frac{\Delta f}{\Delta x} \approx \frac{df}{dg} \cdot \frac{dg}{dx}$$

In the limit that $\Delta x \rightarrow 0$, this result becomes exact and $\frac{\Delta f}{\Delta x} \rightarrow \frac{df}{dx}$.

REMARK

Owing to the notation $\frac{df}{dx}$, it is tempting to think of a derivative as a fraction. The derivative is *not* a fraction (it's a limit!). However, the chain rule allows us to pretend (only for a moment) that what we're dealing with are fractions, and that they cancel like fractions cancel. Indeed, if f is a function of u , and u is a function of x , then the chain rule says

$$\frac{df}{dx} = \frac{df}{du} \frac{du}{dx}.$$

The proof of the chain rule is *not* simply “cancel the ‘ du ’s and you're on your way!”, and as stated above, is quite involved.

The chain rule does tell us that Leibniz notation is good notation though! Sometimes good notation can help us remember statements of theorems, and can help us think more clearly.

Example 7

Let $f(x) = e^{x^2}$.

- Determine $f'(3)$.
- Determine an expression for $f'(x)$ which can be used to find $f'(a)$ for any $a \in \mathbb{R}$.

Solution:

- We first need to describe $f(x)$ as the composition of two functions. If we let $f(g) = e^g$ and $g(x) = x^2$, then $f(g(x)) = e^{g(x)} = e^{x^2}$ as desired.

Let's compute some derivatives.

$$\frac{df}{dg} = \frac{d}{dg}(e^g) = e^g \quad \text{and} \quad \frac{dg}{dx} = \frac{d}{dx}(x^2) = 2x$$

Evaluating these at $x = 3$ gives

$$\left. \frac{df}{dg} \right|_{g(3)=9} = e^9 \quad \text{and} \quad \left. \frac{dg}{dx} \right|_{x=3} = 6$$

Therefore, by the Chain Rule

$$f'(3) = \left. \frac{df}{dx} \right|_{x=3} = \left. \frac{df}{dg} \right|_{g(3)=9} \cdot \left. \frac{dg}{dx} \right|_{x=3} = 6e^9$$

- To get a general expression for $f'(x)$ we stop short in the calculation in part (a) of evaluating our derivatives at a specific x value. However, we do take the extra step of eliminating g from our final expression since we defined this intermediate function only as a tool for applying the Chain Rule.

We therefore have

$$f'(x) = \frac{df}{dx} = \frac{df}{dg} \cdot \frac{dg}{dx} = (e^y) \cdot (2x) = 2xe^{x^2}$$

Notice, we can quickly reproduce the answer to part (a) with this expression by plugging in $x = 3$.

EXERCISE

Let $f(x) = \tan^3(x)$. Determine $f'(x)$ and $f'(\frac{\pi}{4})$. (When applying the Chain Rule, remember that $\tan^3(x)$ just means $(\tan(x))^3$.)

Let's revisit now that we have the Chain Rule, how to compute the derivative of the exponential function a^x with any base $a > 0$.

Let $f(x) = a^x$ and observe that we can rewrite this as

$$f(x) = a^x = (e^{\ln(a)})^x = e^{\ln(a)x}$$

In other words, we can view $f(x) = a^x$ as the composition of $f(g) = e^g$ and $g(x) = \ln(a)x$. Therefore, by the Chain Rule, we get

$$f'(x) = f'(g(x)) \cdot g'(x) = e^{g(x)} \cdot \ln(a) = e^{\ln(a)x} \ln(a) = a^x \ln(a)$$

Let's summarize this result for future reference.

Fact 5

Let $f(x) = a^x$ with $a > 0$, then

$$\frac{df}{dx} = \frac{d}{dx} a^x = a^x \ln(a)$$

Another useful feature of the Chain Rule is that it can be applied iteratively to differentiate functions which are built up from composing multiple simpler functions.

Example 8

If $f(x) = e^{\sin(x^2)}$, then determine $f'(x)$.

Solution: Observe that we can view $f(x)$ as the composition of three functions. That is, $f(x) = f(g(h(x)))$ where $f(g) = e^g$, $g(h) = \sin(h)$, and $h(x) = x^2$.

Applying the Chain Rule once gives

$$\frac{df}{dx} = \frac{df}{dg} \cdot \frac{dg}{dx}$$

Applying it again to $\frac{dg}{dx}$ and then evaluating each derivative gives

$$\begin{aligned} \frac{df}{dx} &= \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{dx} \\ &= (e^{g(h(x))}) \cdot (\cos(h(x))) \cdot (2x) \\ &= (e^{\sin(x^2)}) \cdot (\cos(x^2)) \cdot (2x) \\ &= 2x \cos(x^2) e^{\sin(x^2)} \end{aligned}$$

This looks complicated until you realize we are really just differentiating from the outside in. At each step, we differentiate an outer function while momentarily ignoring the details of its inner function. We then differentiate the inner function. If we can differentiate it easily, we do so and we're done. If not, we treat it as another outer function and peel another layer off our derivative onion.

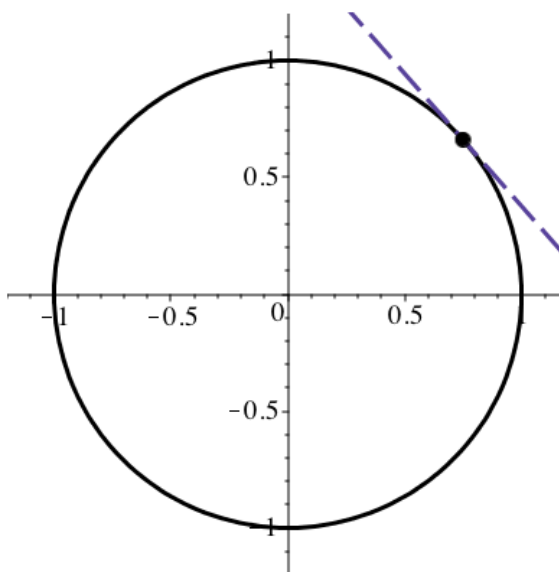
EXERCISE

Determine the derivative of $f(x) = \cos^2(4x)$.

3.8 Implicit Differentiation

Given a function $f(x)$, the derivative $f'(a)$ will tell us the slope of a tangent line to the curve $y = f(x)$ at $x = a$. What if we want to find the slope of a tangent line to a curve which cannot be described in the form $y = f(x)$? A familiar example of such a curve is the unit circle $x^2 + y^2 = 1$. The circle fails the vertical line test, so it is not the graph of a function. That is, we cannot describe the y -coordinates on this curve *explicitly* in terms of a function. Instead, the y -coordinates on this curve are described *implicitly* by an equation relating x and y .

This issue aside, we can still draw a tangent line to the circle anywhere along its circumference. So, how do we find the slope $\frac{dy}{dx}$ of such a tangent line?



A somewhat inefficient way of approaching this particular example is to break this up into two separate problems. That is, we take the equation of the circle and solve for y to get two curves $y = \sqrt{1 - x^2}$ and $y = -\sqrt{1 - x^2}$ corresponding to the upper and lower halves of the circles, respectively. With these functions in hand, we consider only the half of the circle we're interested in and compute a derivative. Or, if we want to investigate both halves, we compute two derivatives and consider both cases. This works for a circle, but it's not always practical or even possible to take an equation relating x and y and partition the curve described by that equation in such a way that we can solve for y in terms of x on each section. Thankfully, there is a much better way of approaching this problem.

Recall, the chain rule tells us how to get started in differentiating the square of a function.

$$\frac{d}{dx} (f(x))^2 = 2f(x) f'(x)$$

If we differentiate the equation $x^2 + y^2 = 1$ and treat y as being implicitly defined in terms of x , then applying the chain rule to deal with the derivative of the y^2 terms gives

$$\frac{d}{dx}(x^2 + y^2) = \frac{d}{dx}(1) \quad \implies \quad 2x + 2y \frac{dy}{dx} = 0$$

Observe that we cannot “complete” taking the derivative because we don’t know how to further simplify the $\frac{dy}{dx}$ term. But that’s all right, we’re going to bypass that step because our ultimate goal is to get an expression for exactly this $\frac{dy}{dx}$. We can now do this by doing a bit of algebra to solve for $\frac{dy}{dx}$.

$$2x + 2y \frac{dy}{dx} = 0 \quad \implies \quad \frac{dy}{dx} = -\frac{x}{y}$$

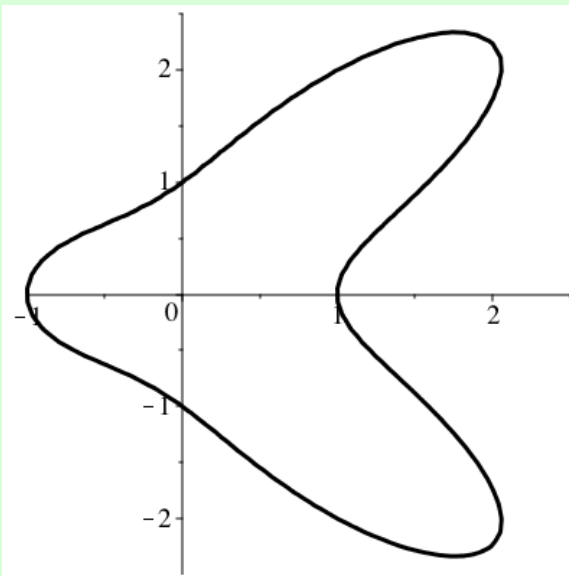
Having the derivative stated in terms of both x and y is new to us. We are used to finding $f'(x)$ as a function in x only. This is actually to be expected here though. For a curve like this, it is not enough to specify an x -coordinate only to refer to a unique point on the circle. For example, there are two points on the circle with $x = 0$. Instead, we need to provide, in general, both an x and a y coordinate to specify a point on the circle. Once we do that, the equation $\frac{dy}{dx} = -\frac{x}{y}$ will tell us the slope of a tangent to the circle at that point. For example, at the point $(x, y) = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right)$, we get $\frac{dy}{dx} = -1$ consistent with what we see in the diagram above.

REMARKS

- In the example above, we call it **implicit differentiation** when we determine $\frac{dy}{dx}$ by differentiating the equation $x^2 + y^2 = 1$ and solving algebraically for the derivative.
- When differentiating an equation in both x and y , we can always make use of the chain rule to get an expression which is linear in $\frac{dy}{dx}$ and can, therefore, be rearranged to solve for $\frac{dy}{dx}$.

Example 9

Consider the curve described by the equation $x^4 - 4xy^2 + y^4 = 1$.



Determine the slope of a tangent line to this curve at the point $(1, 2)$.

Solution: We begin by differentiating the given equation implicitly. We will use the notation $y' = \frac{dy}{dx}$ for simplicity.

$$\frac{d}{dx}(x^4 - 4xy^2 + y^4) = \frac{d}{dx}(1) \quad \implies \quad 4x^3 - 4y^2 - 8xyy' + 4y^3y' = 0$$

Next, we solve for y' and simplify to get

$$y' = \frac{x^3 - y^2}{2xy - y^3}$$

When $x = 1$ and $y = 2$, this gives

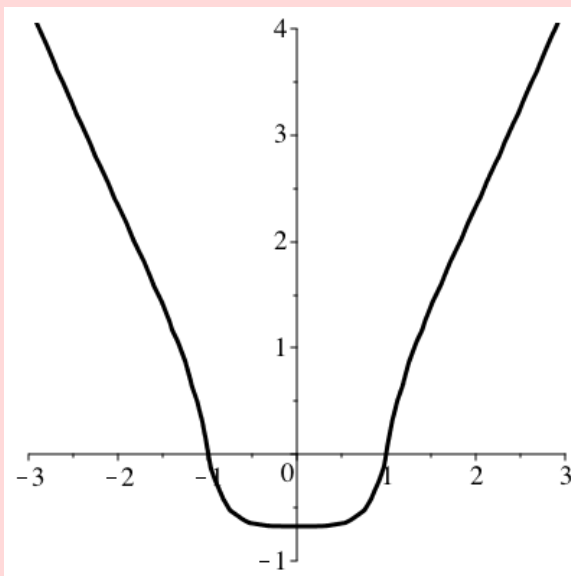
$$y' = \frac{1^2 - 2^2}{2(1)(2) - 2^3} = \frac{3}{4}$$

Therefore, the line tangent to the curve at $(1, 2)$ has slope $\frac{3}{4}$.

EXERCISE

The graph below depicts the curve described by the equation

$$x^4 - y^3 - y = 1.$$



Determine the slope of a tangent to this curve at $(1, 0)$.

(Note: This curve is the graph of a function - that is, there exists an $f(x)$ such that $y = f(x)$ - but this function is very complicated and difficult to find.)

3.8.1 Logarithmic Functions

In previous sections, we discussed how to differentiate a variety of functions. However, there are two important classes which we still have not dealt with. They are logarithmic functions and inverse trigonometric functions. In both cases, we will use implicit differentiation to determine the derivatives of these functions. Let's look at logarithmic functions first.

Fact 6 Let $f(x) = \log_a(x)$ with $a > 0$, then

$$\frac{df}{dx} = \frac{1}{x \ln(a)}$$

REMARK

When $a = e$, we have $\log_e(x) = \ln(x)$ and since $\ln(e) = 1$, we get

$$\frac{d}{dx} \ln(x) = \frac{1}{x}$$

To derive the general result, we first note that

$$y = \log_a(x) \quad \implies \quad x = a^y$$

Differentiating the latter expression implicitly and applying the Chain Rule as needed yields

$$1 = a^y \ln(a) \frac{dy}{dx} \quad \implies \quad \frac{dy}{dx} = \frac{1}{a^y \ln(a)} = \frac{1}{x \ln(a)}$$

Example 10

Use implicit differentiation to prove the Power Rule: $\frac{d}{dx}x^n = nx^{n-1}$ for $n \neq 0$.

Proof: Let $y = x^n$. Taking the natural logarithm of both sides and assuming $n \neq 0$ gives

$$\ln(y) = \ln(x^n) \quad \implies \quad \ln(y) = n \ln(x)$$

Next, we differentiate implicitly to obtain

$$\frac{1}{y} \frac{dy}{dx} = n \frac{1}{x} \quad \implies \quad \frac{dy}{dx} = n \frac{y}{x} = n \frac{x^n}{x} = nx^{n-1}$$

Therefore, when $n \neq 0$, we have $\frac{d}{dx}x^n = nx^{n-1}$.

EXERCISE

Let $f(x) = \ln(x^2 - 1)$.

1. Use the Chain Rule to find $f'(x)$.
2. Factor $x^2 - 1$ and then use log rules to write $f(x)$ as the sum of two separate logarithmic functions. Differentiate this expression and verify that you get the same answer as in part (a).

3.8.2 Inverse Trigonometric Functions

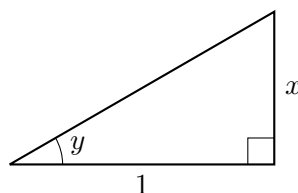
It is not obvious how to differentiate an inverse trigonometric function like $f(x) = \arctan(x)$ (which may also be written as $\tan^{-1}(x)$). However, with implicit differentiation and a bit of geometry, we can work it out.

To begin, let $y = \arctan(x)$. Equivalently, we have $\tan(y) = x$. Differentiating this implicitly (while recalling that $\frac{d}{dx} \tan(x) = \sec^2(x)$) we get

$$\sec^2(y) \frac{dy}{dx} = 1 \quad \implies \quad \frac{dy}{dx} = \cos^2(y)$$

This is great already. If we want to know $f'(x)$ for a specific value of x , we can work out the corresponding value of y (since $y = f(x)$) and plug that into this expression. But we can do even better! Using a bit of geometry, we can get $f'(x)$ entirely in terms of x .

Imagine a right-angled triangle where we'll label one of the acute angles y . Relative to this angle, suppose the opposite side has length x and the adjacent side length has length 1. We've constructed a triangle for which $\tan(y) = x$.



By the Pythagorean theorem, the hypotenuse of this triangle will have length $\sqrt{1+x^2}$. Therefore, we have $\cos(y) = \frac{1}{\sqrt{1+x^2}}$ or $\cos^2(y) = \frac{1}{1+x^2}$. We can substitute this into our result for $\frac{dy}{dx}$ above to get

$$\frac{d}{dx} \arctan(x) = \frac{1}{1+x^2}$$

A similar procedure can be used to find the derivatives of all six inverse trigonometric functions.

Fact 7

The derivatives of the inverse trigonometric functions $\arcsin(x)$, $\arccos(x)$, and $\arctan(x)$ are:

$$\frac{d}{dx} \arcsin(x) = \frac{1}{\sqrt{1-x^2}} \quad \frac{d}{dx} \arccos(x) = -\frac{1}{\sqrt{1-x^2}} \quad \frac{d}{dx} \arctan(x) = \frac{1}{1+x^2}$$

EXERCISE

Use implicit differentiation to show $\frac{d}{dx} \arcsin(x) = \frac{1}{\sqrt{1-x^2}}$.

3.9 Tangent Lines and Linear Approximations

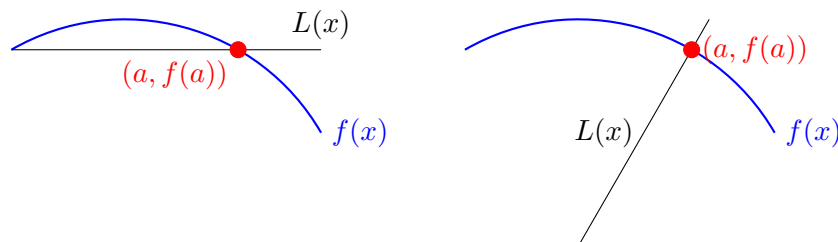
In this section, we will start to see first-hand that calculus is more than just a tool for computing rates of change. Specifically, we can use it to extract from complicated mathematical functions - which may describe real-world phenomena - simple approximations of those functions. In this section, those approximations will be linear functions, so we'll need to be mindful of the accuracy of the approximations. However, it turns out that the ideas we explore in this section can be extended to build arbitrarily accurate polynomial representations of any function and polynomials are generally very nice functions to work with.

Back to the task at hand, the basic idea that we'll be exploiting in this section is that the more you zoom in on a section of the graph of a reasonably well-behaved non-linear function, the flatter the graph looks. (Just like the surface of the Earth looks flat unless you're looking at it from space.) This means that we should be able to draw a line at any point on the graph of the non-linear function and use the linear function describing that line as an approximation of the non-linear function provided we don't stray too far from where the line is constructed.

Now which line should we draw? Here's where differentiability comes in!

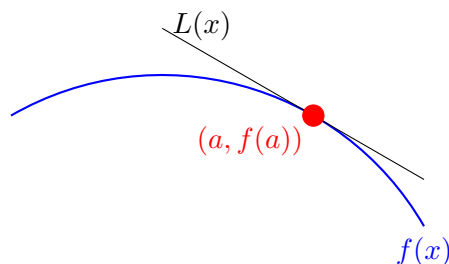
Suppose we have a function $f(x)$ that is differentiable at $x = a$, that is, $f'(a)$ exists. Let's try to find a line that approximates $f(x)$ at $x = a$. What properties should such a line have?

Suppose the line is the graph of a linear function $L(x)$. At the very least, the line should pass through the point $(a, f(a))$ on the graph of $f(x)$. That is, we insist that $L(a) = f(a)$. There are many such lines! For example, here are two of them for a graph of a particular function $f(x)$:



Although both of these linear functions agree with the function at the point $x = a$, they differ from the function $f(x)$ as you move away from $x = a$.

Now, since $f(x)$ is not a linear function, we'll never get a linear function to agree with $f(x)$ exactly. Nonetheless, we can try to find the best approximation (we'll explore what "best" means a little later on) to $f(x)$ by a linear function near $x = a$. Looking at the graph of $f(x)$, a good guess would be a linear function that "follows" the graph of $f(x)$ as closely as it can:



So, the line that we're after is the graph of a linear function whose slope at $x = a$ agrees with the instantaneous slope of $f(x)$ at $x = a$. Such a line is the tangent line to the graph of $f(x)$ at $x = a$.

Let's explore this idea now with a concrete example.

Example 11

Let $f(x) = \sqrt{x}$.

1. Find an equation for the tangent line to the curve $y = f(x)$ at $x = 100$.
2. Compare the approximate decimal value of $\sqrt{101}$ (e.g., computed with a calculator) to the y value of the tangent line at $x = 101$.

Solution:

1. If a line passes through a point (x_0, y_0) and has slope m , then its x and y values satisfy

$$y = y_0 + m(x - x_0)$$

This is sometimes called the **point-slope equation of a line**.

We know that the tangent line will intersect $y = f(x)$ when $x = 100$, so it must contain the point $(x_0, y_0) = (100, 10)$.

Next, its slope will be equal to the instantaneous rate of change of $f(x)$ at $x = 100$. In other words, $m = f'(100)$. Using the power rule to compute $f'(x)$ we find

$$f'(x) = \frac{d}{dx} \left(x^{\frac{1}{2}} \right) = \frac{1}{2} x^{-\frac{1}{2}} \quad \implies \quad f'(100) = \frac{1}{20}$$

Therefore, the desired tangent line equation is

$$y = 10 + \frac{1}{20}(x - 100)$$

2. Using a calculator, we find $\sqrt{101} \approx 10.04988$. Meanwhile, when $x = 101$, the y value on our tangent line is $y = 10 + \frac{1}{20}(101 - 100) = 10.05$.

Observe that the tangent line value gives an excellent approximation. Moreover, since the equation of the tangent line is linear, then only basic arithmetic operations (i.e., addition, subtraction, multiplication, and division) are needed to compute the approximate value.

It is worth graphing out $y = 10 + \frac{1}{20}(x - 100)$ and $y = \sqrt{x}$ on the same axes, and verifying that the former is indeed tangent to the latter at $(100, 10)$.

We generalize the method used in the previous example to define the linear approximation or **linearization** of a function.

Definition 3.9.1

linear
approximation or
linearization,
tangent line

Let $f(x)$ be differentiable at $x = a$. The linearization of $f(x)$ at $x = a$ is denoted $L_a(x)$ and given by

$$L_a(x) = f(a) + f'(a)(x - a).$$

The tangent line to the curve $y = f(x)$ at the point $(a, f(a))$ is the graph of the linearization, $y = L_a(x)$.

EXERCISE

Suppose $f(x)$ is a function differentiable at $x = a$. Let $L_a(x)$ be the linearization of $f(x)$ at $x = a$. Prove that $L_a(a) = f(a)$ and $L'_a(a) = f'(a)$.

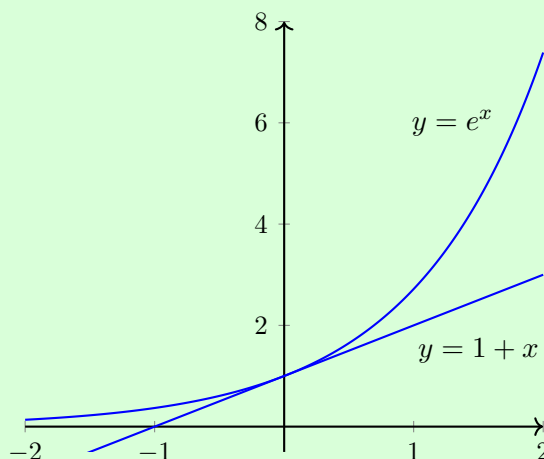
Example 12

Determine the linearization of $f(x) = e^x$ at $x = 0$ and use it to approximate $e^{0.1}$.

Solution: Since $f'(x) = e^x$, we have $f(0) = 1$ and $f'(0) = 1$. Therefore, the linearization of e^x at $x = 0$ is

$$L_0(x) = f(0) + f'(0)(x - 0) = 1 + x$$

This means that $e^x \approx 1 + x$ when x is close to zero. Graphing $y = e^x$ and $y = 1 + x$ confirms this result.



When $x = 0.1$, the linearization gives us $e^{0.1} \approx L_0(0.1) = 1 + 0.1 = 1.1$. Using a calculator, we get $e^{0.1} \approx 1.105$, so the linearization indeed gives a very accurate approximation for not a whole lot of effort.

EXERCISE

Find the linearization of $f(x) = x^{1/3}$ at $x = 27$ and use it to approximate $(26.5)^{1/3}$. Compare your linear approximation to the approximation output by a calculator.

LOOKING AHEAD

What makes a linear function a “good” linear approximation? Why is there no “good” linear approximation to $f(x) = |x|$ at $x = 0$? These are both great questions.

One way to study how good an approximation is, is to look at the error between the approximation and the function. Suppose we wish to approximate the function $f(x)$ at $x = a$ by a linear function $L(x)$. The error in the approximation is the quantity $R(x) = L(x) - f(x)$.

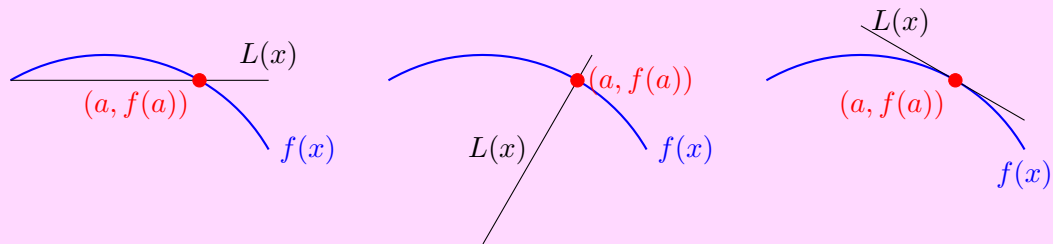
In the examples above, we used a calculator to compute $R(x)$ near $x = a$. However, our goal is not just to do away with calculators. More importantly, we’d like to be able to replace a non-linear function with a linear one and be confident that this replacement is reasonable provided we are working within a particular interval.

Now, unless $f(x)$ is a linear function, you cannot expect the error to be zero everywhere. Also, away from $x = a$, the function could behave completely wildly, curving dramatically away from the linear approximation. So really, it makes sense to study the error as x gets close to a .

A first guess would be to insist that

$$\lim_{x \rightarrow a} R(x) = 0.$$

That is, as we get close to $x = a$, the error goes to 0. While this seems good at first, all three of the following graphs of linear functions satisfy that $R(x) \rightarrow 0$ as $x \rightarrow a$:

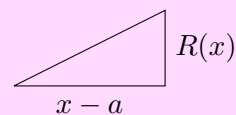


So we need to be a little more careful. It turns out that a good thing to consider is not just the error $R(x)$, but the *ratio* between the error $R(x)$ and the distance away from $x = a$, which is the quantity $x - a$.

One way to characterize $L(x)$ as a “good” linear approximation to $f(x)$ at $x = a$ is to insist that

$$\lim_{x \rightarrow a} \frac{R(x)}{x - a} = 0.$$

Where does this come from? Well, suppose we’re looking at some value x near a . Imagine a right-angled triangle, with base $x - a$ and height $R(x)$.



Then the ratio $\frac{R(x)}{x - a}$ measures how flat that triangle is. If the limit goes to 0 as x approaches a , then the triangle gets flatter and flatter as x approaches a .

In fact, we have the following excellent fact (which is left as an exercise below for you to prove).

Fact 8 (Differentiability and good linear approximations)

Suppose $f(x)$ is continuous at $x = a$.

The function f is differentiable at $x = a$ if and only if there exists a linear function $L(x)$ satisfying $\lim_{x \rightarrow a} \frac{L(x) - f(x)}{x - a} = 0$.

When you work on proving this fact, along the way you will prove that if there is a linear function $L(x)$ satisfying $\lim_{x \rightarrow a} \frac{L(x) - f(x)}{x - a} = 0$, then that linear function must be the linearization of $f(x)$ at $x = a$. Neat!

The take-home message from Fact 8 is the following: Assuming continuity, differentiability is equivalent to the existence of a good linear approximation, and a good linear approximation is one satisfying the limit condition in Fact 8. This is roughly how differentiability is defined for functions of more than one variable, which is something you may encounter if you study multivariable calculus in the future.

We won't worry any more about errors for now because it will be addressed more fully when we study Taylor polynomials. There we will see that we can both create better approximations by using polynomials (with linear functions being a special case) and also determine upper bounds on the error associated with using one of these polynomials.

EXERCISE

Prove Fact 8.

EXERCISE

Let $f(x) = |x|$.

1. Let $L(x) = 0$. Prove that $\lim_{x \rightarrow 0} \frac{L(x) - f(x)}{x - 0} \neq 0$.
2. Let $L(x)$ be any linear function. Prove that $\lim_{x \rightarrow 0} \frac{L(x) - f(x)}{x - 0} \neq 0$.

3.10 Newton's Method

Suppose you are tasked with writing some code to evaluate n -th roots. For concreteness, let's say that includes evaluating something like $\sqrt[5]{3}$. You might try constructing the linearization of the function $f(x) = \sqrt[5]{x}$ at $x = 1$ (because it's easy to evaluate $f(1)$ and $f'(1)$). However, $x = 1$ is not very "close" to $x = 3$ and yields $\sqrt[5]{3} \approx 1.4$. This is not a very accurate approximation! The correct decimal approximation for $\sqrt[5]{3}$ is 1.2457....

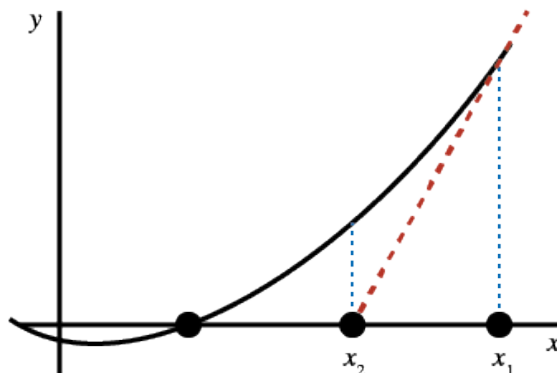
Instead, let's turn this into a root-finding problem and then exploit a clever way of using tangent lines to seek out roots.

To reframe this as a root finding problem, we note that we'd like to find x such that $x = \sqrt[5]{3}$. Equivalently, we'd like to find x such that $x^5 = 3$ or $x^5 - 3 = 0$. If we define $f(x) = x^5 - 3$, then by construction, the root of $f(x)$ (i.e., the value of x such that $f(x) = 0$) will satisfy $x = \sqrt[5]{3}$.

Now we need to come up with a procedure for finding the roots of a function. Geometrically, this is equivalent to finding the x -intercepts of a function. With a bit of thought and maybe

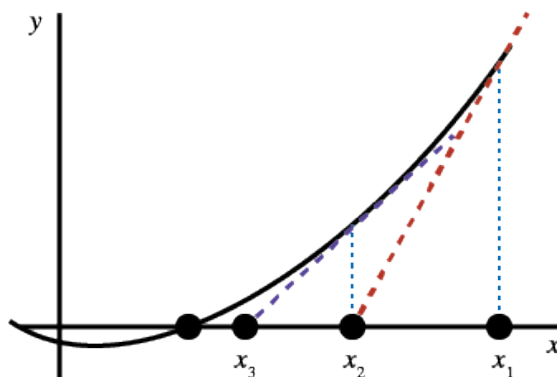
some guidance from the Intermediate Value Theorem, it is not usually too hard to determine an interval in which we can be confident a root exists. But how do we hone in on it with precision? This is where the tangent line comes in.

Consider an initial guess x_1 as an x -value near the root of a function $f(x)$. Next, construct the tangent line to $y = f(x)$ at the point $(x_1, f(x_1))$.



Under most circumstances, this tangent line will intersect the x -axis closer to the actual root of $f(x)$ than x_1 . So, if we call this tangent line intercept x_2 , then this process has generated a value of x which approximates the desired root better than our initial guess of x_1 .

Let's go again! Construct the tangent line to $y = f(x)$ at $(x_2, f(x_2))$, follow that tangent line to the x -axis, and call the intercept x_3 . Typically, we will again be closer to the actual root.



This process can be repeated as many times as we want with each iteration producing a better approximation of the root. It also turns out that we can usually achieve a very high level of precision (e.g., as good as the precision of a calculator) with just a few iterations.

We call this technique **Newton's Method** after - you guessed it - Sir Isaac Newton.

Let's look now at the algebra involved in applying this technique. Recall, we start by constructing the tangent line to $y = f(x)$ at $(x_1, f(x_1))$ and find its x -intercept which we call x_2 . Equivalently, we take the linearization of $f(x)$ at $x = x_1$,

$$L_{x_1}(x) = f(x_1) + f'(x_1)(x - x_1),$$

set $L_{x_1}(x_2) = 0$, and solve for x_2 .

$$L_{x_1}(x_2) = 0 \quad \implies \quad x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$$

To find the next point x_3 , we would simply repeat this process with x_2 in place of x_1 . This would just give

$$x_3 = x_2 - \frac{f(x_2)}{f'(x_2)}$$

This pattern continues, so we can infer that in the i -th iteration, we can determine x_{i+1} from x_i according to the formula

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

Let's summarize what we've found.

Algorithm 1 (Newton's Method)

To approximate a root of the function $f(x)$:

1. Start with an initial guess near a root of $f(x)$. Call it x_1 .
2. Iterate using the equation:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

3. Stop once you have successive values whose matching digits have the desired precision.

To clarify the last point, suppose a few iterations of Newton's method have yielded $x_3 = 3.14785$, $x_4 = 3.14162$, and $x_5 = 3.14159$. By themselves, we do not know how good of an approximation any of these values are. However, comparing successive values allows us to gain some insight. For example, x_3 and x_4 share the same first three digits and we expect further approximations to *increase* the precision, so we can reasonably infer that the root is at 3.14 ± 0.01 . Similarly, x_4 and x_5 share the first four digits, so at this point we can be confident that the root is at 3.141 ± 0.001 .

CAUTION

- For an iteration of Newton's method with input x_i to work, $f'(x_i)$ must exist and be non-zero.
- For a function with more than one root, different starting values will converge to different roots. For some functions, it can be tricky to predict a good starting point to seek out a particular root.
- It is possible for Newton's method to get caught in a cycle (i.e., $x_{i+m} = x_i$ for all $m, i \in \mathbb{Z}$ with $m \geq 2$) and fail to produce good approximations.

- Issues with applying Newton's method can usually be circumvented by choosing a different initial guess, perhaps guided by taking a moment to consider the shape of the graph of the function.

Example 13

Let's revisit the example at the beginning of this section and use Newton's Method to approximate $\sqrt[5]{3}$.

Solution: We first define $f(x) = x^5 - 3$ to turn this into a root-finding problem. Next, as we'll be applying Newton's Method, we calculate $f'(x) = 5x^4$.

Now we need an initial guess to get started. We note that $f(x)$ is a continuous function and also that $f(x)$ changes signs between $x = 1$ and $x = 2$. Therefore, by the Intermediate Value Theorem, there must be a root on the interval $(1, 2)$. It would be reasonable to start in the middle of this interval, but a closer look at the values of $f(x)$ at the endpoints ($f(1) = -2$ and $f(2) = 29$) suggests that the root is probably closer to the left endpoint. So, let's start with $x_1 = 1$.

Applying Newton's method once, we have:

$$\begin{aligned} x_2 &= x_1 - \frac{f(x_1)}{f'(x_1)} \\ &= 1 - \frac{f(1)}{f'(1)} \\ &= 1 - \frac{(1)^5 - 3}{5(1)^4} \\ &= 1.4 \end{aligned}$$

A few more iterations, keeping ten digits as we go, and we get

$$\begin{aligned} x_3 &= 1.276184923 \\ x_4 &= 1.247150132 \\ x_5 &= 1.245734166 \\ x_6 &= 1.245730940 \\ x_7 &= 1.245730940 \end{aligned}$$

The actual value of $\sqrt[5]{3}$ to ten significant digits is 1.245730940. Observe that Newton's method produces this number after just five iterations with x_6 (and confirms it with x_7). The rate at which this method gains precision is impressive!

REMARK

You may be wondering how important your initial guess x_1 is. In the previous example, we reasoned that we should take x_1 near the left endpoint of the interval $[1, 2]$ and so went with $x_1 = 1$. If we instead tried the middle of the interval (i.e., $x_1 = 1.5$), the values would have converged at a similar rate. If we'd been a bit more clever and assumed it would be closer to the left but not all the way and guessed, say, $x_1 = 1.2$, then we'd would have been

able to skip two iterations for the same level of precision. Either way though, the values converge quite quickly.

EXERCISE

Use Newton's method to approximate $\sqrt{2}$.

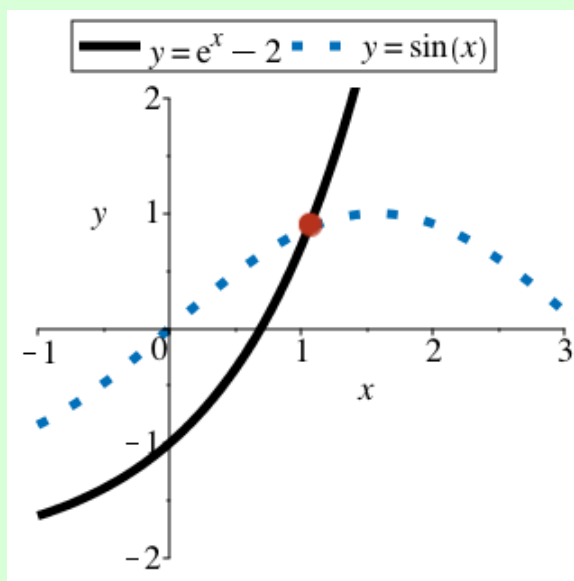
At this point, you could be forgiven for thinking that Newton's method is pretty neat but that you can also just do the same calculations much faster with a calculator. However, the real value in Newton's method is its utility for solving equations.

Example 14

Determine the solution to the following equation to four significant digits.

$$e^x - 2 = \sin(x)$$

Solution: This equation cannot be solved algebraically for x . However, if we plot $y = e^x - 2$ and $y = \sin(x)$, we see that a solution does exist near $x = 1$.



As such, we define $f(x) = e^x - 2 - \sin(x)$ to turn this into a root-finding problem. The solution to $f(x) = 0$ will be the same value that solves the original equation.

Let's apply Newton's method with $x_1 = 1$ (since the two curves in our graph intersect near there). We have $f'(x) = e^x - \cos(x)$, so

$$x_2 = 1 - \frac{e^1 - 2 - \sin(1)}{e^1 - \cos(1)} \approx 1.0566$$

We keep five significant digits since the question ultimately wants us to find the solution to four significant digits. The extra digit will avoid round-off errors affecting the final answer.

Performing two more iterations of Newton's method, we find

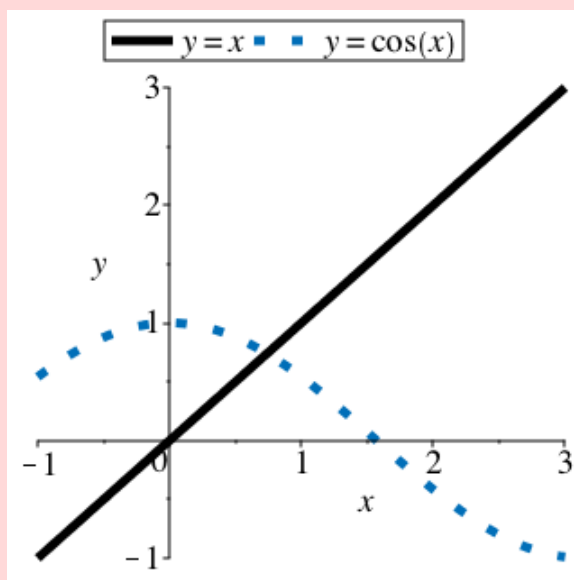
$$x_3 \approx 1.0541$$

$$x_4 \approx 1.0541$$

To the desired precision, our values are not changing after just three iterations. Therefore, we can infer that the solution to the original equation is 1.054 to four significant digits.

EXERCISE

Use Newton's method to find the solution to $\cos(x) = x$ to four significant digits. The following plot including $y = \cos(x)$ and $y = x$ should help you choose an appropriate x value to get started.

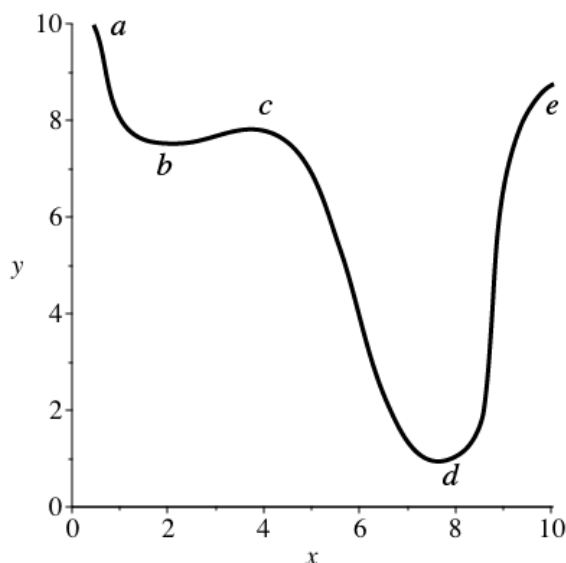


3.11 Local Extreme Values

Have you ever wondered if there is an optimal speed to drive a car so that your vehicle is using as little fuel as possible per kilometre driven? A car uses a certain amount of fuel just to idle, so going at very slow speeds is not ideal. But wind resistance increases non-linearly with speed, so very fast speeds are also not ideal. Maybe there is a sweet spot in the middle somewhere?

In fact, there must be! Consider the function describing your fuel economy in terms of your car's speed. It's reasonable to suppose that this function is continuous and therefore, by the Extreme Value Theorem, there must be a global maximum (i.e., most fuel efficient speed) and a global minimum (i.e., least fuel efficient speed) on the interval of speeds that your car can operate over. In this section, we'll look at how we can use differentiation to locate these global extreme values as well as local extreme values.

To help visualize the different types of extreme values, consider the following graph.



A quick inspection reveals that there is a **global maximum** at the point labelled a and a **global minimum** at a point labelled d . The points labelled b and c are special too, though. The point b has a lesser value than all the points in some neighbourhood around it, so we call it a **local minimum**. Similarly, we call the point c a **local maximum**. (We'll come back to point e momentarily.)

We previously defined global extreme values as follows.

Definition 3.11.1
absolute max/min

Let f be a function with domain D and let $c \in D$.

- If $f(c) \geq f(x)$ for all $x \in D$, then we say $f(c)$ is the absolute maximum of f on D .
- If $f(c) \leq f(x)$ for all $x \in D$, then we say $f(c)$ is the absolute minimum of f on D .

We now also define local extreme values. Recall that an open interval in \mathbb{R} is a set of the form (a, b) , which denotes all numbers x satisfying $a < x < b$. If we say that an open interval (a, b) contains c , then $a < c < b$.

Definition 3.11.2
local max/min

Let f be a function with domain D and let $c \in D$.

- We say $f(c)$ is a local maximum of f on D if there exists an open interval $I \subseteq D$ containing c such that $f(c) \geq f(x)$ for all $x \in I$.
- We say $f(c)$ is a local minimum of f on D if there exists an open interval $I \subseteq D$ containing c such that $f(c) \leq f(x)$ for all $x \in I$.

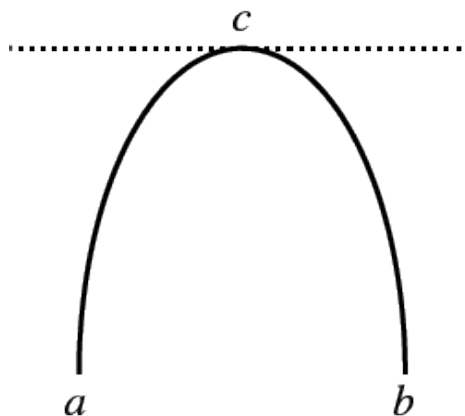
REMARK

Any global extreme value which is not at the left or right-endpoint of the domain D , is also a local extreme value. For example, in the graph above, the point d is both a global and local minimum. However, the definition of a local extreme value (i.e., it exists on an *open* interval) requires that there be points on either side of where that local extreme value is located. Therefore, the endpoints of D cannot be local extreme values. For example, in the graph above, the point a is a global maximum, but it is not a local maximum. Similarly, the point e is not considered an extreme value at all.

Now that we can classify the extreme values of a function, how do we find them? Let's begin with global extrema.

Locating Global Extreme Values

Suppose $f(x)$ has domain $D = [a, b]$. It is possible that the global extreme values occur at the endpoints of the domain, so we'll keep this in mind for later. It is also possible that they occur on the interval (a, b) . In that case, it turns out there are only two conditions we need to consider to locate them. To understand the first condition, consider the following graph.



There is a global extreme value at $x = c$. Observe that the tangent line to $y = f(x)$ is horizontal at $x = c$. This is because the function changes between increasing and decreasing at the exact location it has a local extreme value. In mathematical terms, at this point the derivative of f is zero - that is, $f'(c) = 0$.

Next, suppose that there is an abrupt (i.e., discontinuous) change in the slope of $y = f(x)$.



This situation could also yield an extreme value. Moreover, when the slope changes abruptly like this, the derivative of f does not exist. This gives us another possibility to consider - specifically, $f'(c)$ does not exist.

We group together these two cases with the following definition.

Definition 3.11.3
critical point

Let $f(x)$ have domain D . A point $c \in D$ is called a critical point of f if either

- $f'(c) = 0$ or
- $f'(c)$ does not exist.

We can now summarize our thoughts into an algorithm for determining the global extreme values of a function defined on a closed interval. This algorithm is often referred to as the Closed Interval Method.

Algorithm 2 (Closed Interval Method)

Let $f(x)$ be a continuous function with domain $D = [a, b]$. The global extrema (maximum and minimum values) can be determined as follows:

1. Determine all critical points (c_1, c_2, \dots, c_n) of f .
2. Evaluate f at each critical point.
3. Evaluate f at the endpoints of the domain (i.e., compute $f(a)$ and $f(b)$).
4. Compare the values obtained in steps 2 and 3. The largest value is the global maximum. The smallest value is the global minimum.

Example 15

Determine the global maximum and global minimum of $f(x) = 16x^2 - x^4$ on the interval $D = [-2, 3]$.

Solution: We apply the Closed Interval Method. To start, we compute $f'(x) = 32x - 4x^3$.

The derivative $f'(x)$ exists for all $x \in D$ so we will only possibly find critical points if the derivative is zero.

$$f'(x) = 0 \implies 4x(8 - x^2) = 0 \implies x(x - 2\sqrt{2})(x + 2\sqrt{2}) = 0$$

Since $-2\sqrt{2} < -2$, the only critical points are at $x = 0$ and $x = 2\sqrt{2}$. At these critical points, f takes the values

$$f(0) = 0 \quad \text{and} \quad f(2\sqrt{2}) = 64$$

At the endpoints of the interval, we have

$$f(-2) = 48 \quad \text{and} \quad f(3) = 63$$

Comparing values we see

$$f(0) < f(-2) < f(3) < f(2\sqrt{2})$$

Therefore, $f(2\sqrt{2}) = 64$ is the global maximum and $f(0) = 0$ is the global minimum.

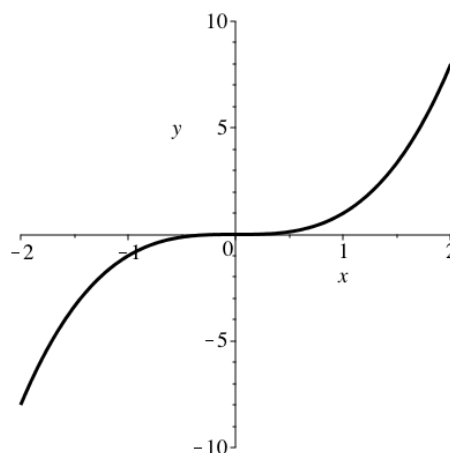
EXERCISE

Determine the global extreme values of $g(x) = \ln(x^2 - 5x + 8)$ on the interval $[0, 6]$.

3.11.1 Locating Local Extreme Values

We can build off of the previous section to help us locate local extrema. In particular, we'll use the idea that any local extrema must be located at critical points. This is because a critical point is the only place where a function can switch between increasing and decreasing.

However, not every critical point is a local extremum. For example, consider the function $f(x) = x^3$. A search for critical points would quickly turn up $f'(0) = 0$, so there is a single critical point at $x = 0$.



Indeed, the graph of $y = f(x)$ has a horizontal tangent line at $x = 0$, but the derivative of f is positive on either side of $x = 0$, so there is no local maximum or minimum.

The lesson here is that we must also consider the sign of $f'(x)$ on either side of a critical point. If it changes signs, then we have a local extremum. As a bonus, knowing the sign of $f'(x)$ on either sign of a local extremum also allows us to identify it as a maximum or minimum. We summarize this idea with the following test.

Theorem 9 (First Derivative Test)

Let $f(x)$ be continuous on interval D and suppose $c \in D$ is a critical point of f , then:

- If there exists an interval $I = (a, b)$ with $c \in I$ such that $f'(x) < 0$ for $x \in (a, c)$ and $f'(x) > 0$ for $x \in (c, b)$, then f has a local minimum at $x = c$.
- If there exists an interval $I = (a, b)$ with $c \in I$ such that $f'(x) > 0$ for $x \in (a, c)$ and $f'(x) < 0$ for $x \in (c, b)$, then f has a local maximum at $x = c$.

To prove this test, we will rely on the all-powerful Mean Value Theorem. We will return to a proof of the First Derivative Test once we encounter the Mean Value Theorem in the future.

In simpler terms, if the derivative of f changes from negative to positive across a critical point, then there is a local minimum there. Similarly, if the derivative changes from positive to negative, then there is a local maximum. If the sign of the derivative does *not* change, then there is neither a local minimum nor a local maximum.

Example 16

Determine the local extrema of $f(x) = x - \sqrt{x}$ on the domain $x > 0$.

Solution: We apply the First Derivative Test. To begin, we compute $f'(x) = 1 - \frac{1}{2\sqrt{x}}$. The derivative exists for all $x > 0$, so critical points will possibly occur only when $f'(x) = 0$.

$$f'(x) = 0 \quad \Longrightarrow \quad x = \frac{1}{4}$$

Observe that when $0 < x < \frac{1}{4}$, we have

$$\begin{aligned} x < \frac{1}{4} &\Longrightarrow \sqrt{x} < \frac{1}{2} \\ &\Longrightarrow \frac{1}{\sqrt{x}} > 2 \\ &\Longrightarrow \frac{1}{2\sqrt{x}} - 1 > 0 \\ &\Longrightarrow 1 - \frac{1}{2\sqrt{x}} < 0 \\ &\Longrightarrow f'(x) < 0 \end{aligned}$$

Similarly, we can show that when $x > \frac{1}{4}$, we have $f'(x) > 0$.

Therefore, by the First Derivative Test, f has a local minimum at $x = \frac{1}{4}$. This local minimum is equal to $f(\frac{1}{4}) = -\frac{1}{4}$.

Fact 10 When applying the First Derivative Test, it is helpful to keep in mind that the sign of $f'(x)$ can only possibly change at a critical point.

This fact gives us a shorter path to establishing the sign of the derivative around a critical point. Let's use the previous example to see why. Since $f(x)$ only had a critical point at $x = \frac{1}{4}$, then the sign of $f'(x)$ must not change on the interval $(0, \frac{1}{4})$. As such, we can simply check the sign of $f'(x)$ at a single point in this interval to determine its sign on the entire interval. For instance, since $\frac{1}{9} \in (0, \frac{1}{4})$ and $f'(\frac{1}{9}) = -\frac{1}{2} < 0$, then $f'(x) < 0$ to the left of $x = \frac{1}{4}$. Similarly, since $1 \in (\frac{1}{4}, \infty)$ and $f'(1) = \frac{1}{2} > 0$, we can infer that $f'(x) > 0$ to the right of $x = \frac{1}{4}$. This approach is generally much more efficient as it avoids considering a sequence of inequalities as we did in the solution above.

EXERCISE

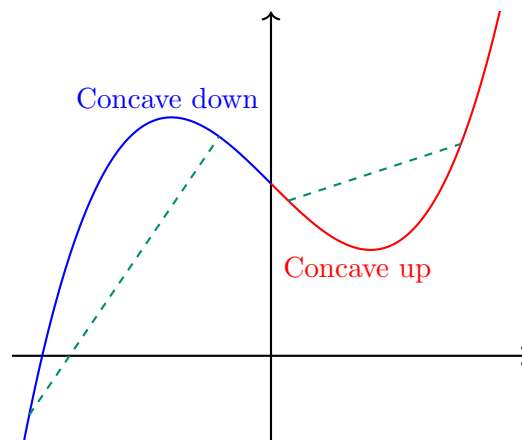
Use the First Derivative Test to locate the local maximum and local minimum points of $f(x) = x^3 - 4x^2 + 4x + 1$ where the domain of f is all real numbers.

It is also possible to use the second derivative to classify critical points. The underlying idea here is to use the second derivative to determine the **concavity** of function in the neighbourhood of a critical point.

Definition 3.11.4

concave up/down

Let f be defined on an interval I . If for every pair of points $a, b \in I$ the secant line connecting $(a, f(a))$ with $(b, f(b))$ lies above (below) the graph of f , then we say f is concave up (down).



Roughly speaking, concave up means the graph opens in the upwards direction and concave down means it opens in the downwards direction.

Fact 11 Given a function $f(x)$, if the second derivative of f exists over some interval I , then the definition of concavity simplifies to:

- If $f''(x) > 0$ over some interval, then f is concave up on that interval.
- If $f''(x) < 0$ over some interval, then f is concave down on that interval.

We will return to prove this fact once we have the Mean Value Theorem in our toolbelt.

This special case definition of concavity turns out to be most useful in practice since we'll typically deal only with functions that are at least twice differentiable.

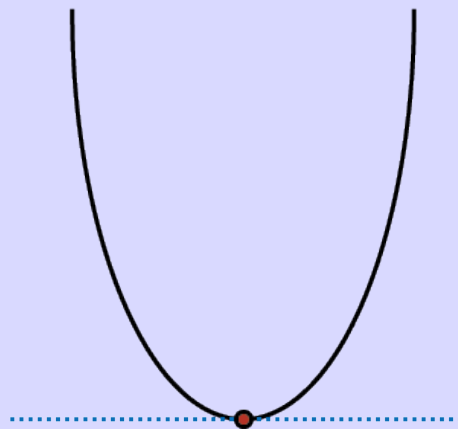
Geometrically, the reason we can connect the second derivative of a function to its concavity is because the second derivative tells us the rate of change of the slope of a tangent line to the graph of the function. If the function is opening upwards around a critical point, then the slopes of a succession of tangent lines are increasing with those slopes being negative to the left of the critical point, zero at the critical point, and positive to the right of the critical point. Most importantly though, if $f'(x)$ is increasing, then $f''(x)$ is positive. Similarly, when a graph is opening downwards, $f'(x)$ is decreasing and $f''(x)$ is negative.

When we use concavity to classify critical points, we say we are using the Second Derivative Test.

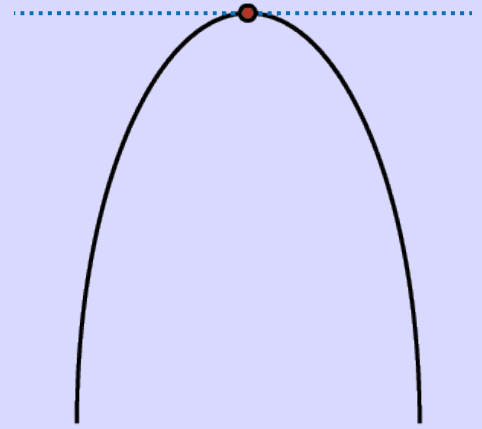
Theorem 12 (Second Derivative Test)

Suppose f has a critical point at $x = c$ (i.e., $f'(c) = 0$) and $f''(c)$ exists, then:

- If $f''(c) > 0$, then f has a local minimum at $x = c$.
- If $f''(c) < 0$, then f has a local maximum at $x = c$.



Concave up with local minimum at critical point



Concave down with local maximum at critical point

The proof of the Second Derivative Test will be deferred until we have the Mean Value Theorem at our disposal.

Example 17

Use the Second Derivative Test to determine the local extrema of $f(x) = x - \sqrt{x}$ on the domain $x > 0$.

Solution: We solved this problem above using the First Derivative Test. There we found $f'(x) = 1 - \frac{1}{2\sqrt{x}}$ and located a critical point at $x = \frac{1}{4}$.

To apply the Second Derivative Test, we next differentiate $f'(x)$ to get $f''(x) = \frac{1}{4x^{3/2}}$. Since $f''(\frac{1}{4}) = 2$ is positive, then the critical point at $x = \frac{1}{4}$ is a local minimum.

When using the First Derivative Test, we needed to consider the sign of $f'(x)$ on either side of the critical point. For the Second Derivative Test, we only need to consider the sign of $f''(x)$ at the critical point, but we do also need to compute $f''(x)$. So, neither method is significantly faster than the other. However, if it will be useful to have the second derivative on hand - perhaps to help gain insight into the behaviour of the function by knowing its concavity at various intervals across its domain - then the Second Derivative Test route may serve a dual purpose.

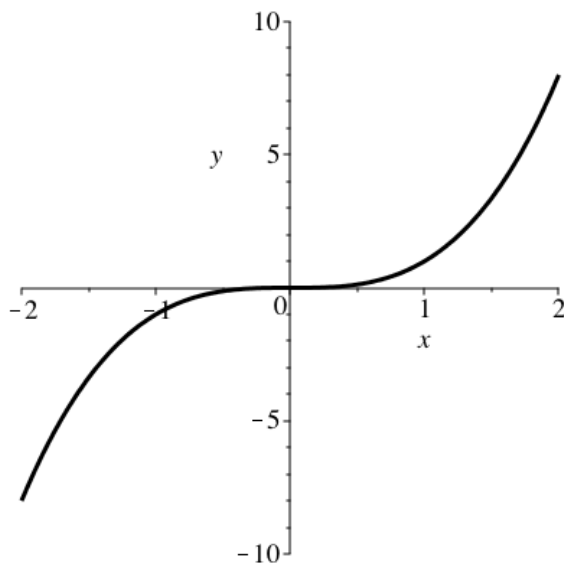
On the topic of using the second derivative to probe the behaviour of a function, we define one more type of point.

Definition 3.11.5
inflection point

For a function $f(x)$ defined on domain D , $a \in D$ is called an inflection point if both of the following are true:

1. $f''(a) = 0$ or $f''(a)$ does not exist and
2. $f''(a)$ changes signs at $x = a$.

A common example of a function with an inflection point is $f(x) = x^3$. We have $f''(0) = 0$ while $f''(x)$ is negative to the left of $x = 0$ and positive to the right of $x = 0$.



Observe that the graph of $f(x) = x^3$ changes from concave down to concave up at the inflection point at $x = 0$.

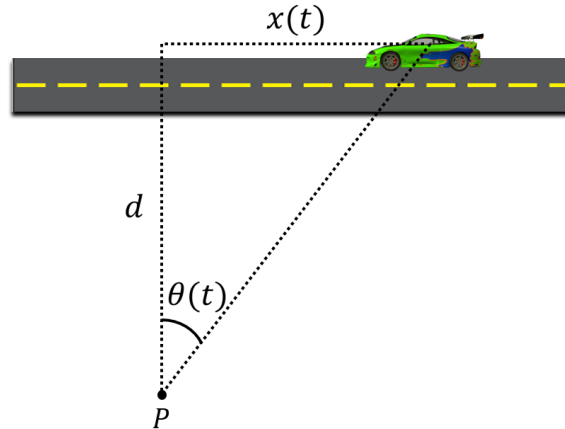
EXERCISE

Let $f(x) = x^3 - x^2$.

1. Locate the critical points of $f(x)$.
2. Use the second derivative test to identify all local extrema of $f(x)$.
3. Locate all points of inflection of $f(x)$.

3.12 Related Rates

If you've ever watched traffic drive along a road from a short distance away, you've probably noticed the closer a car is, the faster it *appears* to be moving. In the instant that it passes directly in front of you, it seems to zoom by. Of course, you know that the cars passing by are (for the most part) driving at a uniform speeds. So what is behind this effect?



When you are watching a car drive by, you are not directly observing its speed. Instead, you are observing the angular rate at which you need to turn your head or rotate your eyeballs to keep your focus on the car. Mathematically, we can say you are measuring $\frac{d\theta}{dt}$ where θ is the angle the car makes with some reference line. For simplicity, let's make this reference perpendicular to the road and passing through your position which we'll label P . We'll also suppose that you are a perpendicular distance d from the road and the car's position along the road is described by a function $x(t)$.

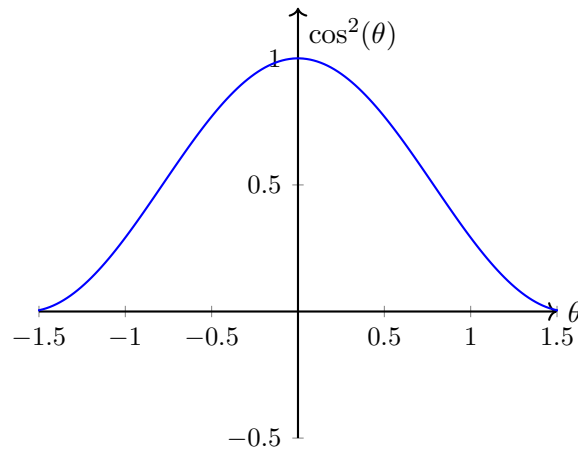
With a bit of trigonometry, we see that

$$\tan(\theta(t)) = \frac{x(t)}{d}$$

Differentiating this implicitly and using the Chain Rule, we get

$$\sec^2(\theta) \frac{d\theta}{dt} = \frac{1}{d} \frac{dx}{dt} \quad \implies \quad \frac{d\theta}{dt} = \frac{\cos^2(\theta)}{d} \frac{dx}{dt}$$

Now we can see what gives rise to the effect we described previously. If a car is driving at a constant speed, then $\frac{dx}{dt}$ is a constant. The distance d is also a constant. However, the angular rate of change, $\frac{d\theta}{dt}$, is moderated by a factor of $\cos^2(\theta)$. When the car is far away, we have $\theta = \pm\frac{\pi}{2}$. Looking at a plot of this term for the relevant angles, we see the angular speed ramps up as θ nears 0 and reaches its maximum there.



Problems like this one where we have two variables whose rates of change are related are, appropriately, called Related Rates problems. It is often the case that we cannot immediately write down a relationship between these rates. However, if we can write down a relationship between the variables, then we can use implicit differentiation and the chain rule to find that relationship.

Let's look at a few more related rates problems.

Example 18

A spherical balloon is inflated at a rate of $500 \text{ cm}^3/\text{s}$. At what rate is the radius increasing when the radius is 10 cm?

Solution: The volume and radius of a sphere are related by the equation $V = \frac{4}{3}\pi r^3$. Treating V and r as functions of time and differentiating this expression with respect to time, we get

$$\frac{dV}{dt} = 4\pi r^2 \frac{dr}{dt} \quad \implies \quad \frac{dr}{dt} = \frac{1}{4\pi r^2} \frac{dV}{dt}$$

When the radius is equal to 10 cm, we get

$$\frac{dr}{dt} = \frac{1}{4\pi(10 \text{ cm})^2} (500 \text{ cm}^3/\text{s}) \approx 0.4 \text{ cm/s}$$

Therefore, at the instant the balloon has a radius of 10 cm, the radius is increasing at a rate of approximately 0.4 cm/s.

EXERCISE

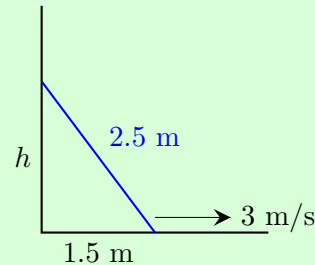
A cylindrical piece of clay is squashed such that its height decreases at a constant rate of 3 mm/s. The piece of clay remains cylindrical and maintains a constant volume. If at a particular instant, the height is 150 mm and the radius 50 mm, then determine the rate at which the radius must be increasing at that instant.

Example 19

A 2.5 m tall ladder is placed on a horizontal floor, leaning against a vertical wall, when tragedy strikes! The ladder slides down the wall, maintaining contact with the wall and the floor at all times. When the base of the ladder is 1.5 m from the wall, the bottom is sliding

away from the wall at 3 m/s. At that time, how fast is the top of the ladder sliding down the wall?

Solution: First, we exercise our artistic freedom and draw a picture.



Let t denote time, $x(t)$ denote the distance from the base of the wall to the base of the ladder at time t , and $h(t)$ denote the height above the floor of the point of contact of the ladder with the wall at time t .

Since the floor is horizontal and the wall is vertical, we know that

$$x(t)^2 + h(t)^2 = 2.5^2. \quad (3.1)$$

Let t_0 be the time at which $x(t_0) = 1.5$. Then from the question we know $x'(t_0) = 3$ (notice the sign is positive since x is increasing as the ladder falls). Here we are denoting x' to mean the derivative of x with respect to t .

Our goal is to compute $h'(t_0)$. Implicitly differentiating Equation 3.1 with respect to t and evaluating at t_0 gives

$$2x(t_0)x'(t_0) + 2h(t_0)h'(t_0) = 0. \quad (3.2)$$

By exploiting the Pythagorean theorem we have $h(t_0) = 2$. Putting everything we know into Equation 3.2 yields

$$2(1.5)(3) + 2(2)h'(t_0) = 0.$$

Solving this equation gives $h'(t_0) = -\frac{9}{4} = -2.25$. Note that the negative sign is consistent with the expectation (and the wording in the question!) that the ladder is sliding down the wall (and not up).

Therefore, the ladder is sliding down the wall at 2.25 m/s.

EXERCISE

A person with height 1.80 m walks towards a street lamp. The light on the street lamp is at a height of 3.10 m. At what rate is the height of the person's shadow changing when the person is 2.20 m from the street lamp?

EXERCISE

A circuit with two resistors is initially operating such that the two resistors have resistances $R_1 = 12 \Omega$ and $R_2 = 18 \Omega$. A change in temperature causes the resistances of the resistors

to increases such that

$$\frac{dR_1}{dt} = 0.1 \text{ } \Omega/\text{s} \quad \text{and} \quad \frac{dR_2}{dt} = 0.2 \text{ } \Omega/\text{s}$$

How does the total resistance of the circuit change if the resistors are connected in series? What if they are connected in parallel?

EXERCISE

The human ear is sensitive to a power levels which range over many orders of magnitudes. For example, the buzzing of a mosquito emits approximately 10^{-10} W of sound energy while a jackhammer emits approximately 1 W. As such, we use a logarithmic scale in units of decibels (dB) to describe “loudness”, L .

The loudness of a sound is related to power via the equation

$$L = 10 \log_{10} \left(\frac{P}{P_0} \right)$$

where $P_0 = 10^{-12}$ W is a constant.

Suppose the volume of a speaker is increased such that the power output increases at a known rate $\frac{dP}{dt}$. Find an expression describing how the rate at which the loudness increases in terms of this rate as well as the current power.

Chapter 4

The Mean Value Theorem

4.1 The Mean Value Theorem

Suppose you are making a trip to Waterloo from Toronto and have arranged for a taxi to drive you there. Your trip involves a 110 km stretch of Highway 401 with a speed limit of 100 km/hr. The taxi pulls onto the highway at 8:00 p.m. and exits an hour later at 9:00 p.m.. Is it possible that your driver obeyed the speed limit throughout the trip? Take a moment and think about it.

If they'd driven at a constant speed, their speed would have been 110 km/hr. Our intuition tells us that, if at any point, they'd driven slower than 110 km/hr, then they must have driven faster at some other point (or points) to make up for this. So, their speed somewhere along the way must have been at least 110 km/hr.

We'll revisit this problem in a little while and see how to prove this result rigorously, but fundamental to that proof will be the Mean Value Theorem. Before we can state that though, we have a bit of work to do.

4.1.1 Fermat's Theorem

Now that we've studied extreme values, the following statement should not feel far-fetched.

Theorem 1 (Fermat's theorem)

If a function f has a local maximum or minimum at $x = c$ and $f'(c)$ exists, then $f'(c) = 0$.

There are many theorems named after Fermat, so this one is also sometimes called the **interior extremum theorem**. This theorem is closely connected to the first and second derivative tests which told us that we should only look for local maximum and minimum points where a function has a critical point. We only have critical points if $f'(x) = 0$ or if $f'(x)$ does not exist. So, if we ignore the case where $f'(x)$ does not exist, then we are only getting critical points and, in turn, possibly local maximum or minimum points when $f'(x) = 0$.

Proof: We break the proof into two cases.

Case (i): Suppose f has a local maximum at $x = c$ and $f'(c)$ exists.

By the definition of a local maximum, there exists an open interval I containing c such that for all $x \in I$, $f(x) \leq f(c)$.

Let $c + h \in I$ with $h \neq 0$, then $f(c + h) - f(c) \leq 0$.

If we take $h > 0$, then

$$\frac{f(c+h) - f(c)}{h} \leq 0 \quad \implies \quad \lim_{h \rightarrow 0^+} \frac{f(c+h) - f(c)}{h} \leq 0$$

Since we assume $f'(c)$ exists, then this implies $f'(c) \leq 0$.

Similarly, if we take $h < 0$, then

$$\frac{f(c+h) - f(c)}{h} \geq 0 \quad \implies \quad \lim_{h \rightarrow 0^-} \frac{f(c+h) - f(c)}{h} \geq 0$$

Again, assuming $f'(c)$ exists, then this implies $f'(c) \geq 0$.

Combining these two results, we have $0 \leq f'(c) \leq 0$ which can only be true if $f'(c) = 0$.

Case (ii): This can be shown using the same approach as in case (i).

Therefore, if f has a local maximum or minimum at $x = c$ and $f'(c)$ exists, then $f'(c) = 0$. □

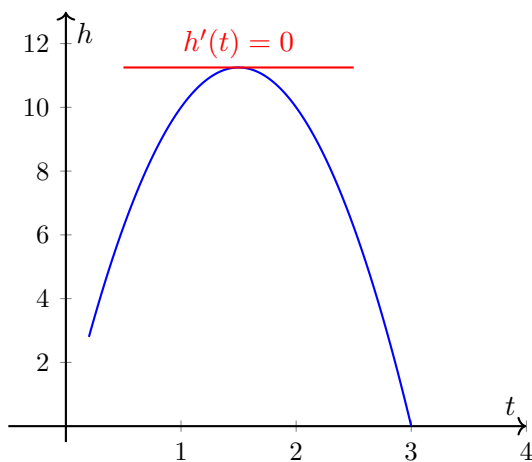
EXERCISE

Is the converse of Fermat's theorem true?

4.1.2 Rolle's Theorem

Let's consider another familiar, real-world problem. If you throw a ball directly up into the air and catch it a moment later, is there an instant at which the speed of the ball is zero? There is! As the ball goes up, gravity works to slow it down and there is an instant in time when the ball reaches its highest point that it stops for an instant and then proceeds to make the return trip back to your hand.

To be more precise, let $h(t)$ be the height of the ball above your hand at time t . Ignoring air resistance, $h(t)$ will be a quadratic function like $h(t) = v_0t - \frac{1}{2}gt^2$ where v_0 is the initial speed of the ball and g is acceleration due to gravity.



At the point where the ball reaches its maximum height, a tangent to the graph is horizontal. Therefore, the derivative of $h(t)$, which is the velocity of the ball, is zero. We can be certain of this thanks to Fermat's theorem!

Now, replace the ball with a bird that you gently launch upwards from your hand, it flies around for a while, and then it returns to your hand. Its path will be a bit more complicated now, but we could still ask the question "Will there be at least one instant in time while the bird is flying around where its velocity in the vertical direction is zero?"

Again, the answer is 'yes'. We can be certain of this because the bird will attain at least one local maximum in its height. By Fermat's theorem, we know that $h'(t) = 0$ when that happens.

Rephrasing this problem more abstractly leads to Rolle's theorem.

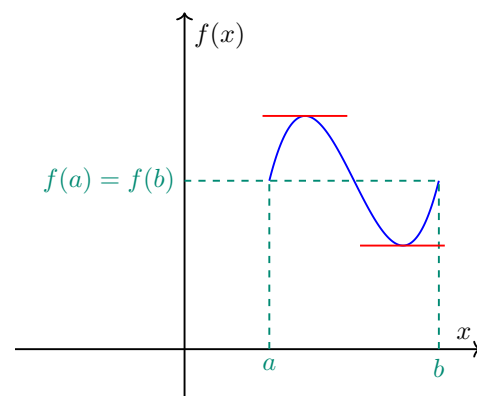
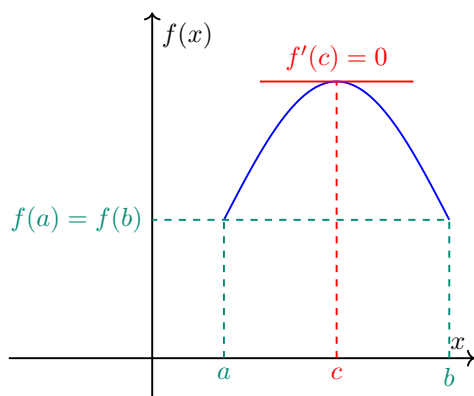
Theorem 2

If f is a function satisfying:

- f is continuous on the interval $[a, b]$,
- f is differentiable on the interval (a, b) , and
- $f(a) = f(b)$,

then there exists $c \in (a, b)$ such that $f'(c) = 0$.

To help digest this statement, consider the following graphs of two functions satisfying the conditions of Rolle's theorem.



Unless f is a constant function (in which case $f'(c) = 0$ everywhere on (a, b)) f takes values *not* equal to the values at the endpoints on (a, b) . Since f is continuous on $[a, b]$, it must have at least one local extremum by the Extreme Value Theorem. Moreover, since f is differentiable on (a, b) , it must be the case that $f'(x) = 0$ at any extrema by Fermat's theorem. Let's take these ideas to prove Rolle's theorem more carefully.

Proof: Let f be continuous on $[a, b]$, differentiable on (a, b) , and satisfy $f(a) = f(b)$. We proceed now by considering three cases:

1. $f(x)$ is a constant function on $[a, b]$
2. There exists $d \in (a, b)$ such that $f(d) > f(a)$
3. There exists $d \in (a, b)$ such that $f(d) < f(a)$

(Note, a function could belong to both cases 2 and 3.)

Case 1. If $f(x)$ is a constant function on $[a, b]$, then $f'(x) = 0$ for all $x \in (a, b)$. As such, we can take c to be anywhere on (a, b) and have $f'(c) = 0$.

Case 2. Since $f(x)$ is continuous on $[a, b]$, then by the Extreme Value Theorem, $f(x)$ must have a global maximum on $[a, b]$. Since there exists $d \in (a, b)$ such that $f(d) > f(a)$, then the global maximum cannot be at either $x = a$ or $x = b$. It must instead be on the interval (a, b) . Let the global maximum be at $c \in (a, b)$, then $f'(c) = 0$ by Fermat's theorem.

Case 3. This case proceeds similarly to case 2.

Therefore, there must exist $c \in (a, b)$ such that $f'(c) = 0$. □

EXERCISE

At the beginning of this section there is a discussion about a bird flying away from your hand and eventually landing back in your hands.

Rolle's theorem guarantees that there is some point where the vertical velocity of the bird is zero. Can you see how Rolle's theorem does this? What assumptions are placed on the movement of the bird for this to be true? Are they reasonable assumptions?

4.1.3 Mean Value Theorem

We are now ready to formally state and prove the Mean Value Theorem.

Theorem 3 (Mean Value Theorem)

If f is a function satisfying:

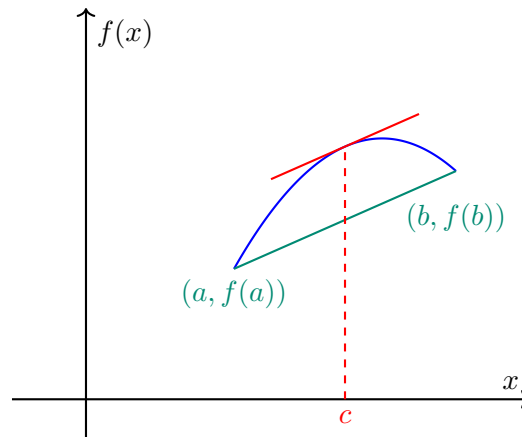
- f is continuous on the interval $[a, b]$ and

- f is differentiable on the interval (a, b) ,

then there exists $c \in (a, b)$ such that

$$f'(c) = \frac{f(b) - f(a)}{b - a}$$

The quotient $\frac{f(b) - f(a)}{b - a}$ gives the slope of the line connecting the points $(a, f(a))$ and $(b, f(b))$.



With this in mind, we see that the Mean Value Theorem is just saying that at some point $c \in (a, b)$, the instantaneous rate of change of $f'(c)$ will be equal to the average (or mean) change in f over the entire interval $[a, b]$. Geometrically, this is equivalent to saying there will exist at least one tangent line to the graph of f on (a, b) whose slope is equal to the slope of the line connecting the endpoints of f on $[a, b]$.

EXERCISE

Sketch out the graph of a function that is continuous on $[0, 1]$, differentiable on $(0, 1)$, and satisfies that there are exactly two different values of $c \in (0, 1)$ so that $f'(c) = f(1) - f(0)$.

Observe that Rolle's theorem is a special case of the Mean Value Theorem. In particular, the first two hypotheses of Rolle's theorem and of the Mean Value Theorem are identical. If we then take the conclusion of the Mean Value Theorem with the additional hypothesis $f(a) = f(b)$, it reduces to $f'(c) = 0$. We will use this connection to Rolle's theorem to prove the Mean Value Theorem by cleverly defining a new function for which Rolle's theorem can be applied and the result mapped back to our original function to prove the general statement.

Proof: Let f be continuous on $[a, b]$ and differentiable on (a, b) .

We define g to be the linear function with $g(a) = f(a)$ and $g(b) = f(b)$. Geometrically, $y = g(x)$ is the line passing through the points $(a, f(a))$ and $(b, f(b))$. Using the point-slope equation of a line, we have

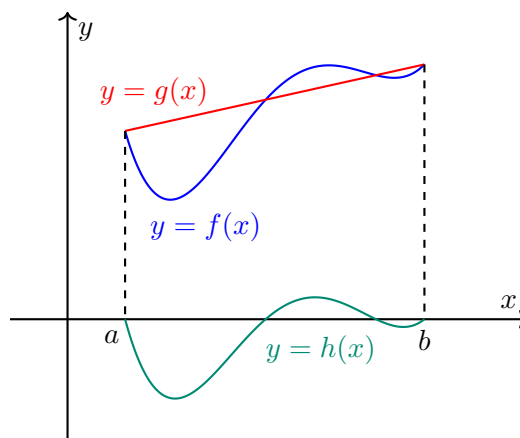
$$g(x) = f(a) + \frac{f(b) - f(a)}{b - a}(x - a)$$

Since $g(x)$ is a linear function, it is continuous and differentiable everywhere. We also note that

$$g'(x) = \frac{f(b) - f(a)}{b - a}$$

We next define h as the difference between f and g . That is,

$$h(x) = f(x) - g(x) = f(x) - f(a) - \frac{f(b) - f(a)}{b - a}(x - a)$$



Observe the following about the function h .

1. Since f and g are both continuous on $[a, b]$, then h is continuous on $[a, b]$.
2. Since f and g are both differentiable on (a, b) , then h is differentiable on (a, b) .
3. By definition of g , we have $g(a) = f(a)$ and $g(b) = f(b)$ and therefore $h(a) = 0$ and $h(b) = 0$.

Therefore, the function h satisfies the conditions for applying Rolle's theorem and we can conclude that there exists $c \in (a, b)$ such that $h'(c) = 0$.

It follows that

$$f'(c) = g'(c) \quad \implies \quad f'(c) = \frac{f(b) - f(a)}{b - a}$$

as required. □

We can now prove that the taxi driver at the beginning of this section was speeding!

Example 1

A car makes a 110 km trip in 1 hour on a highway with a speed limit of 100 km/hr. Prove that the car must have exceeded the speed limit at some point during the trip.

Solution:

Let $x(t)$ be the distance travelled by the car as a function of time with x measured in kilometres and t measured in hours. Let $t = 0$ be the start of the trip. We make the physically reasonable assumptions that x is continuous on $[0, 1]$ and differentiable on $(0, 1)$.

We are given that $x(0) = 0$ and $x(1) = 110$. Therefore, by the Mean Value Theorem, there exists a time $c \in (0, 1)$ such that

$$x'(c) = \frac{x(1) - x(0)}{1 - 0} = 110$$

The derivative of the position is the speed of the car and will have units km/hr. As such, we can conclude that there must have been a time during the trip when the car was driving at 110 km/hr thereby exceeding the speed limit.

EXERCISE

Let $f(x) = x^2 - \sin(\pi x)$. Prove that there exists $c \in (0, 1)$ such that $f'(c) = 1$.

(Once you've done that, see if you can figure out what c is. You should get an equation that you can't solve algebraically, but to which you can probably guess the solution.)

EXERCISE

Suppose $f(x)$ satisfies $f(0) = 5$ and $1 \leq f'(x) \leq 2$ for all $x \in \mathbb{R}$. Use the Mean Value Theorem to argue $7 \leq f(2) \leq 9$.

4.1.4 Extreme Values Revisited

In the previous section on extreme values, we stated three results without proof. We will restate and prove these in this section, but we first need just a couple more ingredients.

Definition 4.1.1

increasing/
decreasing

We say a function f is **increasing** on an interval I if for all $a, b \in I$

$$a < b \quad \implies \quad f(a) < f(b)$$

We say a function f is **decreasing** on an interval I if for all $a, b \in I$

$$a < b \quad \implies \quad f(a) > f(b)$$

With this definition, we can state the following fact.

Fact 4 Let f be differentiable on an interval I .

- i. If $f'(x) > 0$ on I , then f is increasing on I .
- ii. If $f'(x) < 0$ on I , then f is decreasing on I .

Proof: First, let us assume that $f'(x) > 0$ on an interval I . Since f is differentiable, we can apply the Mean Value Theorem using any two points $a, b \in I$ to get that there exists $c \in (a, b)$ such that

$$f'(c) = \frac{f(b) - f(a)}{b - a} \quad \implies \quad f(b) - f(a) = f'(c)(b - a)$$

We know $f'(c) > 0$. If we further assume that $b > a$, it follows that

$$f(b) - f(a) > 0 \quad \implies \quad f(b) > f(a)$$

Since this holds for any $a, b \in I$, then f is increasing on I . This proves part (i). The proof of part (ii) is similarly done. \square

The previous fact is something we know to be true in our hearts, but we need the Mean Value Theorem to prove! Sometimes the most seemingly obvious statements require a serious amount of mathematics to prove.

First Derivative Test

The First Derivative Test provided a tool for classifying a critical point as a local minimum or local maximum when the sign of the first derivative changed at that critical point.

Theorem 5 (First Derivative Test)

Let $f(x)$ be continuous on interval D and suppose $c \in D$ is a critical point of f , then:

- If there exists an interval $I = (a, b)$ with $c \in I$ such that $f'(x) < 0$ for $x \in (a, c)$ and $f'(x) > 0$ for $x \in (c, b)$, then f has a local minimum at $x = c$.
- If there exists an interval $I = (a, b)$ with $c \in I$ such that $f'(x) > 0$ for $x \in (a, c)$ and $f'(x) < 0$ for $x \in (c, b)$, then f has a local maximum at $x = c$.

Proof: Let f be continuous on an interval $I = (a, b)$ and let $c \in I$ be a critical point of f . Furthermore, suppose $f'(x) < 0$ on (a, c) and $f'(x) > 0$ on (c, b) .

Since $f'(x) < 0$ on (a, c) , then f is decreasing on (a, c) . Therefore, $f(c) < f(x)$ for all $x \in (a, c)$.

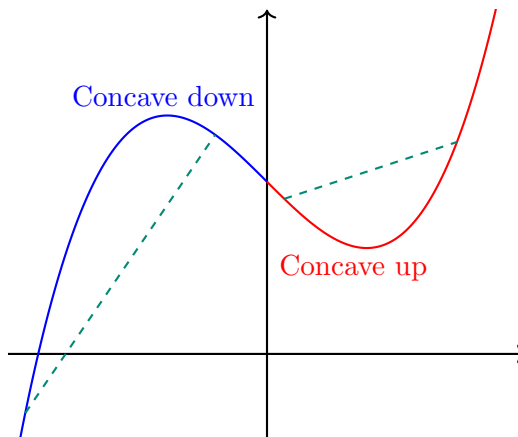
Similarly, since $f'(x) > 0$ on (c, b) , then f is increasing on (c, b) . Therefore, $f(c) < f(x)$ for all $x \in (c, b)$.

It follows that $f(c) \leq f(x)$ for all $x \in I$ and is therefore a local minimum of f .

The proof for the case where f has a local maximum at the critical point follows similarly. \square

Concavity

The definition of concavity says that for a function f defined on an interval I , if for every pair of points $a, b \in I$ the secant line connecting $(a, f(a))$ with $(b, f(b))$ lies above (below) the graph of f , then we say f is concave up (down).



We claimed that when $f''(x)$ exists on the interval I , we can identify concave up and concave down based on the sign of $f''(x)$.

Fact 6

Given a function $f(x)$, if the second derivative of f exists over some interval I , then the definition of concavity simplifies to:

- If $f''(x) > 0$ over some interval, then f is concave up on that interval.
- If $f''(x) < 0$ over some interval, then f is concave down on that interval.

We are now ready to prove this fact.

Proof: We first assume that $f''(x) > 0$ on an interval I . This means that $f'(x)$ is increasing on I .

Let $a, b \in I$ and, without loss of generality, suppose $a < b$.

Define g to describe the secant line connecting $(a, f(a))$ and $(b, f(b))$. That is,

$$g(x) = f(a) + \frac{f(b) - f(a)}{b - a}(x - a).$$

We will show that for any $d \in (a, b)$, $f(d) < g(d)$ or, in other words, the secant line lies above the function.

We note for later that

$$g(d) = f(a) + \frac{f(b) - f(a)}{b - a}(d - a).$$

By the Mean Value Theorem, there exists $c_1 \in (a, d)$ such that

$$f'(c_1) = \frac{f(d) - f(a)}{d - a}.$$

Similarly, there exists $c_2 \in (d, b)$ such that

$$f'(c_2) = \frac{f(b) - f(d)}{b - d}.$$

Since $f'(x)$ is increasing on I and $c_1 < c_2$, we have

$$f'(c_1) < f'(c_2) \quad \implies \quad \frac{f(d) - f(a)}{d - a} < \frac{f(b) - f(d)}{b - d}$$

Solving this inequality for $f(d)$ we get

$$\begin{aligned} f(d) &< \frac{f(b)(d - a) - f(a)(d - b)}{b - a} \\ &= \frac{f(b)(d - a) - f(a)(d - a + a - b)}{b - a} \\ &= \frac{f(b)(d - a) - f(a)(d - a)}{b - a} + f(a) \\ &= f(a) + \frac{f(b) - f(a)}{b - a}(d - a) \\ &= g(d) \end{aligned}$$

Since $f(d) < g(d)$, the secant line lies above the graph of the function.

A similar proof addresses the case when $f''(x) < 0$. □

Second Derivative Test

After identifying the critical points of a function, assuming $f''(x)$ exists at each critical point, we can use the Second Derivative Test to determine which critical points are local minimums or local maximums.

Theorem 7 (Second Derivative Test)

Suppose f has a critical point at $x = c$ (i.e., $f'(c) = 0$) and $f''(c)$ exists, then:

- If $f''(c) > 0$, then f has a local minimum at $x = c$.
- If $f''(c) < 0$, then f has a local maximum at $x = c$.

The proof that this test works also involves a nice application of the Mean Value Theorem.

Proof: Let $f''(x)$ exist on an interval I containing the critical point $x = c$ of f . This means that $f'(c) = 0$.

Now we consider the case when $f''(c) > 0$ on I . This means that $f'(x)$ is increasing on the interval I .

Let $a \in I$ with $a < c$. By the Mean Value Theorem, there exists $d_1 \in (a, c)$ such that

$$f'(d_1) = \frac{f(c) - f(a)}{c - a} \quad \implies \quad f(c) = f(a) + f'(d_1)(c - a)$$

Since $f'(x)$ is increasing on I and $d_1 < c$, then $f'(d_1) < f'(c)$. Building on this, since $(c - a) > 0$, we get

$$f(c) = f(a) + f'(d_1)(c - a) < f(a) + f'(c)(c - a) = f(a)$$

The last equality follows from the fact that $f'(c) = 0$. This establishes that $f(c)$ is less than all values of $f(x)$ to the left of $x = c$ on the interval I .

Next, we consider $b \in I$ with $b > c$. By the Mean Value Theorem, we know that there exists $d_2 \in (c, b)$ such that

$$f'(d_2) = \frac{f(b) - f(c)}{b - c} \implies f(c) = f(b) - f'(d_2)(b - c)$$

Since $f'(x)$ is increasing on I and $c < d_2$, then $f'(c) < f'(d_2)$. Since $(b - c) > 0$, we then have

$$f(c) = f(b) - f'(d_2)(b - c) < f(b) - f'(c)(b - c) = f(b)$$

This establishes that $f(c)$ is less than all values of $f(x)$ to the right of $x = c$ on the interval I .

Since $f(c)$ is less than all other values of f on the interval I , then it must be a local minimum of f .

A similar argument works for the case where $f''(x) < 0$ on I . □

4.2 Indeterminate Forms

Recall, we previously worked out the following limit which we referred to as the fundamental trigonometric limit

$$\lim_{x \rightarrow 0} \frac{\sin(x)}{x} = 1$$

Since $x = 0$ is not in the domain of $\frac{\sin(x)}{x}$, we could not simply substitute $x = 0$ into the numerator and denominator. This would give a ratio that looks like $\frac{0}{0}$ and it's not clear how to interpret this. In fact, it is so unclear that we have a special name for ratios like this; we call them **indeterminate forms**.

To actually evaluate the limit above, we relied on a lengthy geometric argument. In this section, we will develop a method to more efficiently compute limits that yield indeterminate forms like the one above.

4.2.1 L'Hôpital's rule

The technique we'll develop is called l'Hôpital's rule named after the mathematician Guillaume de l'Hôpital. To motivate the rule, let's consider a limit that behaves like the one at the start of this section - that is, suppose we want to compute the following limit.

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)}$$

but $f(a) = 0$ and $g(a) = 0$. Direct substitution would give the indeterminate form $\frac{0}{0}$, so we need to try something else.

In the neighbourhood of $x = a$, each of $f(x)$ and $g(x)$ should be well approximated by their linearizations (so, immediately, we need that f and g are differentiable). Moreover, in the limit as $x \rightarrow a$, we should be able to replace each function with its linearization and the error in doing so should go to zero. Mathematically, this idea yields

$$\begin{aligned} \lim_{x \rightarrow a} \frac{f(x)}{g(x)} &= \lim_{x \rightarrow a} \frac{f(a) + f'(a)(x - a)}{g(a) + g'(a)(x - a)} \\ &= \lim_{x \rightarrow a} \frac{f'(a)(x - a)}{g'(a)(x - a)} \\ &= \lim_{x \rightarrow a} \frac{f'(a)}{g'(a)} \\ &= \frac{f'(a)}{g'(a)} \end{aligned}$$

We make use of the fact that $f(a) = 0$ and $g(a) = 0$ to get rid of these terms from the first line to the second line.

One thing we need to be careful about here is that $g'(a) \neq 0$. One way we could handle this is by using a higher-degree approximation of both $f(x)$ and $g(x)$ than the linearization so that we are left with non-zero polynomials in the numerator and denominator, but a cleaner approach is to rewrite the final expression as the limit of the quotient $\frac{f'(x)}{g'(x)}$. In doing so, l'Hôpital's rule can be used iteratively, as needed, to circumvent this potential issue.

A complete proof of l'Hôpital's rule requires a generalization of the mean value theorem (called Cauchy's mean value theorem) and is beyond the scope of this course, but the idea presented above holds for all of the functions we'll be working with. As we'll see in the precise statement of the rule below, we can also adapt the derivation above to handle indeterminate forms of the type $\frac{\pm\infty}{\pm\infty}$.

Theorem 8 (l'Hôpital's Rule)

Let I be an open interval containing the point $x = a$. Further, let f and g be differentiable on I (except possibly at $x = a$). If either

- $\lim_{x \rightarrow a} f(x) = 0$ and $\lim_{x \rightarrow a} g(x) = 0$ or
- $\lim_{x \rightarrow a} f(x) = \pm\infty$ and $\lim_{x \rightarrow a} g(x) = \pm\infty$, then

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = \lim_{x \rightarrow a} \frac{f'(x)}{g'(x)}.$$

In other words, if the limit of a quotient function gives an indeterminate form, then you can differentiate the numerator and denominator and try computing the limit again. L'Hôpital's rule can be applied repeatedly if needed.

Fact 9 L'Hôpital's rule can be applied to one-sided limits.

Let's try out l'Hôpital's rule on what is sometimes called the *fundamental log limit*.

Example 2 Let's evaluate $\lim_{x \rightarrow \infty} \frac{\ln(x)}{x}$.

As $x \rightarrow \infty$, both $\ln(x) \rightarrow \infty$ and $x \rightarrow \infty$, so we have an indeterminate form. Since both x and $\ln(x)$ are differentiable everywhere, we can apply l'Hôpital's rule. We have

$$\lim_{x \rightarrow \infty} \frac{\ln(x)}{x} = \lim_{x \rightarrow \infty} \frac{\frac{1}{x}}{1} = 0.$$

Looking at the graphs of $\ln(x)$ and x , this is perhaps believable, as x grows to $+\infty$ much faster than $\ln(x)$ does.

REMARK

Let's look at the motivating example, $\lim_{x \rightarrow 0} \frac{\sin(x)}{x}$. Substituting in $x = 0$ gives us $\frac{0}{0}$, so we should be in the clear to use l'Hôpital's rule. Well, not quite. Let's see what happens if we do:

$$\lim_{x \rightarrow 0} \frac{\sin(x)}{x} = \lim_{x \rightarrow 0} \frac{\cos(x)}{1} = 1.$$

Great, right?

Unfortunately not. In the process of using l'Hôpital's rule, we used the fact that $\frac{d}{dx} \sin(x) = \cos(x)$. However, the proof of that fact relies on knowing $\lim_{x \rightarrow 0} \frac{\sin(x)}{x} = 1$.

There is some serious circular reasoning going on here! Although the section on l'Hôpital's rule was motivated by trying to compute this limit, you cannot actually use l'Hôpital's rule to do so! The lengthy geometric argument from Section 3.6.1 is the only way we have computed this particular limit in these notes.

EXERCISE

Can you use l'Hôpital's rule to evaluate

$$\lim_{x \rightarrow 0} \frac{\cos(x) - 1}{x} \quad ?$$

Why or why not?

EXERCISE

Use l'Hôpital's rule to evaluate

$$\lim_{x \rightarrow 2} \frac{x - 2}{x^2 - 4}.$$

EXERCISE

Why can't you use l'Hôpital's rule to evaluate $\lim_{x \rightarrow 0} \frac{\cos(x)}{x^2}$? What is the limit equal to? Just for fun, see what happens if you *do* differentiate the numerator and denominator and then take the limit. (You should get the wrong answer if you try to use l'Hôpital's rule.)

4.2.2 Indeterminate Products

We can use l'Hôpital's rule to evaluate some limits even if they don't, at first, appear to yield the indeterminate forms $\frac{0}{0}$ or $\frac{\pm\infty}{\pm\infty}$. Consider the following limit.

$$\lim_{x \rightarrow \infty} x e^{-x}$$

As x tends to infinity, we have a product consisting of one term which grows without bound and another which tends to zero. It is not clear if such a limit should give 0, some finite non-zero number, or tend to infinity. When we have something that behaves in this way - that is, something which upon direct substitution has the form $\infty \cdot 0$ - we call it an indeterminate product.

There is a nice way of dealing with indeterminate products. We simply rewrite it as a quotient function by moving one of the product functions into a denominator. Let's try this out.

Example 3 Evaluate $\lim_{x \rightarrow \infty} x e^{-x}$.

Solution: We note that $e^{-x} = \frac{1}{e^x}$ so that the limit can be rewritten as

$$\lim_{x \rightarrow \infty} x e^{-x} = \lim_{x \rightarrow \infty} \frac{x}{e^x}$$

Observe that direct substitution now gives the indeterminate product $\frac{\infty}{\infty}$. To this, we can apply l'Hôpital's rule.

$$\begin{aligned} \lim_{x \rightarrow \infty} \frac{x}{e^x} &= \lim_{x \rightarrow \infty} \frac{\frac{d}{dx}(x)}{\frac{d}{dx}(e^x)} \\ &= \lim_{x \rightarrow \infty} \frac{1}{e^x} \\ &= \lim_{x \rightarrow \infty} e^{-x} \\ &= 0 \end{aligned}$$

REMARK

In the previous example, we could have also tried to rewrite the limit as

$$\lim_{x \rightarrow \infty} x e^{-x} = \lim_{x \rightarrow \infty} \frac{e^{-x}}{\frac{1}{x}}$$

This would even give us the indeterminate form $\frac{0}{0}$ to which we could apply l'Hôpital's rule. Unfortunately, applying the rule in this form only makes things worse. Each application yields another indeterminate form. For example, applying a few iterations iteratively would give

$$\begin{aligned}\lim_{x \rightarrow \infty} \frac{e^{-x}}{\frac{1}{x}} &= \lim_{x \rightarrow \infty} \frac{e^{-x}}{\frac{1}{x^2}} \\ &= \lim_{x \rightarrow \infty} \frac{e^{-x}}{\frac{2}{x^3}} \\ &= \lim_{x \rightarrow \infty} \frac{e^{-x}}{\frac{6}{x^4}} \\ &= \dots\end{aligned}$$

There is no end in sight. If you get an indeterminate product, move part of your limit argument to the denominator, and then this happens to you, it is worth checking to see if moving the *other* part of your limit argument to the denominator yields better results.

EXERCISE

Evaluate $\lim_{x \rightarrow 0^+} x \ln(x)$.

4.2.3 Indeterminate Differences

Suppose you have a limit of the form

$$\lim_{x \rightarrow a} (f(x) - g(x))$$

but both $f(x)$ and $g(x)$ tend to infinity in this limit. This gives a limit that looks like $\infty - \infty$. We call this an indeterminate difference. Sometimes it is possible to convert an indeterminate difference into a standard indeterminate form.

Example 4

Evaluate $\lim_{x \rightarrow 0^+} \left(\frac{1}{x} - \frac{1}{\ln(x+1)} \right)$.

Solution: Observe that $\lim_{x \rightarrow 0^+} \frac{1}{x} = \infty$ and $\lim_{x \rightarrow 0^+} \frac{1}{\ln(x+1)} = \infty$, so direct substitution gives an indeterminate difference.

Let's manipulate the argument into a quotient function.

$$\lim_{x \rightarrow 0^+} \left(\frac{1}{x} - \frac{1}{\ln(x+1)} \right) = \lim_{x \rightarrow 0^+} \frac{\ln(x+1) - x}{x \ln(x+1)}$$

Direct substitution in this form gives the indeterminate form $\frac{0}{0}$, so l'Hôpital's rule can be applied.

$$\begin{aligned} \lim_{x \rightarrow 0^+} \left(\frac{1}{x} - \frac{1}{\ln(x+1)} \right) &= \lim_{x \rightarrow 0^+} \frac{\ln(x+1) - x}{x \ln(x+1)} \rightarrow \frac{0}{0}, \quad \text{apply LH rule} \\ &= \lim_{x \rightarrow 0^+} \frac{\frac{1}{x+1} - 1}{\ln(x+1) + \frac{x}{x+1}} \\ &= \lim_{x \rightarrow 0^+} \frac{-x}{(x+1)\ln(x+1) + x} \rightarrow \frac{0}{0}, \quad \text{apply LH rule} \\ &= \lim_{x \rightarrow 0^+} \frac{-1}{\ln(x+1) + 2} \\ &= -\frac{1}{2} \end{aligned}$$

EXERCISE

Evaluate $\lim_{x \rightarrow 0^+} \left(\frac{1}{\sin(x)} - \frac{1}{x^2} \right)$.

4.2.4 Indeterminate Powers

The last type of indeterminate forms we'll consider can occur when we have a limit like

$$\lim_{x \rightarrow a} f(x)^{g(x)}$$

If direct substitution gives either 0^0 , ∞^0 , or 1^∞ , then we have an indeterminate power. In all three cases, the trick we'll use is the same. In particular, we take advantage of continuity and instead consider the limit of the logarithm of the argument and, afterwards, reverse this operation by exponentiating our result.

Example 5

Evaluate $\lim_{x \rightarrow 0^+} x^x$.

Solution: Direct substitution gives the indeterminate power 0^0 , so we let $y = x^x$ and consider the limit as $x \rightarrow 0^+$ of $\ln(y) = x \ln(x)$.

$$\begin{aligned} \lim_{x \rightarrow 0^+} \ln(y) &= \lim_{x \rightarrow 0^+} x \ln(x) \rightarrow 0 \cdot (-\infty), \quad \text{indeterminate power} \\ &= \lim_{x \rightarrow 0^+} \frac{\ln(x)}{\frac{1}{x}} \rightarrow \frac{-\infty}{\infty}, \quad \text{apply LH rule} \\ &= \lim_{x \rightarrow 0^+} \frac{1}{-\frac{1}{x^2}} \\ &= \lim_{x \rightarrow 0^+} -x \\ &= 0 \end{aligned}$$

Now we use this result to evaluate the original limit.

$$\begin{aligned}\lim_{x \rightarrow 0^+} x^x &= \lim_{x \rightarrow 0^+} e^{x \ln(x)} \\ &= e^{\lim_{x \rightarrow 0^+} x \ln(x)} \\ &= e^0 \\ &= 1\end{aligned}$$

Therefore, $\lim_{x \rightarrow 0^+} x^x = 1$.

EXERCISE

Evaluate $\lim_{x \rightarrow \infty} x^{\frac{1}{x}}$.

Chapter 5

Extra Topics

5.1 Optimization

Now that we have developed techniques for locating and classifying the extrema of a function, we can apply these ideas in context to find optimal solutions to real-world problems. The basic idea is that we will construct a function which describes something we'd like to maximize or minimize in terms of a parameter and then use calculus to determine the optimal value of that parameter. We call this process optimization.

Depending on the context of the problem we're trying to solve, the domain of the function representing the quantity we're looking to optimize may be a closed interval or may be an open interval. When it is a closed interval, we know that global extrema will exist by the Extreme Value Theorem. Moreover, we can use the Closed Interval Method to find these extrema. On the other hand, if we have an open interval or we just want to identify all local extreme values, then we can turn to the First and Second Derivative Tests.

Let's look at a simple example to compare these two scenarios.

Example 1

Let P be the point $(p, 0)$ where p can be varied (i.e., P can slide along the x -axis) and let Q be the point $(5, 4)$.

- i. Determine the minimum distance between P and Q if $p \in [1, 2]$.
- ii. Determine the minimum distance between P and Q if $p \in \mathbb{R}$.

Solution:

- i. Let $f(p)$ be the distance between the points P and Q as a function of p . Then,

$$f(p) = \sqrt{(5-p)^2 + 4^2}$$

This function is continuous on \mathbb{R} , so it is continuous on $[1, 2]$. Therefore, we can apply the Closed Interval Method to find the global minimum of f on the given interval.

To do this, we first compute the derivative of f with respect to p . We denote this derivative $f'(p)$.

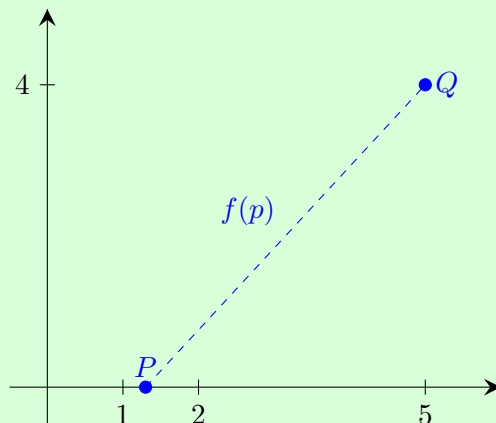
$$f'(p) = \frac{p-5}{\sqrt{(5-p)^2 + 4^2}}$$

Observe that $f'(p)$ exists and is not equal to zero for all $p \in (1, 2)$. Therefore, the extreme values of f on $[1, 2]$ will be attained at the interval endpoints. These are

$$f(1) = 4\sqrt{2} \quad \text{and} \quad f(2) = 5$$

Since $f(2) < f(1)$, the minimum distance between P and Q is 5 and occurs when P is at $(2, 0)$.

A sketch reveals that this result is not at all surprising.



The closer we can get the point P to sitting directly below the point Q , the smaller the distance will be between the two points. But since P can only slide to the right as far as $(2, 0)$, this configuration will give the minimum distance.

- ii. Based on the diagram we drew for the first part, we now know what to expect as our answer. That is, the minimum distance will be attained when P is at $(5, 0)$ but suppose we didn't have a nice diagram to help us solve the problem. In that case, we might proceed as follows.

Let f be defined as in part (i.) but now $p \in \mathbb{R}$. Let's locate and classify the critical points of f .

Using the expression from part (i.), observe that $f'(p)$ exists for all $p \in \mathbb{R}$, so critical points can only possibly occur when the derivative is zero.

$$f'(p) = 0 \quad \implies \quad p = 5$$

There is a single critical point (right where we expected) at $p = 5$. To classify it, let's apply the First Derivative Test. To do this, we recall that a continuous function can only change between increasing and decreasing at a critical point. Therefore, we need only check if f is increasing or decreasing at a single convenient point in each of the intervals $(-\infty, 5)$ and $(5, \infty)$. Let's check $p = 2$ and $p = 8$.

$$f'(2) = -\frac{3}{5} \quad \text{and} \quad f'(8) = \frac{3}{5}$$

Since $f'(2) < 0$, then f must be decreasing on $(-\infty, 5)$. Similarly, since $f'(8) > 0$, then f must be increasing on $(5, \infty)$. Therefore, by the First Derivative Test, $f(5) = 4$ is a local minimum of f . We can further infer based on the known intervals of increasing and decreasing that $f(5)$ is the global minimum of f .

(Note, we could have also used the Second Derivative Test here to show $f''(5) > 0$ and therefore $f(5)$ is local minimum.)

REMARKS

- The First and Second Derivative Tests are designed to classify critical points as being local maxima or local minima (or neither). A bit more justification is therefore required to further show that such an extreme value may be a global extrema.
- Derivatives of distance formulas like $f(p) = \sqrt{(5-p)^2 + 4^2}$ can get cumbersome because of the square root function. However, if we're doing an optimization problem, we can usually work with the square of the distance instead. This is because distance is maximized (minimized) when distance-squared is maximized (minimized).

EXERCISE

Find the global maximum and minimum distances from the curve $y = x^2$ with $x \in [0, 2]$ to the point $(0, 2)$.

EXERCISE

What is the minimum vertical distance between the parabolas $y = x^2 + 1$ and $y = -x^2 + 2x - 1$?

These geometrical examples help us establish the basic mechanics of optimization. Let's look next at applying these mechanics to some application problems.

Example 2

A ball is thrown with speed v_0 at an angle θ measured relative to the horizontal. The only force acting on the ball is gravity. If the ball lands at the same height from which it is thrown, the horizontal distance it covers before first hitting the ground is given as a function of θ by $d(\theta) = \frac{v_0^2}{g} \sin(2\theta)$ where g is acceleration due to gravity. What is the optimal angle for throwing the ball to maximize the horizontal distance covered?

Solution: First, we note that a reasonable domain for this problem is $\theta \in [0, \frac{\pi}{2}]$. Throwing the ball with a negative angle would have it impact the ground immediately and throwing the ball with an angle greater than $\frac{\pi}{2}$ would have it go backwards.

Since we have a closed interval, we will apply the Closed Interval Method. This requires us to calculate the derivative of d with respect to θ .

$$d'(\theta) = \frac{2v_0^2}{g} \cos(2\theta)$$

This derivative function is defined for all $\theta \in [0, \frac{\pi}{2}]$. It is, however, equal to zero on this interval at $\theta = \frac{\pi}{4}$. Therefore, this is the only critical point in our interval.

Next, we compute d at the critical point and at the interval endpoints.

$$d(0) = 0, \quad d\left(\frac{\pi}{4}\right) = \frac{2v_0^2}{g}, \quad d\left(\frac{\pi}{2}\right) = 0$$

Since $d(0) = d\left(\frac{\pi}{2}\right) < d\left(\frac{\pi}{4}\right)$, the optimal angle to throw the ball is $\frac{\pi}{4}$ (or 45°).

EXERCISE

The yield y of a food crop depends on the nitrogen content n of the soil (in appropriate units) according to the equation

$$y(n) = \frac{3n}{1 + 2n^2}$$

Determine the optimal nitrogen level for maximizing the yield.

In many optimization problems, the quantity we're looking to optimize depends on more than one variable. However, if we have additional constraints on these variables, we can often simplify the problem so that we ultimately only need to find extreme values of a function of one variable.

Example 3

Determine real numbers x and y such that the sum of x and y is 10 and the product of x and y is a maximum.

Solution: At first glance, it appears we need to find the maximum of the product xy which depends on both x and y . However, since x and y need to satisfy the constraint $x + y = 10$, then there is really only one degree of freedom in this problem. This is because anywhere that y appears, we can replace it with $10 - x$.

We use this observation to define the following function to optimize.

$$f(x) = x(10 - x) = 10x - x^2$$

Now, let's find the critical points of f . We have

$$f'(x) = 10 - 2x$$

The derivative function is defined for all $x \in \mathbb{R}$, so we solve $f'(x) = 0$ to find any critical points. This yields a single critical point at $x = 5$.

We now apply the Second Derivative Test to classify this point. Observe that $f''(x) = -2$ so f is concave down everywhere. This means that the critical point at $x = 5$ is a maximum. Moreover, since f is continuous and has no other critical points, this is also a global maximum. Therefore, when $x = y = 5$, the product $xy = 25$ is a global maximum with the given constraint that $x + y = 10$.

This problem has a nice geometrical interpretation. If we imagine a rectangle with side lengths x and y , then saying $x + y = 10$ fixes the perimeter of the rectangle as 20. Meanwhile, the product xy represents the area of the rectangle. So, if we want a rectangle with fixed perimeter to have the maximum possible area, then that rectangle should be a square.

EXERCISE

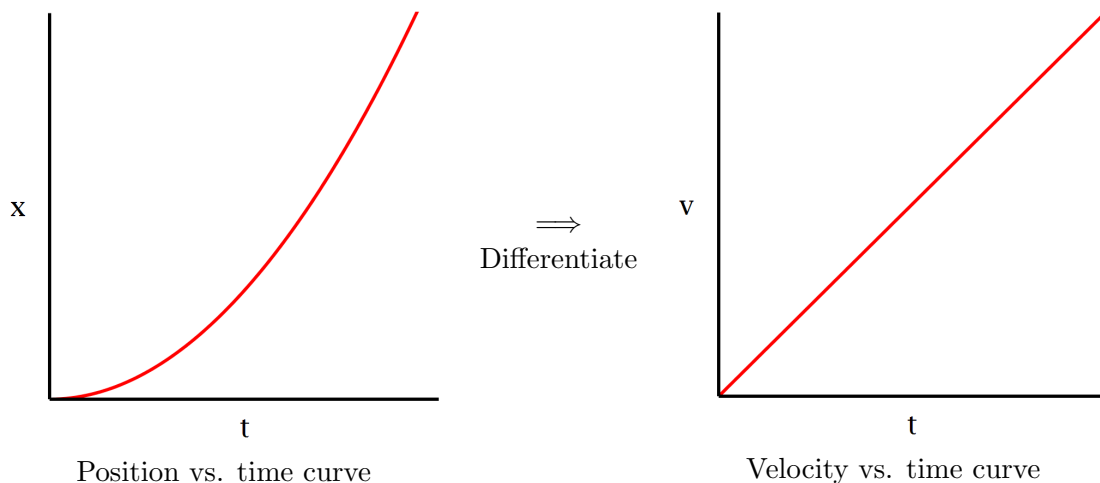
A company wishes to manufacture food cans that are cylindrical in shape (including a circular disk for the top and bottom) with a volume of 400 cm^3 . Find the radius and height for such cans that minimize the surface area of the can.

5.2 Antidifferentiation

Suppose we are given a function $x(t)$ describing the one-dimensional position of an object in terms of time. By differentiating $x(t)$ we get the velocity of the object. That is, in general

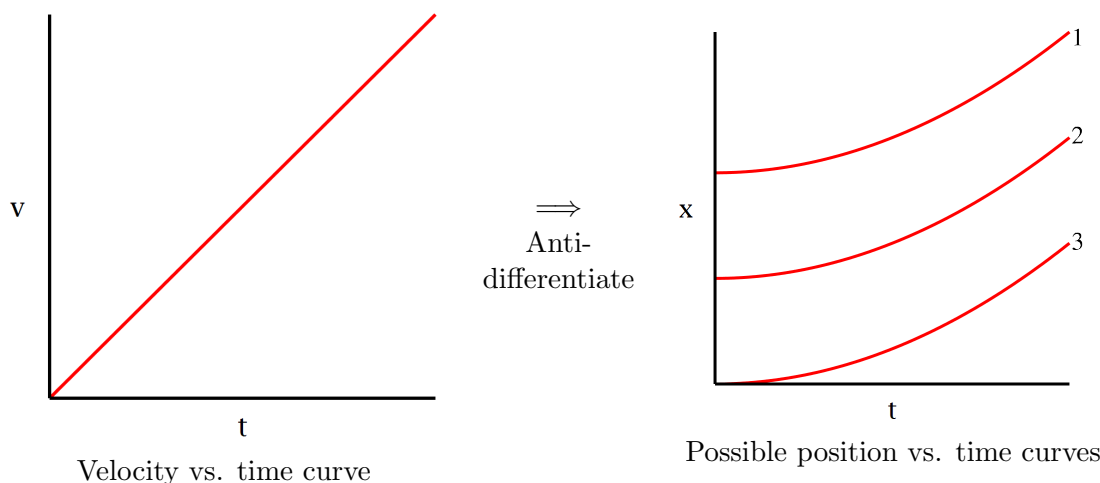
$$v(t) = \frac{dx}{dt}$$

As a concrete example of this relationship, consider an object undergoing constant acceleration, a , and having zero initial velocity. In this case, we would have $x(t) = \frac{1}{2}at^2$ and differentiating this gives $v(t) = at$.



Now, suppose we are given the velocity $v(t)$ of the object as a starting point. Is it possible to determine the position? In other words, could we go from the graph on the right to the graph on the left?

The answer is “almost”. The derivative of a function tells us how that function changes from point-to-point which is another way of saying it tells us the shape of the function. Having the shape of the function correct for all points in the domain is great. It would tell us how to graph the function and situate it in the x -direction, but it wouldn't tell us how to situate it in the y -direction. That is, the function would only be determined up to a vertical shift.



We call this “reversing differentiation” process **antidifferentiation** and we call the family of functions we get as a result of this process **antiderivatives**. In general, the antiderivatives of a function are identical up to an additive constant.

Definition 5.2.1
antiderivative

A function $F(x)$ is an antiderivative of $f(x)$ if $F'(x) = f(x)$.

Fact 1

If $F_1(x)$ and $F_2(x)$ are both antiderivatives of $f(x)$, then $F_2(x) = F_1(x) + C$ for some constant C .

EXERCISE

Use the Mean Value Theorem to prove the previous fact. (Hint: Define a helper function $h(x) = F_1(x) - F_2(x)$.)

Due to the previous fact, when finding the antiderivatives of a function, it is sufficient to find a single antiderivative $F(x)$ and then write the entire family of antiderivatives as $F(x) + C$ for an arbitrary constant C .

Example 4

Determine the antiderivative of $f(x) = x$.

Solution: We need to find a function $F(x)$ such that $F'(x) = x$. Our experience with the Power Rule suggests the antiderivative should be proportional to x^2 . By differentiating x^2 , we can figure out the multiplicative constant we need to get a derivative of just x .

$$\frac{d}{dx}(x^2) = 2x \quad \implies \quad \frac{d}{dx}\left(\frac{1}{2}x^2\right) = x$$

Therefore, the general antiderivative is $F(x) = \frac{1}{2}x^2 + C$ where C is an arbitrary constant.

EXERCISE

Determine the antiderivative of $f(x) = x^{5/2}$.

EXERCISE

Determine the antiderivative of $g(x) = \sin(x)$.

The limit definition of the derivative is the foundation for computing derivatives. It allowed us to determine the derivatives of elementary functions as well as establish differentiation rules for dealing with products, compositions, and so on. In combination, we get a reliable process for computing derivatives of arbitrary functions by “turning a crank”.

In contrast, we do not have a definition from which we can build a process for computing antiderivatives. Instead, we must rely on our knowledge of differentiation to guide us.

For example, what is the antiderivative of $\frac{1}{x}$ when x is positive? After scratching our chins for a minute, we remember that the derivative of $\ln(x)$ is $\frac{1}{x}$ and so the antiderivative of $\frac{1}{x}$ with $x > 0$ is $\ln(x) + C$ for an arbitrary constant C .

Similarly, what is the antiderivative of $3 + 2x$? We know that the derivative of $3x$ is 3 and the derivative of x^2 is $2x$. We also know that derivatives obey a sum formula which could be used here to say that the derivative of $3x + x^2$ is $3 + 2x$. It follows that the antiderivative of $3 + 2x$ is $3x + x^2 + C$ for an arbitrary constant C .

There are at least a couple of lessons here. The first is that we should be mindful of how we can exploit differentiation rules to help us determine antiderivatives. The second is that it will probably be useful to summarize our knowledge of the derivatives of elementary functions in the form of an antiderivative ~~cheat sheet~~ reference table. So, here it is:

Function $f(x)$	Antiderivative $F(x)$
$x^n, \quad n \neq -1$	$\frac{x^{n+1}}{n+1} + C$
$\frac{1}{x}$	$\ln(x) + C$
$\sin(x)$	$-\cos(x) + C$
$\cos(x)$	$\sin(x) + C$
$\sec^2(x)$	$\tan(x) + C$
e^x	$e^x + C$
$b^x, \quad b > 0, b \neq 1$	$\frac{b^x}{\ln(b)} + C$
$\frac{1}{\sqrt{1-x^2}}$	$\arcsin(x) + C$
$\frac{1}{1+x^2}$	$\arctan(x) + C$

REMARK

The reader with a keen eye may notice something a little odd in the second row of the reference table. The antiderivative of $\frac{1}{x}$ is listed as $\ln(|x|) + C$. Why the absolute values?

Well, the first thing to keep in mind is that $\ln(x)$ is not defined if $x < 0$. The second is that

$$\frac{d}{dx} \ln(x) = \frac{1}{x}$$

is only valid for $x > 0$.

So, the question becomes what $F(x)$ has the property that $F'(x) = \frac{1}{x}$ when $x < 0$?

When $x < 0$, applying the chain rule gives that $\ln(-x) = \frac{1}{x}$. Therefore, an antiderivative $F(x)$ of $\frac{1}{x}$ is given by

$$F(x) = \begin{cases} \ln(x) & \text{if } x > 0 \\ \ln(-x) & \text{if } x < 0. \end{cases}$$

We can succinctly write $F(x) = \ln(|x|)$.

Let's look now at how we might find the antiderivative of a function which is not on this list.

Example 5

Find the general antiderivative of $f(x) = x^2e^{x^3}$.

Solution: The antiderivative is not immediately clear. However, we do know that when we differentiate an exponential function, the same exponential function will appear in the derivative. More precisely, by the Chain Rule, we have

$$\frac{d}{dx} \left(e^{g(x)} \right) = e^{g(x)} \frac{dg}{dx}$$

This suggests that our antiderivative will include the function e^{x^3} . Inspecting the derivative of this function, we find

$$\frac{d}{dx} \left(e^{x^3} \right) = 3x^2e^{x^3} \quad \implies \quad \frac{d}{dx} \left(\frac{1}{3}e^{x^3} \right) = x^2e^{x^3} = f(x)$$

Our idea worked. We found an antiderivative and from this we can say that the general antiderivative of $f(x)$ is

$$F(x) = \frac{1}{3}e^{x^3} + C$$

for an arbitrary constant C .

EXERCISE

Find the general antiderivative of $\frac{x}{1+x^4}$. (Hint: What function in the table above most closely resembles this one?)

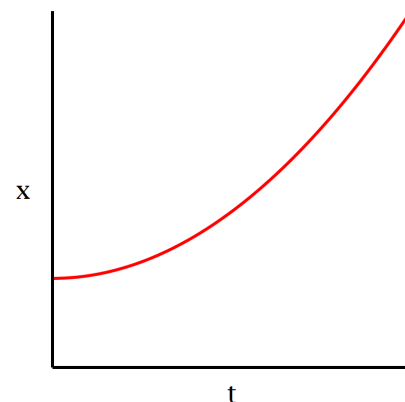
Our general antiderivatives always include an additive arbitrary constant (i.e., a “+C”) to indicate that it is actually a family of solutions. However, in certain contexts we might be able to assign a value to the arbitrary constant. Let's revisit the example of an object undergoing constant acceleration to explore the role this arbitrary constant plays in a bit more detail.

Consider again an object moving in one dimension with velocity $v(t) = at$. Given this velocity, we find the family of antiderivatives describing the position:

$$x(t) = \frac{1}{2}at^2 + C$$

If we set $t = 0$, this equation gives

$$x(0) = C$$



So, if in addition to specifying the velocity of the object, we also specify its location at $t = 0$, we can determine the position function uniquely.

In fact, it is sufficient to specify the location of the object at any time. More generally, if $x(t_0) = x_0$, then

$$x(t_0) = \frac{1}{2}at_0^2 + C = x_0 \quad \implies \quad C = x_0 - \frac{1}{2}at_0^2$$

This makes sense when we recall that a family of antiderivative functions are identical up to a vertical shift. No two such functions intersect, so if we know of just a single point that the graph of our function must contain, this determines our antiderivative function uniquely as the one whose graph passes through that point.

Note, to obtain a unique antiderivative, we are not restricted to dealing with problems of motion. We just need to give an appropriate constraint on the antiderivative we're seeking as in the next example.

Example 6

Find the unique function $f(x)$ such that $f'(x) = x^{2/3}$ and $f(1) = 0$.

Solution: Antidifferentiating $f'(x)$, we have $f(x) = \frac{3}{5}x^{5/3} + C$. Imposing the condition that $f(1) = 0$ gives

$$\frac{3}{5}(1)^{5/3} + C = 0 \quad \implies \quad C = -\frac{3}{5}$$

Therefore, $f(x) = \frac{3}{5}(x^{5/3} - 1)$.

EXERCISE

Determine $g(x)$ if $g'(x) = \frac{1}{2} \cos\left(\frac{x}{2}\right)$ and $g(\pi) = 2$.

5.2.1 Higher-Order Antiderivatives

In certain circumstances, we may need to compute higher-order antiderivatives. This is done iteratively. For example, suppose we are given $f''(x) = 2$ and asked to determine $f(x)$.

Antidifferentiating once gives $f'(x) = 2x + C_1$ for an arbitrary constant C_1 . But now we must antidifferentiate again with this arbitrary constant in place. Doing so, we arrive at $f(x) = x^2 + C_1x + C_2$ where C_2 is another arbitrary constant independent of C_1 .

Observe that, in general, each time we antidifferentiate, we will introduce a new arbitrary constant. Each arbitrary constant represents a degree of freedom in our family of antiderivatives. If we're seeking a unique antiderivative, then we will need one independent piece of information about $f(x)$ per constant.

Continuing the example above, suppose we also know that $f(0) = -1$ and $f(1) = 0$. The first condition allows us to specify C_2 .

$$f(0) = 0^2 + C_1 \cdot 0 + C_2 = -1 \quad \implies \quad C_2 = -1$$

The second condition along with now knowing C_2 lets us determine C_1 .

$$f(1) = 1^2 + C_1 \cdot 1 - 1 = 0 \quad \implies \quad C_1 = 0$$

This uniquely fixes the function as $f(x) = x^2 - 1$.

We often call these constraints that allow us to eliminate degrees of freedom in our solution *initial conditions*. In general, the initial conditions will yield a system of equations for the arbitrary constants. Provided this system of equations yields a unique solution - that is, there are as many linearly independent equations as there are constants - we will get a uniquely specified antiderivative.

Example 7

A projectile is thrown from a hot-air balloon at an altitude of 1,800 m, with an initial velocity of 10 m/s downwards. If the projectile undergoes constant acceleration of -10 m/s², how long does it take to reach the ground?

Solution: We are given an initial velocity, but we are not given the velocity function, so we cannot just antidifferentiate once to get a position function. However, we *are* told that the acceleration is constant and we recall that acceleration is the derivative of velocity. So, we should be able to antidifferentiate the constant acceleration once to get velocity and antidifferentiate again to get position. Let's do this in general terms and plug in the known quantities later.

Let $x(t)$ be the height above the ground, $v(t)$ be the velocity (with upwards velocity being positive), and $a(t)$ be the accelerating (with upwards acceleration being positive). Furthermore, denote the initial position $x(0) = x_0$, the initial velocity $v(0) = v_0$, and the constant acceleration a_0 .

If $v'(t) = a(t) = a_0$, then we can antidifferentiate once to get $v(t) = a_0t + C_1$. Imposing the condition that $v(0) = v_0$ yields $C_1 = v_0$, so

$$v(t) = a_0t + v_0$$

Next, since $x'(t) = v(t)$, antidifferentiating again gives $x(t) = \frac{1}{2}a_0t^2 + v_0t + C_2$. Imposing the condition $x(0) = x_0$ gives $C_2 = x_0$, so

$$x(t) = \frac{1}{2}a_0t^2 + v_0t + x_0$$

(If you spend much time doing physics, these formulas for $x(t)$ and $v(t)$ are no doubt very familiar to you. What's great about this approach is that it can be generalized for when $a(t)$ is not a constant function.)

In this problem, we are given $x_0 = 1,800$ m, $v_0 = -10$ m/s, and $a_0 = -10$ m/s². Then, with units omitted, the projectile has position

$$x(t) = -5t^2 - 10t + 1800$$

The particle reaches the ground when $x(t) = 0$. Solving this equation with the quadratic formula gives

$$t = \frac{-(-10) \pm \sqrt{(-10)^2 - 4(-5)(1800)}}{2(-5)} = -1 \pm 19$$

Since the projectile will hit the ground with $t > 0$, we take the positive solution. Therefore, restoring units, we conclude that the projectile reaches the ground 18 seconds after it is thrown.

EXERCISE

A particle undergoes periodic acceleration $a(t) = \sin(2t)$. If it has initial velocity $v(0) = 0$ and initial position $x(0) = 1$, then determine $x(t)$.

5.3 Modelling

Many real-world processes are modelled by equations relating quantities and their rates of change.

For example, suppose the force acting on a charged particle in an electric field depends on the particle's position in that field. According to Newton's Second Law, force is proportional to the rate of change of the rate of change of position. This yields an equation relating the position of a particle to its second derivative.

As another example from chemistry, the Law of Mass Action is invoked to describe the relationship between the abundance of reactants and the rate of the reactions in which they participate.

In this section, we will explore mathematical models described in terms of quantities and their rates of change. We call these models which describe an unknown function in terms of one or more of its derivatives **differential equations**. Since we have not yet covered integration, we will be limited in what differential equations we can solve, but as we'll see, we can still solve some rather interesting problems.

To begin, let's consider a purely mathematical example.

Example 8

The rate of change of $f(x)$ satisfies the equation $f'(x) = 3x^2 - 2x$ and $f(0) = 1$. Determine $f(x)$.

Solution: We can antidifferentiate the expression for $f'(x)$ to get $f(x) = x^3 - x^2 + C$ for an arbitrary constant C . Furthermore, since $f(0) = 1$, it must be the case that $C = 1$. Therefore, $f(x) = x^3 - x^2 + 1$.

EXERCISE

Find a function $g(x)$ satisfying $g''(x) = e^x$ with the conditions $g(0) = 0$ and $g'(0) = 0$.

Now let's look at some differential equations in context.

Example 9

Suppose the temperature, T , in a lake decreases linearly with depth z such that

$$\frac{dT}{dz} = -k$$

for some constant $k > 0$. If the surface temperature is 20°C and the temperature at a depth of 5 metres is 10°C , then determine the rate constant k and the depth at which this model predicts a temperature of 5°C .

Solution: We can antidifferentiate the given equation to get $T(z) = -kz + C$ for some constant C . Imposing the initial conditions $T(0) = 0$ and $T(5) = 10^\circ\text{C}$ gives $C = 20^\circ\text{C}$ and $k = 2^\circ\text{C}/\text{m}$. This gives (with units omitted)

$$T(z) = -2z + 20$$

To find the depth at which this model predicts a temperature of 5°C , we set $T(z) = 5$ and solve for z .

$$-2z + 20 = 5 \quad \implies \quad z = 7.5$$

Therefore, this model predicts the temperature will be 5°C at a depth of 7.5m.

(A bit of thought reveals that this model is not very realistic and/or would only work down to a depth of 10m at which point the water would start to freeze.)

EXERCISE

Suppose the rate at which a reservoir fills after a storm obeys the differential equation

$$\frac{dV}{dt} = \frac{k}{\sqrt{t}}$$

where V is the volume of water in the reservoir, t is time after the storm, and k is a rate constant. If the volume is zero at $t = 0$ and 8 million litres after 1 hour, determine the volume of water in the reservoir at $t = 2$ hours.

In the previous examples, we had a derivative of our function depending only on some function of the independent variable. In general, differential equations will relate the function and one or more of its derivatives. Solving such differential equations often requires more advanced techniques, but also allows us to model a much richer variety of phenomenon. We won't get into these advanced techniques here, but we can still make some headway without those techniques.

5.3.1 Exponential Processes

For many populations, the rate at which the population grows is proportional to the population size. For example, consider a well-resourced population of rabbits with no predators. The more rabbits there are, the faster the rabbit population will grow. Let's formulate this relationship in mathematical terms.

If the rate of change in time t of a population with size P is proportional to the population size, then P obeys the differential equation

$$\frac{dP}{dt} = kP(t)$$

for some constant k .

If, for a moment, we set $k = 1$, then this says P is a function which is equal to its own derivative. We know such a function - the natural exponential function, e^t . We can even

bring an arbitrary k back into the fold by recalling that the derivative of e^{kt} will give k times e^{kt} by the chain rule. So, the function e^{kt} is a solution to the differential equation. However, we also recognize that we could multiply e^{kt} by an arbitrary constant and it would still satisfy the differential equation. If we write this arbitrary constant as P_0 (for reasons that will become clear in a moment), then we have the following solution to the differential equation

$$P(t) = P_0 e^{kt}$$

This turns out to be the most general solution to this differential equation. Notice, also, that $P(0) = P_0$ so the constant P_0 represents the population at time zero (which is why this notation was chosen).

The differential equation above with k positive therefore describes **exponential growth**. There are also many processes which exhibit exponential growth that do not strictly involve populations. For example, the spread of a virus, inflation, participants in a pyramid scheme, nuclear fission, and more.

When the constant k is negative, the same differential equation describes **exponential decay**. In that case, the rate of change of the size of the system is proportional to the size of the system but negative. A few examples of systems that undergo exponential decay include nuclear decay, transistor size, and chemical reactant concentration.

Example 10

A sample of bacteria initially contains 1,000 cells/L. When observed 3 days later, the sample is found to contain 10,000 cells/L. Assuming the bacteria concentration follows an exponential growth model, what bacteria concentration should we expect after 7 days?

Solution: Let $C(t)$ be the concentration of bacteria. Since $C(t)$ follows the exponential growth model, we know $C(t) = C_0 e^{kt}$ for some constants C_0 and k .

We are given $C(0) = 1,000$, so $C_0 = 1,000$. We are also given $C(3) = 10,000$. We can use this to solve for k .

$$1000e^{3k} = 10000 \quad \implies \quad k = \frac{\ln(10)}{3}$$

Therefore,

$$C(t) = 1000e^{t \ln(10)/3}$$

After 7 days, the concentration would be $C(7) = 1000e^{7 \ln(10)/3} \approx 215443$ cells/L.

EXERCISE

The **half-life** of a system is the time it takes for the size of the system to halve. A particular isotope of Californium (Californium-250) has a half-life of approximately 13.1 years. A sample of Californium containing 50 mg of Californium-250 is stored in a container. How much of the isotope will remain in this sample 50 years later?

Exponential Decay Towards Equilibrium

Consider a beverage sitting out at room temperature. If the beverage is colder than its surroundings, then it will warm up over time. On the other hand, if it is hotter than its surroundings, then it will cool down. Less obvious, but something you've probably still noticed, is that the greater the difference in temperature between the beverage and the surrounding temperature, the faster its temperature changes. For example, at 90°C cup of tea will drop 10°C in temperature in much less time than a 40°C cup of tea would.

We can summarize this observation as “the rate of change of the beverage is proportional to the difference between its temperature and the surrounding temperature”. This idea is known as **Newton’s Law of Cooling**. If we let $T(t)$ be the temperature of the beverage as a function of time and let T_S be the constant surrounding temperature, then this gives the differential equation model

$$\frac{dT}{dt} = k(T(t) - T_S)$$

where k is a constant that depends on the thermal properties of the system. Observe that k will be negative here since either (i) the beverage is hotter than the surroundings ($T(t) > T_S$), so we need $k < 0$ to make $\frac{dT}{dt} < 0$ (i.e., the beverage cools down) or (ii) the beverage is colder than the surroundings ($T(t) < T_S$), so we need $k < 0$ to make $\frac{dT}{dt} > 0$ (i.e., the beverage warms up).

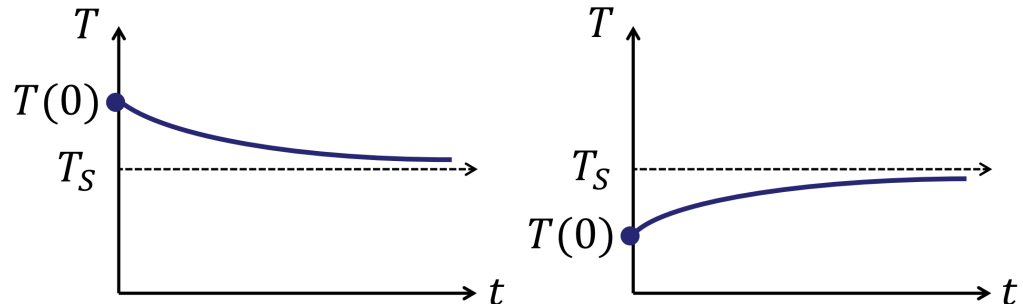
To solve this differential equation, let us introduce a new function $y(t) = T(t) - T_S$. Since T_S is a constant, then $\frac{dy}{dt} = \frac{dT}{dt}$ and the differential equation becomes

$$\frac{dy}{dt} = ky(t)$$

We know this has the general solution $y(t) = y_0 e^{kt}$. Moreover, $y_0 = y(0) = T(0) - T_S$ is the initial temperature difference. Rewriting the solution in terms of $T(t)$ gives

$$T(t) - T_S = (T(0) - T_S)e^{kt} \quad \implies \quad T(t) = (T(0) - T_S)e^{kt} + T_S$$

In this form, the solution $T(t)$ has a nice interpretation. The term $(T(0) - T_S)e^{kt}$ represents the temperature difference between the beverage and its surroundings as a function of time. Since $k < 0$, this difference decreases with time and eventually tends to zero. In that limit, only the term T_S remains on the right-hand side of the expression for $T(t)$ which is what we would expect - that is, the temperature of the beverage will eventually be the same as the temperature of the surroundings.



Many scientific models involve this type of exponential convergence to an equilibrium state. For example, the time it takes for a compressed elastic object to return to its original form, the concentration of proteins in biological systems under certain circumstances, or the voltage across a capacitor in an RC circuit.

Example 11

A skydiver steps out of a plane. Their speed in the vertical direction is initially zero. They accelerate due to gravity but due to air resistance, they will stop accelerating if they reach a speed called terminal velocity. At this moment, the gravitational force $F_g = mg$ is exactly balanced by the force of air resistance at terminal velocity $F_r = bv_t$. Here, m is the mass of the skydiver, g is acceleration due to gravity, b is a drag coefficient, and v_t is terminal velocity.

Suppose that the rate of change with respect to time t of their speed $v(t)$ is proportional to the difference between their speed and terminal velocity.

- i. Write out a differential equation to describe this statement.
- ii. Determine the solution to this differential equation in terms of any rate constants introduced in part (a) and terminal velocity v_t . Impose the condition $v(0) = 0$.
- iii. If the skydiver reaches a velocity of 5 m/s after 1 second of falling and $v_t = 50$ m/s, then how long does it take to reach a velocity of 45 m/s?

Solution:

- i. We have the differential equation.

$$\frac{dv}{dt} = k(v(t) - v_t)$$

for some constant k .

- ii. Using the same techniques we employed in the Newton's Law of Cooling example, we get

$$v(t) = (v_0 - v_t)e^{kt} + v_t$$

Since $v_0 = v(0) = 0$, this further simplifies to

$$v(t) = v_t(1 - e^{kt})$$

- iii. With $v_t = 50$ we have $v(t) = 50(1 - e^{kt})$. Imposing the condition that $v(1) = 5$ allows us to solve for k

$$50(1 - e^k) = 5 \quad \implies \quad k = \ln\left(\frac{9}{10}\right)$$

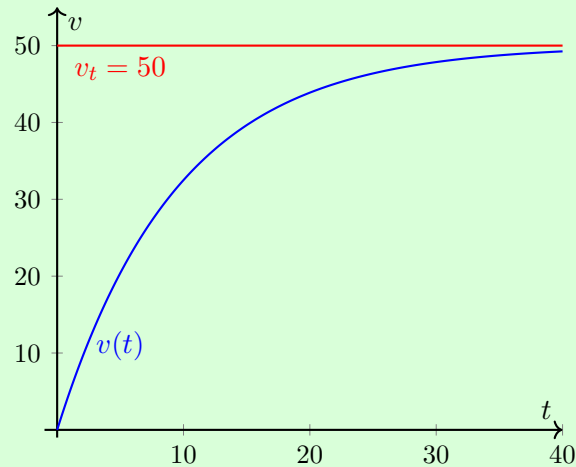
We can now set $v(t) = 45$ and solve for t to find when the velocity is 45 m/s.

$$50\left(1 - e^{\ln\left(\frac{9}{10}\right)t}\right) = 45 \quad \implies \quad t = \frac{\ln\left(\frac{1}{10}\right)}{\ln\left(\frac{9}{10}\right)} \approx 21.9$$

Therefore, it takes about 22 seconds to reach a velocity of 45 m/s.

(For comparison, without air resistance, the velocity would increase linearly with time. This means that if the velocity is 5 m/s after 1 second, then it would only take 9 seconds to reach 45 m/s.)

Here is a graph of $v(t)$, with the terminal velocity in red. Compare the graph to the graphs above relating to Newton's law of cooling.

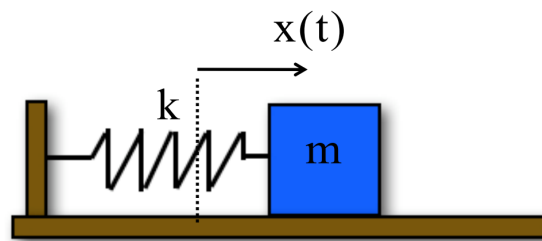


EXERCISE

A cold glass of water initially at 5°C is left outside where the temperature is 30°C . After 20 minutes, the water has warmed up to 15°C . The temperature of the water as a function of time obeys Newton's Law of Cooling. When will the temperature reach 25°C ?

5.3.2 Oscillatory Motion

Consider a block with mass m attached to a spring with spring constant k . The spring constant quantifies the "strength" of the spring and by Hooke's Law we know that if the block is displaced a distance x from equilibrium, the spring will exert a force $F_s = -kx$ on the block. If we ignore friction, then this is the only force acting on the block.



Newton's Second Law tells us that an object with mass m will accelerate when a non-zero net force, F_{net} is present and, moreover, that this acceleration a is determined by the equation

$$F_{\text{net}} = ma$$

For the mass on the spring above, the net force is $F_s = -kx$. We can also write the acceleration as the second time derivative of the displacement. This yields the following differential equation describing the displacement of the block as a function of time.

$$ma = -kx \quad \Longrightarrow \quad m \frac{d^2x}{dt^2} = -kx(t) \quad \Longrightarrow \quad \frac{d^2x}{dt^2} = -\frac{k}{m}x(t)$$

The quantities k and m are constants, so let's omit them for a moment. Doing so gives a simpler differential equation of the form

$$\frac{d^2y}{dt^2} = -y(t)$$

Do we know any functions which would satisfy this relation? That is, do we know any functions which when you differentiate them twice, you just get the function back with a minus sign? We do! The cosine and sine functions have this property.

$$\frac{d^2}{dt^2}(\cos(t)) = \frac{d}{dt}(-\sin(t)) = -\cos(t) \quad \text{and} \quad \frac{d^2}{dt^2}(\sin(t)) = \frac{d}{dt}(\cos(t)) = -\sin(t)$$

Since the derivative operator is linear, we can use this observation to construct the following solution to the simplified differential equation:

$$y(t) = c_1 \cos(t) + c_2 \sin(t)$$

for arbitrary constants c_1 and c_2 . This turns out to be the most general solution to this differential equation. (Our experience with antidifferentiation suggests that we should expect *two* arbitrary constants due to the presence of a second derivative in the differential equation.)

To adapt this solution to the differential equation for a mass on a spring, we note that if the argument of our trigonometric functions has a constant coefficient, then by the chain rule, we will pick up a factor of this constant each time we differentiate. We can then choose the constant appropriately to balance the left and right hand sides of the differential equation. Let's do that.

We will let $x(t) = c_1 \cos(\omega t) + c_2 \sin(\omega t)$ where ω is a constant. With a bit of work, we can show

$$\frac{d^2x}{dt^2} = -\omega^2 x(t)$$

If we set $\omega^2 = \frac{k}{m}$, then we reproduce the differential equation for $x(t)$. Therefore, the displacement of the mass on the spring is described by

$$x(t) = c_1 \cos\left(\sqrt{\frac{k}{m}}t\right) + c_2 \sin\left(\sqrt{\frac{k}{m}}t\right)$$

The quantity $\omega = \sqrt{\frac{k}{m}}$ determines how quickly $x(t)$ oscillates, so it is called the angular frequency. The time it takes for one oscillation is given by the period $T = \frac{2\pi}{\omega} = 2\pi\sqrt{\frac{m}{k}}$. Observe that the greater the mass, the longer the period. In contrast, the stronger the spring, the shorter the period.

The constants c_1 and c_2 allow us to use this solution to describe the displacement of the mass as a function of time regardless of how we set the mass into motion. For example, if we release the mass with initial position $x(0) = A$ and zero initial velocity $x'(0) = 0$, then you can show you'll get $c_1 = A$ and $c_2 = 0$ so $x(t) = A \cos\left(\sqrt{\frac{k}{m}}t\right)$.

EXERCISE

Show that $x(t) = c_1 \cos(\omega t) + c_2 \sin(\omega t)$ can be rewritten in the form $x(t) = A \cos(\omega t + \phi)$. The cosine angle addition formula should be useful for this: $\cos(a + b) = \cos(a)\cos(b) - \sin(a)\sin(b)$.

There are many real-world phenomena described by a differential equation of the form

$$\frac{d^2x}{dt^2} = -\omega^2 x(t)$$

This includes the swinging of a pendulum with small oscillations, the voltage in an LC circuit, molecular vibrations, acoustic waves in a pipe, and more. In other words, once we write down the differential equation for any of these problems, we can solve it by adapting the solution above.

EXERCISE

According to Kirchoff's voltage law, the charge Q in the capacitor of an LC circuit can be modelled as a function of time t by the differential equation

$$L \frac{d^2Q}{dt^2} + \frac{1}{C} Q(t) = 0$$

where L is the inductance and C is the capacitance. Find the general solution to this differential equation and an expression for the period of oscillations in the charge Q in terms of L and C .

This section on modelling only scrapes the surface with regards to how we can use calculus to analyze and describe the behaviour of real-world applications. To take the next steps though, we will need integration, better approximation techniques, calculus of multivariable vector functions, and more sophisticated methods for solving differential equations.

Part II

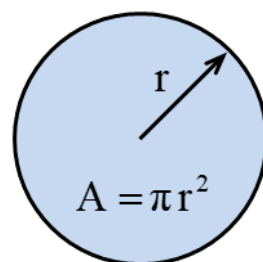
MATH 648

Chapter 6

Sequences and Series

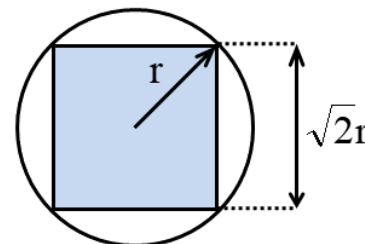
6.1 Sequences

We all know that the area A of a circle of radius r is proportional to the square of its radius. More precisely, $A = \pi r^2$. This fact dates back to at least the days of Archimedes (circa 250 BCE). But such an exact relationship cannot be determined by measurement alone. That is, if we construct a circle and can only measure its area and radius to a certain degree of precision, then we would not be able to rule out that they are related by some awful equation like $A = \pi r^{2.001}$.



To establish that the area of a circle is exactly proportional to the square of its radius, an abstract geometrical argument is needed.

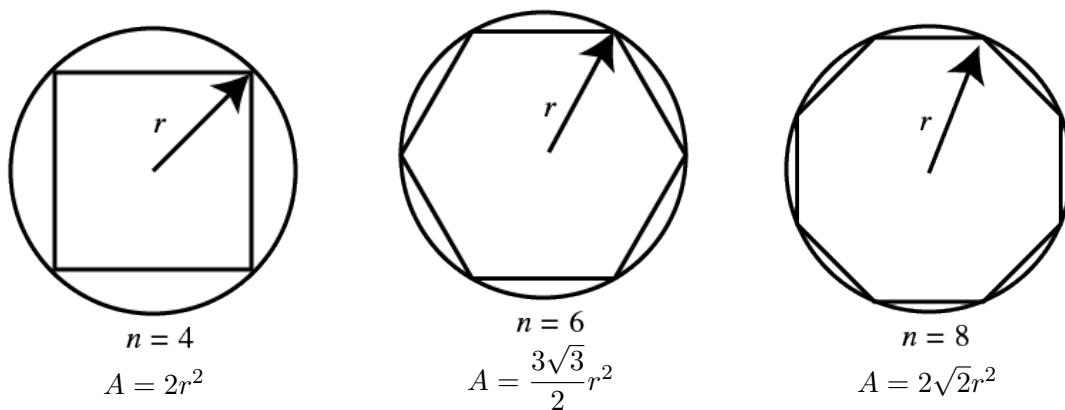
As a sample of how such an argument might work, consider a circle of radius r with a square inscribed in it. A bit of geometry gives the edge length of this square as $\sqrt{2}r$. As such, we know the square will have area $2r^2$. Therefore, we can infer that the area of the circle satisfies $A > 2r^2$.



EXERCISE

Inscribe a circle of radius r *inside* a square to determine an upper bound on the area of the circle.

Now, we repeat this construction but with polygons having an increasing numbers of sides.



We can generalize the process for finding the area of the inscribed polygons by dividing each regular n -gon into n congruent triangles. The area of each triangle will always be proportional to r^2 , so the area of the n -gon will always be proportional to r^2 . As the number of sides tends to infinity, this relation holds while at the same time, the area of the polygon tends to the area of the circle. Therefore, the area of the circle must be exactly proportional to the square of its radius.

The basis of the previous argument is known as the Method of Exhaustion and inspired the more sophisticated technique more commonly used today known as integration. Before we dive into integration, though, we need to lay some groundwork. Specifically, we'll define an integral as the limit of a sequence, where each term in the sequence is a sum which gets closer and closer to the desired value of the integral (whatever that may be!). So, we need to spend some time thinking about sequences, sums of sequences, and the limiting behaviour of sequences and their sums.

6.1.1 Introduction to Sequences

A sequence is just an **ordered** list of numbers. For example, the first five non-zero perfect squares listed in order forms the sequence $\{1, 4, 9, 16, 25\}$. In general, we can denote an arbitrary sequence with k elements by $\{a_1, a_2, \dots, a_k\}$ or, more concisely, $\{a_n\}_{n=1}^k$.

In some instances, we may be able to write down an explicit formula for the terms in a sequence as a function of the index. For example, we could write the sequence of perfect squares above as $\{n^2\}_{n=1}^5$ since

$$\{n^2\}_{n=1}^5 = \{1^2, 2^2, 3^2, 4^2, 5^2\} = \{1, 4, 9, 16, 25\}$$

Observe that we are essentially taking a sample of outputs from the continuous function $f(n) = n^2$ in this case. In other words, the sequence is just a set of outputs of a function over a domain consisting of a set of consecutive integers. It will be helpful at times to keep this connection between functions and sequences in mind.

We call a sequence with infinitely many terms an **infinite sequence**. These can be denoted by

$$\{a_n\}_{n=1}^{\infty} = \{a_1, a_2, \dots, a_n, \dots\}$$

For example, consider the following infinite sequence

$$\{\cos(n\pi)\}_{n=0}^{\infty} = \{\cos(0), \cos(\pi), \cos(2\pi), \dots\} = \{1, -1, 1, -1, 1, -1, \dots\}$$

Observe that our index need not start at 1. This is true of finite and infinite sequences.

In contrast to the examples above, it is not necessary that we be able to write out an explicit formula for the n -th term in a sequence. We could, for example, define a sequence to be made up of the digits of π .

$$\{3, 1, 4, 1, 5, 9, 2, 6, 5, 3, 5, \dots\}$$

Alternatively, we could define a sequence via a recursion relation. As an example you've probably seen, suppose we set $a_1 = a_2 = 1$ and then require that $a_n = a_{n-1} + a_{n-2}$ for $n \geq 3$. This yields the Fibonacci sequence

$$\{1, 1, 2, 3, 5, 8, 13, 21, 34, 55, \dots\}$$

Example 1

Find an expression for the n -th term in the sequence $\{a_n\}$ with $a_1 = 1$ and $a_n = 3a_{n-1}$ for $n \geq 2$.

Solution: Writing out the first few terms, we have

$$\begin{aligned} a_1 &= 1, \\ a_2 &= 3a_1 = 3(1) = 3, \\ a_3 &= 3a_2 = 3(3) = 3^2, \\ a_4 &= 3a_3 = 3(3^2) = 3^3, \end{aligned}$$

We observe that $a_n = 3^{n-1}$ for $n \geq 1$. To be a little more formal, we could prove this using induction, but we'll leave that to you!

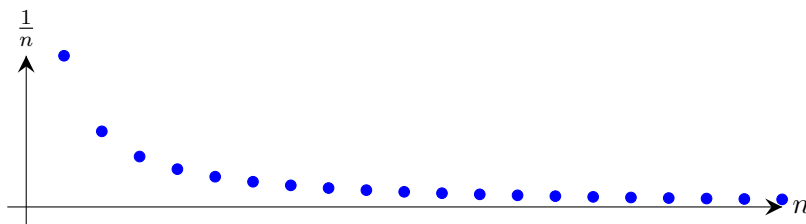
6.1.2 Limit of a Sequence

Often, an important feature of an infinite sequence is the 'eventual' behaviour of the terms in the sequence. Specifically, it can be useful to know if the values of the terms in the sequence are approaching a finite fixed value, growing arbitrarily large (positive or negative), or not settling towards a particular state at all.

Consider the sequence

$$\left\{ \frac{1}{n} \right\}_{n=1}^{\infty} = \left\{ 1, \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \dots \right\}$$

Here is a graph plotting the values of the sequence:



This sequence is known as the Harmonic Sequence and as n tends to infinity, the terms in the sequence tend to zero.

We can state this fact by writing down the following limit.

$$\lim_{n \rightarrow \infty} \frac{1}{n} = 0$$

In analogy with the corresponding limit at infinity of the function $f(x) = \frac{1}{x}$, it seems reasonable to write out such an expression. However, we should pause for a moment because our definition of a limit at infinity was specifically for a function whose domain contains an interval of the form $(a, \infty) \subset \mathbb{R}$. Sequences aren't functions with this property at all! Thankfully, our limit definition can be easily adapted for a sequence.

Definition 6.1.1
limit of a sequence

We say that the limit of the sequence $\{a_n\}$ as n approaches ∞ is L and write

$$\lim_{n \rightarrow \infty} a_n = L$$

if for every $\epsilon > 0$, there exists some positive integer k such that if

$$n > k,$$

then

$$|a_n - L| < \epsilon.$$

REMARK

If for a sequence $\{a_n\}$ we have $\lim_{n \rightarrow \infty} a_n = L$, then we say the sequence is **convergent**, or that the sequence **converges**. On the other hand, if this limit does not exist, then we say the sequence is **divergent**, or that the sequence **diverges**.

Example 2

Use the definition of the limit of a sequence to show $\{a_n\}$ with $n \geq 1$ where $a_n = \frac{n-1}{n}$ converges to 1.

Proof: Let $\epsilon > 0$ and take k be a positive integer such that $k > \frac{1}{\epsilon}$. Then, for all $n > k$, we have

$$\begin{aligned} n > \frac{1}{\epsilon} &\implies \frac{1}{n} < \epsilon \\ &\implies \left| -\frac{1}{n} \right| < \epsilon \\ &\implies \left| 1 - \frac{1}{n} - 1 \right| < \epsilon \\ &\implies \left| \frac{n-1}{n} - 1 \right| < \epsilon \end{aligned}$$

Therefore, we have $\lim_{n \rightarrow \infty} a_n = 1$.

EXERCISE

Prove that the sequence $\{a_n\}_{n=1}^{\infty}$ with $a_n = n$ is divergent. (Hint: Assume that the sequence converges to some finite value L and argue that this leads to a contradiction.)

Since the definition of the limit of a sequence is analogous to the definition of the limit at infinity of a function, an analogous version of the properties that hold for the latter, also hold for the former.

Fact 1 If $\{a_n\}$ and $\{b_n\}$ are convergent sequences and c is a constant, then

1. $\lim_{n \rightarrow \infty} (a_n \pm b_n) = \lim_{n \rightarrow \infty} a_n \pm \lim_{n \rightarrow \infty} b_n$ (sum/difference law)
2. $\lim_{n \rightarrow \infty} (ca_n) = c \lim_{n \rightarrow \infty} a_n$ (scalar multiplication law)
3. $\lim_{n \rightarrow \infty} (a_n \cdot b_n) = \lim_{n \rightarrow \infty} a_n \cdot \lim_{n \rightarrow \infty} b_n$ (product law)
4. $\lim_{n \rightarrow \infty} \left(\frac{a_n}{b_n} \right) = \frac{\lim_{n \rightarrow \infty} a_n}{\lim_{n \rightarrow \infty} b_n}$ provided $\lim_{n \rightarrow \infty} b_n \neq 0$ (quotient law)
5. $\lim_{n \rightarrow \infty} (a_n^c) = \left(\lim_{n \rightarrow \infty} a_n \right)^c$ provided $a_n > 0$ and $c > 0$ (power law)

The Squeeze theorem can also be applied when computing limits of sequences.

Theorem 2 (Squeeze theorem for sequences)

Let $\{a_n\}$, $\{b_n\}$, and $\{c_n\}$ be infinite sequences. If $\lim_{n \rightarrow \infty} b_n = \lim_{n \rightarrow \infty} c_n = L$ and there exists an integer N so that $b_n \leq a_n \leq c_n$ for all $n > N$, then $\lim_{n \rightarrow \infty} a_n = L$.

EXERCISE

Prove the Squeeze theorem for sequences.

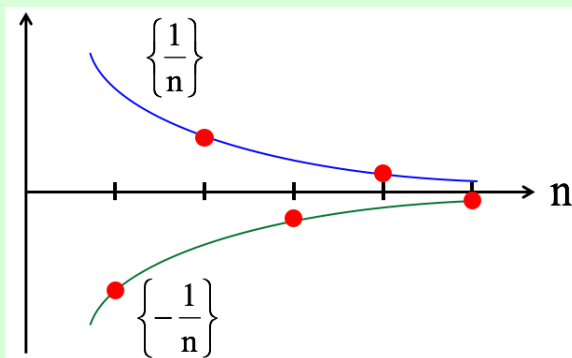
Example 3

Let $a_n = \frac{\cos(n\pi)}{n}$ for $n \geq 1$. Show $\{a_n\}$ converges to zero.

Solution: This is the Harmonic Sequence with alternating signs.

$$\{a_n\} = \left\{ -1, \frac{1}{2}, -\frac{1}{3}, \frac{1}{4}, -\frac{1}{5}, \dots \right\}$$

The magnitude of the values in this sequence tend to zero, but how do we deal with the alternating signs? Observe that we can always bound the terms in this sequence between the sequences $\{-\frac{1}{n}\}$ and $\{\frac{1}{n}\}$.



Let's take advantage of this idea to apply the Squeeze theorem.

Since $-1 \leq \cos(n\pi) \leq 1$ and $n > 0$, we have

$$-\frac{1}{n} \leq \frac{\cos(n\pi)}{n} \leq \frac{1}{n} \implies -\frac{1}{n} \leq a_n \leq \frac{1}{n}$$

Taking the limit as $n \rightarrow \infty$ and noting that $\lim_{n \rightarrow \infty} \frac{1}{n} = 0$ gives

$$0 \leq \lim_{n \rightarrow \infty} a_n \leq 0$$

Therefore, by the Squeeze theorem, the sequence converges to zero.

EXERCISE

Determine whether the sequence $\{a_n\}$ with $n \geq 1$ is convergent or divergent where

$$a_n = n^2 e^{-n^2}$$

Another very useful property we can use when investigating sequence convergence is the composition between a function and the terms in a sequence.

Theorem 3 If $\lim_{n \rightarrow \infty} a_n = L$ and f is a continuous function at L , then

$$\lim_{n \rightarrow \infty} f(a_n) = f(L)$$

We've seen this kind of behaviour before (see Theorem 4 from Chapter 2), where continuity allows us to bring limits inside the argument of the function!

Proof: This proof is an exercise in unwrapping definitions. Let $\epsilon > 0$. Since f is continuous at L , there exists $\delta > 0$ so that $|x - L| < \delta$ implies $|f(x) - f(L)| < \epsilon$. Since $\lim_{n \rightarrow \infty} a_n = L$, there is some integer N so that $n > N$ implies $|a_n - L| < \delta$. Putting these two implications together we have that if $n > N$, $|f(a_n) - f(L)| < \epsilon$. This precisely satisfies the definition of $\lim_{n \rightarrow \infty} f(a_n) = f(L)$. \square

Example 4

Determine what value the sequence $\{a_n\}$ with $n \geq 1$ where $a_n = e^{\frac{1}{\sqrt{n}}}$ converges to.

Solution: By the previous theorem, since $f(x) = e^x$ is continuous on \mathbb{R} (and therefore at zero) and $\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}} = 0$ we have

$$\begin{aligned}\lim_{n \rightarrow \infty} a_n &= e^{\left(\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}}\right)} \\ &= e^0 \\ &= 1\end{aligned}$$

Therefore, the sequence converges to 1.

EXERCISE

Determine to what value the sequence $\{a_n\}$ with $n \geq 1$ converges if

$$a_n = \sqrt{\frac{n^2 - n}{4n^2 + 1}}$$

6.1.3 Monotonicity and Boundedness

Consider the sequence

$$\left\{ \frac{n-1}{n} \right\}_{n=1}^{\infty} = \left\{ 0, \frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \frac{4}{5}, \dots \right\}$$

We previously proved that this sequence converges to 1. With a bit of thought about the behaviour of the sequence, this is not surprising. We just need to make two observations. First, each term is strictly larger than the one before it. Second, the terms can never *exceed* a value of one. This means that the terms keep growing, but since they can't grow arbitrarily large, they must instead be honing in on some finite value. Let's make this idea more precise with a few definitions.

Definition 6.1.2

A sequence $\{a_n\}$ with $n \geq 1$ is:

increasing/decreasing

- increasing if $a_{n+1} > a_n$ for all $n \geq 1$.
- non-decreasing if $a_{n+1} \geq a_n$ for all $n \geq 1$.
- decreasing if $a_{n+1} < a_n$ for all $n \geq 1$.
- non-increasing if $a_{n+1} \leq a_n$ for all $n \geq 1$.

If a sequence is either non-increasing or non-decreasing, we call it **monotonic**.

EXERCISE

Show that $\left\{\frac{1}{\sqrt{n}}\right\}_{n=1}^{\infty}$ is a decreasing sequence.

Definition 6.1.3
bounded

A sequence $\{a_n\}$ with $n \geq 1$ is:

- bounded above if there exist $M \in \mathbb{R}$ such that $a_n \leq M$ for all $n \geq 1$.
- bounded below if there exist $m \in \mathbb{R}$ such that $a_n \geq m$ for all $n \geq 1$.

If a sequence is bounded both above and below, we say it is a **bounded sequence**.

EXERCISE

Show that the sequence $\left\{\frac{4n+1}{2n}\right\}_{n=1}^{\infty}$ is bounded below by 2.

Combining the notions of monotonicity and boundedness leads to the powerful Theorem 4. We call it the monotonic sequence theorem, but it is also known as the monotone convergence theorem.

Theorem 4 (Monotonic sequence theorem)

If a sequence is monotonic and bounded, then it is convergent.

Proof: We will first prove the result for non-decreasing sequences. Let $\{a_n\}$ with $n \geq 1$ be a non-decreasing sequence bounded above. Let L be the least upper bound of $\{a_n\}$, that is, an upper bound for $\{a_n\}$ so that if M is any other upper bound, then $M \geq L$. A least upper bound always exists by the Completeness Axiom for the real numbers.

Given $\epsilon > 0$, we can always find some element of the sequence within ϵ of L . (If we could not, then we would be able to reduce the least upper bound and this cannot be done by definition of the least upper bound.) Let a_k for some positive integer k be such an element. That is,

$$a_k > L - \epsilon$$

Since the sequence is also non-decreasing, it must be true that $a_n \geq a_k$ for all $n > k$. It follows that for all $n > k$

$$a_n > L - \epsilon \quad \implies \quad L - a_n < \epsilon$$

Since $L > a_n$, we can rewrite this as

$$|a_n - L| < \epsilon$$

Therefore, the $\lim_{n \rightarrow \infty} a_n = L$ and the sequence is convergent.

A similar argument can be used to demonstrate that a non-increasing sequence which is bounded below will be convergent. \square

REMARK

Observe that the proof of the monotonic sequence theorem also tells us what the limit should be! In particular, the limit of a non-decreasing sequence is its least upper bound, and the limit of a non-increasing function is its greatest lower bound.

EXERCISE

Let $\{a_n\}$ be the sequence recursively defined by $a_1 = 1$ and $a_{n+1} = \sqrt{3 + 2a_n}$ for $n \geq 1$. Prove that $\{a_n\}$ converges.

6.2 Series

We will make use of infinite sequences to confidently approximate function values when we study Taylor polynomials. We will also use them to properly formulate the notion of an integral. In both cases, we will need to sum the infinitely many terms in a sequence. We call the sum of the terms in an infinite sequence the **series**. Roughly speaking, we can visualize how an infinite sequence is related to a series as follows,

$$\underbrace{\{a_1, a_2, \dots, a_n, \dots\}}_{\text{Infinite Sequence}} \quad \text{versus} \quad \underbrace{a_1 + a_2 + \dots + a_n + \dots}_{\text{Series}}$$

but we can also give a precise definition.

Definition 6.2.1
series

For a sequence $\{a_n\}$ with $n \geq 1$, the sum

$$\sum_{n=1}^{\infty} a_n = a_1 + a_2 + a_3 + \dots + a_n + \dots$$

is called a series.

As we have done in this definition, we will frequently use sigma notation to write series more concisely. When we do this, we refer to n as the **index** of a term a_n in the series.

REMARKS

- In general, the sum defining a series does not need to start with index $n = 1$, but it does need to include infinitely many terms to be a series.
- It is common to hear a series referred to as an “infinite series” with the word “infinite” added to indicate that there are infinite many terms. Technically, this is redundant given the definition.

When the sum of the infinitely many terms in a series is equal to a finite number, we say the series is convergent and that the series converges to this value. If the sum does *not* tend to a specific, finite value, then we say the series is divergent. We'll give a more functional definition of series convergence and divergence shortly, but first let's look at a couple of examples.

Example 5

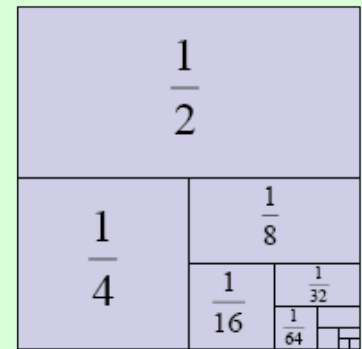
Determine the sum of the series $\sum_{n=1}^{\infty} \left(\frac{1}{2}\right)^n$.

Solution: This series represents the sum of all positive integer powers of $\frac{1}{2}$. That is, if we let s denote the series, then

$$s = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \cdots + \frac{1}{2^n} + \cdots$$

You may recognize these terms as those of a geometric sequence and possibly know how to compute the sum. If not, don't worry because we'll come back to the general case for geometric series in good time. For now, let's instead use a fun visual argument to determine the value of the series.

Let us represent the terms in our series by areas of rectangles. We start by constructing a rectangle with width 1 and height $\frac{1}{2}$ to get a rectangle with area equal to $\frac{1}{2}$. For our next rectangle (which will be a square), we halve the width to get dimensions $\frac{1}{2} \times \frac{1}{2}$ for an area of $\frac{1}{4}$. Next, we halve the height to get a rectangle with dimensions $\frac{1}{2} \times \frac{1}{4}$ and area $\frac{1}{8}$. We continue in this manner alternating between halving the width or the height to construct an infinite array of rectangles whose areas match up with the terms in our series.



Observe that we can piece our rectangles together in a square with dimension 1×1 for a total area equal to 1. Therefore, the series converges to 1.

The result of the previous example can be a bit surprising at first. We are adding together infinitely many strictly positive terms. Should the sum not then be infinite? In this case, the answer is 'no' and the reason is that the terms are approaching zero fast enough to allow for the sum to be finite. We must be careful though, it is *not* true that any series whose terms tend to zero will converge. The rate at which the terms approach zero is crucial.

Example 6

Show that the series $\sum_{n=1}^{\infty} \frac{1}{n}$ is divergent.

Solution: This series is known as the **Harmonic series** and we can argue that it is unbounded

by grouping terms in a clever way. Observe that

$$\begin{aligned} \sum_{n=1}^{\infty} \frac{1}{n} &= 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} + \dots \\ &= 1 + \frac{1}{2} + \left(\frac{1}{3} + \frac{1}{4}\right) + \left(\frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8}\right) + \dots \\ &> 1 + \frac{1}{2} + \left(\frac{1}{4} + \frac{1}{4}\right) + \left(\frac{1}{8} + \frac{1}{8} + \frac{1}{8} + \frac{1}{8}\right) + \dots \\ &> 1 + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \dots \end{aligned}$$

By reaching deeper and deeper into the sequence as needed, we can create infinitely many groupings whose sums are each greater than $\frac{1}{2}$. If we attempt to add together infinitely many copies of $\frac{1}{2}$, we get a sum tending to $+\infty$. Since the sum of the Harmonic series will grow even faster than this, then it is be divergent.

Inspired by this example, we state the following result.

Theorem 5 (Divergence Test)

Given a sequence $\{a_n\}$ with $n \geq 1$, if $\lim_{n \rightarrow \infty} a_n \neq 0$, then the series $\sum_{n=1}^{\infty} a_n$ is divergent.

This theorem is simply saying that if the terms in a sequence tending to some finite non-zero value as we progress further and further down the sequence, then the sum of the series can never settle to a particular finite value.

REMARK

It is worth noting that the Divergence Test does *not* tell us if a series for which $\lim_{n \rightarrow \infty} a_n = 0$ will be convergent or divergent.

6.2.1 Partial Sums

When adding together the infinitely many terms in an infinite sequence, it can be helpful as an intermediate step (or for entirely different reasons) to consider the sum of just the first k terms in a series. We call this a partial sum.

Definition 6.2.2 partial sum

Given a series $\sum_{n=1}^{\infty} a_n$, for each integer $k \geq 1$, we define the k -th **partial sum** of the series as the sum of the first k terms and denote it by

$$s_k = \sum_{n=1}^k a_n$$

Observe that the partial sums of a series form a sequence of partial sums $\{s_k\}$. It is this sequence we use to formally define series convergence.

Definition 6.2.3
series convergence

Given a sequence $\{a_n\}$ with $n \geq 1$, let s_k denote the k -th partial sum of the sequence. We say the series $\sum_{n=1}^{\infty} a_n$ is **convergent** if the sequence $\{s_k\}$ is convergent. Moreover, if the sequence $\{s_k\}$ is convergent, then we define

$$\sum_{n=1}^{\infty} a_n = \lim_{k \rightarrow \infty} s_k$$

If this limit does not exist because the sequence $\{s_k\}$ is divergent, then we say the series is **divergent**.

EXERCISE

Now that we have a formal definition for what it means for a series to converge (or diverge), prove the divergence test.

Example 7

Find an expression for the partial sum of the following series and then determine if the series is convergent or divergent.

$$\sum_{n=1}^{\infty} \frac{1}{n(n+1)} = \frac{1}{2} + \frac{1}{6} + \frac{1}{12} + \dots$$

Solution: Let's look at the first few terms in the sequence of partial sums $\{s_k\}$.

$$\begin{aligned} s_1 &= \frac{1}{2} \\ s_2 &= \frac{1}{2} + \frac{1}{6} = \frac{2}{3} \\ s_3 &= \frac{2}{3} + \frac{1}{12} = \frac{3}{4} \end{aligned}$$

It looks like the partial sum has the form $\frac{k}{k+1}$, but we must justify this. We could use induction, but let's do something a bit different here for fun.

Observe that we should be able to rewrite a_n as follows.

$$\frac{1}{n(n+1)} = \frac{A}{n} + \frac{B}{n+1}$$

for some constants A and B . (We call this a partial fraction expansion.) With a bit of algebra, we find $A = 1$ and $B = -1$.

Writing a_n in this form it becomes more apparent how our suspected formula for the k -th partial sum arises.

$$\begin{aligned} s_k &= \sum_{n=1}^k \left(\frac{1}{n} - \frac{1}{n+1} \right) \\ &= \left(1 - \frac{1}{2} \right) + \left(\frac{1}{2} - \frac{1}{3} \right) + \left(\frac{1}{3} - \frac{1}{4} \right) + \cdots + \left(\frac{1}{k} - \frac{1}{k+1} \right) \\ &= 1 - \frac{1}{k+1} \\ &= \frac{k}{k+1} \end{aligned}$$

Note, we can be more precise in the second step by splitting the sum as $\sum_{n=1}^k \frac{1}{n} - \sum_{n=1}^k \frac{1}{n+1}$, re-indexing the second sum to start at $n=2$ (i.e., so that it is $\sum_{n=2}^{k+1} \frac{1}{n}$), and then cancelling all terms present in both sums.

Back to the problem at hand, since

$$\begin{aligned} \lim_{k \rightarrow \infty} s_k &= \lim_{k \rightarrow \infty} \frac{k}{k+1} \\ &= \lim_{k \rightarrow \infty} \frac{1}{1 + \frac{1}{k}} \\ &= 1 \end{aligned}$$

we conclude that the series is convergent and, moreover, that it converges to a value of 1.

EXERCISE

Determine if the following series is convergent or divergent. If it is convergent, determine to what value the series converges.

$$\sum_{n=1}^{\infty} \ln \left(\frac{n}{n+1} \right)$$

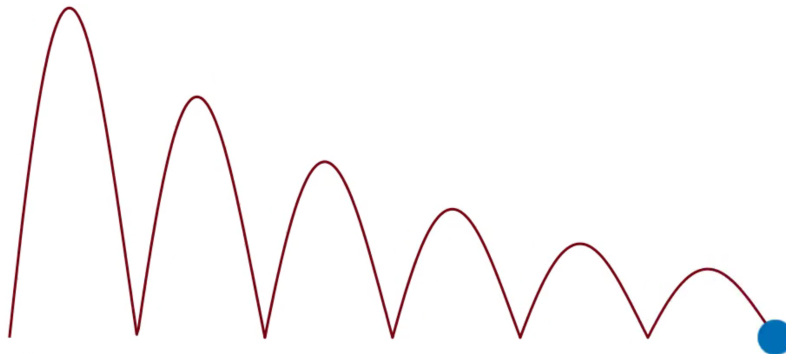
EXERCISE

Prove that the series given by summing all natural numbers is divergent by:

- using the Divergence Test.
- showing the the sequence of partial sums is divergent.

6.2.2 Geometric Series

Imagine we launch a bouncy ball such that it reaches a height of 5 meters before its first bounce. On each bounce, the ball then loses energy and only reaches two-thirds the height it attained on its previous bounce.



What is the total vertical distance d traversed by the ball?

Before the first bounce, the ball will cover a vertical distance of 10 m (5 m up and 5 m down). After the first bounce, the ball will cover two-thirds of 10 m. And after the next bounce, it will cover two-thirds of two-thirds of 10 m.

If we let the distance travelled on each trip be d_n , then we get the sequence $\{d_n\}$ with $n \geq 1$ where

$$d_n = 10 \left(\frac{2}{3}\right)^{n-1}$$

Since each term in the sequence is the same scalar multiple of the previous term (i.e., $a_{k+1} = \frac{2}{3}a_k$ for $k \geq 1$), we call this a **geometric sequence**.

Adding together the infinitely many terms in this geometric sequence gives us a **geometric series** representing the total vertical distance traversed by the ball.

$$d = 10 + 10 \left(\frac{2}{3}\right) + 10 \left(\frac{2}{3}\right)^2 + \dots = 10 \sum_{n=1}^{\infty} \left(\frac{2}{3}\right)^{n-1}$$

Definition 6.2.4 geometric series

A series whose successive terms differ by a multiplicative factor r is called a **geometric series**. If the first term in the series is a , then the series takes the form

$$a \sum_{n=1}^{\infty} r^{n-1} = a + ar + ar^2 + \dots \quad (6.1)$$

We previously encountered the geometric series $\sum_{n=1}^{\infty} \left(\frac{1}{2}\right)^{n-1}$, which we determined is convergent with an argument involving areas. However, not all geometric series will converge. For example, $\sum_{n=1}^{\infty} 2^{n-1} = 1 + 2 + 4 + \dots$ is clearly divergent; we could justify this quickly with

the Divergence Test. Let's see now if we can come up with a general rule for classifying a geometric series as convergent or divergent. To do this, we will take the approach of finding an expression for the partial sum of a geometric series.

Let s_k denote the sum of the first k terms in the geometric series $\sum_{n=1}^{\infty} ar^{n-1}$. That is,

$$s_k = \sum_{n=1}^k ar^{n-1} = a + ar + ar^2 + \cdots + ar^{k-1}$$

If we multiply s_k by r , we have

$$rs_k = \sum_{n=1}^k ar^n = ar + ar^2 + ar^3 + \cdots + ar^k$$

Taking the difference between s_k and rs_k leaves just the first term in the former sum and the last term in the latter sum; every other term cancels out.

$$s_k - rs_k = a - ar^k$$

Solving this expression when $r \neq 1$ for s_k yields

$$s_k = a \frac{1 - r^k}{1 - r}$$

When the limit of s_k as $k \rightarrow \infty$ exists, the series will converge. Moreover, this limit will exist when the limit of the term r^k in the numerator of the partial sum exists. We have a few cases to consider to complete our investigation.

- If $-1 < r < 1$, then as $k \rightarrow \infty$, r^k converges to 0.
- If $r > 1$ or $r < -1$, then as $k \rightarrow \infty$, r^k diverges.
- If $r = -1$, then r^k oscillates between -1 and 1 , so is divergent.
- If $r = 1$, we cannot use the partial sum formula since we divided by $1 - r$ in our derivation. However, a series of the form $\sum_{n=1}^{\infty} a = a + a + a + \cdots$ will clearly diverge.

We can summarize our results as follows.

Fact 6

The geometric series $\sum_{n=1}^{\infty} ar^{n-1}$ with $a \neq 0$ is convergent if and only if $|r| < 1$. Furthermore, when $|r| < 1$, the series converges to $\frac{a}{1 - r}$.

Using this result, we can now determine the total vertical distance traversed by the bouncing ball discussed earlier. In that case, we had $a = 10$ and $r = \frac{2}{3}$. Since $|r| < 1$, the series will converge to $\frac{10}{1 - \frac{2}{3}} = 30$. Therefore, the ball travels a total vertical distance of 30 m.

EXERCISE

Consider the sequence $\{a_n\}$ with $a_1 = 4$ and $a_{n+1} = -\frac{3}{5}a_n$ for $n \geq 2$. Determine the sum of the series

$$\sum_{n=1}^{\infty} a_n = 4 - \frac{12}{5} + \frac{36}{25} - \frac{108}{125} + \dots$$

6.2.3 Series Properties

We finish this section by noting some properties that will help us manipulate series and test them for convergence.

Theorem 7 (series properties)

If $\sum_{n=1}^{\infty} a_n$ and $\sum_{n=1}^{\infty} b_n$ are convergent series, then:

1. The series $\sum_{n=1}^{\infty} ca_n$ is convergent for any $c \in \mathbb{R}$ with

$$\sum_{n=1}^{\infty} ca_n = c \sum_{n=1}^{\infty} a_n$$

2. The series $\sum_{n=1}^{\infty} (a_n + b_n)$ is convergent with

$$\sum_{n=1}^{\infty} (a_n + b_n) = \sum_{n=1}^{\infty} a_n + \sum_{n=1}^{\infty} b_n$$

Our next property encodes the idea that a series will not change from convergent to divergent or vice versa by removing a finite number of terms.

Theorem 8 If $\sum_{n=1}^{\infty} a_n$ is a convergent series, then for any positive integer j , the series $\sum_{n=j}^{\infty} a_n$ is convergent.

We can make sense of this statement by recognizing that

$$\sum_{n=1}^{\infty} a_n - \sum_{n=j}^{\infty} a_n = \sum_{n=1}^{j-1} a_n$$

The quantity on the right-hand side of the previous equation is the sum of a finite number of finite numbers which is therefore a finite number. Therefore, the two series on the left-hand side of this equation must both converge or both diverge. We gain a further lesson from this result. In particular, whether or not a series is convergent depends on what is called the *tail* of the sequence being summed. That is, if the series is adding together elements of a sequence $\{a_n\}$, then all that matters for convergence is the behaviour of the terms in this sequence as $n \rightarrow \infty$.

6.3 Comparison Tests

Consider the following series:

$$\sum_{n=1}^{\infty} \frac{1}{2^n + 1} = \frac{1}{3} + \frac{1}{5} + \frac{1}{9} + \dots$$

Now, compare term-by-term to the following geometric series which converges to 1:

$$\sum_{n=1}^{\infty} \frac{1}{2^n} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = 1$$

Observe that for each term in the first series, the corresponding term in the second series is larger. Thus, we can establish the following inequality.

$$0 < \sum_{n=1}^{\infty} \frac{1}{2^n + 1} < \sum_{n=1}^{\infty} \frac{1}{2^n} = 1$$

It would be tempting to infer at this point that the original series must also be convergent. However, all this really tells us is that the sum is always somewhere between 0 and 1. It is still possible, for example, that the sum bounces around between values on the interval $(0, 1)$ endlessly rather than settling down to a specific value. If this possibility feels wrong to you in this example, then you're in good company. And in case you can't quite put your finger on why, it is because the terms in the series are all positive. This means the sequence of partial sums can only 'bounce' in the positive direction. But since the sequence of partial sums is increasing and bounded, it must converge by the Monotonic Sequence Theorem.

6.3.1 Comparison Test

Let us generalize the idea above for testing series convergence via comparison with a series whose behaviour we know and dub it the Comparison Test.

Theorem 9 (Comparison Test)

Let $\sum_{n=1}^{\infty} a_n$ and $\sum_{n=1}^{\infty} b_n$ be series with positive terms.

- If $\sum_{n=1}^{\infty} b_n$ converges and $a_n \leq b_n$ for all n , then $\sum_{n=1}^{\infty} a_n$ converges.
- If $\sum_{n=1}^{\infty} b_n$ diverges and $a_n \geq b_n$ for all n , then $\sum_{n=1}^{\infty} a_n$ diverges.

Example 8

We previously showed that the series $\sum_{n=2}^{\infty} \frac{1}{n(n-1)}$ is convergent. Use this result to prove that the series $\sum_{n=1}^{\infty} \frac{1}{n^2}$ is convergent.

Solution: Observe that for $n \geq 2$, we have $n^2 - n < n^2$. It follows that for $n \geq 2$,

$$\frac{1}{n^2} < \frac{1}{n(n-1)}$$

Let $a_n = \frac{1}{n^2}$ and $b_n = \frac{1}{n(n-1)}$. Then we have $0 < a_n \leq b_n$ for all $n \geq 2$. It follows by the Comparison Test that the series $\sum_{n=2}^{\infty} \frac{1}{n^2}$ is convergent. Adding in the $n = 1$ term will not change the convergent nature of this series, so the series $\sum_{n=1}^{\infty} \frac{1}{n^2}$ is also convergent.

EXERCISE

Use the result of the previous example to prove that $\sum_{n=1}^{\infty} \frac{1}{n^p}$ converges for all $p \geq 2$.

We call a series of the form $\sum_{n=1}^{\infty} \frac{1}{n^p}$ a ***p*-Series**. Series of this form serve as useful references when applying comparison tests. In the example and exercise above, it is established that these series converge for at least $p \geq 2$. We also previously showed it does not converge when $p = 1$ (i.e., for the Harmonic Series). In fact, it turns out they converge if and only if $p > 1$. We will fully justify this when we introduce the Integral Test for convergence, but state the result now so that we can refer to it as needed.

Fact 10

The series $\sum_{n=1}^{\infty} \frac{1}{n^p}$ is called a ***p*-series** and is convergent if and only if $p > 1$.

Returning to the Comparison Test, observe that it can also be used to argue for divergence.

Example 9

Prove that the series $\sum_{n=2}^{\infty} \frac{1}{n - \sqrt{n}}$ is divergent.

Solution: Let $a_n = \frac{1}{n - \sqrt{n}}$ and $b_n = \frac{1}{n}$. For $n \geq 2$, it holds that $\frac{1}{n - \sqrt{n}} > \frac{1}{n} > 0$.

Therefore, by the Comparison Test, since $\sum_{n=2}^{\infty} \frac{1}{n}$ is a divergent *p*-Series (specifically the Harmonic Series without the first term), then the given series is also divergent.

EXERCISE

Determine whether the following series is convergent or divergent.

$$\sum_{n=1}^{\infty} \frac{2}{3n-2}$$

It is important to note that the comparison test tells us nothing if the inequalities are reversed.

EXERCISE

Come up with sequences $\{a_n\}$, $\{b_n\}$, and $\{c_n\}$ satisfying

- $a_n \leq b_n$ and $a_n \leq c_n$ for all n ,
- $\sum_{n=1}^{\infty} a_n$ converges,
- $\sum_{n=1}^{\infty} b_n$ converges, and
- $\sum_{n=1}^{\infty} c_n$ diverges.

6.3.2 Limit Comparison Test

As we now know, the convergence or divergence of a series depends on the behaviour of the terms in the tail of the series - that is, the behaviour of the sequence $\{a_n\}$ for large values of the index n . Therefore, if two series differ but their terms have the same generic behaviour as $n \rightarrow \infty$, we should expect both series to either converge or diverge.

As an example, consider the following series:

$$\sum_{n=1}^{\infty} \frac{1}{2^n - 1} = 1 + \frac{1}{3} + \frac{1}{7} + \frac{1}{15} + \dots$$

Compare this term-by-term to the convergent geometric series:

$$\sum_{n=1}^{\infty} \frac{1}{2^n} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \dots$$

We cannot use the Comparison Test here because the terms in the first series are term-by-term *greater* than those in the second (convergent) series. However, as n gets really big, the terms in the two series behave nearly identically. For example, when $n = 100$, we have

$$\left(\frac{1}{2^{100} + 1} \right) - \left(\frac{1}{2^{100}} \right) \approx 10^{-61}$$

So, while the terms differ noticeably for small n , in the limit as n tends to infinity, the differences become insignificant. As such, we expect the first series to also converge. This is the idea behind the Limit Comparison Test.

Theorem 11 (Limit Comparison Test)

Let $\sum_{n=1}^{\infty} a_n$ and $\sum_{n=1}^{\infty} b_n$ be series with positive terms. If $\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = c$ where $c > 0$ is finite, then either both series are convergent or both series are divergent.

Let's apply this test to the previous example by defining $a_n = \frac{1}{2^n - 1}$ and $b_n = \frac{1}{2^n}$. This gives the limit:

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{a_n}{b_n} &= \lim_{n \rightarrow \infty} \frac{2^n}{2^n - 1} \\ &= \lim_{n \rightarrow \infty} \frac{1}{1 - 2^{-n}} \\ &= 1 \end{aligned}$$

Since $\sum_{n=1}^{\infty} \frac{1}{2^n}$ is a convergent geometric series and the limit $\lim_{n \rightarrow \infty} \frac{a_n}{b_n}$ is a positive finite number, then by the Limit Comparison Test, the series $\sum_{n=1}^{\infty} \frac{1}{2^n - 1}$ is convergent.

REMARK

When $\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = 1$, then the terms in the two series being compared truly behave identically as $n \rightarrow \infty$. However, notice that the Limit Comparison Test only requires the limit to be a positive finite number c . This is because the convergence or divergence of a series is unaffected by rescaling every term (or, more precisely, every term in the tail of the series) by the same constant. Letting c be any positive number accounts for this fact.

Example 10

Determine whether the following series is convergent or divergent.

$$\sum_{n=1}^{\infty} \frac{\sqrt{n+1}}{3n^2+n}$$

Solution: Let's inspect the large n behaviour of the terms in the series to help us identify another series to compare this one to.

Observe that when n is very large, $\sqrt{n+1} \approx \sqrt{n}$ and $3n^2 + n \approx 3n^2$. Therefore, for large n ,

$$\frac{\sqrt{n+1}}{3n^2+n} \approx \frac{\sqrt{n}}{3n^2} = \frac{1}{3n^{3/2}}$$

So, for large n , we expect the given series to behave like the convergent p -Series $\sum_{n=1}^{\infty} \frac{1}{n^{3/2}}$ up to a constant multiple. With this in mind, let's apply the Limit Comparison Test.

Let $a_n = \frac{\sqrt{n+1}}{3n^2+n}$ and let $b_n = \frac{1}{n^{3/2}}$.

Now, take the limit of the ratio $\frac{a_n}{b_n}$ as $n \rightarrow \infty$.

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{a_n}{b_n} &= \lim_{n \rightarrow \infty} \frac{\sqrt{n+1}}{3n^2+n} \cdot \frac{n^2}{\sqrt{n}} \\ &= \lim_{n \rightarrow \infty} \frac{\sqrt{\frac{n}{n} + \frac{1}{n}}}{\frac{3n^2}{n^2} + \frac{n}{n^2}} \\ &= \lim_{n \rightarrow \infty} \frac{\sqrt{1 + \frac{1}{n}}}{3 + \frac{1}{n}} \\ &= \frac{1}{3}\end{aligned}$$

Since this limit is a positive finite number and the series $\sum_{n=1}^{\infty} \frac{1}{n^{3/2}}$ is convergent, then the given series is also convergent by the Limit Comparison Test.

EXERCISE

Prove that the following series is divergent.

$$\sum_{n=1}^{\infty} \frac{3^n}{2^n + 1}$$

(Hint: Begin by identifying a divergent geometric series whose terms behave similarly for large n .)

6.4 Alternating Series

Demonstrating convergence of a series with a procedure like the Comparison Test can be very helpful. However, it is somewhat frustrating that in many cases we cannot readily determine to what value a series converges when we use a Comparison Test.

For example, we can argue that since

$$\sum_{n=1}^{\infty} \frac{\sin^2(n)}{2^n} \leq \sum_{n=1}^{\infty} \frac{1}{2^n},$$

both series have positive terms, and the geometric series on the right-hand side converges, then by the Comparison Test, the series on the left-hand side converges. But this doesn't tell us what the series on the left-hand side converges to.

Computing a few sample partial sums for the series $\sum_{n=1}^{\infty} \frac{\sin^2(n)}{2^n}$, we find $s_5 \approx 0.6278$, $s_{10} \approx 0.6368$, $s_{20} \approx 0.6375$, and $s_{100} \approx 0.6375$. It looks like the series is converging to some number near 0.6375, but we can't truly be sure of this just by inspecting a few partial sums because there are still infinitely many terms left in the series!

However, since we know the geometric series above converges to 1 and we could also argue that the sequence of partial sums of the series $\sum_{n=1}^{\infty} \frac{\sin^2(n)}{2^n}$ is increasing, then the sum of this series must be somewhere in the interval $(s_k, 1]$ for all k . This is not a fantastic bound on the series sum, but it is a start.

In this section, we'll explore a particular type of series called an alternating series. We'll come up with a test to check if an alternating series is convergent, but we'll also use alternating series as a test bed for thinking more about how to approximate the sum of a series using a partial sum and being able to assign a precise upper bound to the error on that approximation.

6.4.1 Alternating Series Test

To get us started, let's define an alternating series.

Definition 6.4.1 alternating series

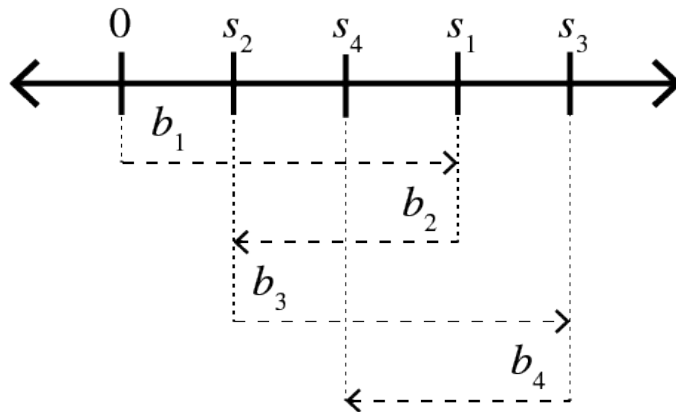
A series $\sum_{n=1}^{\infty} a_n$ where $a_n = (-1)^n b_n$ or $a_n = (-1)^{n+1} b_n$ with $b_n > 0$ for all n is called an **alternating series**.

Put more simply, an alternating series is a series where the signs of the terms alternate between positive and negative.

An example of an alternating series is the **alternating harmonic series**.

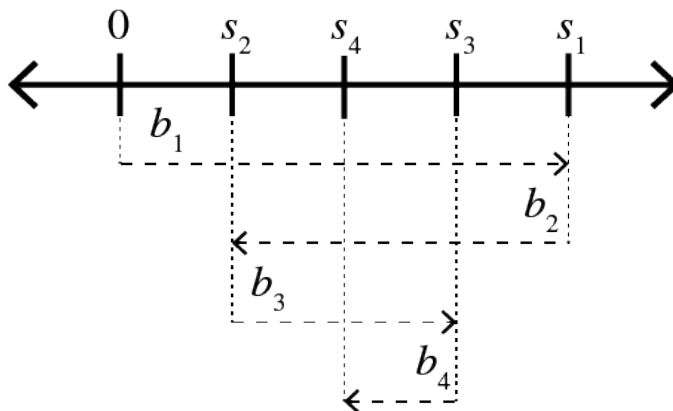
$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots$$

It can be helpful to think of the terms in an alternating series as representing steps right (for positive terms) and left (for negative terms) on the real number line. The position after n steps gives the partial sum s_n . If this process of stepping back-and-forth tends towards a specific location on the number line, then the alternating series is convergent.



In the scenario depicted above, the eventual behaviour (i.e., $n \rightarrow \infty$) of the sequence of partial sums is unclear. However, under some circumstances it is possible to quickly determine that an alternating series is convergent.

Consider a series for which the step sizes get progressively smaller (i.e., $b_{n+1} < b_n$ for all n). In this case, the n -th step will always land between the previous two steps. Put another way, s_n will always lie somewhere between s_{n-2} and s_{n-1} for $n \geq 3$. We illustrate what this might look like below.



Observe that since the step-sizes are decreasing, the interval between consecutive partial sums narrows as more steps are taken. As n tends to infinity, this interval converges to a single point. Therefore, the sequence of partial sums converges and, in turn, the alternating series is convergent.

We summarize this argument with the following convergence test.

Theorem 12 (Alternating Series Test (AST))

An alternating series $\sum_{n=1}^{\infty} (-1)^{n+1} b_n$ with $b_n > 0$ is convergent if it satisfies:

1. $b_{n+1} \leq b_n$ for all n (i.e., the sequence $\{b_n\}$ is non-increasing) and
2. $\lim_{n \rightarrow \infty} b_n = 0$

Proof: Consider the sequence of partial sums with an even number of terms, $\{s_{2k}\}$. Observe that for $n \geq 2$, we have

$$s_{2k} = s_{2k-2} + (b_{2k-1} - b_{2k})$$

Since $b_{2k-1} \geq b_{2k}$, it follows that $s_{2k} \geq s_{2k-2}$. In other words, this sequence is monotonic non-decreasing.

We can also write s_{2k} as follows.

$$\begin{aligned} s_{2k} &= \sum_{n=1}^{2k} (-1)^{n+1} b_n \\ &= b_1 - b_{2k} - \sum_{n=2}^{2k-1} (-1)^n b_n \\ &= b_1 - b_{2k} - \sum_{n=1}^{k-1} (b_{2n} - b_{2n+1}) \end{aligned}$$

Since $b_{2n} \geq b_{2n+1}$, then each difference $(b_{2n} - b_{2n+1})$ is non-negative. This means the right-hand side of this equation is equal to b_1 minus a bunch of non-negative terms. It follows that $s_{2k} \leq b_1$ which means s_{2k} is bounded above.

Since the sequence $\{s_{2k}\}$ is monotonic non-decreasing and bounded above, it must be convergent. Suppose then that the limit of this sequence is L . That is,

$$\lim_{k \rightarrow \infty} s_{2k} = L$$

Next, consider the sequence of partial sums with an odd number of terms, $\{s_{2k+1}\}$. Observe that

$$\begin{aligned} \lim_{k \rightarrow \infty} s_{2k+1} &= \lim_{k \rightarrow \infty} (s_{2k} + b_{2k+1}) \\ &= L + 0 \\ &= L \end{aligned}$$

Since the sequences $\{s_{2k}\}$ and $\{s_{2k+1}\}$ both converge to the same value L , then the sequence $\{s_k\}$ must also converge to this same value. Therefore, the series is convergent. \square

At the end of the proof we use a fact about sequences converging based on the behaviour of two different subsequences (i.e., sequences obtained from the original one by leaving some terms out). Here is an exercise to prove that what we did is above board!

EXERCISE

Let $\{a_n\}$ be a sequence. Consider the subsequences

$$\begin{aligned} \{a_{2n}\} &= \{a_2, a_4, a_6, \dots\} \\ \{a_{2n-1}\} &= \{a_1, a_3, a_5, \dots\}. \end{aligned}$$

1. Prove that if $\lim_{n \rightarrow \infty} a_n = L$, then $\lim_{n \rightarrow \infty} a_{2n} = L$ and $\lim_{n \rightarrow \infty} a_{2n-1} = L$.
2. Prove that if $\lim_{n \rightarrow \infty} a_{2n} = L$ and $\lim_{n \rightarrow \infty} a_{2n-1} = L$, then $\lim_{n \rightarrow \infty} a_n = L$.

REMARKS

- We only state and prove the Alternating Series Test for series whose first term is positive, but it works equally well when the first term is negative.
- Observe that if the second condition of the Alternating Series Test is *not* met (i.e., if $\lim_{n \rightarrow \infty} b_n \neq 0$), then we can instead conclude that the series is *divergent* by the Divergence Test.

Example 11

Prove that the alternating harmonic series $\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n}$ is convergent.

Proof: Let $b_n = \frac{1}{n}$ and apply the Alternating Series Test (AST).

Since $n + 1 \geq n$ and n is positive, we have

$$\frac{1}{n+1} \leq \frac{1}{n} \quad \implies \quad b_{n+1} \leq b_n$$

satisfying the first condition of the test.

Next, since $\lim_{n \rightarrow \infty} b_n = \lim_{n \rightarrow \infty} \frac{1}{n} = 0$, the second condition of the test is also met.

Therefore, the alternating harmonic series is convergent by the Alternating Series Test.

EXERCISE

Prove that the series $\sum_{n=1}^{\infty} \sec(n\pi)e^{-n}$ is convergent.

EXERCISE

Can the Alternating Series Test be applied to show that the following series is convergent?

$$\sum_{n=1}^{\infty} \frac{\cos(n)}{n^2}$$

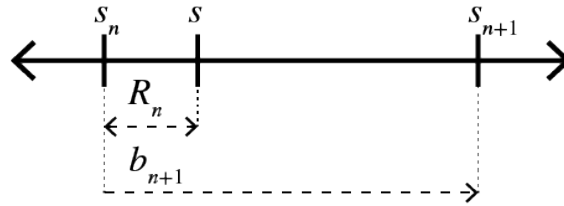
6.4.2 Alternating Series Approximations

Similar to the Comparison Test, the Alternating Series Test (AST) can check for convergence, but does not determine the value of a convergent series. However, for an alternating series satisfying the AST, we can with relative ease establish an estimate of how close a partial sum is to the series sum.

To see how this works, let's first get some notation down. Suppose we have an alternating series $\sum_{n=1}^{\infty} (-1)^{n+1} b_n$ satisfying the AST. Let's denote the n -th partial sum by s_n and the

series sum by s . We will denote the difference between the series sum and the n -th partial sum by $R_n = s - s_n$. We will call R_n the **remainder** after computing the n -th partial sum.

Recall that the series sum s will always lie in the interval between s_n and s_{n+1} , so the remainder cannot be larger than the width of this interval.



The width of this interval is given by magnitude of the next term in the series, b_{n+1} . This gives us the following theorem.

Theorem 13 (Alternating Series Remainder Theorem)

If $\sum_{n=1}^{\infty} (-1)^{n+1} b_n$ is an alternating series satisfying the Alternating Series Test, then the remainder R_n which is the difference between the series sum and the n -th partial sum satisfies

$$|R_n| \leq b_{n+1}$$

In practice, this result is very useful. Series which can be shown to converge by the AST are very common in applications. In the same applications, the sum of the series is needed but usually only to some fixed level of precision. According to this remainder theorem then, it is sufficient to find for what value of n we first get b_{n+1} to be less than the required precision. Once that is known, we can then be assured that the partial sum s_n approximates the series sum within that level of precision.

Example 12

Determine the sum of the following series within a precision of ± 0.001 .

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n^3}$$

Solution: We'd like to apply the Alternating Series Remainder Theorem but to do that, we must first show that this series converges by the Alternating Series Test (AST).

Let $b_n = \frac{1}{n^3}$. For $n \geq 1$, we'll have $b_{n+1} \leq b_n$. Also, we have $\lim_{n \rightarrow \infty} b_n = 0$. Therefore, the series converges by the AST.

With that out of the way, let's find the value of n such that the remainder R_n on the partial sum s_n satisfies

$$|R_n| \leq 0.001 = \frac{1}{1000}$$

By the Alternating Series Remainder Theorem, if we can find a value of n such that $b_{n+1} \leq \frac{1}{1000}$, then we can guarantee that the remainder R_n has the desired precision. Let's solve this inequality.

$$\begin{aligned} b_{n+1} \leq \frac{1}{1000} &\implies \frac{1}{(n+1)^3} \leq \frac{1}{1000} \\ &\implies \frac{1}{n+1} \leq \frac{1}{10} \\ &\implies n+1 \geq 10 \\ &\implies n \geq 9 \end{aligned}$$

Therefore, every partial sum s_n with $n \geq 9$ will be within ± 0.001 of the actual series sum s . With s_9 , we get

$$s_9 = 1 - \frac{1}{8} + \frac{1}{27} - \frac{1}{64} + \frac{1}{125} - \frac{1}{216} + \frac{1}{343} - \frac{1}{512} + \frac{1}{729} \approx 0.902$$

Therefore the sum of the series rounded to three decimal places is equal to 0.902 ± 0.001 .

We could even go slightly further than we did with the information we obtained in the previous example. Since the next term in the series would be $\frac{-1}{1000}$, then the partial sum *decreases* in the next step. This means the series sum rounded to three decimal places must be either 0.901 or 0.902 (i.e., it can't be 0.903).

REMARK

You can probably imagine that a computer comes in very handy for problems like these. In particular, it would be simple to write an algorithm which skips over solving any inequalities and just uses a 'while' loop to iteratively compute partial sums with the condition to stop when it encounters a term whose magnitude is less than the desired precision.

EXERCISE

Determine the sum of the series $\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n!}$ correct to three decimal places. How many terms would you need to get the sum correct to ten decimal places?

6.5 Ratio Test

6.5.1 Absolute Convergence

Consider the following series

$$\sum_{n=1}^{\infty} \frac{\cos(n)}{n^2}$$

We know that the series $\sum_{n=1}^{\infty} \frac{1}{n^2}$ is a convergent p -series, so it is tempting to try to apply the Comparison Test. Unfortunately, we can't because the terms in the series above are not all positive. We might then ask if we could try the Alternating Series Test, but again we hit a roadblock because the terms do not alternate regularly between positive and negative.

But it feels like the $\cos(n)$ really shouldn't matter too much. It is always between -1 and 1 , so really just serves to make the magnitudes of the terms in the convergent series $\sum_{n=1}^{\infty} \frac{1}{n^2}$ smaller. If anything, this reduction in magnitude of the terms and mix of positive and negative terms should make the sequence of partial sums for the series converge faster. This leads to an idea. Let's just focus on the magnitude of the terms. More precisely, let's look at a new series given by adding together the absolute value of every term in the series above:

$$\sum_{n=1}^{\infty} \left| \frac{\cos(n)}{n^2} \right|$$

The terms in this series are now positive by construction. Moreover, since $0 \leq |\cos(n)| \leq 1$ and $n^2 > 0$, we can establish the following inequality:

$$0 \leq \sum_{n=1}^{\infty} \left| \frac{\cos(n)}{n^2} \right| = \sum_{n=1}^{\infty} \frac{|\cos(n)|}{n^2} \leq \sum_{n=1}^{\infty} \frac{1}{n^2}$$

Now we *can* use the Comparison Test. Since the p -series on the right is convergent and all series are made up of positive terms, then the series $\sum_{n=1}^{\infty} \left| \frac{\cos(n)}{n^2} \right|$ is convergent.

This idea of looking at the sum of the absolute values of the terms in a series will be useful, so we give it a name.

Definition 6.5.1

absolute
convergence

If the series $\sum |a_n|$ is convergent, then we say the series $\sum a_n$ is **absolutely convergent**.

The main reason that absolute convergence is useful is the following theorem.

Theorem 14

If a series $\sum a_n$ is absolutely convergent, then it is convergent.

Here's the main idea of the proof, and you will be walked through the formal details in the next exercise.

Let s_+ be the sum of all the positive terms in the series $\sum a_n$. Similarly, let s_- be the sum of the absolute values of all the negative terms in the series $\sum a_n$. Since the series is absolutely convergent, then we know

$$\sum |a_n| = s_+ + s_-$$

is finite. Since s_+ and s_- are both positive and sum to a finite number, they must each be finite on their own. It follows that

$$\sum a_n = s_+ - s_-$$

is finite since it is the difference of two finite numbers.

Now, each of the values s_+ and s_- are limits of partial sums, so there are some details that need checking to make sure the argument above is valid. Those details are in this exercise:

EXERCISE (Proof of Theorem 14)

Let $\{a_n\}$ be a sequence with the property that the series $\sum_{n=1}^{\infty} |a_n|$ converges. Define two sequences $\{p_n\}$ and $\{q_n\}$ by

$$p_n = \begin{cases} a_n & \text{if } a_n \geq 0 \\ 0 & \text{otherwise.} \end{cases} \quad \text{and} \quad q_n = \begin{cases} -a_n & \text{if } a_n < 0 \\ 0 & \text{otherwise.} \end{cases}$$

So, we can think of $\{p_n\}$ as the positive part of $\{a_n\}$, and $\{q_n\}$ as the negative part.

1. Prove that $\sum_{n=1}^{\infty} p_n$ and $\sum_{n=1}^{\infty} q_n$ both converge.
2. Show that $a_n = p_n - q_n$ for all n .
3. Let $s_+ = \sum_{n=1}^{\infty} p_n$ and $s_- = \sum_{n=1}^{\infty} q_n$. Prove that $\sum_{n=1}^{\infty} a_n$ converges, and that it converges to $s_+ - s_-$.

REMARK

It is important to keep in mind that the converse statement is not true in general. That is, convergence does **not** imply absolute convergence. This is because some series converge not because their terms tend to zero sufficiently quickly as $n \rightarrow \infty$, but rather the convergence relies on a certain degree of cancellation between positive and negative terms. The alternating harmonic series is an example of a series which is convergent but not absolutely convergent. When this happens, we say a series is **conditionally convergent**.

EXERCISE

Show that the series $\sum_{n=1}^{\infty} \frac{\sin(\sqrt{2}n)}{n^3}$ is convergent.

6.5.2 The Ratio Test

When possible, showing a series is absolutely convergent is a great way to prove that the series is also just convergent. However, the fact that absolute convergence implies convergence is also a key ingredient for proving another very powerful method of convergence testing. Inspiration for this comes from looking at the behaviour of a geometric series.

Observe that for a geometric series $\sum_{n=1}^{\infty} ar^n$, consecutive terms satisfy the following equality.

$$\left| \frac{a_{n+1}}{a_n} \right| = \left| \frac{ar^{n+1}}{ar^n} \right| = |r|$$

Recall, a geometric series converges if and only if $|r| < 1$. So, with the above equality in mind, we could equivalently say that a geometric series is convergent if and only if

$$\left| \frac{a_{n+1}}{a_n} \right| < 1$$

It is natural to wonder if this condition can be generalized to non-geometric series. But let's give ourselves the best chance of success by only worrying about this ratio as $n \rightarrow \infty$. After all, the convergence of a series only actually depends on how the terms in the series behave in the large n limit. It turns out this idea works quite nicely and is what we call the Ratio Test.

Theorem 15 (Ratio Test)

For a series $\sum_{n=1}^{\infty} a_n$, determine the limit $L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right|$.

1. If $L < 1$, then the series is absolutely convergent, which implies the series is convergent.
2. If $L > 1$ (including $L \rightarrow \infty$), then the series is divergent.
3. If $L = 1$, then the test is inconclusive.

Proof: We will handle each case separately.

1. Suppose $L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| < 1$. Take r to be any number satisfying $L < r < 1$. It follows that there exists some positive integer k such that for all $n \geq k$, we have

$$\left| \frac{a_{n+1}}{a_n} \right| < r \quad \implies \quad |a_{n+1}| < |a_n|r$$

Since this inequality is satisfied for all $n \geq k$, one can show by induction (the details are left as an exercise) that

$$|a_{k+j}| < r^j |a_k| \quad \text{for } j \geq 1$$

Next, consider the series $\sum_{n=k+1}^{\infty} |a_n|$. Making an appropriate index change, we have

$$\sum_{n=k+1}^{\infty} |a_n| = \sum_{j=1}^{\infty} |a_{k+j}| < \sum_{j=1}^{\infty} |a_k| r^j$$

The series $\sum_{j=1}^{\infty} |a_k| r^j$ is a convergent geometric series since $|r| < 1$, so by the Comparison Test, the series $\sum_{n=k+1}^{\infty} |a_n|$ is convergent.

The series $\sum_{n=k+1}^{\infty} |a_n|$ is convergent.

Since $\sum_{n=k+1}^{\infty} |a_n|$ is convergent, then the series $\sum_{n=1}^{\infty} |a_n|$ will also be convergent as the two series only differ by a finite number of finite terms.

Therefore, the series $\sum_{n=k+1}^{\infty} a_n$ is absolutely convergent and thus convergent.

2. Suppose $L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| > 1$ or that this limit diverges to ∞ . In either case, there will exist some positive integer k such that for all $n \geq k$, we have

$$\left| \frac{a_{n+1}}{a_n} \right| > 1 \quad \implies \quad |a_{n+1}| > |a_n|$$

It follows that $\lim_{n \rightarrow \infty} a_n \neq 0$ and therefore the series is divergent by the Divergence Test.

3. We demonstrate that the test is inconclusive when $L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| = 1$ by providing an example of a convergent series and an example of a divergent series both satisfying this limit.

First, consider the convergent p -series $\sum_{n=1}^{\infty} \frac{1}{n^2}$ so that $a_n = \frac{1}{n^2}$. We have

$$L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| = \lim_{n \rightarrow \infty} \frac{\frac{1}{(n+1)^2}}{\frac{1}{n^2}} = \lim_{n \rightarrow \infty} \frac{1}{\left(1 + \frac{1}{n}\right)^2} = 1$$

Next, consider the divergent p -series $\sum_{n=1}^{\infty} \frac{1}{n}$ so that $a_n = \frac{1}{n}$. We have

$$L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| = \lim_{n \rightarrow \infty} \frac{\frac{1}{n+1}}{\frac{1}{n}} = \lim_{n \rightarrow \infty} \frac{1}{1 + \frac{1}{n}} = 1$$

Since $L = 1$ can occur for both convergent and divergent series, the test is inconclusive.

□

Example 13

Determine whether the following series is convergent or divergent.

$$\sum_{n=1}^{\infty} \frac{n!}{n^4}$$

Solution: Apply the Ratio Test with $a_n = \frac{n!}{n^4}$.

$$L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| = \lim_{n \rightarrow \infty} \frac{\frac{(n+1)!}{(n+1)^4}}{\frac{n!}{n^4}} = \lim_{n \rightarrow \infty} \frac{n+1}{\left(1 + \frac{1}{n}\right)^4} \rightarrow \infty$$

Since $L \rightarrow \infty$, the series is divergent by the Ratio Test.

EXERCISE

Determine whether the following series is convergent or divergent.

$$\sum_{n=1}^{\infty} \frac{n^2}{2^n} \cos(n\pi)$$

EXERCISE

Show that the Ratio Test is always inconclusive for a p -series. (The lesson here is that the Ratio Test will not work well for any series whose term behaves like those of a p -series in the large n limit.)

6.6 Power Series

6.6.1 Introduction to Power Series

We've spent a lot of time discussing various types of series and tests we can use to check for convergence. In this section, we will start putting what we've learned to work for a new purpose. Specifically, we're going to discover that we can represent functions using series. These series, however, need to be something a bit more general than an infinite sum of numbers. To represent a function, they'll need to include a variable. A natural and, as we'll see, useful thing to consider is a sum of integer powers of that variable. We refer to such a construction as a power series.

Definition 6.6.1

power series

A power series centred at $x = a$ is a series of the form

$$\sum_{n=0}^{\infty} c_n(x-a)^n = c_0 + c_1(x-a) + c_2(x-a)^2 + \dots$$

where the c_n 's are constants called the coefficients of the series. The quantity x is considered to be a variable in this context.

Note that when a power series is centred at $x = 0$, it simplifies to

$$\sum_{n=0}^{\infty} c_n x^n = c_0 + c_1 x + c_2 x^2 + \dots$$

Observe that a power series is essentially an infinite degree polynomial. For a given sequence of coefficients $\{c_n\}$, a power series may converge for some values of x and diverge for other values of x .

As an example, consider the power series centred at $x = 0$ with $c_n = 1$ for all n . This gives the power series

$$\sum_{n=0}^{\infty} x^n = 1 + x + x^2 + \dots$$

After staring at this for a moment, we realize this is essentially just the general form of a geometric series with $a = 1$ and r replaced by x . So we can infer that this power series is convergent for $|x| < 1$ and divergent for $|x| \geq 1$. Moreover, when it is convergent, we know what it converges to: $\frac{1}{1-x}$.

This is where the breakthrough happens. Reversing this argument, we realize that for all $x \in (-1, 1)$, we can replace the rational function $\frac{1}{1-x}$ with the power series above. That is,

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n, \quad \text{whenever } x \in (-1, 1)$$

This means that for any calculation involving the function $\frac{1}{1-x}$, we could use this power series representation in its place. Depending on the problem at hand, this could be very useful because polynomials are generally very nice to work with.

6.6.2 Power Series Convergence

We'll discover more functions which can be represented by power series, but as in the example above, it only makes sense to do this for values of x such that the power series is convergent. Thankfully, the following theorem can be used to handle most or all of that work.

Theorem 16

For a power series $\sum_{n=0}^{\infty} c_n(x-a)^n$, exactly one of the following statements is true:

1. The series converges only for $x = a$.
2. There is a positive number, R , such that the series is convergent for $|x - a| < R$ and divergent for $|x - a| > R$.
3. The series is convergent for all x .

Proof: First, observe that when $x = a$, every term in the series is zero. Therefore, the series will always be convergent when $x = a$.

Now, assume that $x \neq a$ and apply the Ratio Test. So, let $a_n = c_n(x-a)^n$ and consider

the following limit:

$$\begin{aligned} L &= \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| \\ &= \lim_{n \rightarrow \infty} \left| \frac{c_{n+1}(x-a)^{n+1}}{c_n(x-a)^n} \right| \\ &= |x-a| \lim_{n \rightarrow \infty} \left| \frac{c_{n+1}}{c_n} \right| \end{aligned}$$

We now consider three cases:

- If $\lim_{n \rightarrow \infty} \left| \frac{c_{n+1}}{c_n} \right| = 0$, then $L = 0$. By the Ratio Test, the series is convergent for all x . This outcome corresponds to statement 3 in the theorem.
- If $\lim_{n \rightarrow \infty} \left| \frac{c_{n+1}}{c_n} \right| \rightarrow \infty$, then $L \rightarrow \infty$. By the Ratio Test, the series is divergent for all x except at $x = a$ where we previously determined the series would always be convergent. This outcome corresponds to statement 1 in the theorem.
- If $\lim_{n \rightarrow \infty} \left| \frac{c_{n+1}}{c_n} \right| = \frac{1}{R}$ for some positive number R (i.e., this limit is a finite but non-zero number), then $L = \frac{|x-a|}{R}$. In this case, we have by the Ratio Test convergence when

$$L < 1 \quad \implies \quad \frac{|x-a|}{R} < 1 \quad \implies \quad |x-a| < R$$

and divergence when

$$L > 1 \quad \implies \quad \frac{|x-a|}{R} > 1 \quad \implies \quad |x-a| > R$$

This outcome corresponds to statement 2 in the theorem.

□

In each case in the previous theorem, the power series is convergent on some open interval which we call the **interval of convergence**. We call R the **radius of convergence** of the power series.

- When a power series only converges at $x = a$, we say the radius of convergence is $R = 0$ and the interval of convergence is just the point where it converges.
- When a power series converges for all x , we say the radius of convergence is $R = \infty$ and the interval of convergence is $(-\infty, \infty)$.
- When the radius of convergence R of a power series is finite and non-zero, the interval of convergence has endpoints $x = a - R$ and $x = a + R$. The power series will converge everywhere on the open interval $(a - R, a + R)$ but may also be convergent at the one or both endpoints. As such, we must check the endpoints explicitly to determine the interval of convergence in this case.

REMARK

The success of the Ratio Test in establishing the general convergence properties of a power series suggests that it will also be useful for analyzing specific examples. However, when we find a finite, non-zero radius of convergence, then we should also be prepared to apply other convergence tests at the endpoints of the interval of convergence.

Example 14

Determine the radius of convergence and interval of convergence of the power series

$$\sum_{n=0}^{\infty} \frac{(x+2)^n}{n}$$

Solution: Let's perform the Ratio Test with $a_n = \frac{(x+2)^n}{n}$.

$$\begin{aligned} L &= \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| \\ &= \lim_{n \rightarrow \infty} \left| \frac{\frac{(x+2)^{n+1}}{n+1}}{\frac{(x+2)^n}{n}} \right| \\ &= |x+2| \lim_{n \rightarrow \infty} \frac{1}{1 + \frac{1}{n}} \\ &= |x+2| \end{aligned}$$

By the Ratio Test, the power series will converge when $|x+2| < 1$ and diverge when $|x+2| > 1$. Therefore, the radius of convergence is $R = 1$.

To determine the interval of convergence, we first note that $|x+2| < 1$ is equivalent to $x \in (-3, -1)$. So, the endpoints of the interval of convergence are at $x = -3$ and $x = -1$. We now need to explicitly check if the power series is convergent at these points to see if they should be included in the interval of convergence.

When $x = -3$, the power series is

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{n}$$

We recognize this as the alternating harmonic series which can be shown to be convergent by the Alternating Series Test. Thus, $x = -3$ is included in the interval of convergence.

When $x = -1$, the power series is

$$\sum_{n=0}^{\infty} \frac{1}{n}$$

We recognize this as the harmonic series which is a divergent p -series. As such, $x = -1$ is *not* included in the interval of convergence.

In summary, the interval of convergence is $[-3, -1)$.

EXERCISE

Determine the radius of convergence and interval of convergence of the power series

$$\sum_{n=0}^{\infty} n!(x-1)^n$$

6.6.3 Manipulating Power Series Representations of Functions

It won't be long before we are swimming in power series representation of functions, but at the moment, we just have one:

$$\frac{1}{1-x} = 1 + x + x^2 + x^3 + \cdots = \sum_{n=0}^{\infty} x^n, \quad x \in (-1, 1)$$

This relation serves as a great test bed though for exploring how we might be able to manipulate the power series representation for one function to get a power series representation for a different function.

For example, if we replace every occurrence of x with $-x^2$ in the expression above, we get

$$\frac{1}{1-(-x^2)} = 1 + (-x^2) + (-x^2)^2 + (-x^2)^3 + \cdots = \sum_{n=0}^{\infty} (-x^2)^n$$

This is effectively a geometric series with common ratio $(-x^2)$, so converges only when $|-x^2| < 1 \implies x \in (-1, 1)$. Tidying things up a bit, we get the following power series representation of a new function:

$$\frac{1}{1+x^2} = \sum_{n=0}^{\infty} (-1)^n x^{2n}, \quad x \in (-1, 1)$$

Other ways we can manipulate power series representations of functions include multiplying by a scalar and multiplying by an integer power of $(x-a)$. For example,

$$\frac{3x^3}{1+x^2} = 3x^3 \sum_{n=0}^{\infty} (-1)^n x^{2n} = \sum_{n=0}^{\infty} 3(-1)^n x^{2n+3}$$

EXERCISE

Determine a power series representation along with an interval of convergence for $\frac{2x^2}{1+3x}$.

Differentiation and Anti-differentiation of Power Series

Observe that we can also obtain new power series representations of functions through differentiation. For example, differentiating the power series representation for $\frac{1}{1-x}$ we have

$$\frac{d}{dx} \left(\frac{1}{1-x} \right) = \frac{d}{dx} \left(\sum_{n=0}^{\infty} x^n \right) \quad \implies \quad \frac{1}{(1-x)^2} = \sum_{n=0}^{\infty} nx^{n-1}$$

Observe that the first term in the new power series is going to be zero, so we could either exclude the $n = 0$ term from our sum or re-index the sum to give

$$\frac{1}{(1-x)^2} = \sum_{n=0}^{\infty} (n+1)x^n$$

Using the Ratio Test, we can show that this new power series has radius of convergence $R = 1$. The Divergence Test can then be used to show that the interval of convergence is $(-1, 1)$. These results are consistent with the following fact, the proof of which we will omit from this course.

Fact 17

When differentiating a power series, the radius of convergence remains the same. However, differentiating may **remove endpoints** from the interval of convergence.

The intuition behind this fact goes as follows (assuming the power series is centred at $x = 0$ for simplicity). Before differentiating, the power series converges for the values of x (excluding the endpoints) such that x^n shrinks sufficiently fast to overpower whatever growth may be occurring in the coefficients c_n as $n \rightarrow \infty$. Differentiating, roughly speaking, introduces a multiplicative factor of n in our power series, but this is still not going to overcome the x^n because an exponential function will always eventually win a battle with a polynomial function. Thus, the radius of convergence stays the same.

However, if the interval of convergence included one or both endpoints before differentiating, these may be lost because it is precisely at the endpoints that the x^n terms cease to behave as decreasing exponential functions. So, if there *was* endpoint convergence before differentiating, it *could* be lost with the introduction of a factor of n in c_n .

With all this talk of differentiating power series, you might be wondering if we can anti-differentiate power series. Of course we can! We can even infer from the previous argument that anti-differentiating a power series will not change the radius of convergence. However, it could *introduce* endpoints to the interval of convergence (if they weren't already included before anti-differentiating).

Fact 18

When anti-differentiating a power series, the radius of convergence remains the same. However, anti-differentiating may **add endpoints** to the interval of convergence.

As an example, let's revisit the following power series:

$$\frac{1}{1+x^2} = \sum_{n=0}^{\infty} (-1)^n x^{2n}, \quad x \in (-1, 1)$$

You may recognize the function on the left-hand side of this expression as the derivative of $\arctan(x)$. That is, $\frac{d}{dx}(\arctan(x)) = \frac{1}{1+x^2}$. So, if we can figure out what power series we need to differentiate to get the power series on the right-hand side of the same expression, then we will have a power series representation for $\arctan(x)$.

This is actually not hard to do because the power series behaves like a polynomial (which, remember, is why we love power series!). So, we really just need to know the antiderivative of x^{2n} . Which we do! It's $\frac{1}{2n+1}x^{2n+1}$.

We have established now that

$$\frac{d}{dx}(\arctan(x)) = \frac{d}{dx} \left(\sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1} \right)$$

We can drop the derivative operators if we introduce an arbitrary constant.

$$\arctan(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1} + C$$

Conveniently, we can determine this constant by requiring that this equality hold when $x = 0$. Doing so gives

$$\arctan(0) = 0 + C \quad \implies \quad C = 0$$

Finally, we have a power series representation for $\arctan(x)$.

$$\arctan(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1}, \quad x \in [-1, 1]$$

Remember, we said that anti-differentiating a power series could add endpoints to the interval of convergence. When $x = \pm 1$, the power series is $\sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} (\pm 1)$ which is convergent by the Alternating Series Test. That's why we've written the interval of convergence now with the endpoints included.

EXERCISE

Verify that substituting $x = 1$ into the power series for $\arctan(x)$ yields the Leibniz formula for π .

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

EXERCISE

Using the power series representation (which you can derive as a separate exercise)

$$\frac{2x}{1-x^2} = \sum_{n=0}^{\infty} 2x^{2n+1}, \quad x \in (-1, 1)$$

determine a power series for $F(x) = \ln(1-x^2)$ with interval of convergence $(-1, 1)$.

(Hint: Differentiate $F(x)$ to discover how it is related to the given power series.)

6.6.4 Working with Divergent Power Series

We're going to do something bizarre now just for fun. We start with the power series representation of $\frac{1}{1-x}$ and replace x with $-x$ to get

$$\frac{1}{1+x} = \sum_{n=0}^{\infty} (-1)^n x^n, \quad x \in (-1, 1)$$

Next, we differentiate this expression to get

$$\frac{1}{(1+x)^2} = \sum_{n=1}^{\infty} (-1)^{n+1} n x^{n-1}, \quad x \in (-1, 1)$$

This power series is only convergent for $x \in (-1, 1)$ but what happens if we set $x = 1$? If we do, it says

$$\sum_{n=1}^{\infty} (-1)^{n+1} n = 1 - 2 + 3 - 4 + \dots = \frac{1}{(1+1)^2} = \frac{1}{4}$$

Nonsense, of course. The series is clearly divergent since the terms do not tend to zero. But, more disturbingly, we are adding and subtracting integers and somehow ending up with a non-integer.

But let's keep going. Next, consider the following series:

$$\sum_{n=1}^{\infty} n = 1 + 2 + 3 + 4 + \dots$$

Another clearly divergent series, but let's suppose for a moment that it converged to some finite value s . In that case, we can do the following trickery:

$$\begin{array}{r} s = 1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + \dots \\ -4s = 0 - 4 + 0 - 8 + 0 - 12 + 0 - 16 + \dots \\ \hline -3s = 1 - 2 + 3 - 4 + 5 - 6 + 7 - 8 + \dots \end{array}$$

We've reproduced the series of alternating integers which earlier "summed" to $\frac{1}{4}$. It follows that if $-3s = \frac{1}{4}$, then $s = -\frac{1}{12}$.

To summarize, we have the following result:

$$1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + \dots = -\frac{1}{12}$$

We've really gone and made a mess of things. Infinitely many positive integers summing to a negative non-integer. Preposterous!

This is where things get unsettling for a moment. There is a phenomenon in quantum physics called the Casimir effect for which the theory underlying the effect involves computing the sum of the positive integers. If you just accept that the sum is infinite, then theory does not agree with experimental results. However, if you take this sum to be equal to $-\frac{1}{12}$, then theory and experiment agree perfectly!

Clearly, there's more to the story. The sum of all positive integers truly is divergent in the standard sense. However, this result does come out more rigorously through a process

called analytic continuation. There, we work with functions on the complex plane and if you know how something called an analytic function behaves in some region of the complex plane, then it can be uniquely extended throughout the rest of the complex plane. One such function is the Riemann-Zeta function

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

Here s is a complex-valued number and the function is defined for $\operatorname{Re}(s) > 1$, but can be analytically continued elsewhere including $s = -1$ where the series becomes $\sum_{n=1}^{\infty} n$ and takes the value $-\frac{1}{12}$.

Chapter 7

Taylor Polynomials

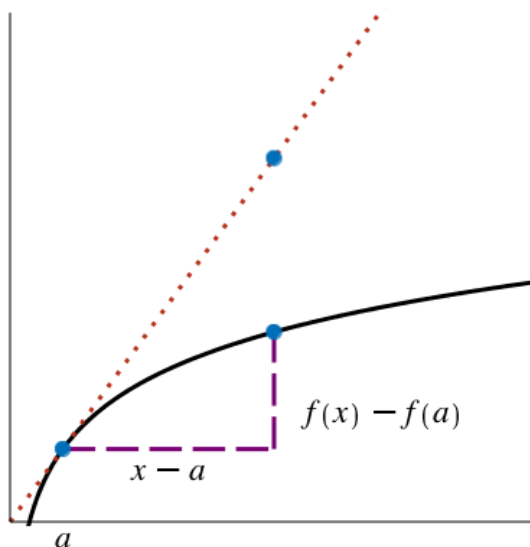
7.1 Introduction to Taylor Polynomials

We can use the linearization of a function to approximate values of that function near the point $x = a$ where we construct the linearization. Geometrically, the graph of the linearization is the line tangent to the graph of the function at $x = a$.

Recall, we denote the linearization of $f(x)$ at $x = a$ by $L_a(x)$ and it is given by

$$L_a(x) = f(a) + f'(a)(x - a)$$

Observe in the graph below that the linearization may provide good approximations of a function near the point $x = a$, but those approximations can quickly become a very poor as you move farther away.



The reason a linearization can only work so well is simple. The function we're trying to approximate is non-linear; we can't guarantee that its graph won't eventually curve away from any tangent line we draw. So, if we want a better approximation, we're going to need

to go beyond using linear functions. But let's at least see how well we can do if we stick with polynomials.

As a trial run of this idea, let's construct a quadratic approximation of $f(x) = e^x$ centred at $x = 0$. For reference, the linearization was $L_0(x) = 1 + x$.

We will denote our quadratic approximation by $Q_0(x)$ and let's allow it to be an arbitrary quadratic function for the moment. That is,

$$Q_0(x) = c_0 + c_1x + c_2x^2$$

where c_0 , c_1 , and c_2 are constants that we'll fix based on features we'd like our quadratic approximation function to exhibit.

To start, we'd like $Q_0(x)$ to take the same value as $f(x)$ at $x = 0$. This gives

$$Q_0(0) = f(0) \quad \implies \quad c_0 + c_1(0) + c_2(0)^2 = e^0 \quad \implies \quad c_0 = 1$$

Next, we'd like the derivatives of $Q_0(x)$ and $f(x)$ to be the same at $x = 0$ because it would probably be nice if their tangent lines coincided at this point.

$$Q_0'(0) = f'(0) \quad \implies \quad c_1 + 2c_2(0) = e^0 \quad \implies \quad c_1 = 1$$

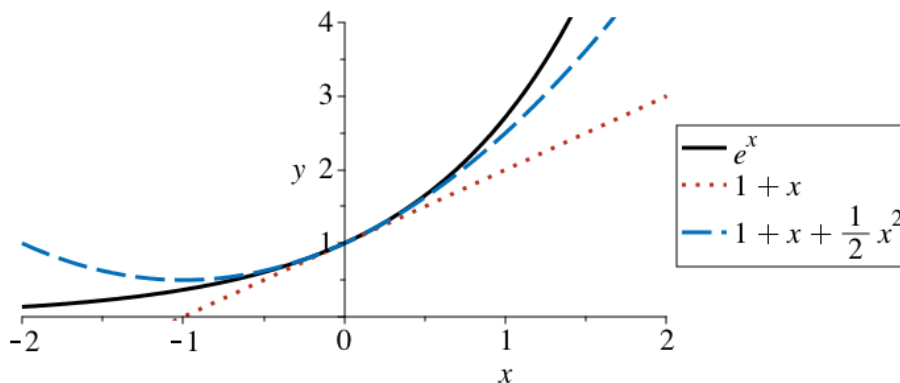
Notice, our quadratic approximation is turning out to just be the linear approximation plus a quadratic "adjustment" term.

$$Q_0(x) = 1 + x + c_2x^2 = L_0(x) + c_2x^2$$

Now, how do we determine the coefficient of the quadratic piece, c_2 ? The natural extension of what we did for the linearization is to require the approximation and the function to have the same *second* derivatives at $x = 0$. This is also perfectly sensible geometrically. The second derivative determines how the function curves at a point. This is why we use the second derivative to determine concavity. Imposing this constraint gives

$$Q_0''(0) = f''(0) \quad \implies \quad 2c_2 = e^0 \quad \implies \quad c_2 = \frac{1}{2}$$

Now we have our quadratic approximation, $Q_0(x) = 1 + x + \frac{1}{2}x^2$. Plotting both the linear and quadratic approximations against our function, we see that the quadratic approximation, indeed, provides a better fit to $y = e^x$ near $x = 0$.



It seems we're on to something here. In particular, if we want an even better approximation, we can extend this procedure and construct higher-order polynomial approximations for an arbitrary function. The key insight being that we can fix the coefficients of a polynomial approximation by requiring the non-zero derivatives of the polynomial match the derivatives of the function at each order at the point where the approximation is centred.

EXERCISE

Find the cubic approximation of $f(x) = e^x$ centred at $x = 0$ by adding an arbitrary cubic term to the quadratic approximation and then requiring the third derivative of the approximation and $f^{(3)}(x)$ take the same value at $x = 0$.

But why stop with a finite degree polynomial? What happens if we use a power series? If a power series and a function have the same derivatives at every order, are they not then just two different ways of describing the same function? Let's explore this idea now.

7.2 Taylor Series

Suppose $f(x)$ is a function which can be represented by a power series in $(x - a)$ with radius of convergence R (keeping in mind that R could be infinite). That is, there exists coefficients $\{c_n\}$ such that

$$f(x) = \sum_{n=0}^{\infty} c_n(x - a)^n, \quad |x - a| < R$$

Observe that if we set $x = a$, we get

$$f(a) = c_0 + c_1(0) + c_2(0)^2 + \cdots = c_0$$

This means that coefficient c_0 must be equal to $f(a)$. Let's see if we can determine more coefficients.

Differentiating the power series representation, we have

$$f'(x) = \sum_{n=0}^{\infty} n c_n(x - a)^{n-1}, \quad |x - a| < R$$

Again, setting $x = a$ gives

$$f'(a) = c_1 + 2c_2(0) + 3c_3(0)^2 + \cdots = c_1$$

We now have $c_1 = f'(a)$. This is working well. Let's differentiate again.

$$f''(x) = \sum_{n=0}^{\infty} n(n-1) c_n(x - a)^{n-2}, \quad |x - a| < R$$

It follows that

$$f''(a) = 2c_2 + 3 \cdot 2c_3(0) + 4 \cdot 3c_4(0)^2 \cdots = 2c_2$$

so $c_2 = \frac{1}{2}f''(a)$.

Each time we take a derivative and set $x = a$, we determine a new coefficient!

Repeating this procedure a few more times gives

$$\begin{aligned} f^{(3)}(a) = 3 \cdot 2c_3 &\implies c_3 = \frac{1}{2 \cdot 3} f^{(3)}(a) \\ f^{(4)}(a) = 4 \cdot 3 \cdot 2c_4 &\implies c_4 = \frac{1}{2 \cdot 3 \cdot 4} f^{(4)}(a) \\ f^{(5)}(a) = 5 \cdot 4 \cdot 3 \cdot 2c_5 &\implies c_5 = \frac{1}{2 \cdot 3 \cdot 4 \cdot 5} f^{(5)}(a) \end{aligned}$$

A pattern is emerging from which we can infer the following general formula. (As an exercise, you can also prove this by induction.)

$$c_n = \frac{f^{(n)}(a)}{n!}$$

Let's summarize our findings.

Theorem 1 (Taylor Series of $f(x)$)

If coefficients c_n exist such that

$$f(x) = \sum_{n=0}^{\infty} c_n (x-a)^n, \quad |x-a| < R$$

then for all n

$$c_n = \frac{f^{(n)}(a)}{n!}$$

We call a power series representation of f a **Taylor series**.

REMARK

When $x = 0$, we get the special case $f(x) = \sum_{n=0}^{\infty} c_n x^n$ with $c_n = \frac{f^{(n)}(0)}{n!}$. This special case power series representation of f - that is, a Taylor series centred at $x = 0$ - is often referred to as a Maclaurin series.

Example 1

Determine the Taylor series for $f(x) = e^x$ centred at $x = 0$.

Solution: Since e^x is unchanged by differentiation, we have $f^{(n)}(x) = e^x$. This gives the coefficients

$$c_n = \frac{f^{(n)}(0)}{n!} = \frac{1}{n!}$$

Therefore, the Taylor series for e^x centred at $x = 0$ is

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

It is important to keep in mind that we can only use a power series for values of x for which it is convergent. To find the radius of convergence of this Taylor series, we can apply the Ratio Test. We let $a_n = \frac{|x|^n}{n!}$ and compute

$$L = \lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| = \lim_{n \rightarrow \infty} \left| \frac{\frac{x^{n+1}}{(n+1)!}}{\frac{x^n}{n!}} \right| = |x| \lim_{n \rightarrow \infty} \frac{1}{n+1} = 0$$

Since $L < 1$, the power series is convergent regardless of the value of x . We will prove this shortly but, incredibly, this Taylor series serves as an exact equivalent of the function e^x for all x .

Let's look at another example.

Example 2

Determine the Taylor series centred at $x = 0$ for $f(x) = \sin(x)$.

Solution: Let's start by computing a handful of derivatives of $f(x)$ and evaluating them at $x = 0$.

$$\begin{aligned} f(x) = \sin(x) &\implies f(0) = 0 \\ f'(x) = \cos(x) &\implies f'(0) = 1 \\ f''(x) = -\sin(x) &\implies f''(0) = 0 \\ f^{(3)}(x) = -\cos(x) &\implies f^{(3)}(0) = -1 \\ f^{(4)}(x) = \sin(x) &\implies f^{(4)}(0) = 0 \end{aligned}$$

Observe that since $f^{(4)}(x) = f(x)$, the derivatives begin to repeat and we can infer the following pattern for odd terms and even terms, respectively:

$$\begin{aligned} f^{(2n+1)}(0) &= (-1)^n \\ f^{(2n)}(0) &= 0 \end{aligned}$$

Therefore, the Taylor series for $\sin(x)$ centred at $x = 0$ is

$$\begin{aligned} \sin(x) &= \sum_{n=0}^{\infty} \frac{f^{(n)}(0)}{n!} x^n \\ &= \sum_{n=0}^{\infty} \frac{\cancel{f^{(2n)}(0)}}{(2n)!} x^{2n} + \sum_{n=0}^{\infty} \frac{f^{(2n+1)}(0)}{(2n+1)!} x^{2n+1} \\ &= \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1} \\ &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \end{aligned}$$

Note, in the second step of this calculation, we split the series into one sum involving only terms with even powers of x and another involving only terms with odd powers of x .

EXERCISE

Show that the Taylor series for $f(x) = \cos(x)$ centred at $x = 0$ is

$$\cos(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} x^{2n} = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$$

EXERCISE

Show that the Taylor series centred at $x = 0$ for both $\sin(x)$ and $\cos(x)$ converge for all x .

It is worth pausing for a moment to soak in these past few results. It appears that we can represent the functions e^x , $\sin(x)$, and $\cos(x)$ with infinite degree polynomials for any value of x . Polynomials are often much easier to work with than transcendental functions, particularly when we start looking at integration, so we can expect that these Taylor series representations will be very useful.

Of course, we can use the same approach to find Taylor series for other functions too. In fact, we already did in the previous section on power series; we just didn't call them Taylor series at the time. The following table compiles the Taylor series representations centred at $x = 0$ for some common functions along with their radii of convergence.

Function	Taylor series (centred at $x = 0$)	Radius of Convergence
e^x	$\sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots$	$R = \infty$
$\sin(x)$	$\sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$	$R = \infty$
$\cos(x)$	$\sum_{n=0}^{\infty} \frac{(-1)^n x^{2n}}{(2n)!} = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots$	$R = \infty$
$\frac{1}{1-x}$	$\sum_{n=1}^{\infty} x^n = 1 + x + x^2 + x^3 + \dots$	$R = 1$
$\ln(1+x)$	$\sum_{n=1}^{\infty} \frac{(-1)^{n-1} x^n}{n} = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$	$R = 1$
$\arctan(x)$	$\sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{2n+1} = x - \frac{x^3}{3} + \frac{x^5}{5} - \frac{x^7}{7} + \dots$	$R = 1$

Of course, not all Taylor series need to be centred at $x = 0$. The reasons that we might choose a non-zero centre when computing the Taylor series of a function will become more clear when we talk about Taylor polynomials, but one reason we can give now is that $x = 0$ may not even be in the domain of the function as in the next example.

Example 3

Determine the Taylor series of $f(x) = \ln(x)$ centred at $x = 1$.

Solution: We begin as usual by computing derivatives of $f(x)$, but now we will evaluate them at $x = 1$.

$$\begin{aligned} f(x) = \ln(x) &\implies f(1) = 0 \\ f'(x) = \frac{1}{x} &\implies f'(1) = 1 \\ f''(x) = -\frac{1}{x^2} &\implies f''(1) = -1 \\ f^{(3)}(x) = \frac{2}{x^3} &\implies f^{(3)}(1) = 2 \\ f^{(4)}(x) = -\frac{3 \cdot 2}{x^4} &\implies f^{(4)}(1) = -3 \cdot 2 \\ f^{(5)}(x) = \frac{4 \cdot 3 \cdot 2}{x^5} &\implies f^{(5)}(1) = 4 \cdot 3 \cdot 2 \end{aligned}$$

Observe that the zeroth-order term in the Taylor series will be zero and the remaining terms will follow the pattern

$$f^{(n)}(1) = (-1)^{n-1}(n-1)!, \quad n \geq 1$$

Therefore, the Taylor series for $\ln(x)$ centred at $x = 1$ is

$$\begin{aligned} \ln(x) &= \sum_{n=1}^{\infty} \frac{(-1)^{n-1}(n-1)!}{n!} (x-1)^n \\ &= \sum_{n=1}^{\infty} \frac{(-1)^{n-1}}{n} (x-1)^n \\ &= (x-1) - \frac{(x-1)^2}{2} + \frac{(x-1)^3}{3} - \frac{(x-1)^4}{4} + \dots \end{aligned}$$

EXERCISE

Show that the Taylor series for $f(x) = \sin(x)$ centred at $x = \frac{\pi}{2}$ is

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} \left(x - \frac{\pi}{2}\right)^{2n}$$

If you stared long enough at the Taylor series for $\sin(x)$ centred at $x = \pi/2$, you may have noticed that it looks just like the Taylor series for $\cos(x)$ centred at $x = 0$ in that they have the exact same coefficients. This is no coincidence. To see why, consider the Taylor series for $\cos(y)$ centred at $y = 0$ and then make the substitution $y = x - \frac{\pi}{2}$, you'll get the expression

$$\cos\left(x - \frac{\pi}{2}\right) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} \left(x - \frac{\pi}{2}\right)^{2n}$$

but $\cos\left(x - \frac{\pi}{2}\right) = \sin(x)$, so this is in perfect agreement with the result of the previous exercise.

Similarly, if you take the Taylor series for $\ln(1+y)$ centred at $y = 0$ and make the substitution $y = x - 1$, you'll get the Taylor series for $\ln(x)$ centred at $x = 1$ that we found in the previous example. The point is not that the latter calculations were not worthwhile, but rather a reminder that since a Taylor series is just a power series representation of a function, we can determine new Taylor series by modifying known ones.

Example 4 Determine the Taylor series for $x \sin(x^2)$ centred at $x = 0$.

Solution: Recall, the Taylor series for $\sin(y)$ centred at $y = 0$ is

$$\sin(y) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} y^{2n+1}$$

Now make the substitution $y = x^2$.

$$\sin(x^2) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} (x^2)^{2n+1} = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{4n+2}$$

Finally, multiply both sides of the previous expression by x .

$$x \sin(x^2) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{4n+3}$$

EXERCISE

Determine the Taylor series for e^{-x^2} centred at $x = 0$ by starting with the Taylor series for e^x centred at $x = 0$.

7.3 Taylor Polynomials

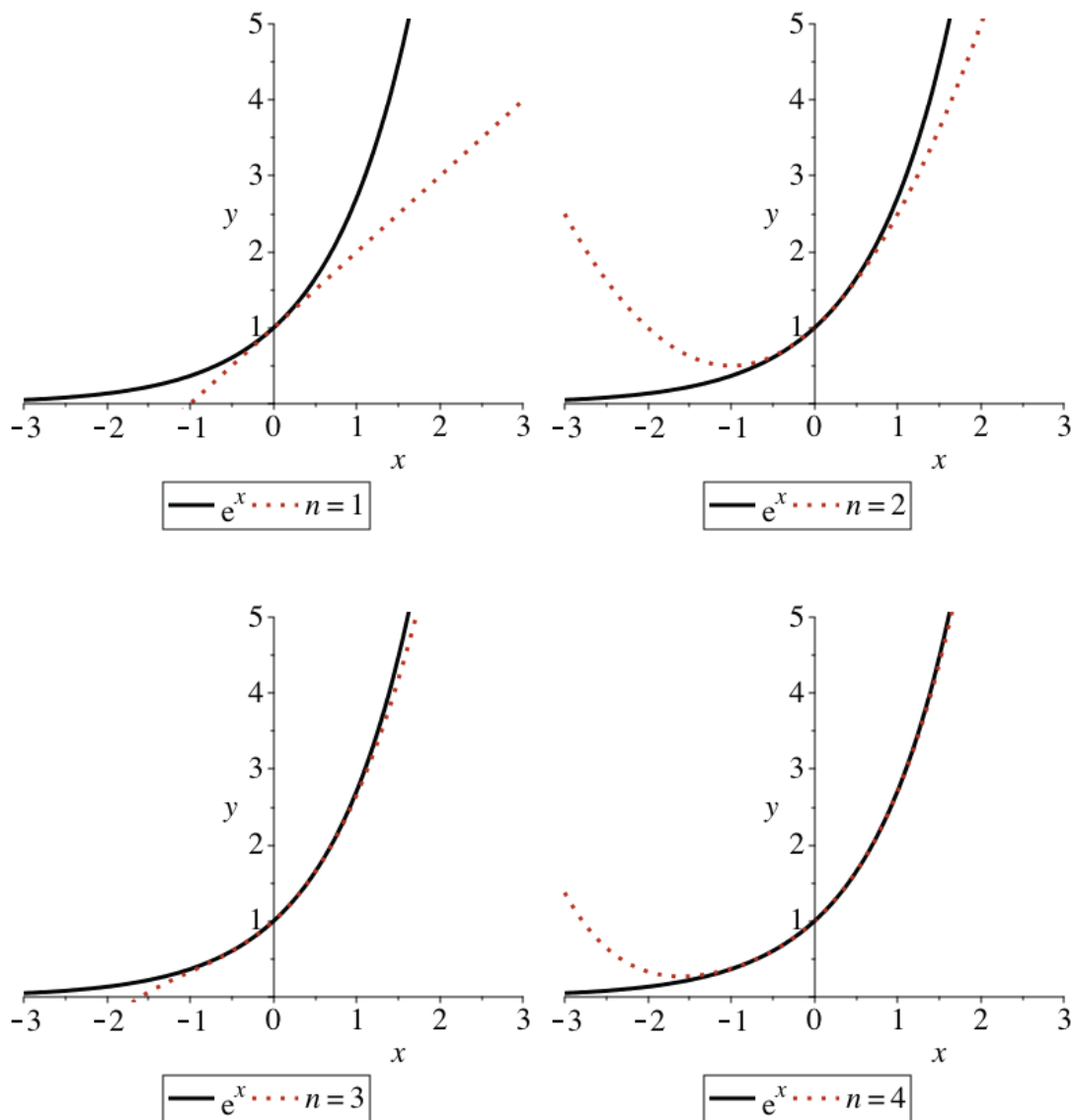
Let's return to the idea of approximating functions. We've defined a Taylor series so that its derivatives match that of the function it is representing at every order at the point where the Taylor series is centred. Therefore, if we simply truncate a Taylor series - that is, consider just a partial sum of the power series - we will get a finite-degree polynomial which should serve to approximate a function in some neighbourhood of where the Taylor series is centred.

For example, the Taylor series for $f(x) = e^x$ centred at $x = 0$ is

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots$$

If we truncate this sum after a finite number of terms, we should get increasingly better approximations for e^x . We'll justify this more rigorously soon, but it's easily enough to

plot a few graphs and see if this idea passes a visual inspection. Below we plot $y = e^x$ against the first, second, third, and fourth-degree polynomials obtained by truncating the Taylor series. Observe also that these truncations correspond to the linear approximation (linearization), quadratic approximation, cubic approximation, and quartic approximation of e^x centred at $x = 0$, respectively.



Observe that each additional term in the Taylor series improves the approximation. At least from this example, the idea of truncating a Taylor series to get an approximation seems worthwhile, so let's give these truncated polynomials a name.

Definition 7.3.1
Taylor Polynomial

Let the function $f(x)$ be at least n times differentiable at $x = a$, then we call $T_n(x)$ the n -th degree Taylor polynomial for f centred at $x = a$ where

$$T_n(x) = \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (x - a)^k$$

REMARK

If a function can be differentiated over and over again to no end (that is, $f^{(n)}(a)$ exists for all n), then we can write down a Taylor series centered at a as

$$T(x) = \sum_{k=0}^{\infty} c_k(x-a)^k$$

with $c_k = \frac{f^{(k)}(a)}{k!}$. Notice that the Taylor polynomials play the role of a partial sum, since

$$T_n(x) = \sum_{k=0}^n c_k(x-a)^k.$$

So, we ought to be able to think of the Taylor polynomials as better and better approximations of the original function, which we can, as we will see soon. This is analogous to how the partial sums of a convergent series provide better and better approximations of the value to which the series converges.

Example 5

Determine the Taylor polynomials with odd degree up to degree $n = 7$ of $f(x) = \sin(x)$ centred at $x = 0$.

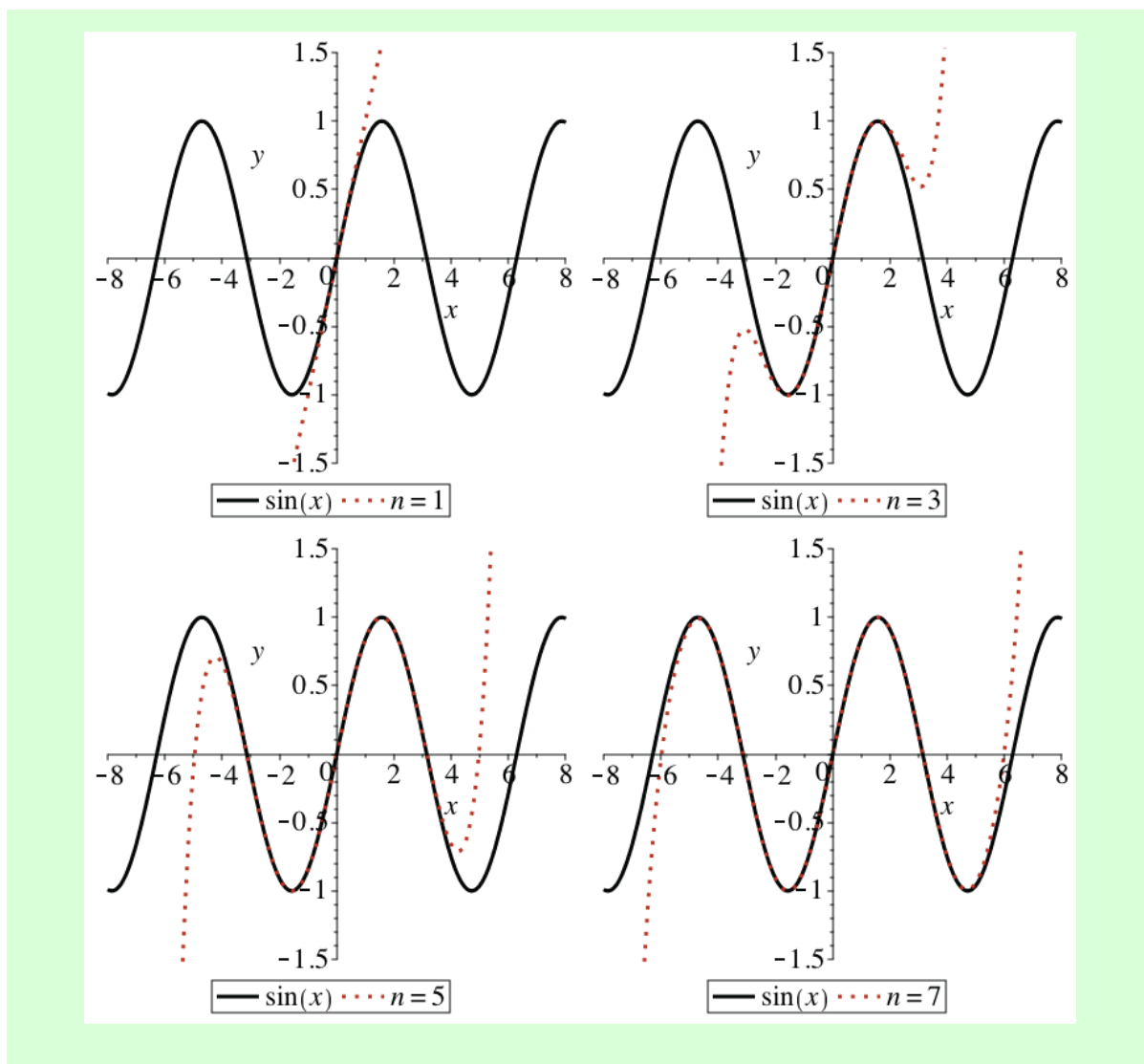
Solution: Recall the Taylor series centred at $x = 0$ for $\sin(x)$.

$$\sin(x) = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} x^{2n+1} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

Rather than compute the Taylor polynomials from scratch, we can truncate the Taylor series as needed.

$$\begin{aligned} T_1(x) &= x \\ T_3(x) &= x - \frac{x^3}{3!} \\ T_5(x) &= x - \frac{x^3}{3!} + \frac{x^5}{5!} \\ T_7(x) &= x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} \end{aligned}$$

Observe in the following plots that the higher the degree of the Taylor polynomial, the better it works as an approximation for $\sin(x)$.



EXERCISE

Determine the 2nd-degree Taylor polynomial, $T_2(x)$, centred at $x = 4$ for $f(x) = \sqrt{x}$.

7.4 Taylor's Inequality

It certainly seems like these Taylor polynomials do a good job of approximating functions, but we should really be more rigorous about this. For that purpose, we make the following definition.

Definition 7.4.1 Taylor Polynomial Remainder

Let the function $f(x)$ be at least n times differentiable at $x = a$, then we call $R_n(x)$ the n -th degree **Taylor polynomial remainder** for the Taylor polynomial $T_n(x)$ of f centred at $x = a$ and, specifically,

$$R_n(x) = f(x) - T_n(x)$$

REMARK

The remainder $R_n(x)$ gives the signed difference between the function $f(x)$ and the approximation $T_n(x)$. If we are only interested in the magnitude of the error on an approximation, then we can use $|R_n(x)|$.

The following theorem establishes an upper bound on the size of the remainder.

Theorem 2 (Taylor's Inequality)

Let I be an interval containing $x = a$. If M is the maximum value of $|f^{(n+1)}(x)|$ on the interval I , then the remainder, $R_n(x)$, of the n -th degree Taylor polynomial, $T_n(x)$, centred at $x = a$ satisfies

$$|R_n(x)| \leq \frac{M}{(n+1)!} |x-a|^{n+1}$$

everywhere on the interval I .

REMARK

When applying Taylor's inequality, it is ideal to take $x = a$ to be one endpoint of the interval I and the other endpoint to be the point where the Taylor polynomial is being used. We cannot take it to be any smaller than this and taking the interval to be any larger can only serve to increase the bound on the remainder.

We will hold off on giving a proof of Taylor's inequality until we have covered integration. Having the Fundamental Theorem of Calculus at our disposal makes the proof cleaner and easier to understand. Instead, let's proceed to see how the theorem can be used to establish an error bound on an approximation.

Example 6

Use the second-degree Taylor polynomial for $f(x) = \ln(x)$ centred at $x = 1$ to approximate $\ln(1.5)$ and provide an upper bound on the error.

Solution: We previously found the Taylor series for $\ln(x)$ centred at $x = 1$. Truncating the series appropriately gives

$$T_2(x) = (x-1) - \frac{1}{2}(x-1)^2$$

We also note for later that $f^{(3)}(x) = -\frac{2}{x^3}$.

Using the second-degree Taylor polynomial, we find

$$\ln(1.5) \approx T_2(1.5) = (1.5-1) - \frac{1}{2}(1.5-1)^2 = 0.375$$

Now let's apply Taylor's inequality to determine an upper bound on the error of this approximation. Since the Taylor polynomial is centred at $x = 1$ and we have evaluated it at $x = 1.5$, we must find the maximum value, M , of $|f^{(3)}(x)|$ on the interval $[1, 1.5]$.

On this interval, $|f^{(3)}(x)| = \frac{2}{x^3}$ is a decreasing function, so it attains its maximum value at $x = 1$. Therefore, if we take $M = f^{(3)}(1) = 2$, the condition $M \leq |f^{(3)}(x)|$ holds for $x \in [1, 1.5]$. Then, by Taylor's inequality

$$|R_2(1.5)| \leq \frac{2}{3!}|1.5 - 1|^3 = \frac{1}{24} = 0.041\bar{6}$$

Rounding this value up at the third decimal place, we get

$$\ln(1.5) = 0.375 \pm 0.042$$

For comparison, the actual value is $\ln(1.5) \approx 0.405$, so the actual remainder is approximately 0.030 in agreement with the upper bound determined by Taylor's inequality.

EXERCISE

The second-degree Taylor polynomial for $f(x) = \sqrt{x}$ centred at $x = 4$ is

$$T_2(x) = 2 + \frac{1}{4}(x - 4) - \frac{1}{64}(x - 4)^2$$

Use this polynomial to approximate $\sqrt{5}$. Next, apply Taylor's inequality to determine an upper bound on $|R_2(5)|$.

The notion of a Taylor polynomial remainder also gives us a focus for proving that the Taylor series of a function - which we'll denote here by $T(x)$ - is equivalent to the function. In particular, since $f(x) = T_n(x) + R_n(x)$, it follows that

$$\lim_{n \rightarrow \infty} R_n(x) = 0 \quad \implies \quad f(x) = \lim_{n \rightarrow \infty} T_n(x) = T(x)$$

Let's test this idea on a concrete example.

Example 7

Prove that $f(x) = e^x$ is equal to its Taylor series centred at $x = 0$.

Proof: Let $I = [a, b]$ be an interval containing $x = 0$. Since $|f^{(n+1)}(x)| = e^x \leq e^b$ on I for all $n \geq 0$, then if we let $M = e^b$, we have $|f^{(n+1)}(x)| \leq M$ on I .

It follows by Taylor's inequality that

$$|R_n(x)| \leq \frac{e^b}{(n+1)!}|x|^{n+1}$$

Let $a_n = \frac{|x|^{n+1}}{(n+1)!}$ and observe that

$$\lim_{n \rightarrow \infty} \left| \frac{a_{n+1}}{a_n} \right| = |x| \lim_{n \rightarrow \infty} \frac{1}{n+2} = 0$$

Therefore, the sequence $\{a_n\}$ is convergent by the Ratio Test. By the contrapositive of the Divergence Test, this is only possible if the terms in the sequence $\{a_n\}$ tend to zero. Taking the limit of the inequality above, we can therefore say

$$\lim_{n \rightarrow \infty} |R_n(x)| \leq \lim_{n \rightarrow \infty} \frac{e^b}{(n+1)!} |x|^{n+1} = e^b \lim_{n \rightarrow \infty} \frac{|x|^{n+1}}{(n+1)!} = 0$$

Therefore, by the Squeeze theorem, $\lim_{n \rightarrow \infty} |R_n(x)| = 0$. We conclude that the Taylor series of e^x centred at $x = 0$ is equal to e^x on every interval containing $x = 0$ which, equivalently, means for all $x \in \mathbb{R}$.

Based on this example, we can also see how to extend this proof more generally. In particular, for a function $f(x)$, if on an interval I containing $x = a$, the derivatives $f^{(n)}(x)$ for $n \geq 0$ are all bounded, then there exists some number M such that $|f^{(n+1)}(x)| \leq M$ for all $n \geq 0$. We can then show that the limit as $n \rightarrow \infty$ of $|R_n(x)|$ must be zero in the same way we did in the previous example. (The only differences are that we would not replace M with e^b and we would have $|x - a|$ instead of $|x|$.)

7.5 Taylor Polynomial Approximations in Science

Many scientific laws are stated in terms of relationships between variables (e.g., the ideal gas law, Hooke's law, Ohm's law). When a relationship is difficult to work with, we can Taylor expand the functions involved to create an approximate form of the relationship. For example, Snell's law of refraction relates the indices of refraction between two media, n_1 and n_2 , to the angles of incidence θ_1 and θ_2 with the equation $n_1 \sin(\theta_1) = n_2 \sin(\theta_2)$. If we assume the angles of incidence are small such that $\sin(\theta_1) \approx \theta_1$ and $\sin(\theta_2) \approx \theta_2$, then we get the simplified relation $n_1 \theta_1 \approx n_2 \theta_2$.

On the flip side of this coin is also the idea that many of the laws we work with are, in fact, already approximating more intricate dynamics. For example, in Newtonian ("everyday") physics, the energy of a mass m with speed v is given by the familiar equation

$$E = \frac{1}{2}mv^2$$

However, in special relativity, the energy of a mass m moving with speed v is given by

$$E = \frac{mc^2}{\sqrt{1 - \left(\frac{v}{c}\right)^2}}$$

where c is the speed of light. How do we reconcile these two seemingly different formulas for the same thing?

First, we make an important observation. When $v = 0$, the Newtonian formula gives $E = 0$ while the relativistic formula gives $E = mc^2$. So, relativity says an object of mass m has a "rest mass energy" of mc^2 while in Newtonian physics, it has zero energy. We could address this discrepancy by just saying "let's just add mc^2 to the Newtonian value"; doing so won't break anything in the Newtonian theory, but the general formulas would still disagree.

Newtonian physics breaks down at relativistic speeds (i.e., speeds near the speed of light) while special relativity should work for any possible speed. With this in mind, let's inspect

what happens to the relativistic formula for energy if we Taylor expand it around $v = 0$. To simplify this process, let's make the substitution $x = \left(\frac{v}{c}\right)^2$. Then we can rewrite the relativistic formula as

$$E(x) = mc^2(1 - x)^{-1/2}$$

Now let's work out the second-degree Taylor polynomial of $f(x) = (1 - x)^{-1/2}$ centred at $x = 0$ (which, remember, corresponds to $v = 0$).

With a bit of work, we find

$$\begin{aligned} f(x) &= (1 - x)^{-1/2} & \implies & f(0) = 1 \\ f'(x) &= \frac{1}{2}(1 - x)^{-3/2} & \implies & f'(0) = \frac{1}{2} \\ f''(x) &= \frac{3}{4}(1 - x)^{-5/2} & \implies & f''(0) = \frac{3}{4} \end{aligned}$$

This gives the Taylor polynomial

$$T_2(x) = 1 + \frac{1}{2}x + \frac{3}{8}x^2$$

Substituting back in $x = \left(\frac{v}{c}\right)^2$, we obtain a Taylor approximation of E in terms of speed v .

$$E \approx mc^2 \left(1 + \frac{1}{2} \left(\frac{v}{c}\right)^2 + \frac{3}{8} \left(\frac{v}{c}\right)^4 \right) = mc^2 + \frac{1}{2}mv^2 + \frac{3}{8} \frac{mv^4}{c^2}$$

We can interpret each term as follows:

- mc^2 is the rest mass energy
- $\frac{1}{2}mv^2$ is the Newtonian kinetic energy
- $\frac{3}{8} \frac{mv^4}{c^2}$ is the first correction to the energy due to relativistic effects

So we see, the Newtonian formula is just an approximation - specifically, the first-order Taylor polynomial - of the relativistic formula (without the rest mass energy term). This approximation works so well for “everyday” physics because when a particle is moving at non-relativistic speeds (i.e., $v \ll c$), the relativistic corrections which include powers of $\left(\frac{v}{c}\right)^2$ will be very small.

EXERCISE

The atmospheric pressure P at altitude z (where $z = 0$ corresponds to sea level) obeys the relation

$$P(z) = P_0 e^{-z/H}$$

where P_0 is atmospheric pressure at sea level and $H \approx 8.5\text{km}$ is a constant. Determine the first-degree Taylor polynomial of $P(z)$ centred at $z = 0$. Such a formula would be useful for approximating atmospheric pressure for altitudes $z \ll H$.

7.6 Binomial Series

In the earlier example when we computed a Taylor polynomial for the relativistic energy of a moving mass, we encountered a function of the form $f(x) = (1+x)^k$. Functions of this form with $|x| < 1$ occur frequently in practical applications and it is often desirable to replace the function with a Taylor polynomial centred at $x = 0$. This can be done in general.

To determine the coefficients of the Taylor series, we compute the first few to identify a pattern:

$$\begin{aligned} f(x) &= (1+x)^k & \implies & f(0) = 1 \\ f'(x) &= k(1+x)^{k-1} & \implies & f'(0) = k \\ f''(x) &= k(k-1)(1+x)^{k-2} & \implies & f''(0) = k(k-1) \end{aligned}$$

The derivatives evaluated at $x = 0$ behave the same as those of the function x^k so

$$f^{(n)}(0) = k(k-1)\dots(k-n+1)$$

This gives coefficients

$$c_n = \frac{f^{(n)}(0)}{n!} = \frac{k(k-1)\dots(k-n+1)}{n!}$$

When k is a non-negative integer, $c_n = 0$ for $n > k$ (since the coefficients all contain a $(k-k)$ term) and for $n \leq k$, we can simplify c_n to

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

Here $\binom{k}{n}$ is a binomial coefficient (i.e., “ k choose n ”). This means, the series terminates after a finite number of terms. In particular, when k is a non-negative integer

$$(1+x)^k = \sum_{n=0}^k \binom{k}{n} x^n \quad (7.1)$$

This is, of course, just what we should expect. After all, how else could we write $(1+x)^2$ as a polynomial other than just expanding it to get $1+2x+x^2$?

In contrast, if k is any real number which is *not* a positive integer, the Taylor Series has infinitely many terms. To be able to walk away with a nice tidy formula applicable to all real numbers, we generalize our binomial coefficients to be defined for any real-valued k to be

$$\binom{k}{0} = 1 \quad \text{and} \quad \binom{k}{n} = \frac{k(k-1)\dots(k-n+1)}{n!}$$

This allows us to write

$$(1+x)^k = \sum_{n=0}^{\infty} \binom{k}{n} x^n = 1 + kx + \frac{k(k-1)}{2!}x^2 + \frac{k(k-1)(k-2)}{3!}x^3 + \dots$$

We call this the **binomial series**.

To determine the radius of convergence of the binomial series, we apply the Ratio Test which means computing the following limit.

$$\begin{aligned}
 L &= \lim_{n \rightarrow \infty} \left| \frac{\binom{k}{n+1} x^{n+1}}{\binom{k}{n} x^n} \right| \\
 &= |x| \lim_{n \rightarrow \infty} \left| \frac{n!(k-n)!}{(n+1)!(k-(n+1))!} \right| \\
 &= |x| \lim_{n \rightarrow \infty} \left| \frac{(k-n)}{(n+1)} \right| \\
 &= |x|
 \end{aligned}$$

By the Ratio Test, we have convergence when $L < 1$ and divergence when $L > 1$. This means the binomial series is convergent for $|x| < 1$ and has radius of convergence equal to 1. Convergence at the endpoints of the interval $(-1, 1)$ depends on the value of k .

7.6.1 The Binomial Approximation

When $|x| \ll 1$, the terms in the binomial series decreasing in magnitude very quickly. This means we can get a good, but non-trivial, approximation by using the first-degree Taylor polynomial (or linearization). Doing so gives what is known as the **binomial approximation**

$$(1+x)^k \approx 1+kx \quad \text{for } |x| \ll 1$$

The simplicity of this approximation makes it very useful in a wide range of applications.

Example 8

Apply the binomial approximation to $f(x) = (1+x)^{1/2}$ assuming $|x| \ll 1$. Use this result to approximate $\sqrt{231}$.

Solution: The binomial approximation says for $|x| \ll 1$

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

To use this to approximate $\sqrt{231}$, we note that the perfect square nearest to 231 is $225 = 15^2$ and then perform the following manipulation.

$$\begin{aligned}
 \sqrt{231} &= \sqrt{225+6} \\
 &= \sqrt{225} \sqrt{1 + \frac{6}{225}} \\
 &= 15 \left(1 + \frac{2}{75} \right)^{1/2} \\
 &\approx 15 \left(1 + \frac{1}{2} \cdot \frac{2}{75} \right) \quad \text{since } \left| \frac{2}{75} \right| \ll 1 \\
 &= \frac{76}{5} \quad \text{or} \quad 15.2
 \end{aligned}$$

For comparison, a calculator gives $\sqrt{231} \approx 15.1987$.

EXERCISE

Use the binomial approximation to approximate $(67)^{2/3}$.

Let's finish our discussion with an application-based example where the binomial approximation proves useful.

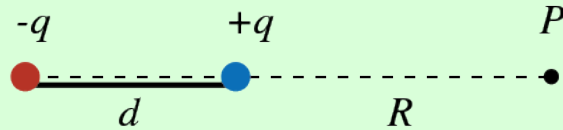
Example 9

According to Coulomb's law, the magnitude of the electric field at a distance r from a charge Q is given by

$$E = \frac{kQ}{r^2}$$

where k is known as Coulomb's constant.

Consider two charges, $+q$ and $-q$, separated by a distance d as pictured below. We call this configuration an **electric dipole**. The point P is aligned with the two charges a distance R from the positive charge and $R + d$ from the negative charge.



Use the binomial approximation to approximate the electric field strength at point P due to the electric dipole assuming $d \ll R$.

Solution: At point P , the magnitude of the electric field is the sum of $\frac{kq}{R^2}$ from the positive charge and $-\frac{kq}{(R+d)^2}$ from the negative charge.

$$E = \frac{kq}{R^2} - \frac{kq}{(R+d)^2} = \frac{kq}{R^2} \left(1 - \frac{1}{\left(1 + \frac{d}{R}\right)^2} \right) = \frac{kq}{R^2} \left(1 - \left(1 + \frac{d}{R}\right)^{-2} \right)$$

Since $d \ll R$, we can apply the binomial approximation to rewrite

$$\left(1 + \frac{d}{R}\right)^{-2} \simeq 1 - \frac{2d}{R}$$

With this, we obtain the following approximation for the magnitude of the electric field

$$E \approx \frac{kq}{R^2} \left(1 - \left(1 - \frac{2d}{R}\right) \right) = \frac{2kqd}{R^3}$$

The field from a single charge varies inversely proportional to the square of the distance. In contrast, the field here for an electric dipole varies at lowest order proportional to the cube of the distance. This means that far away from the dipole, the fields due to the individual charges largely cancel out. What remains are higher-order corrections which decay more quickly with distance from the source.

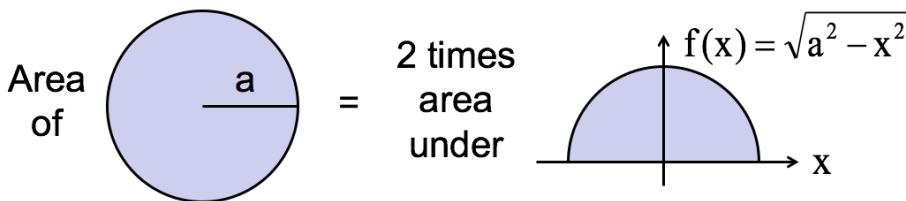
Chapter 8

Integration

Calculus found its roots in solving two problems which, aside from being geometrical in nature, are seemingly unrelated. Those problems are the “tangent problem” and the “area problem”. In the tangent problem, the goal is to find the slope of a line tangent to an arbitrary curve at a specific point. As we have seen, this problem can be solved using the tools of differentiation. For the area problem, the goal is to find the area of a region with a curved boundary. In this chapter, we’ll see that this problem is solved using integration and, incredibly, the processes which solve the tangent and area problems - that is, differentiation and integration - are opposite sides of the same coin.

8.1 Riemann Sums

We can frame the problem of finding the area of a region with a curved boundary in the language of functions. For example, suppose we would like to know the area of a circle of radius a . An equivalent problem is to find the area between the graph of the function $y = \sqrt{a^2 - x^2}$ and the x -axis on the interval $-a \leq x \leq a$.



We’re not going to calculate this area now, but we’ll adopt this convention for describing the areas. That is, we’ll define a function $f(x)$ such that the region between the graph $y = f(x)$ and the x -axis (i.e., the line $y = 0$) over a specified interval for x has the area we’re interested in. You may worry that this is restricting us to only being able to find areas of regions that have one “flat” edge on their boundary, but it is always possible to decompose any region into these simpler instances - just like we did for the area of a circle above. Moreover, we’ll eventually see how to generalize the area finding process for areas between curves rather than between an axis and a curve.

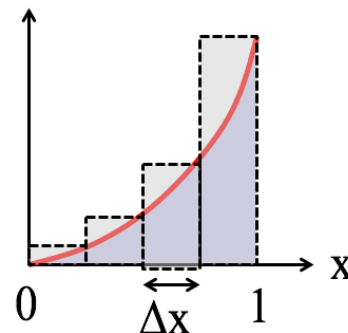
Now, it turns out that finding the area of a circle with the approach we’re about to go through is quite difficult because the function $\sqrt{a^2 - x^2}$ is not particularly nice to work

with (at least not in the Cartesian coordinate system). So let's pick another curve for our working example. We want something as simple as possible algebraically while still describing a curve. The natural choice is the parabola $y = x^2$. Let's also keep the interval as simple as possible and consider the area under this parabola over the interval $[0, 1]$.

Our goal now is to calculate the *exact* area between the curve $y = f(x) = x^2$ and the x -axis on the interval $[0, 1]$. However, let's start by approximating this area with a sequence of four equal-width rectangles.

Letting Δx be the width of each rectangle, then Δx is equal to the full interval width divided by the number of rectangles, so

$$\Delta x = \frac{1}{4}$$



The height of each rectangle is determined by the height of the graph at the right side of each rectangle. (We'll consider later on what may change if we use another point within each sub-interval to set the height of the corresponding rectangle.) Then, the heights of the rectangles from left-to-right are:

$$\left\{ \left(\frac{1}{4}\right)^2, \left(\frac{1}{2}\right)^2, \left(\frac{3}{4}\right)^2, (1)^2 \right\}$$

Summing the areas of all four rectangles, we get the following approximation for the area under the parabola:

$$A \simeq \frac{1}{4} \left(\frac{1}{4}\right)^2 + \frac{1}{4} \left(\frac{1}{2}\right)^2 + \frac{1}{4} \left(\frac{3}{4}\right)^2 + \frac{1}{4} (1)^2 = \frac{15}{32}$$

We call this type of approximation a **Riemann sum**. As we'll see soon, the exact area under the parabola on this interval is $\frac{1}{3}$. So, the approximate result $\frac{15}{32} \approx 0.47$ is not terrible, but leaves a lot of room for improvement. Observe that the reason for the error in a Riemann sum is that a finite-width rectangle will capture more or less area than is actually between the curve and the x -axis. But this error can be reduced by using a larger number of thinner rectangles.

Rather than choosing a specific number of rectangles, let's generalize the previous analysis for an arbitrary number of rectangles n . So, we need to find expressions for the width and height of each rectangle in the set.

The total width of the region is still 1 and we are using equal-width rectangles. Therefore, we have

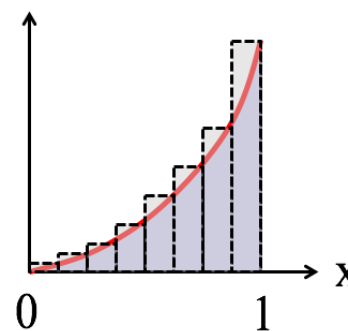
$$\Delta x = \frac{1}{n}$$

The height of each rectangle will be $f(x_i)$ where x_i is the location of the right-endpoint of the i -th rectangle. This point will follow i blocks of width Δx relative to $x = 0$, so will be located at

$$x_i = i\Delta x = \frac{i}{n}$$

Therefore, the i -th block has height

$$f\left(\frac{i}{n}\right) = \frac{i^2}{n^2}$$



With the width and height of each block determined, we can compute the total area of all the rectangles by summing over the areas of each block from $i = 1$ to $i = n$. Let us denote this quantity by R_n and call it the **right-hand Riemann sum** since it is a Riemann sum using n blocks constructed such that the right-edge of the block connects the x -axis to the function.

$$A \simeq R_n = \sum_{i=1}^n f\left(\frac{i}{n}\right) \cdot \Delta x = \sum_{i=1}^n \frac{i^2}{n^2} \cdot \frac{1}{n} = \frac{1}{n^3} \sum_{i=1}^n i^2$$

Recall that there is an elegant formula for the sum of the first n squares, which is

$$\sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6}$$

This allows us to further simplify our expression for R_n .

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

Observe that when $n = 4$, we get $R_4 = \frac{15}{32}$ in agreement with our previous calculation. But a nice thing about having this expression for R_n is we can compute its value for as large a value of n as we like to get increasingly better approximations. For example, with $n = 50$, we get $R_{50} = \frac{1717}{5000} = 0.3434$ which is getting very close to the value $\frac{1}{3}$ which I promised was the *exact* answer.

So, where do we go from here? Well, we could take even bigger values of n to get even better approximations, but doing so won't ever give us the exact area. If you're thinking then that we should take a limit, you're right!

In the limit that $n \rightarrow \infty$, we get infinitely many, infinitesimally thin rectangles and the errors associated with each rectangle vanish. The result will therefore be the *exact* area under the parabola from $x = 0$ to $x = 1$.

$$\begin{aligned} A &= \lim_{n \rightarrow \infty} R_n \\ &= \lim_{n \rightarrow \infty} \frac{2n^2 + 3n + 1}{6n^2} \\ &= \lim_{n \rightarrow \infty} \frac{2 + \frac{3}{n} + \frac{1}{n^2}}{6} \\ &= \frac{1}{3} \end{aligned}$$

EXERCISE

Determine a formula for the Riemann sum for $f(x) = x^2$ on $[0, 1]$ constructed using the left-endpoints of each block to determine the corresponding block height. We call this the **left-hand Riemann sum** and denote it L_n . Next, take the limit as $n \rightarrow \infty$ and check that you get the same answer as above. With this result, how might you use the Squeeze theorem to more rigorously justify that the area under the parabola on this interval is $\frac{1}{3}$?

We can generalize the procedure above to an arbitrary function $f(x)$ over an arbitrary closed interval $[a, b]$ as follows.

We first divide the interval into n subintervals of equal width Δx . Since the interval length is $(b-a)$, this gives

$$\Delta x = \frac{b-a}{n}$$

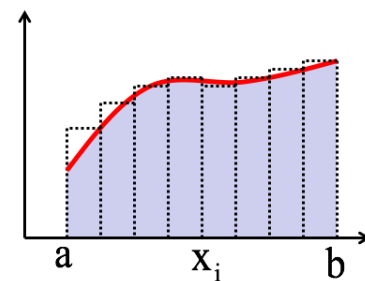
Next, the right-endpoint x_i of each rectangle will be a distance $i\Delta x$ to the right of $x = a$, so

$$x_i = a + i\Delta x = a + \frac{i}{n}(b-a)$$

Then, the height of the i -th block for the right-hand Riemann sum is $f(x_i)$. Let R_n be the right-hand Riemann sum and then take the limit as $n \rightarrow \infty$ to get an expression for the exact area under the curve over the full interval.

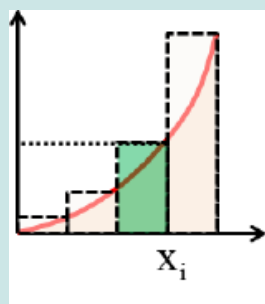
$$A = \lim_{n \rightarrow \infty} R_n = \lim_{n \rightarrow \infty} \left[\sum_{i=1}^n f(x_i) \cdot \Delta x \right]$$

A similar result holds for the left-hand Riemann sum.

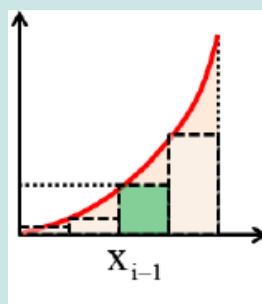


REMARK

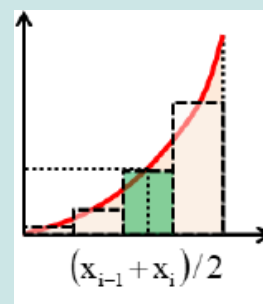
When constructing a Riemann sum approximation, the result depends on whether we use the right-endpoint, x_i , or the left-endpoint, x_{i-1} , to determine the height of the i -th block. And, in fact, there's nothing stopping us from using some other point in the interval like the midpoint $\frac{x_{i-1} + x_i}{2}$.



Right-hand Riemann sum



Left-hand Riemann sum



Midpoint Riemann sum

If using the Riemann sum to obtain an approximate result, then it may be worthwhile thinking about what point $x^* \in [x_{i-1}, x_i]$ to use to determine the block height. However, in the limit that $n \rightarrow \infty$, a Riemann sum will give the *exact* area under the curve regardless of which point in each subinterval is used to determine the block heights provided the function is continuous. In fact, it need not even be continuous under certain circumstances. We explore this idea a bit more later.

8.2 Definite Integrals

The limit of a Riemann sum as a tool for computing areas is a promising idea, so let's define some new notation to represent it.

Definition 8.2.1

definite integral

Let f be a function whose domain includes the interval $[a, b]$ and let us partition this interval into n subintervals of equal size $\Delta x = \frac{b-a}{n}$. Let $\{x_i\}_{i=0}^n$ be the endpoints of the subintervals and let $x_i^* \in [x_{i-1}, x_i]$ for $1 \leq i \leq n$. Then the **definite integral** of f from a to b is

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \sum_{i=1}^n f(x_i^*) \Delta x$$

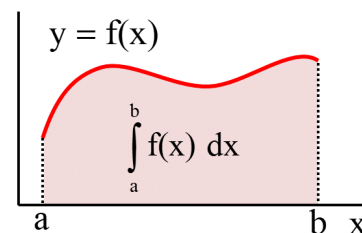
provided that this limit exists and gives the same value for all possible choices of sample points x_i^* . If this is the case, then we say f is **integrable** on $[a, b]$.

REMARK

It is not always the case that a function f satisfies the conditions of being integrable. That is, it's possible that a function f is defined on a closed interval $[a, b]$, but the limit defining the integral doesn't exist, or does not give the same value for all possible choices of sample points x_i^* . Such functions will be investigated more in Section 8.2.2.

Unless explicitly stated otherwise, if we are working with a definite integral of a function, we are assuming the function is integrable on the closed interval over which we are working.

Geometrically, we can think of the definite integral as the area between the curve $y = f(x)$ and the x -axis from $x = a$ to $x = b$.



The symbol \int is called an **integral sign**. It is essentially an elongated letter 'S', for summation, in reference to its relation to a Riemann sum. The numbers a and b attached to the integral sign are called the **limits of integration**. In this context, we refer to the function $f(x)$ as the **integrand**. The symbol dx is called the **differential** (or differential length element) and takes the place of the Riemann sum subinterval width Δx in the limit that $n \rightarrow \infty$.

When evaluating definite integrals using the Riemann sum definition, the following formulas may be useful.

Fact 1

Let n be a positive integer, then

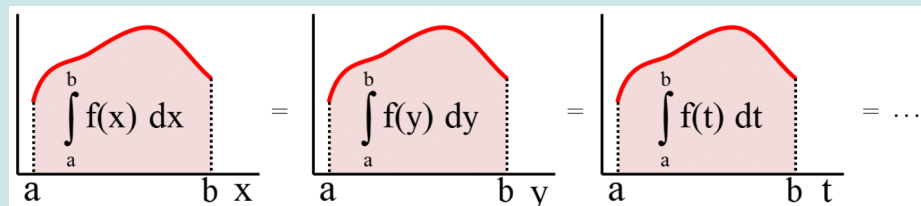
$$\sum_{i=1}^n i = \frac{n(n+1)}{2} \quad \sum_{i=1}^n i^2 = \frac{n(n+1)(2n+1)}{6} \quad \sum_{i=1}^n i^3 = \left(\frac{n(n+1)}{2} \right)^2$$

EXERCISE

Evaluate $\int_0^2 x^3 dx$ using the Riemann sum definition of the definite integral.

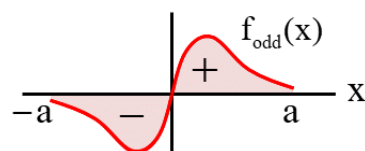
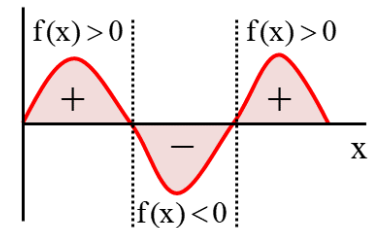
REMARK

It is important to keep in mind with all this notation floating around that given a function f and fixed values of a and b , the definite integral $\int_a^b f(x) dx$, represents the area of a fixed region and is therefore a *number*. The independent variable x is used to connect the components of the integral (bounds, integrand, and differential) and may have some contextual meaning, but computationally it is a “dummy” variable a could be replaced with another symbol without changing the result.



A point worth addressing now is how should we interpret the area “under” a curve which sits *below* the x -axis?

When the curve sits below the x -axis, the corresponding function is negative. Going back to the definition of the integral in terms of a Riemann sum, such blocks will contribute negatively to the total. Therefore, areas below the x -axis will contribute negatively to a definite integral. For this reason, we say the definite integral measures the **signed area** between a curve and the x -axis.



This property of the definite integral can be exploited to take advantage of symmetries of a function to simplify integral calculations. For example, the integral of an odd function (i.e., $f(-x) = -f(x)$) over a symmetric interval is always zero.

$$\int_{-a}^a f_{\text{odd}}(x) dx = 0$$

EXERCISE

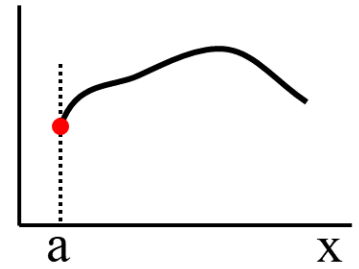
Explain why the definite integral $\int_0^{2\pi} \sin(x) dx$ is zero.

8.2.1 Properties of Definite Integrals

Let's explore some properties of definite integrals.

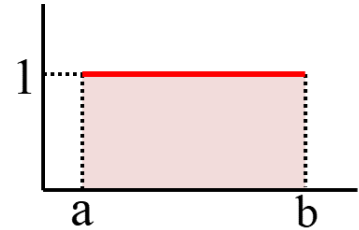
If the lower and upper limits of integration are equal, then the definite integral is zero. This is because it describes a region with zero area.

$$\int_a^a f(x) dx = 0$$



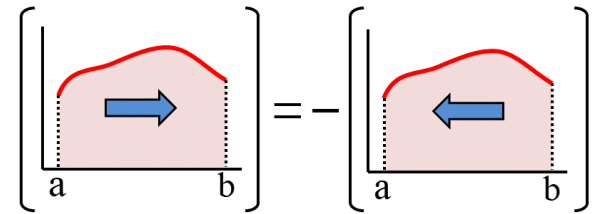
When the integrand is the constant function $f(x) = 1$, then the definite integral is equal to the upper limit of integration minus the lower limit of integration. This is because the integral describes a rectangle with width $(b - a)$ and height 1.

$$\int_a^b 1 dx = b - a$$



Swapping the limits of integration (i.e., integrating from right-to-left) reverses the sign of the definite integral. This sign change can be traced back to the definition of the width of a rectangle in the Riemann sum, Δx .

$$\int_a^b f(x) dx = - \int_b^a f(x) dx$$



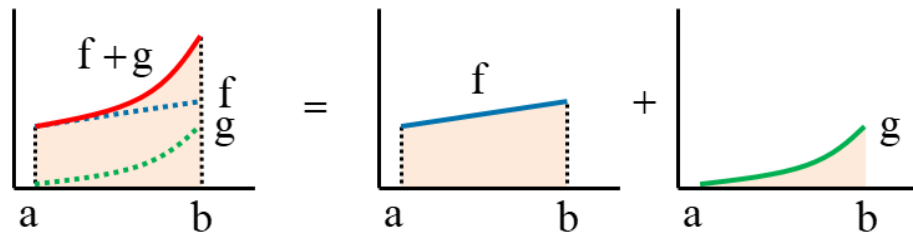
The following three properties also follow from the definition of the definite integral in terms of a Riemann sum.

A constant c multiplying the integrand can be factored out of the integral.

$$\int_a^b c f(x) dx = c \int_a^b f(x) dx$$

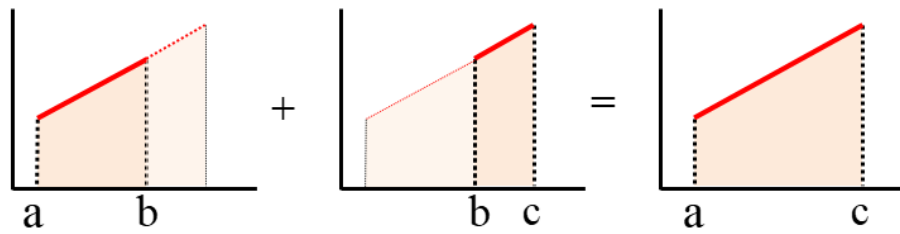
The integral of a sum is equal to the sum of the integrals.

$$\int_a^b (f(x) + g(x)) dx = \int_a^b f(x) dx + \int_a^b g(x) dx$$



Integrals can be combined or split apart if the limits of integration match up appropriately.

$$\int_a^b f(x) dx + \int_b^c f(x) dx = \int_a^c f(x) dx$$



We also have the following useful comparison properties.

- If $f(x) \geq 0$ on $[a, b]$, then $\int_a^b f(x) dx \geq 0$.
- If $f(x) \geq g(x)$ on $[a, b]$, then $\int_a^b f(x) dx \geq \int_a^b g(x) dx$.
- If $\ell \leq f(x) \leq m$ on $[a, b]$, then $\ell(b - a) \leq \int_a^b f(x) dx \leq m(b - a)$.

Example 1

Prove that $2 \leq \int_1^2 2^x dx \leq 4$.

Solution: This follows from the last comparison property. Since $f(x) = 2^x$ is an increasing function on $[1, 2]$, we have $f(1) \leq f(x) \leq f(2)$ on $[1, 2]$. Therefore,

$$f(1)(2 - 1) \leq \int_1^2 2^x dx \leq f(2)(2 - 1) \quad \implies \quad 2 \leq \int_1^2 2^x dx \leq 4$$

EXERCISE

Suppose $a < b < c$ and $\int_a^c f(x) dx < \int_a^c g(x) dx$. Is it true that $\int_a^b f(x) dx < \int_a^b g(x) dx$? Explain your reasoning.

8.2.2 Integrability

Recall the definition of the definite integral as

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \sum_{i=1}^n f(x_i^*) \Delta x$$

provided that this limit exists and gives the same value for all possible choices of sample points x_i^* .

There is a natural question to ask here: Is it always the case that the limit exists and is independent of the choice of sample points? Let's investigate this with a few examples of functions that don't behave very nicely at all.

Example 2 Let f be the function defined on the closed interval $[0, 1]$ by

$$f(x) = \begin{cases} 0 & \text{if } x = 0, \\ \frac{1}{x} & \text{otherwise.} \end{cases}$$

Let's try to compute $\int_0^1 f(x)dx$.

We will do this by selecting the right-most point of each interval, and see what the limit gives us.

If we split up the interval into n rectangles, the corresponding Riemann sum is

$$\begin{aligned} R_n &= \Delta x \sum_{i=1}^n f\left(\frac{i}{n}\right) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{n}{i} \\ &= \sum_{i=1}^n \frac{1}{i}. \end{aligned}$$

The astute reader may recognise R_n as the n th partial sum of the harmonic series! So we know $\lim_{n \rightarrow \infty} R_n$ does not exist. Therefore, the function $f(x)$ is *not* integrable on the closed interval $[0, 1]$.

The previous example is one where the limit defining the definite integral does not exist. The next example exhibits non-integrability in a slightly different way.

Example 3 Consider the function

$$f(x) = \begin{cases} 0 & \text{if } x \in \mathbb{Q} \\ 1 & \text{otherwise.} \end{cases}$$

Let's focus on the domain $[0, 1]$.

This is a perfectly good function, although it's not one we can easily visualize. Strangely, this is an example of a function that is defined on the closed interval $[0, 1]$, but is not continuous anywhere on the interval.

Nonetheless, we can still try and integrate it!

Let's look at the right-hand and left-hand Riemann sums. We have

$$R_n = \frac{1}{n} \sum_{i=1}^n f\left(\frac{i}{n}\right) = \frac{1}{n} \sum_{i=1}^n 0 = 0$$

and

$$L_n = \frac{1}{n} \sum_{i=1}^n f\left(\frac{i-1}{n}\right) = \frac{1}{n} \sum_{i=1}^n 0 = 0.$$

Note that $f\left(\frac{i}{n}\right) = 0$ since $\frac{i}{n}$ is rational. Therefore, $\lim_{n \rightarrow \infty} R_n = \lim_{n \rightarrow \infty} L_n = 0$. Everything is fine, and we have $\int_0^1 f(x)dx = 0$, right?

Well, not quite. Remember, the definition of a definite integral requires that the limit exists and is the same regardless of which sample points within each subinterval (each Δx) we choose. So far we have shown that the limit is zero whether or not we choose the left- or right-hand endpoints of each subinterval. What about other points?

Here's a fact about the real numbers that we will use without proof: Between any two distinct rational numbers, there is an irrational number.

With that in mind, let's choose our i th sample point x_i^* to be any irrational number in the closed interval $[\frac{i-1}{n}, \frac{i}{n}]$.

In this case, the Riemann sum S_n takes the form:

$$S_n = \frac{1}{n} \sum_{i=1}^n f(x_i^*) = \frac{1}{n} \sum_{i=1}^n 1 = 1.$$

Therefore, $\lim_{n \rightarrow \infty} S_n = 1$, which is not equal to the limits of the left- or right-hand Riemann sums computed earlier.

The fact that the limit is *not* independent of our choice of sample points means that the function $f(x)$ is *not* integrable on the closed interval $[0, 1]$.

EXERCISE

Suppose $f(x)$ is a function defined, but not bounded, on the closed interval $[a, b]$, and suppose $f(x) \geq 0$ for all $x \in [a, b]$. Prove that $f(x)$ is not integrable on $[a, b]$.

These two examples (and the exercise) illustrate that not every function is integrable! The natural question arises, are there situations where we can guarantee that a function is integrable?

Here is a theorem which requires a few technical facts about the real numbers to prove, and we will not prove it in this course. It is important, and tells us that most commonly occurring functions we encounter in a mathematics course are indeed integrable.

Fact 2 Let $f(x)$ be a continuous function on the closed interval $[a, b]$. Then $f(x)$ is integrable on $[a, b]$.

So, polynomials, trigonometric functions, exponential functions, logarithmic functions, among many other examples we have seen in these notes are integrable on closed intervals contained in the domains of the functions.

However, these are not the only integrable functions!

Example 4 Let $f(x)$ be the function defined on $[-1, 1]$ by

$$f(x) = \begin{cases} 1 & \text{if } x = 0 \\ 0 & \text{otherwise.} \end{cases}$$

If you were forced to produce a number that was the area under the graph of $f(x)$, hopefully every bone in your body screams “0”, since there is only one part of the function that’s above the x -axis. That part is just one value, it shouldn’t contribute any area! You would be correct, but let’s prove it!

This function is not continuous on $[-1, 1]$ (it has a jump discontinuity at $x = 0$) so we can’t automatically conclude that it’s integrable. We’ll have to proceed using Riemann sums and the definition of the definite integral.

Suppose we divide the interval $[-1, 1]$ into n subintervals of even width. We want to investigate the Riemann sum

$$S_n = \Delta x \sum_{i=1}^n f(x_i^*) = \frac{2}{n} \sum_{i=1}^n f(x_i^*)$$

where x_i^* is some sample point in the i th subinterval.

If n is odd, then 0 appears in exactly one subinterval (the $\frac{n+1}{2}$ -th interval). If n is even, then 0 is in two subintervals (the $\frac{n}{2}$ -th and $(\frac{n}{2} + 1)$ -th subintervals).

Therefore, either zero, one, or two of our sample points x_i^* are equal to 0. If none of our sample points are 0, then

$$S_n = \frac{2}{n} \sum_{i=1}^n f(x_i^*) = \frac{2}{n} \sum_{i=1}^n 0 = 0.$$

If exactly one of our sample points is 0, then

$$S_n = \frac{2}{n} \sum_{i=1}^n f(x_i^*) = \frac{2}{n} f(0) = \frac{2}{n}.$$

If two of our sample points are 0, then

$$S_n = \frac{2}{n} \sum_{i=1}^n f(x_i^*) = \frac{2}{n} (f(0) + f(0)) = \frac{4}{n}.$$

Therefore, regardless of our choice of sample points, we have $0 \leq S_n \leq \frac{4}{n}$. Since $\lim_{n \rightarrow \infty} \frac{4}{n} = 0$, the squeeze theorem permits us to conclude that $\lim_{n \rightarrow \infty} S_n = 0$.

Therefore,

$$\int_{-1}^1 f(x) dx = 0.$$

EXERCISE

Let $f(x)$ be a function defined on $[a, b]$ so that $f(x) = 0$ at all but finitely many points. Prove that $f(x)$ is integrable and that

$$\int_a^b f(x) dx = 0.$$

The most recent example suggests that as long as a function is continuous almost everywhere, then it should be integrable. The following theorem is true, but not proved here.

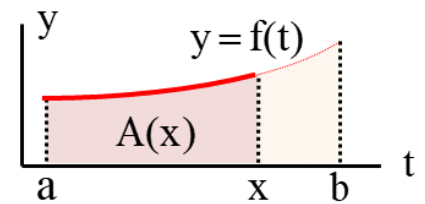
Fact 3 Let $f(x)$ be continuous on the closed interval $[a, b]$ except for finitely many jump discontinuities. Then $f(x)$ is integrable on $[a, b]$.

8.3 The Fundamental Theorem of Calculus

We claimed at the beginning of this chapter that differentiation and integration are two sides of the same coin. More precisely, they are inverse processes of one another which means that instead of computing integrals by evaluating limits of Riemann sums, we can compute them through antidifferentiation. Such an extraordinary connection, if true, would be so valuable that it would be worthy of being called “The Fundamental Theorem of Calculus”. But let’s see for ourselves how this connection arises.

Let us define what we will call an **area function** $A(x)$ that represents the area under a curve $y = f(t)$ from $t = a$ to $t = x$.

$$A(x) = \int_a^x f(t) dt$$



We call this an *area function* because we are allowing the upper limit of integration x to be a variable. This is in contrast to the lower limit of integration a which is fixed.

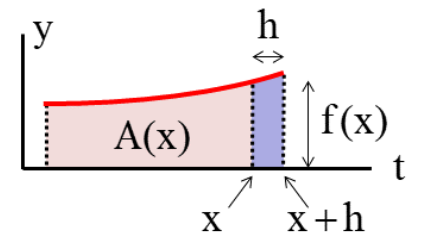
Now suppose we shift the upper limit of integration by a small amount h .

The area under the curve $y = f(t)$ will change by a small amount.

$$A(x+h) - A(x) \approx f(x) \cdot h$$

Rearranging this expression gives

$$f(x) \approx \frac{A(x+h) - A(x)}{h}$$



In the limit that $h \rightarrow 0$, this relation becomes exact and we recognize the right-hand side as the definition of the derivative of $A(x)$ with respect to x .

$$f(x) = \lim_{h \rightarrow 0} \frac{A(x+h) - A(x)}{h} = \frac{dA}{dx}$$

Let’s think about what this is saying for a moment. Remember, $A(x)$ represents the area under the function $y = f(t)$ from a fixed left-endpoint to a variable right-endpoint x . The value of $f(x)$ is the height of the curve $y = f(t)$ at the right-endpoint of the measured area. The equation we’ve just derived says that the rate of change of the area is equal to the height of the function. This is quite reasonable when put in simple terms - that is, we expect the area under a “tall function” to grow faster than the area under a “short function”. In the equation above, we have captured this notion precisely and this is what we call the Fundamental Theorem of Calculus.

Theorem 4 (The Fundamental Theorem of Calculus (Part 1))

Let f be a function continuous on $[a, b]$ and define

$$A(x) = \int_a^x f(t) dt$$

where $x \in [a, b]$. Then $A(x)$ is continuous on $[a, b]$, differentiable on (a, b) , and $A'(x) = f(x)$.

Proof: Let's first investigate the derivative of $A(x)$ on (a, b) . Let $x \in (a, b)$. For h small enough so that $x + h \in (a, b)$ we have

$$\begin{aligned} A'(x) &= \lim_{h \rightarrow 0} \frac{A(x+h) - A(x)}{h} \\ &= \frac{1}{h} \left(\int_a^{x+h} f(t) dt - \int_a^x f(t) dt \right) \\ &= \frac{1}{h} \left(\int_a^x f(t) dt + \int_x^{x+h} f(t) dt - \int_a^x f(t) dt \right) \\ &= \frac{1}{h} \int_x^{x+h} f(t) dt. \end{aligned}$$

Let's assume that $h > 0$, and when $h < 0$ the argument follows similarly.

Since $f(t)$ is continuous on the closed interval $[x, x+h]$, the Extreme Value Theorem tells us there exists $u, v \in [x, x+h]$ so that $f(u)$ and $f(v)$ are the maximum and minimum values of $f(t)$ on $[x, x+h]$. Therefore, by some comparison properties of the definite integral we have

$$f(u)h \geq \int_x^{x+h} f(t) dt \geq f(v)h$$

which implies

$$f(u) \geq \frac{1}{h} \int_x^{x+h} f(t) dt \geq f(v).$$

We are interested in the limit of the middle term as $h \rightarrow 0$. Since $u, v \in [x, x+h]$, as $h \rightarrow 0$, we know that u and v both tend to x . Therefore, taking the limit as $h \rightarrow 0$ of all parts of the above inequality gives

$$\lim_{u \rightarrow x} f(u) \geq \lim_{h \rightarrow 0} \frac{1}{h} \int_x^{x+h} f(t) dt \geq \lim_{v \rightarrow x} f(v).$$

A similar inequality holds if $h < 0$ (can you write it down precisely?).

Since $f(t)$ is continuous at x we have $\lim_{v \rightarrow x} f(v) = \lim_{u \rightarrow x} f(u) = f(x)$. Therefore, by the squeeze theorem, $A'(x) = f(x)$.

We have shown that $A'(x) = f(x)$ on (a, b) and as a consequence, $A(x)$ is differentiable on (a, b) .

A similar argument works to show $A(x)$ is continuous from the right at $x = a$ and from the left at $x = b$. We now prove the former and leave the latter as an exercise.

Consider

$$\lim_{x \rightarrow a^+} A(x) = \lim_{x \rightarrow a^+} \int_a^x f(t) dt.$$

As above, since $f(t)$ is continuous on $[a, x]$, there exist $u, v \in [a, x]$ so that $f(u)$ and $f(v)$ are maximum and minimum values of $f(t)$ on $[a, x]$. Then

$$(x - a)f(u) \geq \int_a^x f(t)dt \geq (x - a)f(v).$$

Consider $\lim_{x \rightarrow a^+} (x - a)f(u)$. Since $u \in [a, x]$, as $x \rightarrow a^+$, $u \rightarrow a^+$. Since $f(t)$ is continuous on $[a, b]$, $\lim_{x \rightarrow a^+} f(u) = \lim_{u \rightarrow a^+} f(u) = f(a)$. Therefore,

$$\lim_{x \rightarrow a^+} (x - a)f(u) = \left(\lim_{x \rightarrow a^+} (x - a) \right) \left(\lim_{x \rightarrow a^+} f(u) \right) = 0f(a) = 0.$$

Similarly, $\lim_{x \rightarrow a^+} (x - a)f(v) = 0$. Therefore by the squeeze theorem,

$$\lim_{x \rightarrow a^+} A(x) = \lim_{x \rightarrow a^+} \int_a^x f(t)dt = 0.$$

Since

$$A(a) = \int_a^a f(t)dt = 0$$

we have $\lim_{x \rightarrow a^+} A(x) = A(a)$ and A is continuous from the right at $x = a$. \square

At first glance, it is not immediately obvious why this theorem is so wonderful and that's why this is just "Part 1". This result *does* have its uses, but let's see where the relationship $A'(x) = f(x)$ can take us. If the derivative of $A(x)$ is $f(x)$, then that is equivalent to saying $A(x)$ is an *anti-derivative* of $f(x)$.

Recall, the general antiderivative, $F(x)$, of a function $f(x)$ is a class of functions that differ by an additive constant. If $A(x)$ belong to this class of functions, then it must be true that

$$F(x) = A(x) + C$$

for some constant C .

Now observe how we can rewrite the following difference in terms of the area function.

$$\begin{aligned} F(b) - F(a) &= (A(b) + C) - (A(a) + C) \\ &= A(b) - A(a) \\ &= \int_a^b f(x) dx - \int_a^a f(x) dx \\ &= \int_a^b f(x) dx \end{aligned}$$

Turning this around, we get part 2 of the Fundamental Theorem of Calculus.

Theorem 5 (The Fundamental Theorem of Calculus (Part 2))

If f is continuous on the interval $[a, b]$, then

$$\int_a^b f(x) dx = F(b) - F(a)$$

where $F(x)$ is any antiderivative of $f(x)$ (i.e., $F'(x) = f(x)$).

This gives us a new way of evaluating definite integrals. We no longer need to take limits of Riemann sums. Instead, we just need to identify an antiderivative of the integrand, evaluate it at the limits of integration, and take the difference.

Example 5

Evaluate the definite integral $\int_0^1 x^2 dx$.

Solution: To begin, we need any antiderivative of $f(x) = x^2$. We could use $\frac{1}{3}x^3$ or $\frac{1}{3}x^3 + 7$ or $\frac{1}{3}x^3 - 99\pi$. Let's just use the general antiderivative $F(x) = \frac{1}{3}x^3 + C$ where C is an arbitrary constant so that we can see why it doesn't matter which antiderivative we choose.

By part 2 of the Fundamental Theorem of Calculus, we have

$$\begin{aligned}\int_0^1 x^2 dx &= F(1) - F(0) \\ &= \left(\frac{1}{3}(1)^3 + C\right) - \left(\frac{1}{3}(0)^3 + C\right) \\ &= \frac{1}{3}\end{aligned}$$

Compare the amount of work involved in this calculation to constructing and simplifying the corresponding Riemann sum and then taking a limit. This way is much faster and simpler.

Notice, also, the reason it doesn't matter which antiderivative we used is because the additive constant will always cancel out. Moving forward with this knowledge, we will usually then pick whatever antiderivative makes for the least amount of writing (e.g., in this case, that would have been $\frac{1}{3}x^3$).

REMARK

By convention, the difference $F(b) - F(a)$ can be written using any of the following notations:

$$F(b) - F(a) = F(x)|_a^b = [F(x)]_a^b$$

EXERCISE

Evaluate $\int_0^{2\pi} \sin(x) dx$ using the Fundamental Theorem of Calculus.

8.4 Indefinite Integrals

The Fundamental Theorem of Calculus has established for us that a definite integral with integrand $f(x)$ can be evaluated quickly in terms of the antiderivative, $F(x)$, of the integrand. Finding an antiderivatives is generally more efficient than working with Riemann sums, so this makes the antiderivative a very useful object in the world of integration. For that reason, we make the following definition.

Definition 8.4.1
indefinite integral

The **indefinite integral** of a function $f(x)$ is denoted $\int f(x) dx$ and is the general antiderivative of $f(x)$. That is,

$$\int f(x) dx = F(x) \iff F'(x) = f(x)$$

Note that there is nothing groundbreaking going on here; this is just a relabelling of the general antiderivative. For example, if $f(x) = x^2$, then $\int f(x) dx = \frac{1}{3}x^3 + C$ where C is an arbitrary constant. This new notation is purely for convenience.

Example 6

Evaluate the indefinite integral $\int (y^4 - 2y^2 - 3) dy$.

Solution: To evaluate this indefinite integral (i.e., find the antiderivative), we need to determine what function, when differentiated with respect to y , yields the integrand. We can do this term-by-term.

$$\begin{aligned} \frac{d}{dy} (y^5) = 5y^4 &\implies \frac{d}{dy} \left(\frac{1}{5}y^5 \right) = y^4 \\ \frac{d}{dy} (y^3) = 3y^2 &\implies \frac{d}{dy} \left(-\frac{2}{3}y^3 \right) = -2y^2 \\ \frac{d}{dy} (y) = 1 &\implies \frac{d}{dy} (-3y) = -3 \end{aligned}$$

Since differentiation is a linear operation, $\frac{d}{dy} \left(\frac{1}{5}y^5 - \frac{2}{3}y^3 - 3y \right) = y^4 - 2y^2 - 3$.

Therefore, the general antiderivative or indefinite integral is

$$\int (y^4 - 2y^2 - 3) dy = \frac{1}{5}y^5 - \frac{2}{3}y^3 - 3y + C$$

where C is an arbitrary constant.

REMARK

It is important to keep in mind the difference between an indefinite integral and a definite integral.

An indefinite integral represents a family of functions

$$\int f(x) dx = F(x) + C$$

since C is an arbitrary constant.

A definite integral is a number

$$\int_a^b f(x) dx = F(b) - F(a)$$

since it can be interpreted as representing the area of a fixed region.

Also, the indefinite integral of $f(x)$ can be used to evaluate a definite integral of $f(x)$. In contrast, by simply knowing the value of a definite integral of a function $f(x)$, we cannot immediately infer what the indefinite integral of $f(x)$ is. (However, in the latter case, we usually work out that indefinite integral implicitly along the way.)

Since a general algorithm for evaluating indefinite integrals does not exist, we state the following commonly occurring indefinite integrals for reference:

Monomials:

$$\int x^n dx = \frac{1}{n+1} x^{n+1} + C \quad \text{for } n \neq -1 \qquad \int \frac{1}{x} dx = \ln(|x|) + C$$

Exponentials:

$$\int e^x dx = e^x + C \qquad \int b^x dx = \frac{1}{\ln(b)} b^x + C$$

Trigonometric:

$$\int \sin(x) dx = -\cos(x) + C \qquad \int \cos(x) dx = \sin(x) + C \qquad \int \sec^2(x) dx = \tan(x) + C$$

EXERCISE

Evaluate the indefinite integral $\int 2^x dx$.

8.4.1 Net Change

We can restate the Fundamental Theorem of Calculus as follows:

$$F(b) - F(a) = \int_a^b F'(t) dt$$

In this form, we can interpret the theorem as saying that the **net change** in a function $F(t)$ over an interval $a \leq t \leq b$ is equal to the integral of its rate of change over that interval.

This interpretation is particularly useful for physical applications where we know the instantaneous rate of change of a quantity (e.g., with respect to time) over an interval and want to know the corresponding net change in that quantity.

Example 7

A particle moving along the x -axis has velocity $v(t) = 2t + 3$ m/s where t is time measured in seconds. What is the net change in position (i.e., displacement) from $t = 0$ to $t = 4$?

Solution: Let $x(t)$ be the position as a function of time. Our goal is to calculate

$$\Delta x = x(4) - x(0)$$

The velocity is the rate of change in the position with respect to time (i.e., $v(t) = x'(t)$) so

$$\begin{aligned} \Delta x = x(4) - x(0) &= \int_0^4 x'(t) dt \\ &= \int_0^4 (2t + 3) dt \\ &= [t^2 + 3t]_0^4 \\ &= (4^2 + 3 \cdot 4) - (0^2 + 3 \cdot 0) \\ &= 28 \end{aligned}$$

Therefore, the particle's displacement over the four-second interval is 28 metres.

EXERCISE

A pool is draining such that the rate of change in its volume of water is given by

$$V'(t) = -20e^{-t} \text{ m}^3/\text{hr}$$

What is the net change in water volume in the pool between $t = 1$ hr and $t = 2$ hr?

8.5 Integration Techniques

Given an integral, if you can identify the antiderivative of the integrand, then you're basically done. However, unlike differentiation which is a process for which you can just 'turn a crank', we only antidifferentiate by recognizing functions as being derivatives of other functions. For example, we know an antiderivative of $\cos(x)$ is $\sin(x)$ because we proved and now know that the derivative of $\sin(x)$ is $\cos(x)$. So, when we *don't* recognize an integrand as the derivative of some other function, what are we to do? In this section, we will work through a handful of common integration techniques. This will certainly not be an exhaustive list - there are *many* more integration techniques out there - but these few should serve you well for many of the integrals you'll encounter.

8.5.1 Integration by Substitution

As inspiration for our first technique, let's try to "reverse" the Chain Rule.

$$\frac{d}{dx} (f(g(x))) = f'(g(x)) \cdot g'(x)$$

Integrating both sides of this equations we get

$$f(g(x)) + C = \int f'(g(x)) \cdot g'(x) dx$$

Reading this equation right-to-left, it is saying that if we can recognize an integrand as the derivative of a composite function $f'(g(x)) \cdot g'(x)$, then the indefinite integral is the composite function plus an arbitrary constant. This is easier said than done, but making a clever substitution in which we try to make an educated guess at what the “inner” function of the composition might be helps identify the form of the composition. Moreover, if the guess is a good one, we are left with a simpler integrand.

To be more precise, we must first pick an “inner” function $g(x)$ and then make the substitution $u = g(x)$. This changes the variable of integration from x to u , so we must rewrite the entire integral in terms of the new variable u . This includes replacing the differential dx using the relation $du = g'(x)dx$ (which is just the Chain Rule). What remains in the integrand after this change of variables will be the “outer” function in terms of u . Let’s see an example.

Example 8

Evaluate the indefinite integral $\int \frac{x}{\sqrt{1+x^2}} dx$.

Solution: Having the expression $1+x^2$ inside a square root makes this challenging to evaluate directly, but it also hints at what we might try as a substitution. In particular, if $u = g(x) = 1+x^2$ were the “inner” function, then the “outer” function would be something like $\frac{1}{\sqrt{u}}$ which would be readily integrated. Let’s try this out and see what happens. That is, make the following substitution

$$u = g(x) = 1 + x^2$$

Our task now is to rewrite the integral in terms of the new variable u . This means rewriting both the integrand and differential.

For the differential, we have

$$du = g'(x) dx = 2x dx \quad \implies \quad dx = \frac{1}{2x} du$$

The factor of $\frac{1}{2x}$ will now become part of the integrand. We could try to rewrite this x in terms of u , but that would be a bit messy (i.e., $x = \pm\sqrt{u-1}$) and, because our substitution is a good one, things clean up nicely without doing that.

$$\int \frac{x}{\sqrt{1+x^2}} dx = \int \frac{x}{\sqrt{1+x^2}} \left(\frac{1}{2x} du \right) = \int \frac{1}{2} \frac{1}{\sqrt{1+x^2}} du$$

To complete the substitution process, we know use $u = g(x) = 1+x^2$ to get rid of any remaining occurrences of the variable x .

$$\int \frac{1}{2} \frac{1}{\sqrt{1+x^2}} du = \int \frac{1}{2} \frac{1}{\sqrt{u}} du = \int \frac{1}{2} u^{-1/2} du$$

We are left with an integral of a power of u for which we know the antiderivative.

$$\int \frac{1}{2} \frac{1}{\sqrt{u}} du = \sqrt{u} + C$$

where C is an arbitrary constant.

We're not done yet though because the variable u is something we invented to help us evaluate the integral. We need to use the relation $u = 1 + x^2$ now to get our answer in terms of the original variable x . This gives us our final result:

$$\int \frac{x}{\sqrt{1+x^2}} dx = \sqrt{1+x^2} + C$$

Let's now summarize the procedure we used in the previous example.

Fact 6 The Substitution Rule

Given a differentiable function $u = g(x)$, then

$$\int f'(g(x)) \cdot g'(x) dx = \int f'(u) du$$

For certain integrals, this rule can be enormously helpful because with a “good” substitution, a complicated integral turns into a much simpler one. The tricky part is finding a “good” substitution. Often, this is achieved through trial and error. We pick a substitution $u = g(x)$ that is a candidate for the “inner” part of the composite function. We then change variables and see if we are left with a simpler integral in the new integration variable u . If not, we start over with a different substitution (or try something else completely). However, with practice, it gets easier to look ahead and see how a particular choice of substitution will play out.

EXERCISE

Make the substitution $u = g(x) = x^5 + 3x^3$ to help you evaluate the indefinite integral

$$\int (5x^4 + 9x^2) \sec^2(x^5 + 3x^3) dx$$

EXERCISE

Evaluate the indefinite integral $\int \frac{2x^3 + x}{x^4 + x^2} dx$ assuming $x \neq 0$.

Substitution with Definite Integrals

Substitution can be used to solve definite integral by solving the corresponding indefinite integral and then applying the Fundamental Theorem of Calculus. However, it is often more convenient to instead apply the substitution directly to the definite integral. The only

catch is that when we do that, we also have to determine limits of integration in terms of the new variable.

Fact 7 **The Substitution Rule for Definite Integrals** Given a function $u = g(x)$ which is differentiable for all $x \in [a, b]$, then

$$\int_{x=a}^b f'(g(x)) \cdot g'(x) dx = \int_{u=g(a)}^{g(b)} f'(u) du$$

REMARKS

- When we use substitution to solve an indefinite integral, we must revert back to the original variable as a final step. This is not necessary for a definite integral as it will work out to the same number either way (provided you remember to rewrite the limits of integration in terms of the new variable).
- We are now labelling our lower limits of integration as being defined with respect to a particular variable (i.e., $x = a$ vs. $u = g(a)$) with the implication that the upper variable is also defined with respect to the same variable. It is not necessary to label our limits in this way, but you may find it helpful as a reminder to change your limits of integration after you make a substitution.

Example 9

Evaluate $\int_{x=e}^{e^2} \frac{(\ln(x))^2}{x} dx$.

Solution: The presence of a $\ln(x)$ makes this integral a bit daunting. So, let's try to replace $\ln(x)$ with a new variable u .

$$u = g(x) = \ln(x) \quad \implies \quad du = g'(x)dx = \frac{1}{x}dx$$

Since this is a definite integral, we must also change the limits of integration. The lower limit $x = e$ becomes $u = g(e) = \ln(e) = 1$ and the upper limit $x = e^2$ becomes $u = g(e^2) = \ln(e^2) = 2$. Therefore, we have

$$\begin{aligned} \int_{x=e}^{e^2} \frac{(\ln(x))^2}{x} dx &= \int_{u=1}^2 u^2 du \\ &= \left[\frac{1}{3}u^3 \right]_1^2 \\ &= \frac{8}{3} - \frac{1}{3} \\ &= \frac{7}{3} \end{aligned}$$

EXERCISE

Use integration by substitution to evaluate $\int_{x=0}^1 \frac{x}{(x^2 + 1)^3} dx$.

8.5.2 Integration by Parts

Recall the product rule for differentiation

$$\frac{d}{dx} [f(x)g(x)] = f'(x)g(x) + f(x)g'(x)$$

Integrating both sides gives

$$\int \frac{d}{dx} [f(x)g(x)] dx = \int f'(x)g(x) dx + \int f(x)g'(x) dx$$

or, after applying the Fundamental Theorem of Calculus on the left-hand side

$$f(x)g(x) = \int f'(x)g(x) dx + \int f(x)g'(x) dx$$

It's not clear how this might help us until we rearrange things a bit.

$$\int f(x)g'(x) dx = f(x)g(x) - \int f'(x)g(x) dx$$

Now we're onto something. Specifically, if we can identify an integrand as a product of two functions or "parts" and we know the antiderivative of one of the two parts, then we can effectively swap the integral on the left-hand side with the one on the right-hand side (plus the product term).

Fact 8 Integration by Parts

Given $u = u(x)$ and $v = v(x)$ with $du = u'(x) dx$ and $dv = v'(x) dx$, then

$$\int u dv = uv - \int v du$$

To use this approach, we need to identify part of the integrand as what we're calling the function $u(x)$ and the rest of the integrand along with the differential as dv .

Let's look at an example to see this in action.

Example 10

Use integration by parts to evaluate $\int x \sin(x) dx$.

Solution: We split the integrand by letting

$$u = x \quad \text{and} \quad dv = \sin(x)dx$$

With these definitions we have

$$du = dx \quad \text{and} \quad v = -\cos(x)$$

Using the rule for integration by parts, we get

$$\begin{aligned} \int x \sin(x) dx &= -x \cos(x) - \int (-\cos(x)) dx \\ &= -x \cos(x) + \int \cos(x) dx \\ &= -x \cos(x) + \sin(x) + C \end{aligned}$$

You may be wondering what would have happened if we had chosen u and dv differently. For example, we could have tried $u = \sin(x)$ and $dv = x dx$. For integration by parts, we typically require that u be easily differentiated and dv be easily anti-differentiated. That is true for these alternate definitions. However, we also want the resulting vdu to be simpler to integrate than the udv we started with. With $v du = \frac{1}{2}x^2 \cos(x) dx$, that would *not* be true.

EXERCISE

Use integration by parts to evaluate $\int x \ln(x) dx$.

(Hint: Which part of the integrand is hardest to integrate? You should define that to be u .)

Integration by parts can also be used to evaluate definite integrals by using the following modified rule.

$$\int_a^b u dv = [uv]_a^b - \int_a^b v du$$

Example 11 Use integration by parts to evaluate $\int_1^e \ln(x) dx$.

Solution: This does not initially appear to be a good candidate for integration by parts. The integrand is not obviously a product, but we have a trick up our sleeve. We'll take $u = \ln(x)$ and just set $dv = dx$. Then we can proceed as follows:

$$\begin{aligned} u = \ln(x) &\implies du = \frac{1}{x} dx \\ dv = dx &\implies v = x \end{aligned}$$

Using the formula for integration by parts for a definite integral, we then have

$$\begin{aligned}\int_1^e \ln(x) dx &= [x \ln(x)]_1^e - \int_1^e x \left(\frac{1}{x}\right) dx \\ &= [x \ln(x)]_1^e - [x]_1^e \\ &= e \ln(e) - 1 \ln(1) - (e - 1) \\ &= e - 0 - e + 1 \\ &= 1\end{aligned}$$

EXERCISE

Evaluate $\int_0^1 x e^{-x} dx$.

Sequential Applications of Integration by Parts

Consider the following indefinite integral:

$$\int x^2 e^x dx$$

This looks like a good candidate for integration by parts (IBP). Let's try it with $u = x^2$ and $dv = e^x dx$. After a bit of work, we find

$$\int x^2 e^x dx = x^2 e^x - 2 \int x e^x dx$$

The antiderivative of $x e^x$ is not immediately obvious, but we've clearly made progress because we started with $x^2 e^x$. We just need to apply (IBP) again; this time with $u = x$ and $dv = e^x dx$.

$$\begin{aligned}\int x^2 e^x dx &= x^2 e^x - 2 \int x e^x dx \\ &= x^2 e^x - 2 \left(x e^x - \int e^x dx \right) \\ &= e^x (x^2 - 2x + 2) + C\end{aligned}$$

The lesson here is that we should not necessarily be deterred when one application of IBP does not produce an antiderivative right away. The more important question is whether or not it has left us in a position where we can still see a path forward.

EXERCISE

Apply integration by parts twice - both times with $u = e^x$ - to the integral $\int e^x \cos(x) dx$. At the end, you won't have an explicit expression for the antiderivative, but you should be able to solve for it with a bit of ingenuity.

(You could also set $dv = e^x dx$ for both applications of IBP and everything should work out similarly well. But what happens if you set $u = e^x$ the first time and then $dv = e^x dx$ the second time or vice versa?)

8.5.3 Trigonometric Integrals

Consider the integral

$$\int \sin^7(x) \cos(x) dx$$

and make the substitution $u = \sin(x)$ which gives $du = \cos(x)dx$ so that

$$\int \sin^7(x) \cos(x) dx = \int u^7 du = \frac{1}{8}u^8 + C = \frac{1}{8} \sin^8(x) + C$$

Similarly, consider the integral

$$\int \sin(x) \cos^4(x) dx$$

and make the substitution $u = \cos(x)$ which gives $du = -\sin(x)dx$ so that

$$\int \sin(x) \cos^4(x) dx = - \int u^4 du = -\frac{1}{5}u^5 + C = -\frac{1}{5} \cos^5(x) + C$$

The lesson here is that powers of sine can be readily integrated by substitution if there is a single factor of cosine present and vice versa.

EXERCISE

Evaluate $\int \sin(x) \cos(x) dx$ two different ways. First, by making the substitution $u = \sin(x)$ and next by making the substitution $u = \cos(x)$. Your antiderivatives from the two approaches may look different at first glance, but you should be able to argue that they are equivalent.

Let's now build upon this result. Consider the integral

$$\int \sin^2(x) \cos^3(x) dx$$

We don't have a single power of either sine or cosine. However, we can use the Pythagorean trigonometric identity $\sin^2(x) + \cos^2(x) = 1$ to remove two powers of cosine.

$$\begin{aligned} \int \sin^2(x) \cos^3(x) dx &= \int \sin^2(x) (1 - \sin^2(x)) \cos(x) dx \\ &= \int \sin^2(x) \cos(x) dx - \int \sin^4(x) \cos(x) dx \end{aligned}$$

Now we have two integrals, but the integrands of both are a power of $\sin(x)$ and a single factor of $\cos(x)$. These can be finished off with the substitution $u = \sin(x)$.

This strategy of using the Pythagorean trigonometric identity can be applied whenever there is an odd power of either $\sin(x)$ or $\cos(x)$.

EXERCISE

Evaluate $\int \sin^5(x) \cos^2(x) dx$ by using $\sin^4(x) = (1 - \cos^2(x))^2$.

The approach above does not work if sine and cosine both appear with even powers. However, in this case, we can make use of the following half-angle identities.

Fact 9 Half-Angle Identities

$$\sin^2(x) = \frac{1 - \cos(2x)}{2} \qquad \cos^2(x) = \frac{1 + \cos(2x)}{2}$$

These identities allow us to rewrite the integrand in powers of $\cos(2x)$ while cutting the degree of the integrand (as a polynomial in trigonometric functions) in half. This trick can be applied iteratively until we are left with a constant plus odd powers of cosine functions. The latter can be integrated using the approach discussed above.

Example 12

Evaluate $\int \sin^2(x) \cos^2(x) dx$.

Solution: We apply the half-angle identities.

$$\begin{aligned} \int \sin^2(x) \cos^2(x) dx &= \frac{1}{2^2} \int (1 - \cos(2x))(1 + \cos(2x)) dx \\ &= \frac{1}{4} \int [1 - \cos^2(2x)] dx \\ &= \frac{1}{4} \int \left[1 - \frac{1}{2}(1 + \cos(4x)) \right] dx \\ &= \frac{1}{8} \int [1 - \cos(4x)] dx \\ &= \frac{1}{8}x - \frac{1}{32} \sin(4x) + C \end{aligned}$$

Observe that in the third line, we applied the cosine half-angle identity a second time.

REMARK

A similar set of techniques exist for certain integrals of the form $\int \tan^m(x) \sec^n(x) dx$ making use of the identity $\sec^2(x) = \tan^2(x) + 1$. However, the authors of these notes expect that the likelihood you will come across such an integral (and can't solve it by rewriting the integrand in terms of sines and cosines) is low enough that our time is better spent on other topics.

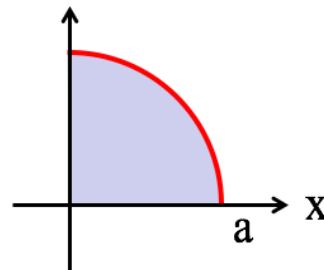
8.5.4 Trigonometric Substitutions

A trigonometric substitution is a type of *inverse substitution*. We call it this because it involves defining the given integration variable to be a function of a new integration variable (e.g., starting with the integration variable x , we might let $x = x(u) = \ln(u)$). This is in contrast to the standard substitutions we've already seen where a new integration variable is defined as a function of the given one (e.g., starting with integration variable x , we might let $u = x^2$).

Let's see how an inverse substitution involving a trigonometric function might be helpful by calculating the area of a circle of radius a . To this end, let $f(x) = \sqrt{a^2 - x^2}$.

From $x = 0$ to $x = a$, the region under the graph of $f(x)$ corresponds to one quarter of the area of a circle of radius a . Therefore, the total area of a circle of radius a is given by

$$A = 4 \int_{x=0}^a \sqrt{a^2 - x^2} dx$$



Let's solve this integral.

It is tempting to try the regular substitution $u = a^2 - x^2$. This will turn the square root function in the integrand into a $u^{1/2}$. On it's own, that would be easy to integrate. However, we must also swap out the differential using the relation $du = -2x dx$. The factor of x here ends up making things worse (try it out). So, we need to be more clever.

Notice that if we could write $a^2 - x^2$ as a perfect square, then we could resolve the square root. If only we knew of some functions $x(\theta)$ and $y(\theta)$ such that $(x(\theta))^2 + (y(\theta))^2 = a^2$. Of course, we *do* know of such functions $x(\theta) = a \sin(\theta)$ and $y(\theta) = a \cos(\theta)$ (or, if you prefer, the other way around). This suggests we make the inverse substitution

$$x(\theta) = a \sin(\theta) \quad \implies \quad dx = a \cos(\theta) d\theta$$

For the limits of integration, we have

$$x = \{0, a\} \quad \implies \quad \theta = \left\{0, \frac{\pi}{2}\right\}$$

By design then, this construction simplifies the integrand as follows.

$$\begin{aligned} \sqrt{a^2 - x^2} &= \sqrt{a^2 - a^2 \sin^2(\theta)} \\ &= a \sqrt{\cos^2(\theta)} \\ &= a |\cos(\theta)| \end{aligned}$$

Observe also that with $\theta \in \left[0, \frac{\pi}{2}\right]$, we have $|\cos(\theta)| = \cos(\theta)$. Therefore, the integral

becomes

$$\begin{aligned}
 A &= 4 \int_{x=0}^a \sqrt{a^2 - x^2} \, dx \\
 &= 4 \int_{\theta=0}^{\pi/2} (a \cos(\theta)) (a \cos(\theta) d\theta) \\
 &= 4a^2 \int_{\theta=0}^{\pi/2} \cos^2(\theta) \, d\theta \\
 &= 4a^2 \int_{\theta=0}^{\pi/2} \left(\frac{1 + \cos(2\theta)}{2} \right) \, d\theta \\
 &= 4a^2 \left[\frac{\theta}{2} + \frac{\sin(2\theta)}{4} \right]_{\theta=0}^{\pi/2} \\
 &= 4a^2 \left[\left(\frac{\pi}{4} + \frac{\sin(\pi)}{4} \right) - \left(\frac{0}{2} + \frac{\sin(0)}{4} \right) \right] \\
 &= \pi a^2
 \end{aligned}$$

As expected, the area of a circle of radius a is given by the well-known formula $A = \pi a^2$.

More importantly, the example above suggests more generally that an integrand including a function of the form $\sqrt{a^2 - x^2}$ may resolve to something simpler with the substitution $x = a \sin(\theta)$.

Similar arguments, which exploit the identity $1 + \tan^2(\theta) = \sec^2(\theta)$, can be made to show that integrands involving functions of the form $\sqrt{a^2 + x^2}$ or $\sqrt{x^2 - a^2}$ may simplify with the substitutions $x = a \tan(\theta)$ and $x = a \sec(\theta)$, respectively.

We summarize these common trigonometric substitutions in the following table.

Function	Substitution	Relevant Identity
$\sqrt{a^2 - x^2}$	$x = a \sin(\theta)$	$\sin^2(\theta) + \cos^2(\theta) = 1$
$\sqrt{a^2 + x^2}$	$x = a \tan(\theta)$	$1 + \tan^2(\theta) = \sec^2(\theta)$
$\sqrt{x^2 - a^2}$	$x = a \sec(\theta)$	$\sec^2(\theta) - 1 = \tan^2(\theta)$

EXERCISE

Use a trigonometric substitution to evaluate $\int \frac{1}{1+x^2} \, dx$.

Hint: The denominator can be viewed as $(\sqrt{1+x^2})^2$.

8.5.5 Integration by Partial Fractions

In this section, we will discuss a general technique for integrating rational functions. Recall that a rational function is one of the form

$$f(x) = \frac{P(x)}{Q(x)}$$

where $P(x)$ and $Q(x)$ are polynomials.

We can classify rational functions based on the relative degrees of the numerator and denominator polynomials.

Definition 8.5.1

proper/improper
rational function

Let $f(x) = \frac{P(x)}{Q(x)}$ where $P(x)$ and $Q(x)$ are polynomials.

- If $\deg(P) < \deg(Q)$, then we say f is a **proper** rational function.
- If $\deg(P) \geq \deg(Q)$, then we say f is an **improper** rational function.

Let's first look at an example where we need to integrate a proper rational function.

Example 13

Evaluate $\int \frac{5x - 1}{x^2 - x - 2} dx$.

Solution: Observe that the denominator polynomial $Q(x) = x^2 - x - 2$ can be factored as $Q(x) = (x + 1)(x - 2)$. This suggests we might be able to find constants A and B such that

$$\frac{5x - 1}{(x + 1)(x - 2)} = \frac{A}{x + 1} + \frac{B}{x - 2}$$

Indeed, with a bit of work, which we leave as an exercise to verify, the unique solution is $A = 2$ and $B = 3$.

We can use this decomposition of the integrand to evaluate the integral.

$$\begin{aligned} \int \frac{5x - 1}{x^2 - x - 2} dx &= \int \frac{2}{x + 1} dx + \int \frac{3}{x - 2} dx \\ &= 2 \ln |x + 1| + 3 \ln |x - 2| + C \end{aligned}$$

The key to solving the previous example was recognizing that the integrand - through what is called the **method of partial fractions** - could be decomposed into readily integrated proper rational functions.

$$\underbrace{\frac{5x - 1}{x^2 - x - 2}}_{\text{Hard to integrate}} = \underbrace{\frac{2}{x + 1} + \frac{3}{x - 2}}_{\text{Easy to integrate}}$$

It turns out that any proper rational function can be decomposed into a sum of rational functions which are “easy to integrate”. The following algorithm can be applied to determine such a partial fraction decomposition.

Algorithm 3**(Method of Partial Fractions)**

Given a proper rational function $f(x) = \frac{P(x)}{Q(x)}$:

1. To determine the denominators of the “target” rational functions, first factor $Q(x)$. Recall that any polynomial can be written as a product of linear factors $(ax + b)$ and irreducible quadratic factors $(ax^2 + bx + c)$ with no real roots). For example:

$$\underbrace{x^4 - 2x^3 + x^2 - 2x}_{\text{Factor out an } x} = x \underbrace{(x^3 - 2x^2 + x - 2)}_{\text{Factor out } (x-2)} = x(x - 2) \underbrace{(x^2 + 1)}_{\text{Irreducible}}$$

2. Write down a “trial” partial fraction expansion according to the following rule:

i. For each linear factor $(ax + b)$, include a term of the form:

$$\frac{A}{ax + b}$$

ii. A linear factor repeated n times, $(ax + b)^n$, gives rise to n terms:

$$\frac{A_1}{ax + b} + \frac{A_2}{(ax + b)^2} + \cdots + \frac{A_n}{(ax + b)^n}$$

iii. For each irreducible quadratic factor $(ax^2 + bx + c)$, we get a term that has the following form:

$$\frac{Ax + B}{ax^2 + bx + c}$$

iii. If an irreducible quadratic factor is repeated n times, $(ax^2 + bx + c)^n$, we get n terms:

$$\frac{A_1x + B_1}{ax^2 + bx + c} + \frac{A_2x + B_2}{(ax^2 + bx + c)^2} + \cdots + \frac{A_nx + B_n}{(ax^2 + bx + c)^n}$$

3. Expand the partial fraction expansion and compare to the desired rational function. This typically involves solving a system of equations.

EXERCISE

Evaluate $\int \frac{3x^2 - 1}{x^3 - x^2} dx$.

Now, what about *improper* rational functions $f(x) = \frac{P(x)}{Q(x)}$ where the degree of the numerator polynomial $P(x)$ is greater than or equal to that of the denominator polynomial $Q(x)$? Here we simply add the extra first step of performing polynomial division. By doing so, we can rewrite the improper rational function as a polynomial plus a proper rational function (both of which we now know how to integrate).

Fact 10

Performing polynomial division on an improper rational function will always yield a polynomial plus a proper rational function. That is, if $P(x)$ and $Q(x)$ are polynomials with $\deg(P) \geq \deg(Q)$, then there exist polynomials $R(x)$ and $S(x)$ such that

$$\frac{P(x)}{Q(x)} = S(x) + \frac{R(x)}{Q(x)}$$

with $\deg(R) < \deg(Q)$.

Example 14

Evaluate $\int \frac{x^2 + 3x + 3}{x + 1} dx$.

Solution: Performing polynomial division, we have

$$\begin{array}{r} x + 2 \\ x + 1 \overline{) x^2 + 3x + 3} \\ \underline{x^2 + x} \\ 2x + 3 \\ \underline{2x + 2} \\ 1 \end{array}$$

We can now evaluate the given integral.

$$\begin{aligned} \int \frac{x^2 + 3x + 3}{x + 1} dx &= \int \left(x + 2 + \frac{1}{x + 1} \right) dx \\ &= \frac{1}{2}x^2 + 2x + \ln|x + 1| + C \end{aligned}$$

where C is an arbitrary constant.

EXERCISE

Evaluate $\int \frac{x^3 + 5x^2 + 6x + 3}{x^2 + 3x} dx$.

8.6 Improper Integrals

In this section, we will investigate areas of regions which are infinite in extent. Perhaps surprisingly, such regions can have finite area. Moreover, they are not simply a mathematical curiosity, but rather have real-world applications. In fact, let's use an application as motivation. Let's determine how much fuel (or energy converted to work) is required to launch a satellite into distant space.

For a constant force, work is equal to force times distance. But the gravitational force acting on a satellite (which our rocket thrusters must work against) is not constant. Instead, it decreases with height according to Newton's Law of Universal Gravitation.

$$F(r) = \frac{GMm}{r^2}$$

In this equation, r is the radial distance from the centre of the Earth to the satellite, G is Newton's gravitational constant, M is the mass of the Earth, and m is the mass of the satellite. To get the total energy required to move a satellite from the surface of the Earth which has radial distance from the centre of the Earth $r = R_E$ to an arbitrary height $r = R$,

we must integrate this force.

$$\begin{aligned} E &= \int_{R_E}^R \frac{GMm}{r^2} dr \\ &= GMm \int_{R_E}^R \frac{1}{r^2} dr \\ &= GMm \left[\frac{1}{R_E} - \frac{1}{R} \right] \end{aligned}$$

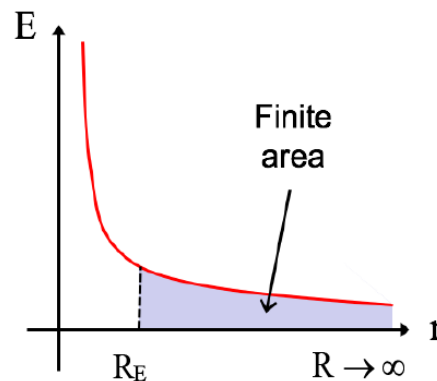
To answer the question of how much energy is required to reach ‘distant space’, we take the limit as $R \rightarrow \infty$ to capture the idea that we must effectively escape the gravitational field of the Earth.

$$\begin{aligned} E &= \lim_{R \rightarrow \infty} \int_{R_E}^R \frac{GMm}{r^2} dr \\ &= GMm \lim_{R \rightarrow \infty} \left[\frac{1}{R_E} - \frac{1}{R} \right] \\ &= \frac{GMm}{R_E} \end{aligned}$$

Therefore, we need just a finite amount of energy to send a satellite arbitrarily far into distant space. Intuitively, this is not unreasonable because the gravitational force gets increasingly weaker as distance increases and eventually tends to zero.

8.6.1 Improper Integrals: Type 1

Let’s step back for a moment and think about the mathematics of what we did in the example above. We evaluated a definite integral for which we allowed one of the limits of integration to tend to infinity. This definite integral then corresponds to finding the ‘area’ of a region which is infinite in extent. But the definite integral was finite. This is quite remarkable when we observe that the region has infinite width and non-zero height!



The fact that this integral gives a finite value motivates us to define integration over intervals which are taken to be infinite in size. We call these improper integrals.

Definition 8.6.1**Improper Integral
(Type 1)**

Let f be integrable on $[a, b]$ for every $b \geq a$, then we write

$$\int_a^\infty f(x)dx = \lim_{b \rightarrow \infty} \int_a^b f(x)dx$$

If the limit exists, we say $\int_a^\infty f(x)dx$ converges. Otherwise, we say it diverges.

Similarly, if f is integrable on $[a, b]$ for every $a \leq b$, then we write

$$\int_{-\infty}^b f(x)dx = \lim_{a \rightarrow -\infty} \int_a^b f(x)dx$$

Again, if the limit exists, we say $\int_{-\infty}^b f(x)dx$ converges. Otherwise, we say it diverges.

Note, we can combine two improper integrals as defined above to write

$$\int_{-\infty}^\infty f(x)dx = \lim_{a \rightarrow -\infty} \int_a^c f(x)dx + \lim_{b \rightarrow \infty} \int_c^b f(x)dx$$

for any $c \in \mathbb{R}$.

Example 15

Evaluate $\int_0^\infty e^{-x} dx$

Solution: We have

$$\begin{aligned} \int_0^\infty e^{-x} dx &= \lim_{b \rightarrow \infty} \int_0^b e^{-x} dx \\ &= \lim_{b \rightarrow \infty} [-e^{-x}]_0^b \\ &= \lim_{b \rightarrow \infty} [(-e^{-b}) - (-e^0)] \\ &= 0 - (-1) \\ &= 1 \end{aligned}$$

REMARK

When evaluating improper integrals where one or both limits of integration are infinite, it is common practice to not actually write the limit explicitly. For example, in the previous example, we could instead write

$$\int_0^\infty e^{-x} dx = [(-e^{-x}) - (-e^0)]_0^\infty = 0 - (-1) = 1$$

EXERCISE

Evaluate $\int_1^{\infty} \frac{1}{x^2} dx$.

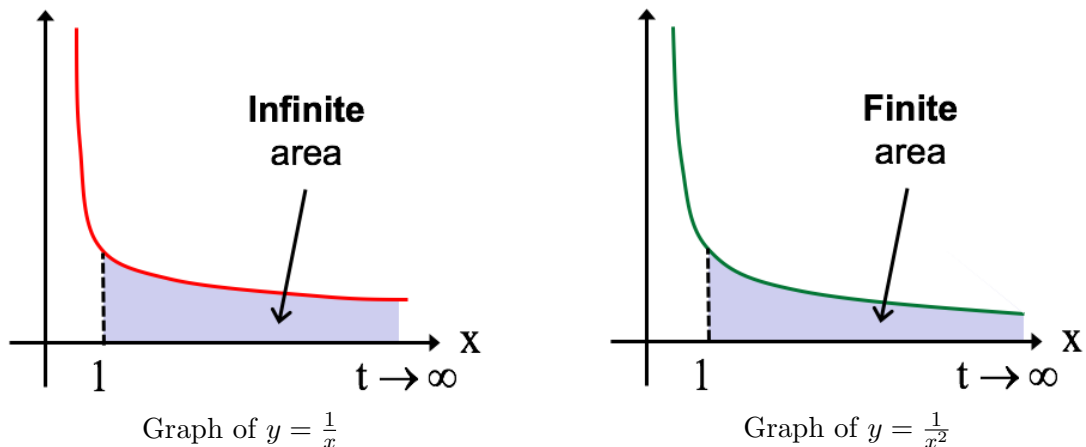
8.6.2 Convergence of Power Functions: Part 1

Consider the following improper integral:

$$\int_1^{\infty} \frac{1}{x} dx = \lim_{t \rightarrow \infty} \int_1^t \frac{1}{x} dx = \lim_{t \rightarrow \infty} [\ln(t) - \ln(1)] = \lim_{t \rightarrow \infty} [\ln(t)] = \infty$$

This integral is divergent and, more precisely, tends towards positive infinity. We can interpret this geometrically to mean that the area under the curve $y = \frac{1}{x}$ on the interval $[1, \infty)$ is infinite.

In contrast, in the previous exercise you should have found $\int_1^{\infty} \frac{1}{x^2} dx = 1$. This means that over the same interval the area of the region under the graph of $y = \frac{1}{x^2}$ is finite.



So $\frac{1}{x^2}$ decreases quickly enough that its graph bounds a finite area on $[1, \infty)$ while $\frac{1}{x}$ does *not* decrease quickly for its graph to bound a finite area. It turns out that $f(x) = \frac{1}{x}$ sits at the boundary between the power functions which are convergent and the power functions which are divergent on the interval $[1, \infty)$.

Theorem 11

The integral $\int_1^{\infty} \frac{1}{x^p} dx$ is convergent for $p > 1$ and divergent for $p \leq 1$.

Proof: In the discussion above we showed that this integral is divergent for $p = 1$. For $p \neq 1$, we have

$$\begin{aligned} \int_1^{\infty} \frac{1}{x^p} dx &= \lim_{t \rightarrow \infty} \int_1^t x^{-p} dx \\ &= \lim_{t \rightarrow \infty} \left[\frac{1}{-p+1} x^{-p+1} \right]_1^t \\ &= \frac{1}{1-p} \left(\lim_{t \rightarrow \infty} \frac{1}{t^{p-1}} - 1 \right) \end{aligned}$$

For $p > 1$, we have $p - 1 > 0$ so $\lim_{t \rightarrow \infty} \frac{1}{t^{p-1}} = 0$. On the other hand, for $p < 1$ we have $p - 1 < 0$ so this limit does not exist. Therefore, the integral $\int_1^{\infty} \frac{1}{x^p} dx$ is convergent for $p > 1$ and divergent for $p \leq 1$. \square

REMARK

We can shift the lower limit of integration in the previous theorem to be any finite positive real number and the result still holds. This is because such an integral will differ from the one in the theorem by a finite amount and a finite difference will never change an integral from convergent to divergent or vice versa.

For example, $\int_4^{\infty} \frac{1}{\sqrt{x}} dx$ is divergent because

$$\begin{aligned} \int_4^{\infty} \frac{1}{\sqrt{x}} dx &= \int_1^{\infty} \frac{1}{\sqrt{x}} dx - \int_1^4 \frac{1}{\sqrt{x}} dx \\ &= \int_1^{\infty} \frac{1}{x^{1/2}} dx - 2 \end{aligned}$$

The integral on the last line is divergent since it corresponds to the case of $p = \frac{1}{2}$ in the previous theorem and subtracting 2 is not going to magically make that finite.

EXERCISE

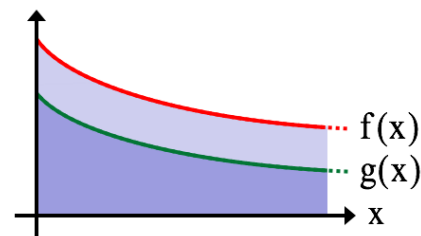
Determine if $\int_1^{\infty} \frac{x-1}{x^2} dx$ is convergent or divergent.

8.6.3 Improper Integral Comparison Test

Some improper integrals cannot be evaluated directly. However, we can often determine if such integrals are convergent or divergent by making an appropriate comparison.

Recall, if $f(x) \geq g(x) \geq 0$ on some interval, then for any subinterval $[a, b]$ (where a or b could be infinite), we have

$$\int_a^b f(x) dx \geq \int_a^b g(x) dx \geq 0$$



From this inequality, we can establish the following test, which should remind you of the Comparison Test for series (from Section 6.3.1).

Theorem 12 (Improper Integral Comparison Test)

Let f and g be continuous on $[a, \infty)$ and suppose $0 \leq g(x) \leq f(x)$.

- If $\int_a^{\infty} f(x) dx$ is convergent, then $\int_a^{\infty} g(x) dx$ is convergent.

- If $\int_a^\infty g(x) dx$ is divergent, then $\int_a^\infty f(x) dx$ is divergent.

To prove this, we will imitate the discussion before the statement of the Comparison Test from Section 6.3.1, which sketches out a proof of the Comparison Test for series.

To do this, we will rely on the following lemma (which we will not prove) that is an analogue of the Monotone Sequence Theorem (Theorem 4).

Recall that a function $f(x)$ is non-decreasing if $a < b$ implies $f(a) \leq f(b)$.

Lemma 13

Suppose $f(x)$ is a function that is non-decreasing on an interval (a, ∞) . If there exists $M \in \mathbb{R}$ so that $f(x) \leq M$ for all $x \in (a, \infty)$, then $\lim_{x \rightarrow \infty} f(x)$ exists.

EXERCISE

Prove Lemma 13. It may be useful to go over the proof of the Monotone Sequence Theorem (Theorem 4) and take inspiration from that proof.

We can now prove Theorem 12

Proof: Let

$$F(t) = \int_a^t f(x) dx \quad \text{and} \quad G(t) = \int_a^t g(x) dx.$$

Note that

$$\lim_{t \rightarrow \infty} F(t) = \int_a^\infty f(x) dx \quad \text{and} \quad \lim_{t \rightarrow \infty} G(t) = \int_a^\infty g(x) dx.$$

By the Fundamental Theorem of Calculus (Theorem 4), $F'(t) = f(t) \geq 0$ and $G'(t) = g(t) \geq 0$. Therefore, both $F(t)$ and $G(t)$ are non-decreasing functions on (a, ∞) .

The second point is the converse of the first point, so we shall only prove the first.

Suppose $\lim_{t \rightarrow \infty} F(t) = L$. Since $F(t)$ is non-decreasing, $F(t) \leq L$ for all $t \in (a, \infty)$ (can you see why?). Since $0 \leq g(x) \leq f(x)$, $G(t) \leq F(t) \leq L$ for all $t \in (a, \infty)$. Since $G(t)$ is non-decreasing and bounded above by L , the limit $\lim_{t \rightarrow \infty} G(t)$ exists by Lemma 13 and

$$\int_a^\infty g(x) dx$$

is convergent. □

EXERCISE

Inspired by the proof you just witnessed, go back and prove the Comparison Test for series.

Example 16

Prove that $\int_1^{\infty} e^{-x^2} dx$ is convergent.

Solution: For $x \geq 1$, we have

$$x^2 \geq x \implies -x^2 \leq -x \implies 0 \leq e^{-x^2} \leq e^{-x}$$

We can also show

$$\begin{aligned} \int_1^{\infty} e^{-x} dx &= [-e^{-x}]_1^{\infty} \\ &= \lim_{t \rightarrow \infty} (-e^{-t}) - (-e^{-1}) \\ &= 0 + e^{-1} \\ &= \frac{1}{e} \end{aligned}$$

Since e^{-x^2} and e^{-x} are both continuous on $[1, \infty)$, it follows by the Improper Integral Comparison Test that $\int_1^{\infty} e^{-x^2} dx$ is convergent.

EXERCISE

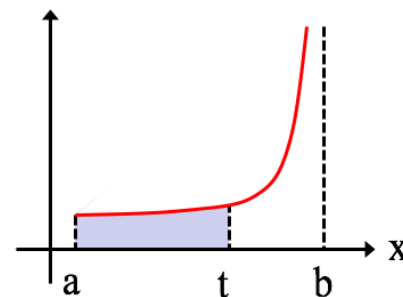
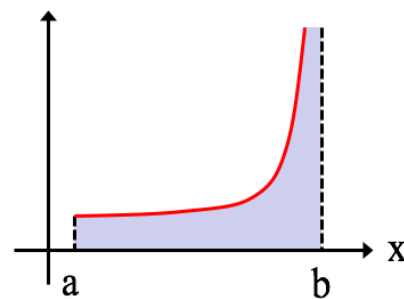
Determine whether $\int_1^{\infty} \frac{1 + \sin^2(x)}{\sqrt{x}} dx$ is convergent or divergent.

8.6.4 Improper Integrals: Type 2

A second type of improper integral occurs when an endpoint of the interval over which we're integrating coincides with a vertical asymptote of the integrand.

In this case, the values of the function diverge to either positive or negative infinity as the asymptote is approached. However, it is possible that the area of the unbounded region may still be finite.

As with Type 1 improper integrals, limits are our friend here. In particular, we can set up the integral with a variable endpoint, perform the integration, and then take the limit as our variable endpoint approaches the location of the vertical asymptote. For example, in the image to the right, we let $t \rightarrow b^-$ to approach the asymptote at $x = b$.



Definition 8.6.2**Improper Integral
(Type 2)**

Let f be integrable on $[a, b)$ and have a vertical asymptote at $x = b$, then

$$\int_a^b f(x) dx = \lim_{t \rightarrow b^-} \int_a^t f(x) dx$$

Similarly, if f is integrable on $(a, b]$ and has a vertical asymptote at $x = a$, then

$$\int_a^b f(x) dx = \lim_{t \rightarrow a^+} \int_t^b f(x) dx$$

When a vertical asymptote exists strictly inside the interval over which we're integrating, we can combine these definitions to evaluate it. For example, suppose $f(x)$ has a vertical asymptote at $x = c$ where $a < c < b$, then

$$\int_a^b f(x) dx = \lim_{t \rightarrow c^-} \int_a^t f(x) dx + \lim_{t \rightarrow c^+} \int_t^b f(x) dx$$

Example 17

Evaluate $\int_0^1 \frac{1}{\sqrt{x}} dx$.

Solution: The function $\frac{1}{\sqrt{x}}$ has a vertical asymptote at $x = 0$ so we must be careful. Let's first write our integral as a limit and then integrate.

$$\begin{aligned} \int_0^1 \frac{1}{\sqrt{x}} dx &= \lim_{t \rightarrow 0^+} \int_t^1 \frac{1}{\sqrt{x}} dx \\ &= \lim_{t \rightarrow 0^+} [2\sqrt{x}]_t^1 \\ &= 2\sqrt{1} - \lim_{t \rightarrow 0^+} 2\sqrt{t} \\ &= 2 - 0 \\ &= 2 \end{aligned}$$

EXERCISE

Determine if $\int_0^1 \frac{1}{x^3} dx$ is convergent or divergent.

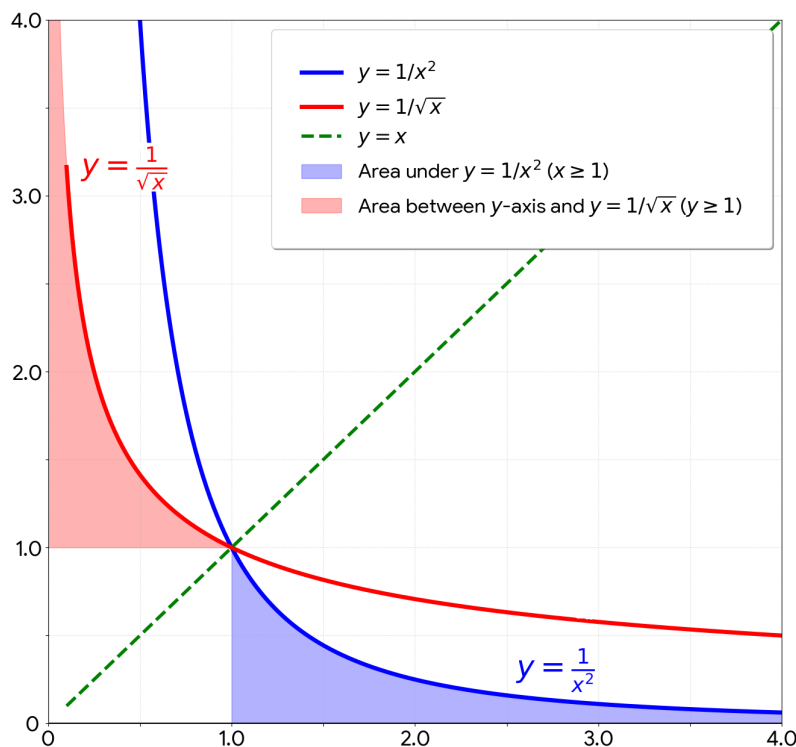
8.6.5 Convergence of Power Functions: Part 2

We have seen that $\int_1^\infty \frac{1}{\sqrt{x}} dx$ is divergent while $\int_0^1 \frac{1}{\sqrt{x}} dx$ is convergent. This means that the rules for convergence of integrals of power functions on the interval $[0, 1]$ are not the same as those on the interval $[1, \infty)$. However, as before, the power function $\frac{1}{x^p}$ marks a boundary between convergence and divergence.

Theorem 14

The integral $\int_0^1 \frac{1}{x^p} dx$ is convergent for $p < 1$ and divergent for $p \geq 1$.

We leave the proof as an exercise, but observe that this is essentially the opposite of what happens for integrals of power functions on the interval $[1, \infty)$ aside from $\frac{1}{x}$ whose integral is divergent on both intervals. Geometrically, this makes sense once you realize that the reflection of $y = \frac{1}{x^p}$ across the line $y = x$ is the curve $y = \frac{1}{x^{1/p}}$.



For example, if the area under $y = \frac{1}{x^2}$ on $[1, \infty)$ is finite, then the area between the y -axis and $y = \frac{1}{x^{1/2}} = \frac{1}{\sqrt{x}}$ for $y \geq 1$ will be finite. It follows that the area under $y = \frac{1}{\sqrt{x}}$ on $[0, 1]$ (which includes an extra square with area 1 adjacent to the origin) will be finite.

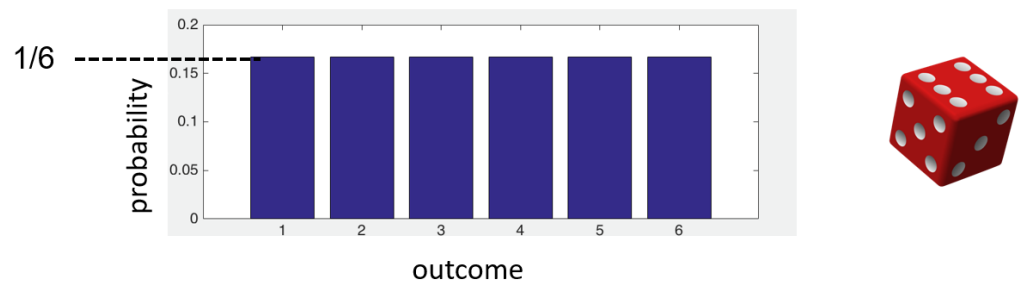
EXERCISE

Determine if $\int_0^2 \frac{1}{(x-1)^2} dx$ is convergent or divergent.

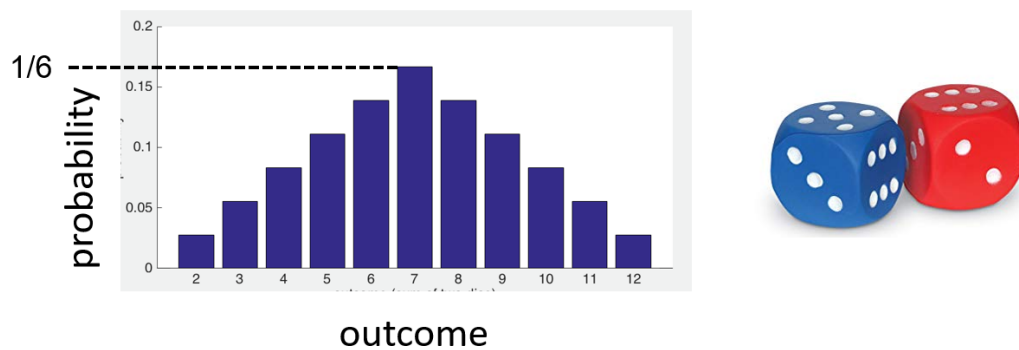
(Warning: Be mindful of the vertical asymptote at $x = 1$.)

8.7 Probability

Consider a standard six-sided die. We can represent the probability of each possible result of the die roll with a bar plot of probability vs. outcome.



In this case, each outcome is equally likely. However, if we roll two dice and take the sum of the numbers rolled things get more interesting. Again, we can capture this with a bar plot.



There are 36 equally likely states which can occur when we roll two dice. However, some states give rise to the same sum, so we no longer have a uniform probability distribution.

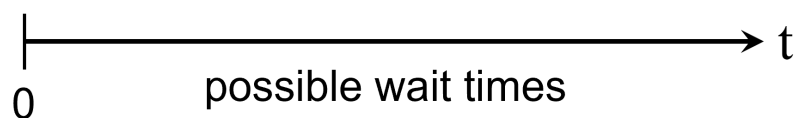
In both of these scenarios, there are a finite number of possible outcomes. For example, the possible outcomes when rolling a single die are the elements of the set $\{1, 2, 3, 4, 5, 6\}$. For distributions like this one, if we denote the probability of outcome i by p_i where there are n possible outcomes, then $0 \leq p_i \leq 1$ for $1 \leq i \leq n$. Moreover, the sum of the probabilities of all possible outcomes will equal one

$$\sum_{i=1}^n p_i = 1$$

When a random process has finitely many outcomes, we say the outcomes follow a **discrete** probability distribution.

8.7.1 Probability Density

In contrast, consider the possible outcomes for someone keeping track of how long it takes for a bus to arrive. The time spent waiting can, hypothetically, be any non-negative number. This means the outcomes cannot be discretely separated.

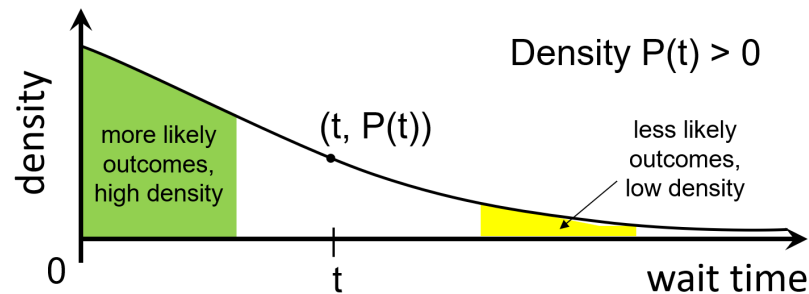


In fact, the probability of any particular outcome (e.g., it takes *exactly* 3 minutes for the bus to arrive) is zero!

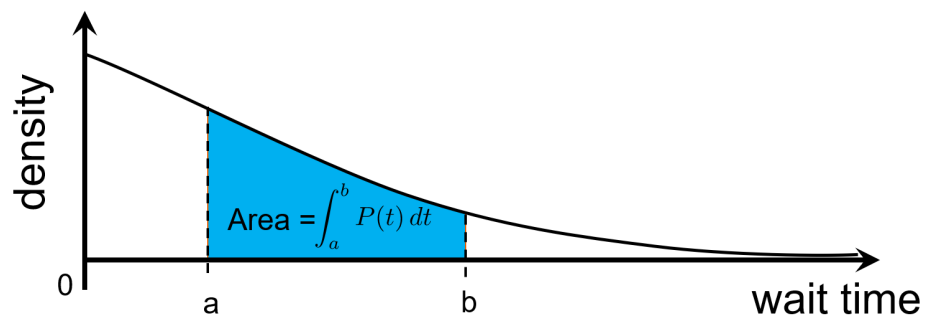
How then can we express, for example, the fact that short waits are more likely than very long wait times?



We assign probabilities to *intervals* of possible outcomes rather than individual outcomes. Mathematically, we achieve this by defining a **probability density function**, $P(t)$.



Each possible wait time t has an associated positive density $P(t)$, but $P(t)$ is *not* the probability of a wait time of t . (Remember, the probability of any individual wait time is zero.) Instead, the area under the probability density curve over an interval $[a, b]$ gives the probability that the wait time lies somewhere in that interval. This means the probability that the wait time is between $t = a$ and $t = b$ is equal to the definite integral $\int_{t=a}^b P(t) dt$.



Example 18

The waiting times when phoning a call centre are described by the probability density function $P(t) = 0.2e^{-0.2t}$. Determine the probability that a caller will wait less than ten minutes on hold rounded to the nearest percent.

Solution: The probability that a caller will wait between $t = 0$ and $t = 10$ minutes is equal

to

$$\begin{aligned}\int_0^{10} 0.2e^{-0.2t} dt &= [-e^{-0.2t}]_0^{10} \\ &= ((-e^{-2}) - (-e^0)) \\ &= 1 - e^{-2} \\ &\approx 0.8647\end{aligned}$$

Therefore, the probability that the caller waits less than ten minutes to the nearest percent is 86%.

In the previous example, we used the probability distribution $P(t) = 0.2e^{-0.2t}$. For this to be a valid probability distribution, the “sum” of all possible outcomes must be equal to one. For a density function $P(t)$ with $t \geq 0$, this means we require

$$\int_{t=0}^{\infty} P(t) dt = 1$$

EXERCISE

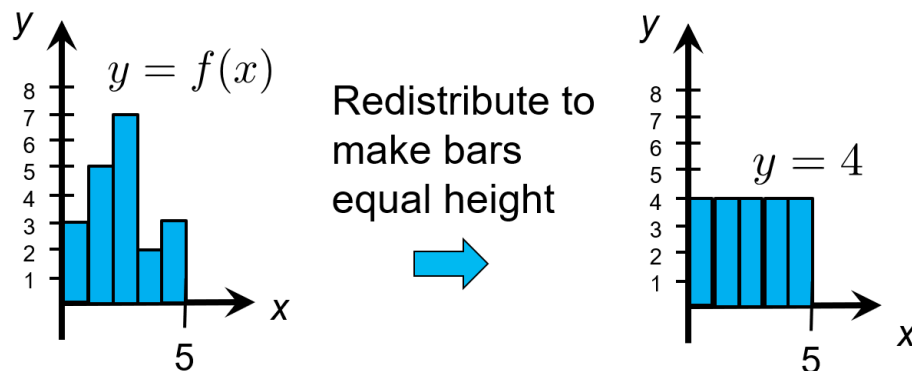
Verify that $\int_{t=0}^{\infty} P(t) dt = 1$ for the probability density in the previous example.

8.7.2 Average Value

Given a finite set of numbers like $\{3, 5, 7, 2, 3\}$ we can compute the average (or mean) by summing the numbers in the set and dividing by the number of numbers in the set.

$$\frac{3 + 5 + 7 + 2 + 3}{5} = 4$$

We can interpret the average geometrically as follows. Let each number in the set be represented by a rectangle of width 1 and area (or height) equal to the number. Next, we stack the rectangles side-by-side and redistribute area among the rectangles until they're all the same height.



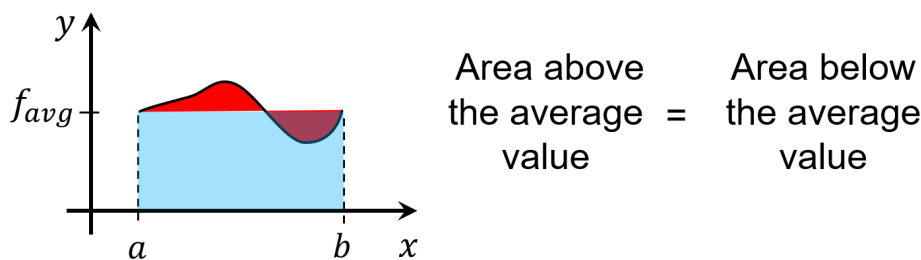
Once we're done redistributing area, each rectangle will have height equal to the average. If we piecewise-define a function $f(x)$ to describe our graph comprised of rectangles before averaging, then we can capture this process in integral notation as follows

$$\int_0^5 f(x) dx = \int_0^5 f_{\text{avg}} dx$$

where f_{avg} is the average value of the function which, being a constant, allows us to factor it out of the integral on the right-hand side to get

$$f_{\text{avg}} = \frac{\int_0^5 f(x) dx}{\int_0^5 dx} = \frac{20}{5} = 4$$

This seems unnecessarily complicated for finding the average of five integers - and it is - but it suggests how we can define the average value for a function defined on a continuous interval. We still want the average value to be the height of the function after redistributing area under the curve $y = f(x)$ to be f_{avg} . Another way of thinking about this redistributing process is to imagine a waterbed; if you push down on the bed somewhere, it goes up somewhere else, but the average height is always the same.



Since total area is conserved during the redistribution process, we have

$$\int_a^b f(x) dx = \int_a^b f_{\text{avg}} dx \quad \implies \quad f_{\text{avg}} = \frac{\int_a^b f(x) dx}{b - a}$$

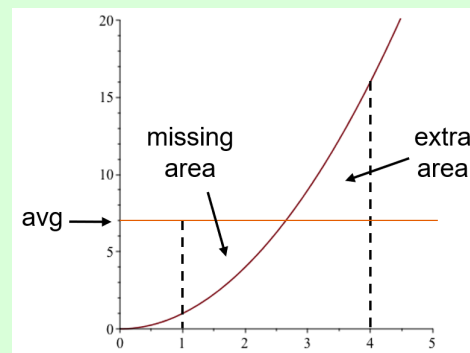
Example 19

Determine the average value of $f(x) = x^2$ over the interval $[1, 4]$.

Solution: We have

$$\begin{aligned} f_{\text{avg}} &= \frac{\int_1^4 x^2 dx}{4 - 1} \\ &= \frac{1}{3} \left[\frac{1}{3} x^3 \right]_1^4 \\ &= 7 \end{aligned}$$

We can visualize this with the diagram on the right. The areas of the regions labelled 'extra area' and 'missing area' are equal.



EXERCISE

The length L (in cm) of a certain species of fish is related to its age t (in months) by

$$L(t) = 1 + \frac{29t}{8+t}$$

If these fish live for 24 months, what is the average length of a fish in the population to the nearest cm?

8.7.3 Expected Value

It is sometimes the case that we have a set of numbers and want to compute something like an average but one where we give different numbers in the set different weights. For example, suppose you're taking a course with quizzes and assignments. You've earned an average grade of 80% on the quizzes and 60% on the assignments. According to the course syllabus, the quizzes are worth 30% of your final grade and the assignments are worth 70% of your final grade. To compute your final grade you don't just average your quiz and assignment grades. Instead, you take a weighted average

$$\frac{30(80\%) + 70(60\%)}{30 + 70} = 66\%$$

In general, the **weighted average** of a set of numbers $\{a_1, a_2, \dots, a_n\}$ with respective weights $\{w_1, w_2, \dots, w_n\}$ is given by

$$\frac{w_1a_1 + w_2a_2 + \dots + w_na_n}{w_1 + w_2 + \dots + w_n}$$

It is common to scale the weights so that they sum to one (e.g., we could have taken the quiz weight above to be 0.3 and the assignment weight to be 0.7). In this case, the weighted average simplifies to a weighted sum.

$$w_1a_1 + w_2a_2 + \dots + w_na_n$$

Weights which sum to one occur naturally when we discuss probability. For example, when rolling a pair of six-sided dice, if we sum the numbers rolled on both dice, we can get any whole number from 2 to 12. The probabilities are not equal for all possible outcomes, but all together they will sum to one.

For random events, it is often useful to determine the average over many instances. We call this the **expected value**. For example, if we roll the two dice many times and average the numbers we get, what should we "expect" that average to be after many rolls? This expected value is given by the weighted sum. For the dice, we would find

$$\frac{1}{36}(2) + \frac{2}{36}(3) + \frac{3}{36}(4) + \dots + \frac{1}{36}(12) = 7$$

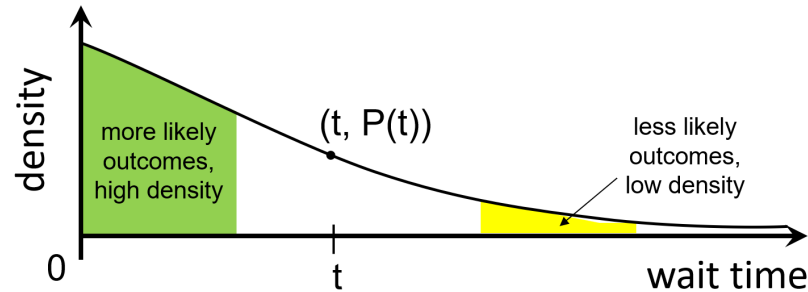
More generally, if a set of outcomes $\{a_1, a_2, \dots, a_n\}$ occur with probabilities $\{p_1, p_2, \dots, p_n\}$, then the expected value is

$$p_1a_1 + p_2a_2 + \dots + p_na_n$$

EXERCISE

A standard six-sided die is been modified so that one face is four times more likely to be rolled than any other face. If the expected value for rolls with this die is 4, then which number is most likely to be rolled?

How do we now define expected value for outcomes defined over a continuum? For example, going back to the problem of wait times, how could we predict the *expected* wait time? Of course, we will need to use the probability density distribution, $P(t)$, but how?



We use the value of the probability density function, $P(t)$, to weight each possible wait time, t . Instead of adding, we integrate! The recurring theme of these applications.

In general, if outcomes t range over an interval $[a, b]$ with probability density $P(t)$, then the expected value, denoted $\langle t \rangle$, is

$$\langle t \rangle = \int_a^b t P(t) dt$$

Note, we do not need to divide by a “total weight” since the “sum” of the probabilities of all possible outcomes which is the integral $\int_a^b P(t) dt$ will be equal to 1 since $P(t)$ is a probability density function.

Example 20

Recall, we determined that for the call centre whose waiting times were described by the probability density function $P(t) = 0.2e^{-0.2t}$, there is about an 86% chance that our call would be answered in the first ten minutes. Now let's determine how long, on average, we should expect to wait for our call to be answered. In other words, let's find the expected value $\langle t \rangle$.

Solution: Assuming that t can range over $[0, \infty)$ and applying integration by parts with

$$u = 0.2t \quad \text{and} \quad dv = e^{-0.2t} dt$$

gives

$$\begin{aligned}
 \langle t \rangle &= \int_0^{\infty} 0.2te^{-0.2t} dt \\
 &= [-te^{-0.2t}]_0^{\infty} + \int_0^{\infty} e^{-0.2t} dt \\
 &= [0 - 0] + \left[\frac{-1}{0.2} e^{-0.2t} \right]_0^{\infty} \\
 &= 0 + \left[0 - \frac{-1}{0.2} \right] \\
 &= 5
 \end{aligned}$$

Therefore, the expected wait time is 5 minutes.

EXERCISE

The time t (in seconds) it takes for blood to clot can be described by the probability density function $P(t) = \frac{a}{t}$ for $t \in [1, 20]$ where a is a constant.

- Determine a so that $P(t)$ is a valid probability density function (i.e., $\int_1^{20} P(t) dt = 1$).
- Determine the expected value for the time it takes for blood to clot.

8.8 Series and Taylor Polynomials Revisited

With integration now in our toolbox, we can expand on some previous topics as well as take care of a little housekeeping.

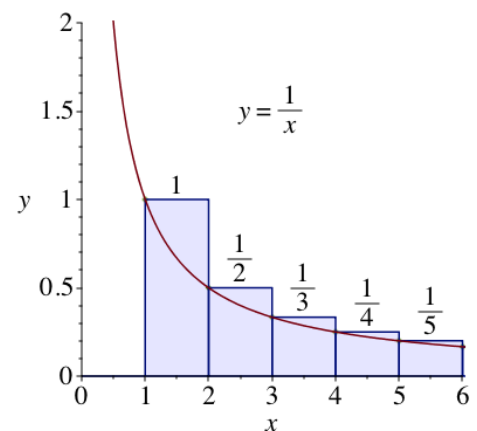
8.8.1 Integral Test

To start, we'll introduce another clever method of testing a series for convergence by relating the series sum to an improper integral.

To see how this works, recall the harmonic series

$$\sum_{n=1}^{\infty} \frac{1}{n} = 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \dots$$

Let us represent this series as the total area of an infinite sequence of rectangular blocks with unit width and height defined using the curve $y = \frac{1}{x}$ as shown in the diagram.



Observe that this total area is greater than the area under the curve $y = \frac{1}{x}$ on the interval $[1, \infty)$. This means

$$\sum_{n=1}^{\infty} \frac{1}{n} \geq \int_1^{\infty} \frac{1}{x} dx = \infty$$

Since this integral is infinite, the series must be divergent. This is the idea underlying the Integral Test.

Theorem 15 (Integral Test)

Let f be a non-negative, continuous, monotonic decreasing function such that $f(n) = a_n$ for each positive integer n . Then $\int_1^{\infty} f(x) dx$ is convergent if and only if $\sum_{n=1}^{\infty} a_n$ is convergent.

REMARKS

- Recall that the convergence or divergence of a series ultimately only depends on the large- n behaviour of the terms in the series. This means in practice that the conditions for the integral test do not actually have to hold for all $n \geq 1$, but rather there just needs to exist some positive integer k such that they hold for all $n \geq k$.
- The integral test cannot be used directly to determine the value of a convergent series. It is only used to determine if a series is convergent or divergent.

We previously encountered series of the form $\sum_{n=1}^{\infty} \frac{1}{n^p}$, which are called p -series. These series are convergent for $p > 1$ and divergent for $p \leq 1$. We proved this using the comparison test but only for $p \geq 2$ and $p \leq 1$. Observe that the integral test can be used to establish this result for all values of p with minimal effort. In particular, since $f(x) = \frac{1}{x^p}$ is a non-negative, continuous, decreasing function and $\int_1^{\infty} \frac{1}{x^p} dx$ is convergent for $p > 1$ and divergent for $p \leq 1$, then the convergence properties for p -series follow by the integral test. Of course, the integral test can be applied to other series as well. Consider the following example.

Example 21

Determine if the following series is convergent or divergent.

$$\sum_{n=1}^{\infty} \frac{\ln(n)}{n^2}$$

Solution: Let $f(x) = \frac{\ln(x)}{x^2}$. For $x \geq 2$, $f(x)$ is a continuous, decreasing, non-negative function and $f(n) = a_n$. Now let's evaluate the integral $\int_1^{\infty} f(x) dx$.

First, we apply integration by parts with:

$$\begin{aligned} u = \ln(x) &\implies du = \frac{1}{x} dx \\ dv = \frac{1}{x^2} dx &\implies v = -\frac{1}{x} \end{aligned}$$

Now, we can integrate as follows

$$\begin{aligned} \int_1^\infty \frac{\ln(x)}{x^2} dx &= \left[\ln(x) \left(-\frac{1}{x} \right) \right]_1^\infty - \int_1^\infty \left(-\frac{1}{x} \right) \left(\frac{1}{x} \right) dx \\ &= \left[-\frac{\ln(x)}{x} \right]_1^\infty + \int_1^\infty \left(\frac{1}{x^2} \right) dx \\ &= 0 - \lim_{x \rightarrow \infty} \left(\frac{\ln(x)}{x} \right) + \left[-\frac{1}{x} \right]_1^\infty \\ &= - \lim_{x \rightarrow \infty} \left(\frac{\ln(x)}{x} \right) - \lim_{x \rightarrow \infty} \left(\frac{1}{x} \right) + 1 \\ &= - \lim_{x \rightarrow \infty} \left(\frac{\ln(x)}{x} \right) + 1 \end{aligned}$$

Finally, we use l'Hôpital's Rule to evaluate the remaining limit to get

$$\int_{x=1}^\infty \frac{\ln(x)}{x^2} dx = - \lim_{x \rightarrow \infty} \left(\frac{\frac{1}{x}}{\frac{1}{1}} \right) + 1 = 1$$

Since this integral converges, then the series $\sum_{n=1}^\infty \frac{\ln(n)}{n^2}$ converges by the integral test.

In the previous example, we dove straight into applying the integral test. This led to a rather lengthy computation. In practice, it is worth considering what other avenues might be available before going down that road. For example, we might ask ourselves if the series is geometric or if the divergence test could be applied. Those checks can both be implemented quickly. However, in this case, they are not helpful (i.e., the series is not geometric and $\lim_{n \rightarrow \infty} a_n = 0$ so the divergence test is inconclusive). On the other hand, we *could* have used the limit comparison test here quite effectively. Nevertheless, the point is that we should try to keep in mind all of the convergence testing tools we have available to us when trying to demonstrate series convergence out in the wild.

EXERCISE

Determine if the series $\sum_{n=1}^\infty n e^{-n^2}$ is convergent or divergent.

8.8.2 Taylor's Inequality Revisited

Now for some housekeeping. We previously stated Taylor's inequality without proof. The general proof of Taylor's theorem is a neat application of the Fundamental Theorem of Cal-

culus, integration by parts, the triangle inequality for integrals, and mathematical induction, so is a great exercise to work through. Let's get started by laying some groundwork.

Theorem 16 (Triangle Inequality for Integrals)

Let f be integrable on $[a, b]$, then

$$\left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx$$

Proof: We have for all x that

$$-|f(x)| \leq f(x) \leq |f(x)|$$

We can integrate this inequality over $[a, b]$ to get

$$-\int_a^b |f(x)| dx \leq \int_a^b f(x) dx \leq \int_a^b |f(x)| dx$$

which is equivalent to

$$\left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx$$

□

We can make sense of this inequality intuitively as follows. On the left-hand side, we compute the integral $\int_a^b f(x) dx$ first and then take the absolute value. As such, if the sign of f changes on $[a, b]$, some cancellation of the areas will occur before we then take the absolute value. In contrast, on the right-hand side, we take the absolute value of $f(x)$ first. This means anywhere $f(x)$ may have been negative, it becomes positive and, in turn, the integral only includes positive contributions (i.e., no cancellation of areas occurs).

REMARK

This is called the Triangle Inequality for Integrals because - while it may not look like it at first glance - it is the generalization of the usual triangle inequality $|x + y| \leq |x| + |y|$ to a Riemann sum rather than just a sum of two elements.

Next, we will prove a theorem which is equivalent to Taylor's inequality and provides a formula for the n -th degree Taylor polynomial remainder $R_n(x)$.

Theorem 17 (Taylor's Integral Remainder Theorem)

If $f^{(n+1)}$ is continuous on an open interval I containing a and $x \in I$, then

$$R_n(x) = \frac{1}{n!} \int_{t=a}^x (x-t)^n f^{(n+1)}(t) dt$$

Proof: We will prove this using the Fundamental Theorem of Calculus and mathematical induction.

First we establish the base case ($n = 0$):

By the Fundamental Theorem of Calculus

$$\int_{t=a}^x f'(t)dt = f(x) - f(a)$$

Since the zeroth-degree Taylor polynomial of $f(x)$ is $f(a)$, we immediately get

$$R_0(x) = f(x) - T_0(x) = f(x) - f(a) = \int_{t=a}^x f'(t)dt$$

Next, we assume the inductive hypothesis that for $n = k$ we have

$$R_k(x) = \frac{1}{k!} \int_{t=a}^x (x-t)^k f^{(k+1)}(t)dt$$

and to complete the inductive argument we want to show

$$R_{k+1}(x) = \frac{1}{(k+1)!} \int_{t=a}^x (x-t)^{k+1} f^{(k+2)}(t)dt$$

Consider just the integral on the right-hand side and apply integration by parts with

$$\begin{aligned} u = (x-t)^{k+1} &\implies du = -(k+1)(x-t)^k dt \\ dv = f^{(k+2)}(t)dt &\implies v = f^{(k+1)}(t) \end{aligned}$$

This gives

$$\begin{aligned} \int_{t=a}^x (x-t)^{k+1} f^{(k+2)}(t)dt &= \left[(x-t)^{k+1} f^{(k+1)}(t) \right]_{t=a}^x + (k+1) \int_{t=a}^x (x-t)^k f^{(k+1)}(t)dt \\ &= -(x-a)^{k+1} f^{(k+1)}(a) + ((k+1)!) R_k(x) \end{aligned}$$

This means that the entire right-hand side above (i.e., with the factor of $\frac{1}{(k+1)!}$ in place) is equal to

$$\frac{1}{(k+1)!} \int_{t=a}^x (x-t)^{k+1} f^{(k+2)}(t)dt = -\frac{(x-a)^{k+1}}{(k+1)!} f^{(k+1)}(a) + R_k(x)$$

Since $R_k(x) = f(x) - T_k(x)$, we then have

$$\begin{aligned} \frac{1}{(k+1)!} \int_{t=a}^x (x-t)^{k+1} f^{(k+2)}(t)dt &= f(x) - T_k(x) - \frac{(x-a)^{k+1}}{(k+1)!} f^{(k+1)}(a) \\ &= f(x) - T_{k+1}(x) \\ &= R_{k+1}(x) \end{aligned}$$

Therefore, since we have established a base case and the $n = k + 1$ case follows from the $n = k$ case, then the statement is true for all n by mathematical induction. \square

We're now ready to prove Taylor's inequality.

Theorem 18 (Taylor's Inequality)

Let I be an interval containing $x = a$. If M is the maximum value of $|f^{(n+1)}(x)|$ on the interval I , then the remainder, $R_n(x)$, of the n -th degree Taylor polynomial, $T_n(x)$, centred at $x = a$ satisfies

$$|R_n(x)| \leq \frac{M}{(n+1)!} |x-a|^{n+1}$$

everywhere on the interval I .

Proof: Assume that there exists M such that for all t between x and a

$$|f^{(n+1)}(t)| \leq M$$

By Taylor's inequality for integrals, we have for all n

$$|R_n(x)| = \left| \frac{1}{n!} \int_{t=a}^x (x-t)^n f^{(n+1)}(t) dt \right|$$

By the triangle inequality for integrals, this implies

$$\begin{aligned} |R_n(x)| &\leq \frac{1}{n!} \int_{t=a}^x |(x-t)^n f^{(n+1)}(t)| dt \\ &= \frac{1}{n!} \int_{t=a}^x |x-t|^n |f^{(n+1)}(t)| dt \\ &\leq \frac{1}{n!} \int_{t=a}^x |x-t|^n M dt \\ &\leq \frac{M}{n!} \int_{t=a}^x |x-t|^n dt \end{aligned}$$

Notice we used our assumption in the third step.

If $x > a$, we have $|x-t|^n = (x-t)^n$ so

$$\begin{aligned} |R_n(x)| &\leq \frac{M}{n!} \int_{t=a}^x (x-t)^n dt \\ &= \frac{M}{n!} \left[\frac{-1}{n+1} (x-t)^{n+1} \right]_{t=a}^x \\ &= \frac{M}{(n+1)!} (x-a)^{n+1} \end{aligned}$$

The $x < a$ case follows similarly and can be accounted for by replacing $(x-a)^{n+1}$ in the last line with $|x-a|^{n+1}$.

Therefore,

$$|R_n(x)| \leq \frac{M}{(n+1)!} |x-a|^{n+1}$$

□

8.8.3 Integral Approximations

Consider the definite integral

$$\int_0^1 \frac{\sin(x)}{x} dx$$

Try as you might, you will not be able to solve this integral exactly. According to a theorem by the mathematician Joseph Liouville, the function $\frac{\sin(x)}{x}$ does not have an antiderivative which can be expressed in terms of elementary functions.

This is quite frustrating. The function $\frac{\sin(x)}{x}$ seems simple enough and it is an enormously important function in science (e.g., optics, signal processing, imaging). What are we then to do with an integral like the one above?

We can replace the integrand with a Taylor polynomial approximation!

Recall, the Taylor series for $\sin(x)$ centred at $x = 0$ is

$$\sin(x) = \sum_{n=0}^{\infty} \frac{x^{2n+1}}{(2n+1)!} = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots$$

This means the Taylor series for $\frac{\sin(x)}{x}$ centred at $x = 0$ is

$$\frac{\sin(x)}{x} = \sum_{n=0}^{\infty} \frac{x^{2n}}{(2n+1)!} = 1 - \frac{x^2}{3!} + \frac{x^4}{5!} - \frac{x^6}{7!} + \dots$$

Since this is a polynomial, integration is easy.

$$\begin{aligned} \int_0^1 \frac{\sin(x)}{x} dx &= \int_0^1 \left(\sum_{n=0}^{\infty} \frac{x^{2n}}{(2n+1)!} \right) dx \\ &= \sum_{n=0}^{\infty} \left(\int_0^1 \frac{x^{2n}}{(2n+1)!} dx \right) \\ &= \sum_{n=0}^{\infty} \left[\frac{x^{2n+1}}{(2n+1)((2n+1)!)} \right]_0^1 \\ &= \sum_{n=0}^{\infty} \frac{1}{(2n+1)((2n+1)!)} \\ &= 1 - \frac{1}{18} + \frac{1}{600} - \frac{1}{35280} + \dots \end{aligned}$$

At this point, we can readily extract an approximation by truncating the series. For example, if we keep three terms, then we have

$$\int_0^1 \frac{\sin(x)}{x} dx \approx 1 - \frac{1}{18} + \frac{1}{600} = \frac{1703}{1800} \approx 0.94611$$

And how good is this as an approximation? Observe that our answer is in the form of an alternating series which we could show is convergent using the Alternating Series Test. This means that we can use the Alternating Series Remainder theorem (ASRT) to bound the error as being less than or equal to the magnitude of the first term in the series which we *don't* use in our approximation. The next term is $\frac{-1}{35280}$ so

$$\left| \int_0^1 \frac{\sin(x)}{x} dx - \frac{1703}{1800} \right| \leq \frac{1}{35280} \approx 0.00003$$

So, I'd say our approximation is excellent especially considering all we did was make a small tweak to a well-known Taylor series, integrate a polynomial, and do a bit of arithmetic.

EXERCISE

Determine the value of $\int_0^{\frac{1}{10}} \frac{1}{1+x^3} dx$ with an error bound of at most $\frac{1}{10^{11}}$.

(This integral can actually be computed exactly, but it is a lengthy calculation involving integration by partial fractions followed by a trigonometric substitution. In contrast, you should be able to determine its value with the precision required in the question with just a few lines of work by Taylor expanding the integrand.)

In the previous example and exercise, we lucked out a bit (by design) in that we could express our definite integral in the form of an alternating series (which can be shown to be convergent by the Alternating Series Test). This let us bound the error using the ASRT. However, we could have also used Taylor's inequality. To do that for the example, we would note that we are effectively using the degree-five Taylor polynomial, $T_5(x)$, centred at $x = 0$ for $f(x) = \frac{\sin(x)}{x}$. Next, we would calculate $f^{(6)}(x)$, note it is decreasing on $[0, 1]$, and so in Taylor's inequality take $M = |f^{(6)}(0)| = \frac{1}{7}$. This then gives $|R_5(x)| \leq \frac{x^6}{7!}$. Upon integrating, we would have

$$\left| \int_0^1 \frac{\sin(x)}{x} dx - \int_0^1 T_5(x) dx \right| = \left| \int_0^1 R_5(x) dx \right| = \frac{1}{7(7!)} = \frac{1}{35280}$$

In this case Taylor's inequality and the ASRT give the same upper bound on the remainder. In general, Taylor's inequality will never give a *tighter* bound than the ASRT, but its advantage is that it can always be applied. The catch is that the ASRT can only be applied to an alternating series which can be shown to be convergent by the Alternating Series Test.

Let's look at an example where we use Taylor's inequality in the context of an integral approximation.

Example 22

Use the second-degree Taylor polynomial for $f(x) = \sqrt{1 + \sin(x)}$ centred at $x = 0$ to approximate

$$\int_0^{1/2} \sqrt{1 + \sin(x)} dx$$

and then use Taylor's inequality to determine an upper bound on the error.

Solution: With a bit of work, we can show

$$\begin{aligned} f(x) &= \sqrt{1 + \sin(x)} & \implies & f(0) = 1 \\ f'(x) &= \frac{\cos(x)}{2\sqrt{1 + \sin(x)}} & \implies & f'(0) = \frac{1}{2} \\ f''(x) &= -\frac{\sqrt{1 + \sin(x)}}{4} & \implies & f''(0) = -\frac{1}{4} \end{aligned}$$

We can construct the second-degree Taylor polynomial of $f(x)$ as follows

$$T_2(x) = f(0) + f'(0)x + \frac{1}{2!}f''(0)x^2 = 1 + \frac{1}{2}x - \frac{1}{8}x^2$$

This gives us the following approximation for the given integral

$$\begin{aligned} \int_0^{1/2} \sqrt{1 + \sin(x)} \, dx &\approx \int_0^{1/2} \left(1 + \frac{1}{2}x - \frac{1}{2}x^2\right) dx \\ &= \left[x + \frac{1}{4}x^2 - \frac{1}{24}x^3\right]_0^{1/2} \\ &= \frac{1}{2} + \frac{1}{16} - \frac{1}{192} \\ &= \frac{107}{192} \\ &\approx 0.5573 \end{aligned}$$

Now let's use Taylor's inequality to determine an upper bound on the error of this approximation.

We first compute that on $[0, \frac{1}{2}]$ we will have

$$|f^{(3)}(x)| = \frac{\cos(x)}{8\sqrt{1 + \sin(x)}}$$

and since on this interval $\cos(x)$ is decreasing while the denominator is increasing, then $|f^{(3)}(x)|$ takes its maximum value at $x = 0$. This means if we take $M = |f^{(3)}(0)| = \frac{1}{8}$ then we ensure $|f^{(3)}(x)| \leq M$ on $[0, \frac{1}{2}]$.

Therefore, by Taylor's inequality, we have for $x \in [0, \frac{1}{2}]$

$$|R_2(x)| \leq \frac{M}{3!}|x|^3 = \frac{1}{48}x^3$$

It follows that

$$\begin{aligned} \left| \int_0^{1/2} \sqrt{1 + \sin(x)} \, dx - \int_0^{1/2} T_2(x) \, dx \right| &\leq \int_0^{1/2} |R_2(x)| \, dx \\ &= \frac{1}{48} \int_0^{1/2} x^3 \, dx \\ &= \frac{1}{48} \left(\frac{1}{4} \left(\frac{1}{2} \right)^4 \right) \\ &= \frac{1}{3072} \\ &\approx 0.0003 \end{aligned}$$

Therefore, rounded to the fourth decimal place, $\int_0^{1/2} \sqrt{1 + \sin(x)} \, dx = 0.5573 \pm 0.0003$.

EXERCISE

Use the second-degree Taylor polynomial for $g(x) = e^{\sin(x)}$ centred at $x = 0$ to approximate

$$\int_0^{1/3} e^{\sin(x)} dx$$

and then use Taylor's inequality to determine an upper bound on the error.

Chapter 9

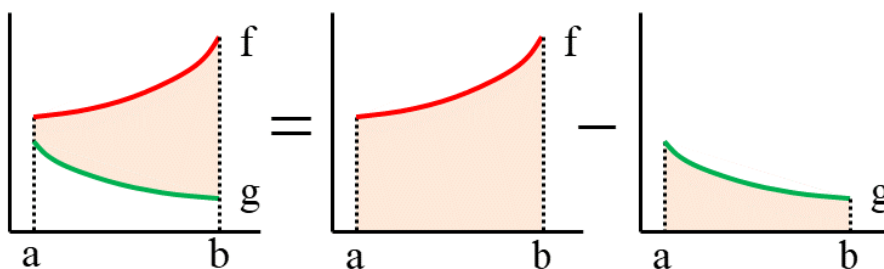
Geometrical Applications of Integration

In this chapter, we explore integration as a tool for calculating geometrical quantities like length, area, and volume.

9.1 Areas Between Curves

We begin with some low-hanging (but essential) fruit: finding the area between two curves. Suppose we have curves $y = f(x)$ and $y = g(x)$ and we'd like to find the area between these two curves on the interval $[a, b]$. For the moment, let's assume that $f(x) \geq g(x) \geq 0$ on $[a, b]$.

We can get the area between these curves by simply computing the definite integral of the difference $f(x) - g(x)$ on the interval $[a, b]$ since it will be equal to the area under $y = f(x)$ minus the area under $y = g(x)$.



$$\begin{aligned} \left(\begin{array}{c} \text{Area between} \\ y = f(x) \text{ and } y = g(x) \end{array} \right) &= \int_a^b f(x) dx - \int_a^b g(x) dx \\ &= \int_a^b [f(x) - g(x)] dx \end{aligned}$$

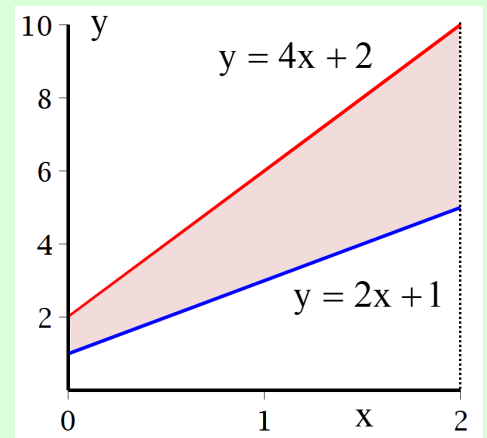
Example 1 Determine the area enclosed by the lines $y = 4x + 2$, $y = 2x + 1$, $x = 0$, and $x = 2$.

Solution: We can reframe this as an area between curves problem by recognizing that the vertical lines $x = 0$ and $x = 2$ give us an interval over which we'd like to find the area between the curves $y = 4x + 2$ and $y = 2x + 1$.

Next, we note that the curve $y = 4x + 2$ sits above the curve $y = 2x + 1$ which sits above the x -axis on the interval $[0, 2]$. So, if we let $f(x) = 4x + 2$ and $g(x) = 2x + 1$, then $f(x) \geq g(x) \geq 0$.

The area between the two curves is then

$$\begin{aligned} \int_0^2 [(4x + 2) - (2x + 1)] dx &= \int_0^2 [2x + 1] dx \\ &= [x^2 + x]_0^2 \\ &= (2^2 + 2) - (0^2 + 0) \\ &= 6 \text{ units}^2 \end{aligned}$$



Integrating the difference of functions to find the area between two curves still works if one or both curves lies below the x -axis. That is, if $f(x) \geq g(x)$ on the interval $[a, b]$, then $\int_a^b [f(x) - g(x)] dx$ will always give a positive value equal to the area between the curves $y = f(x)$ and $y = g(x)$. Let's see this in action.

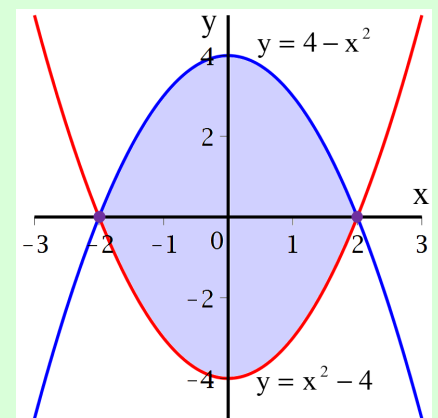
Example 2 Determine the area enclosed by the parabolas $y = x^2 - 4$ and $y = 4 - x^2$.

Solution: We are given the equations of two curves, but we are not told over what interval we should be finding the area between these two curves. However, the wording of the question (i.e., that the curves *enclose* an area) suggests that the curves intersect.

A quick sketch or a quick bit of algebra reveals that the two curves intersect on the x -axis at $x = -2$ and $x = 2$. We can also see that on this interval $y = 4 - x^2$ sits above $y = x^2 - 4$.

Therefore, the area enclosed by the parabolas is

$$\begin{aligned} \int_{-2}^2 [(4 - x^2) - (x^2 - 4)] dx &= \int_{-2}^2 [8 - 2x^2] dx \\ &= \left[8x - \frac{2}{3}x^3 \right]_{-2}^2 \end{aligned}$$



$$\begin{aligned}
 &= \left(8(2) - \frac{2}{3}(2)^3\right) - \left(8(-2) - \frac{2}{3}(-2)^3\right) \\
 &= \frac{64}{3} \text{ units}^2
 \end{aligned}$$

EXERCISE

Determine the area enclosed by the line $y = x + 2$ and the parabola $y = \frac{1}{2}x^2 - 2$.

In some instances, we may be interested in the area between two curves $y = f(x)$ and $y = g(x)$ on an interval where the two curves cross inside the interval. When this happens, we need to swap which curve we identify as being above the other curve. Mathematically, we can do this by modifying our formula for the area between two curves with an absolute value.

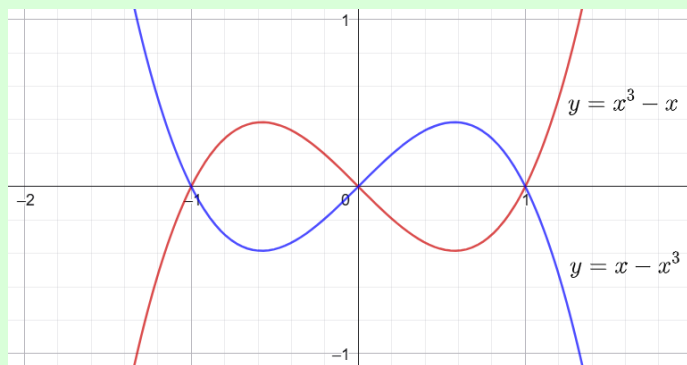
Fact 1 The area between the curves $y = f(x)$ and $y = g(x)$ on the interval $[a, b]$ is given by the definite integral

$$\int_a^b |f(x) - g(x)| dx$$

When evaluating such an integral, we manage the absolute value sign by splitting the integral wherever the integrand changes sign.

Example 3 Determine the area enclosed by $y = x^3 - x$ and $y = x - x^3$.

Solution: We first observe that these curves intersect at $x = -1$, $x = 0$, and $x = 1$. Therefore, we will need to integrate from $x = -1$ to $x = 1$.



We also observe that $y = x^3 - x$ sits above $y = x - x^3$ on $(-1, 0)$ and below it on $(0, 1)$. Therefore, the area enclosed by the two curves is

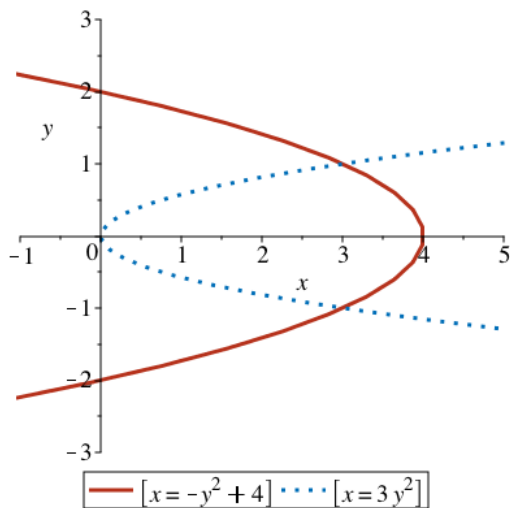
$$\begin{aligned} \int_{-1}^1 |(x^3 - x) - (x - x^3)| dx &= \int_{-1}^0 ((x^3 - x) - (x - x^3)) dx + \int_0^1 ((x - x^3) - (x^3 - x)) dx \\ &= \int_{-1}^0 (2x^3 - 2x) dx + \int_0^1 (2x - 2x^3) dx \\ &= \left[\frac{1}{2}x^4 - x^2 \right]_{-1}^0 + \left[x^2 - \frac{1}{2}x^4 \right]_0^1 \\ &= \left[0 - \left(\frac{1}{2} \right) \right] + \left[\left(\frac{1}{2} \right) - 0 \right] \\ &= 1 \end{aligned}$$

(Note, we could have also made use of symmetry in this problem. In particular, the area enclosed on $[-1, 0]$ was equal to the area enclosed from $[0, 1]$. Making this observation early would have meant that we only needed to determine the area enclosed on one of these intervals and then double it.)

Another complication that we can encounter is having curves which are not oriented in an ideal way. For example, suppose we wish to find the area enclosed by the curves:

$$\begin{aligned} x &= 4 - y^2 \\ x &= 3y^2 \end{aligned}$$

These are just parabolas but opening in the negative and positive x -directions, respectively.



We want to find the area between these two curves from $x = 0$ to $x = 4$ but we can't just integrate along x immediately. The curves are described as functions of y and neither function is always "above" or "below" the other.

There are two (ultimately, equivalent) approaches we can take to handle this problem:

1. Integrate along the y -direction.
2. Swap the roles of x and y and integrate along the x -direction.

Approach 1: Integrate Along y

First, we find that the curves intersect at $y = -1$ and $y = 1$. These will serve as our limits of integration. Next, we observe that if we turn our heads so that the positive x -direction is "up", then $x = 4 - y^2$ sits above $x = 3y^2$.

Now we're ready to integrate and let's use symmetry across the line $x = 0$ to turn an integral over $y \in [-1, 1]$ to two times the same integral over $y \in [0, 1]$.

$$\begin{aligned}
 A &= \int_{y=-1}^1 [(4 - y^2) - (3y^2)] dy \\
 &= 2 \int_{y=0}^1 [(4 - y^2) - (3y^2)] dy \\
 &= 8 \int_{y=0}^1 [1 - y^2] dy \\
 &= 8 [1 - y^2]_{y=0}^1 \\
 &= \frac{16}{3}
 \end{aligned}$$

Approach 2: Swap x and y

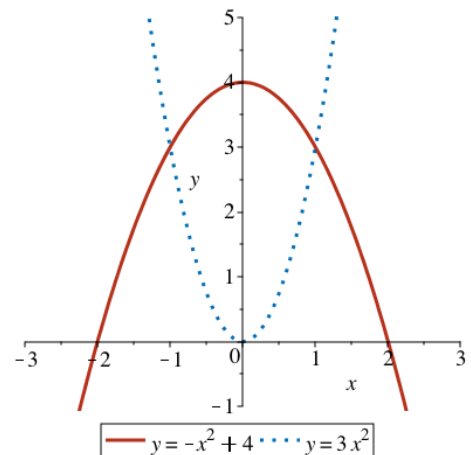
First, we reframe the problem as finding the area between the curves

$$\begin{aligned}
 y &= 4 - x^2 \\
 y &= 3x^2
 \end{aligned}$$

This does not change the area enclosed.

Intersections occur at $x = -1$ and $x = 1$ and $y = 4 - x^2$ sits above $y = 3x^2$.

We leave it as an exercise to fill in the details, but this approach also yields an enclosed area of $\frac{16}{3}$.



EXERCISE

Determine the area enclosed by $x = 2y^2$ and $x = 4 + y^2$.

9.2 Volumes

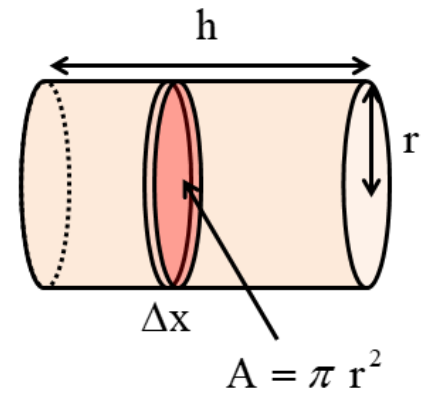
We have seen how to use integration to determine the area of a planar (two-dimensional) region. In this section, we will develop approaches for determining the volumes of a particular kind of solid using integration techniques we've already established.

The basic strategy we'll employ is to imagine dividing the solid up into slices perpendicular to an axis. Without loss of generality, let's suppose we're using the x -axis. If we can write down a function describing the area of a slice at an arbitrary value of x , then we can integrate this function to get the volume. Let's see how this works with a familiar object.

Consider a right cylinder with radius r and height h . Let's orient this cylinder so that the axis of rotational symmetry lies on the x -axis.

Next, imagine we divide the cylinder up into n slices of thickness $\Delta x = \frac{h}{n}$ and area $A = \pi r^2$. The total volume will be the sum of the volumes of the n slices.

$$V = \sum_{i=1}^n A \Delta x = \frac{\pi r^2 h}{n} \sum_{i=1}^n 1 = \pi r^2 h$$



In the limit that the number of slices n tends to infinity so that the thickness of each slice becomes infinitesimal, we could instead write this as

$$V = \int_{x=0}^h A dx = \pi r^2 \int_{x=0}^h dx = \pi r^2 h$$

The Riemann sum and integrals are both trivial here because the areas of the slices do not vary with x . However, you could imagine a more interesting solid where a non-constant function $A(x)$ describes the area of a slice as a function of x . If the object extends from $x = a$ to $x = b$, then its volume would be given by

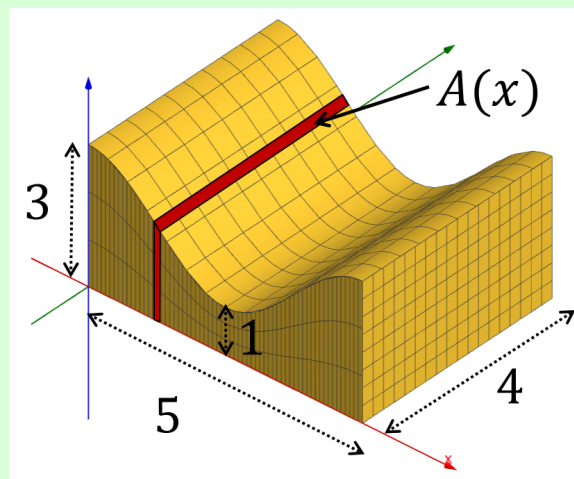
$$V = \int_{x=a}^b A(x) dx$$

Example 4

A wavelike solid extends from $x = 0$ to $x = 5$ and has cross-sectional area

$$A(x) = 8 + 4 \cos\left(\frac{2\pi}{5}x\right)$$

Determine the total volume of the solid.



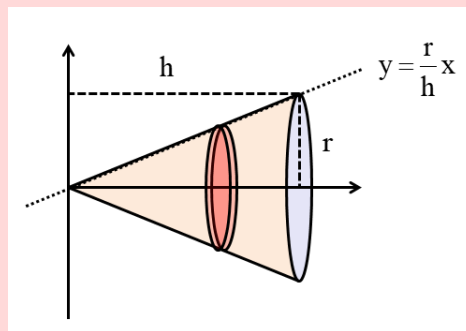
Solution: The volume is

$$\begin{aligned} V &= \int_{x=0}^5 \left(8 + 4 \cos \left(\frac{2\pi}{5} x \right) \right) dx \\ &= \left[8x + \frac{10}{\pi} \sin \left(\frac{2\pi}{5} x \right) \right]_{x=0}^5 \\ &= \left[8(5) + \frac{10}{\pi} \sin(2\pi) \right] - \left[0 + \frac{10}{\pi} \sin(0) \right] \\ &= 40 \end{aligned}$$

EXERCISE

Consider a cone with height h and base radius r . Orient the cone as shown in the diagram so that its axis lies on the x -axis with the cone's vertex at $x = 0$ and base at $x = h$.

- What is the area, A , of a circular slice through this cone perpendicular to the x -axis as a function of x ?
- Integrate $A(x)$ from $x = 0$ to $x = h$ to reproduce the well-known formula for the volume of a cone, $V_{\text{cone}} = \frac{1}{3}\pi r^2 h$.



9.2.1 Volumes by Disks

Setting up an integral to find the volume of a solid given a function describing its cross-sectional area is straightforward. However, finding a function to describe the cross-sectional area can be a whole separate geometry problem on its own. In general, it requires doing a separate integral where the limits of integration are functions rather than constants. This is one way to motivate the need for multiple integrals (i.e., nested integrals). However, in the special case that the volume has rotational symmetry about an axis, we can circumvent needing to do these additional integrals. It is this case we'll explore in this section.

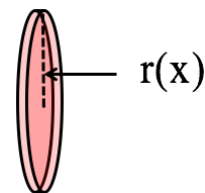
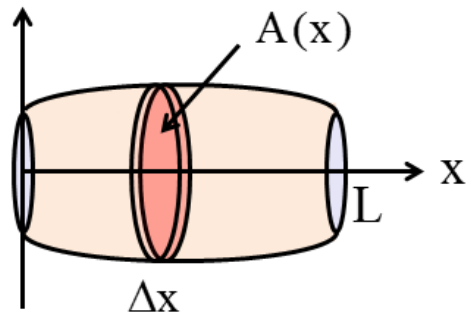
To see how this works, let's position our solid with its axis of rotational symmetry on the x -axis.

We know from above that if we can find a function $A(x)$ to describe the cross-sectional area, then the volume is

$$V = \int_{x=a}^b A(x) dx$$

Now, suppose that the solid has rotational symmetry about the x -axis. This means that there exists a function $r(x)$ such that rotating $y = r(x)$ about the x -axis generates the surface of this solid. It also means that every slice of the solid perpendicular to the x -axis will be a circular disk. Moreover, the radius of a disk at x is equal to $r(x)$ which means that the disk will have area $A(x) = \pi (r(x))^2$.

We summarize this argument with the following fact.



Fact 2

Volume of Revolution by Disks:

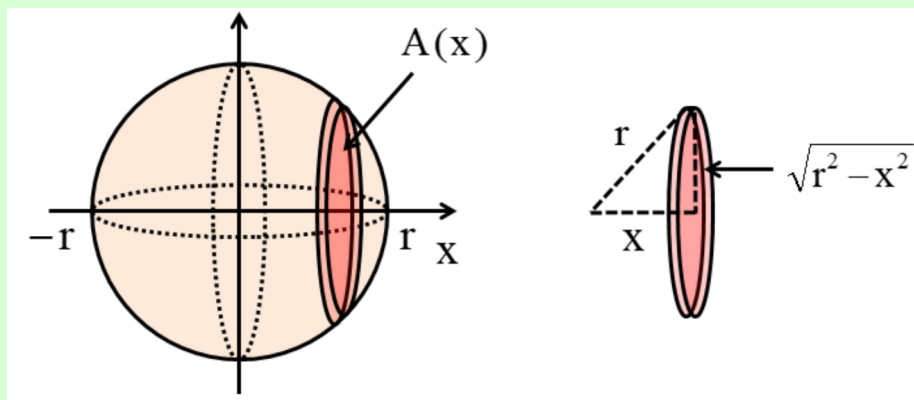
If the area below the curve $y = f(x)$ (where $f(x) \geq 0$) from $x = a$ to $x = b$ is rotated about the x -axis, then the volume of the resulting solid is equal to

$$V = \int_{x=a}^b \pi (f(x))^2 dx$$

Example 5

Determine the volume of a sphere of radius r .

Solution: Observe that we can construct a solid sphere by taking the area below the curve $y = \sqrt{r^2 - x^2}$ from $x = -r$ to $x = r$ and rotating it about the x -axis.



A cross-section of this sphere will be a disk with radius $\sqrt{r^2 - x^2}$. The area of the corresponding disk will be $A(x) = \pi(r^2 - x^2)$.

Integrating these areas from $x = -r$ to $x = r$ gives the volume of the sphere.

$$V_{\text{sph}} = \int_{-r}^r A(x) dx = \int_{-r}^r \pi(r^2 - x^2) dx = \pi \left[r^2x - \frac{1}{3}x^3 \right]_{-r}^r = \frac{4}{3}\pi r^3$$

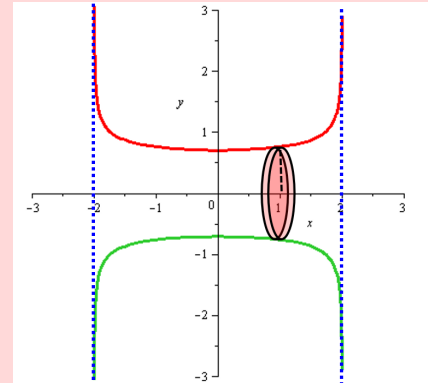
As expected, we have reproduced the well-known formula for the volume of a sphere of radius r .

EXERCISE

Consider the region in the xy -plane between the curve $y = f(x) = (4 - x^2)^{-1/4}$ and the x -axis from $x = -2$ to $x = 2$.

Determine the volume of the solid obtained by rotating this region around the x -axis.

(Note that $f(x)$ has vertical asymptotes at $x = \pm 2$ so this finite volume will be infinite in extent.)



9.2.2 Volumes by Washers

We can adapt our procedure above to also find the volumes of solids which have rotational symmetry about an axis and also have one or more cavities. We do this by recognizing that such volumes are formed by rotating the area between *two* curves around an axis. Therefore, the volume of such a solid can be found by taking the difference between the volume generated by rotating the outer curve around the axis and the volume generated by rotating the inner curve around the axis. In other words, we will end up integrating over a set of infinitesimally thin washers instead of disks.

Fact 3 Volume of Revolution by Washers:

If the area between the curves $y = f(x)$ and $y = g(x)$ (where $f(x) \geq g(x) \geq 0$) from $x = a$ to $x = b$ is rotated about the x -axis, then the volume of the resulting solid is equal to

$$V = \int_{x=a}^b \pi \left((f(x))^2 - (g(x))^2 \right) dx$$

Example 6

Determine the volume of the solid obtained by rotating the area enclosed by the curves $y = x^2$ and $y = \sqrt{x}$ around the x -axis.

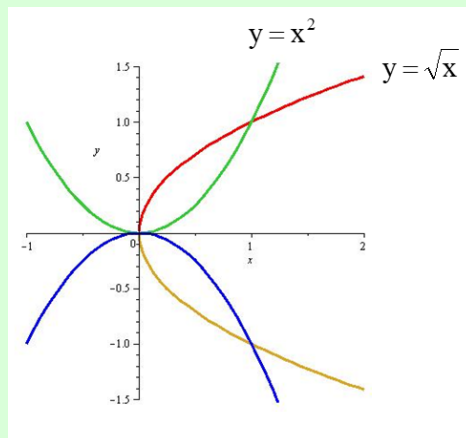
Solution: The area enclosed by the two given curves lies in the first quadrant of the xy -plane between $x = 0$ and $x = 1$. Inside this interval, $y = \sqrt{x}$ lies above $y = x^2$.

Therefore, we can model the volume using an infinite set of washers with outer radius $f(x) = \sqrt{x}$ and inner radius $g(x) = x^2$. The resulting area of each washer is then

$$\begin{aligned} A(x) &= \pi \left((f(x))^2 - (g(x))^2 \right) \\ &= \pi(x - x^4) \end{aligned}$$

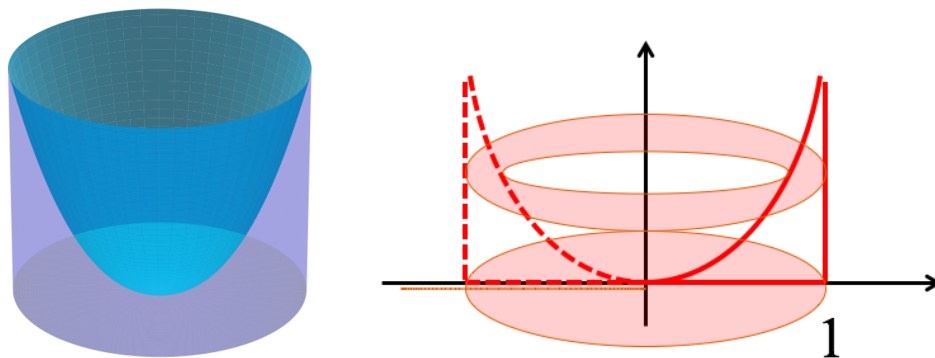
Integrating $A(x)$ from $x = 0$ to $x = 1$ gives the total volume

$$\begin{aligned} V &= \int_{x=0}^1 A(x) dx \\ &= \pi \int_{x=0}^1 (x - x^4) dx \\ &= \pi \left[\frac{1}{2}x^2 - \frac{1}{5}x^5 \right]_{x=0}^1 \\ &= \pi \left(\frac{1}{2} - \frac{1}{5} \right) \\ &= \frac{3\pi}{10} \end{aligned}$$



9.2.3 Volumes by Shells

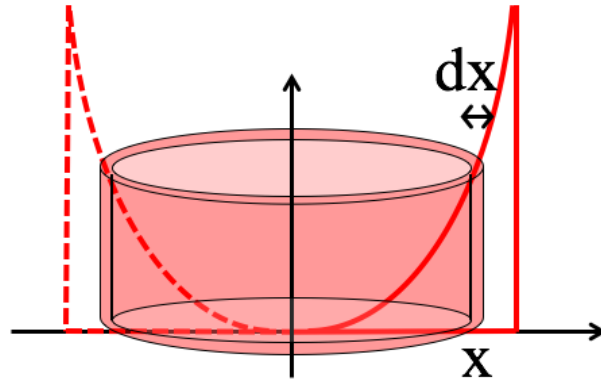
In some cases, volumes of rotation cannot conveniently be computed by integrating over a set of disk or washer-like slices. For example, consider the solid produced by rotating the area enclosed by the curves $y = x^3 + x^2$, $y = 0$, $x = 0$, and $x = 1$ about the y -axis.



The resulting solid is similar to that of a solid cylinder with a bowl-like cavity. This suggests that we could determine the volume by integrating over a set of washers in the y -direction.

Following this approach, the first step is to determine the area of a washer as a function of y . The outer radius of such a washer will be $f(y) = 1$. However, to find the inner radius we need to solve $y = x^3 + x^2$ for x which is not even remotely straightforward to do. Effectively, we've hit a roadblock. Thankfully, we are clever enough to come up with a better way to approach this particular problem.

Instead of dividing the volume up into washers, let's divide it up into concentric cylindrical shells. If we can determine the lateral surface area of a shell of thickness dx in terms of x , then we can integrate in the x -direction to get the total volume of the solid.



To find the surface area of a cylindrical shell $A(x)$, we need its height and radius. The radius is just the coordinate x and the height is given by the curve $y = x^3 + x^2$. In other words, $r(x) = x$ and $h(x) = x^3 + x^2$. Therefore, we have

$$A(x) = 2\pi r(x) h(x) = 2\pi x(x^3 + x^2) = 2\pi(x^4 + x^3)$$

Observe that we now have an expression we can integrate and we did not need to invert a cubic equation to get it.

The volume of the solid is

$$\begin{aligned} V &= \int_{x=0}^1 2\pi x h(x) dx \\ &= 2\pi \int_{x=0}^1 (x^4 + x^3) dx \\ &= 2\pi \left[\frac{1}{5}x^5 + \frac{1}{4}x^4 \right]_{x=0}^1 \\ &= 2\pi \left(\frac{1}{5} + \frac{1}{4} \right) \\ &= \frac{9\pi}{10} \end{aligned}$$

We can summarize this procedure as follows:

Fact 4 Volume of Revolution by Cylindrical Shells:

If the area between the curves $y = f(x)$ and $y = g(x)$ (with $f(x) \geq g(x)$) from $x = a$ to $x = b$ is rotated about the y -axis, then the volume of the resulting solid is equal to

$$V = \int_{x=a}^b 2\pi r(x) h(x) dx$$

where $r(x) = x$ and $h(x) = f(x) - g(x)$.

REMARKS

- The formula above for a volume of revolution by shells assumes the rotation is done around the y -axis, but can be restated as needed to centre the cylindrical shells on any axis.

This brings us to an important related point. While we have provided formulas for each of the methods under certain conditions, we encourage you to focus on the geometry of the underlying construction which led to that formula. This way, if you need to describe, for example, washers centred on the $y = 2$ axis, you're able to do this.

- In the example above, the cylindrical shell method proved to be more effective than the washer method (which is a generalization of the disk method). In general, it could be the case that either method is preferable over the other. Ultimately, it depends on the geometry of the solid and whether it is easier to integrate the function describing the areas of the cylindrical shells or the functions describing the areas of the disks/washers which can be stacked to make up the solid.

Example 7

Determine the volume of the solid obtained by rotating the area enclosed by $y = \sqrt{x}$ and $y = x^2$ about the y -axis.

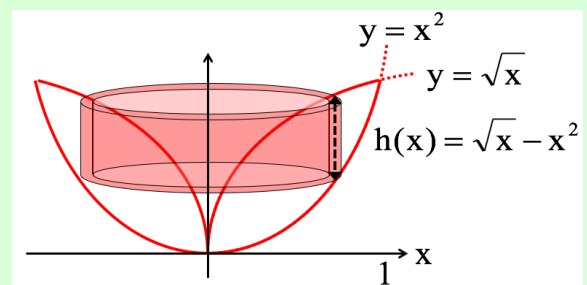
Solution: Let's use cylindrical shells centred on the y -axis. The height $h(x)$ of a cylindrical shell at x is then given by

$$h(x) = \sqrt{x} - x^2$$

The radius of a shell is $r(x) = x$ and the shells range from $x = 0$ to $x = 1$.

Therefore, the volume is

$$\begin{aligned} V &= \int_{x=0}^1 2\pi x(\sqrt{x} - x^2) dx \\ &= 2\pi \int_{x=0}^1 (x^{3/2} - x^3) dx \end{aligned}$$

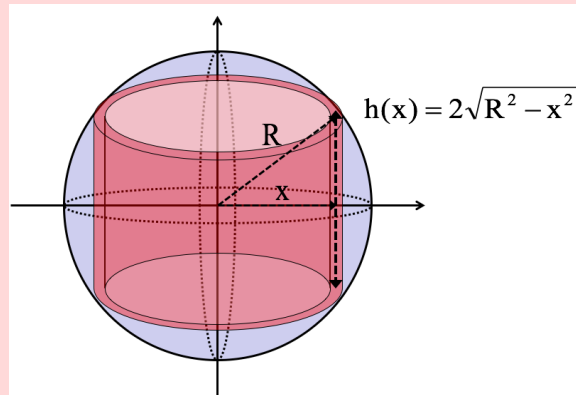


$$\begin{aligned}
 &= 2\pi \left[\frac{2}{5}x^{5/2} - \frac{1}{4}x^4 \right]_{x=0}^1 \\
 &= 2\pi \left(\frac{2}{5} - \frac{1}{4} \right) \\
 &= \frac{3\pi}{10}
 \end{aligned}$$

Observe that we get the same answer as we did when rotating the same area about the x -axis and using the washer method. Can you see why we get the same solid both ways?

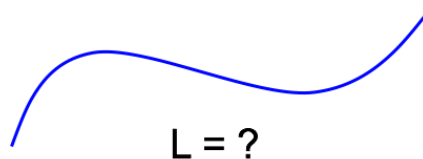
EXERCISE

Show that the volume of a solid sphere of radius R is $\frac{4}{3}\pi R^3$ using cylindrical shells.



9.3 Arc Length

Suppose we have a curve $y = f(x)$ and we'd like to know the length of that curve over some interval $a \leq x \leq b$. First of all, how do we even define length for a curve that is not straight?



A natural definition might be to imagine laying a string along the curve, then straightening the string, and measuring its length on a flat ruler. But how do we translate that into mathematics?

Here is another idea, we could imagine walking along the curve taking very small steps and then adding up the lengths of all the steps taken. This would only give us an approximation because our steps won't perfectly match up with the curve, but the smaller the steps we

take, the better the approximation becomes. This is starting to sound like something we could do with a bit of calculus. Let's explore this idea a bit more carefully.

Suppose the curve $y = f(x)$ starts at (x_1, y_1) and ends at (x_2, y_2) . A very crude approximation we could make for the length L of the curve is to take the straight line distance between the endpoints.

$$L \approx \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Unless the curve is just a line, we can improve this approximation by adding some intermediate points.

The more intermediate points we add, the better we expect the approximation to get.

If we divide the curve up into n segments, find an expression for the length of the curve using n segments, and take the limit as $n \rightarrow \infty$, then we should get the exact length of the curve. We are constructing a Riemann sum!

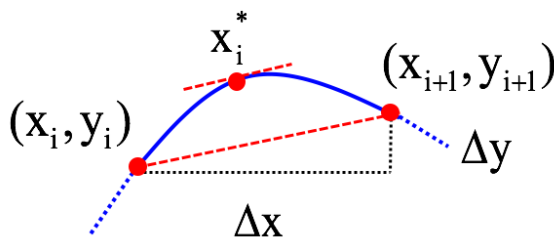
The i -th term in such a Riemann sum will give us the length L_i of the i -th segment. Let's zoom in on the consecutive pair of points: (x_i, y_i) and (x_{i+1}, y_{i+1}) to analyze L_i .

By the Pythagorean theorem

$$L_i = \sqrt{(x_{i+1} - x_i)^2 + (y_{i+1} - y_i)^2}$$

We will assume that the points are evenly spaced in the x direction. This means we can write $x_{i+1} - x_i = \Delta x$ for all i where Δx is a constant.

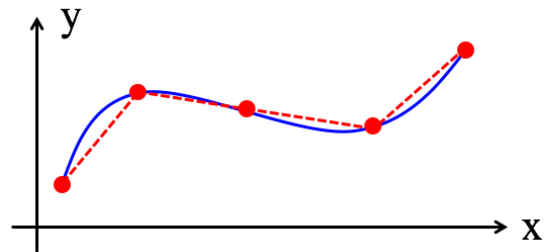
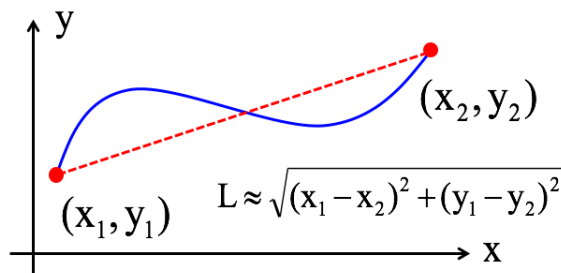
What about the difference $y_{i+1} - y_i$? By the Mean Value Theorem, we know there exists a point $x_i^* \in [x_i, x_{i+1}]$ such that $y_{i+1} - y_i = f'(x_i^*)(x_{i+1} - x_i)$.



Alternatively, we could appeal to the linearization of $f(x)$ at $x = x_i$ to get $y_{i+1} \approx y_i + f'(x_i)(x_{i+1} - x_i)$ and reason that this approximation becomes exact in the limit that $n \rightarrow \infty$, but let's proceed with the more rigorous MVT argument which gives $y_{i+1} - y_i = f'(x_i^*)\Delta x$.

We now have

$$L_i = \sqrt{(\Delta x)^2 + (f'(x_i^*)\Delta x)^2} = \sqrt{1 + (f'(x_i^*))^2} \Delta x$$



Therefore, the total length L of the curve is

$$L = \lim_{n \rightarrow \infty} \sum_{i=1}^n \sqrt{1 + (f'(x_i^*))^2} \Delta x$$

We recognize this as the definition of the definite integral of $\sqrt{1 + (f'(x))^2}$.

Fact 5 Arc Length:

The arc length of the curve $y = f(x)$ between $x = a$ and $x = b$ is equal to

$$\int_{x=a}^b \sqrt{1 + (f'(x))^2} dx$$

Let's test this formula out by computing the circumference of a circle.

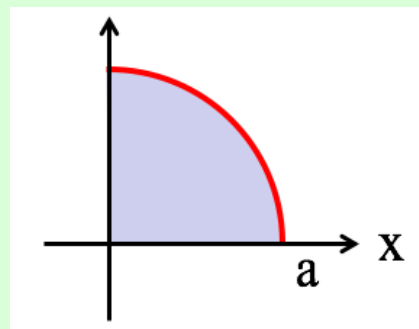
Example 8

Determine the circumference of a circle of radius R .

Since the graph of a circle cannot be described in the form $y = f(x)$, let's consider just the quarter of a circle in the first quadrant which can be described by letting $f(x) = \sqrt{R^2 - x^2}$ with $0 \leq x \leq R$.

For the arc length formula, we need

$$f'(x) = -\frac{x}{\sqrt{R^2 - x^2}}$$



Before we attempt to integrate anything, let's simplify the expression $\sqrt{1 + (f'(x))^2}$ as much as we can.

$$\begin{aligned} \sqrt{1 + (f'(x))^2} &= \sqrt{1 + \frac{x^2}{R^2 - x^2}} \\ &= \sqrt{\frac{R^2}{R^2 - x^2}} \\ &= \frac{R}{\sqrt{R^2 - x^2}} \end{aligned}$$

Now we're ready to integrate to find the circumference and let's not forget that factor of four since we're only computing the length of a quarter-arc.

$$\begin{aligned} C &= 4 \int_{x=0}^R \frac{R}{\sqrt{R^2 - x^2}} dx \quad \text{let } x = R \sin(\theta) \\ &= 4 \int_{\theta=0}^{\pi/2} \frac{R}{\sqrt{R^2 - R^2 \sin^2(\theta)}} (R \cos(\theta) d\theta) \\ &= 4R \int_{\theta=0}^{\pi/2} \frac{\cos(\theta)}{\sqrt{1 - \sin^2(\theta)}} d\theta \\ &= 4R \int_{\theta=0}^{\pi/2} d\theta \end{aligned}$$

$$\begin{aligned}
 &= 4R [\theta]_{\theta=0}^{\pi/2} \\
 &= 2\pi R
 \end{aligned}$$

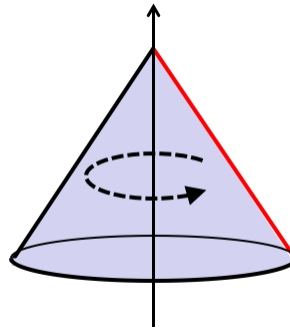
As expected, the circumference of a circle of radius R is $2\pi R$.

EXERCISE

Determine the arc length of the curve $y = x^{3/2}$ from $x = 0$ to $x = 4$.

9.3.1 Surfaces of Revolution

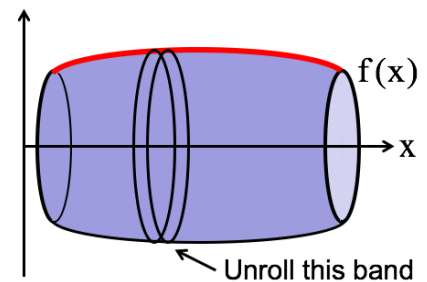
A volume of revolution is a three-dimensional object generated by rotating an area about an axis. A surface of revolution is a two-dimensional object generated by rotating a curve around an axis. For example, a line segment rotated about an axis intersecting one of its endpoints generates a conical surface (or a disk if the line segment and axis are perpendicular).



In this section, we will explore how to determine the area of surfaces generated in this way.

Suppose the curve $y = f(x)$ over some interval is rotated about the x -axis to generate a surface as shown.

Let's focus on a thin "band" generated by rotating a section of the curve $y = f(x)$ with infinitesimal width about the y -axis. The radius of this band will be $f(x)$ so its circumference is $2\pi f(x)$.



The width is a little trickier since the curve $y = f(x)$ is not horizontal and straight, but we can use our knowledge of arc length to work it out. In particular, if the band had an infinitesimal width projected in the x -direction of dx , then its arc length is $\sqrt{1 + (f'(x))^2} dx$. Therefore, the surface area of the band is $2\pi f(x) \sqrt{1 + (f'(x))^2} dx$.

Integrating this quantity to sum up the surface areas of all bands gives an expression for the total surface area.

Fact 6 Surface Area of Revolution:

If the curve $y = f(x)$ from $x = a$ to $x = b$ is rotated about the x -axis to generate a surface, the area of that surface is

$$\int_{x=a}^b 2\pi f(x) \sqrt{1 + (f'(x))^2} dx$$

EXERCISE

Show that if the curve $x = g(y)$ with $c \leq y \leq d$ is rotated about the x -axis to generate a surface, then the area of this surface can be found by evaluating

$$\int_{y=c}^d 2\pi y \sqrt{1 + (g'(y))^2} dy$$

(Note, this provides an alternate method for computing a surface area of revolution in case the curve that is rotated is more conveniently expressed in the form $x = g(y)$ rather than $y = f(x)$.)

Example 9

Suppose that the surface of a car tire is designed to be inscribed in a spherical shell of radius 3 units centred at the origin with $-1 \leq x \leq 1$. Determine the surface area of the tire?

Solution: A sphere of radius 3 can be generated by rotating the curve $y = f(x) = \sqrt{9 - x^2}$ about the x -axis.

We only need to find the area of the portion of such a sphere with $-1 \leq x \leq 1$ and we can do this by dividing up this area into infinitesimal circular strips with circumference $2\pi f(x)$ and width $\sqrt{1 + (f'(x))^2} dx$. This means we will need

$$f'(x) = \frac{-x}{\sqrt{9 - x^2}}$$

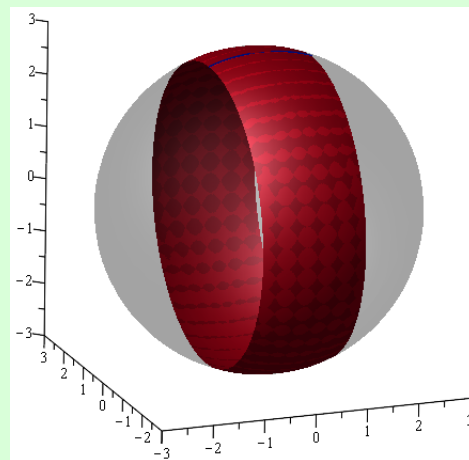
which in turn gives us

$$\sqrt{1 + (f'(x))^2} = \sqrt{1 + \frac{x^2}{9 - x^2}} = \frac{3}{\sqrt{9 - x^2}}$$

We can now determine the area of the surface.

$$\begin{aligned} A &= \int_{x=-1}^1 2\pi \left(\sqrt{9 - x^2}\right) \left(\frac{3}{\sqrt{9 - x^2}}\right) dx \\ &= 6\pi \int_{x=-1}^1 dx \\ &= 12\pi \end{aligned}$$

Note: If we had used a sphere of arbitrary radius R and integrated across the full diameter, we would have derived the formula for the surface area of a sphere, $A = 4\pi R^2$.



EXERCISE

Determine the area of the lateral (side) surface of a right circular cone with height h and base radius r by rotating the line $y = f(x) = h - \frac{h}{r}x$ about the y -axis.

(You should find $A = \pi r \sqrt{h^2 + r^2}$.)

