
MULTIVARIABLE CALCULUS

an introduction

Samuel S. Watson

Textbooks customarily include an abundance of material and leave you to figure out what to focus on. While there are compelling reasons to present many examples in each section, such an approach can be unsatisfactory for those who feel they need a more streamlined source. The present text aims to support a different workflow: (i) read the book for a minimalistic presentation of the core ideas, and (ii) seek alternative sources for additional examples when needed. I hope that you will be able to read each section carefully and work all the examples and exercises, thereby wasting no effort on figuring out what to skip. I recommend the following sources for more exercises, discussion, or examples:

1. Appendix A.6, which lists additional exercises by section. Some of these are original, and others are modified from exercises in James Stewart's *Multivariable Calculus* or Susan J. Colley's *Vector Calculus*.
2. The community calculus textbook, available at communitycalculus.org. This is an open source textbook available (for free) online, and it has *a lot* of exercises.
3. MathInsight (mathinsight.org) has webpages for many multivariable topics with nice exposition and some beautiful applets for exploring the ideas developed in this course.
4. A standard multivariable calculus textbook.
5. 3Blue1Brown, a math video creator with an excellent series on linear algebra. Unfortunately he hasn't done multivariable calculus yet, so these videos will only be available for the vector topics.

The following margin note icons in the text are clickable:

1.  links to a CoCalc worksheet with a relevant calculation (see Appendix A.5 for more discussion)
2.  links to a relevant 3Blue1Brown YouTube video.
3.  links to a relevant page at mathinsight.org

All of the **3D graphics in this PDF may be interactively manipulated**, but that feature requires that the PDF be viewed with Adobe's free Acrobat Reader (<https://get.adobe.com/reader/>). The references in this text are hyperlinked; for example, you can click on Theorem 6.3.1 to navigate directly to Green's theorem.

Please do not hesitate to contact me via email (sswatson@brown.edu) about any mistakes you find in these notes, no matter how minor.



This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/> or send a letter to Creative Commons, 543 Howard Street, 5th Floor, San Francisco, California, 94105, USA. If you distribute this work or a derivative, include the history of the document.

Last updated: August 5, 2019

Contents

1 Transformations of Euclidean space	5
1.1 n -dimensional space	5
1.2 Functions from \mathbb{R}^n to \mathbb{R}^m	6
1.2.1 Visualizing functions	6
1.2.2 Linear transformations	7
1.3 The determinant	9
2 Vectors	13
2.1 Introduction to vectors	13
2.2 The dot product	15
2.3 The cross product	17
3 Three-dimensional Geometry	19
3.1 Lines and planes	19
3.2 Vector-valued functions	22
3.2.1 Paths in space	22
3.2.2 Arclength*	25
3.2.3 Curvature*	26
3.3 Quadric surfaces	27
3.4 Polar, cylindrical, and spherical coordinates	29
4 Multivariable Differentiation	32
4.1 Limits	32
4.2 Partial derivatives	37
4.3 Linear approximation	39
4.4 Taylor's theorem*	42
4.5 Multivariable optimization	43
4.6 Second derivative test	45
4.7 Directional derivative and gradient	47
4.8 The multivariable chain rule	50
4.9 Optimization with Lagrange multipliers	51

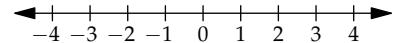
5 Multivariable Integration	55
5.1 Double integration	55
5.2 Triple integration	58
5.3 Polar, cylindrical, and spherical integration	62
5.4 Integration in custom coordinates	64
5.5 Applications of integration*	66
5.5.1 Average value	66
5.5.2 Center of mass	68
5.5.3 Moments of inertia	69
5.5.4 Probability	71
6 Vector Calculus	75
6.1 Vector fields and line integrals	75
6.2 The fundamental theorem of vector calculus	77
6.3 Green's theorem	79
6.4 Surface integrals and flow	81
6.4.1 Surface integrals	81
6.4.2 Flow	85
6.5 Divergence and curl	87
6.5.1 Divergence	88
6.5.2 Curl	89
6.6 Divergence theorem	91
6.7 Stokes' theorem	93
A Appendix	97
A.1 Review	97
A.1.1 Sets and functions	97
A.1.2 Trig review	99
A.1.3 Summation notation	100
A.2 Reference	101
A.2.1 Visualizing functions	101
A.2.2 Polar, cylindrical, and spherical coordinate reference	103
A.3 Technical Appendix	104
A.3.1 The conventional definition of a limit	104
A.3.2 A rigorous formulation of arclength	105
A.3.3 Proof of the second derivative test	108
A.3.4 Wait, is the gradient normal or tangent?	110
A.3.5 The Riemann integral	111
A.4 Big picture	112

A.4.1	How to solve math problems	112
A.4.2	How to write a math solution	113
A.4.3	Thinking about functions	115
A.4.4	The central idea of integral calculus	117
A.5	SageMath	119
A.6	Additional exercises	120

1 Transformations of Euclidean space

1.1 *n*-dimensional space

We visualize the set of **real numbers** (denoted by \mathbb{R} or \mathbb{R}^1) as a line, called the *real number line* (Figure 1.1).



The location of a number x on this line is the point whose *signed* distance from 0 is x . The word *signed* means that distances measured from 0 to a point which is left of 0 count as negative.

The set \mathbb{R}^2 of ordered pairs (x, y) where x and y are real numbers can be thought of as a *plane*, since a point in a plane can be specified by two real numbers: its signed distances from two lines which meet at a right angle. These two lines are called the *x-axis* and the *y-axis* (Figure 1.2).

Figure 1.1 The real number line

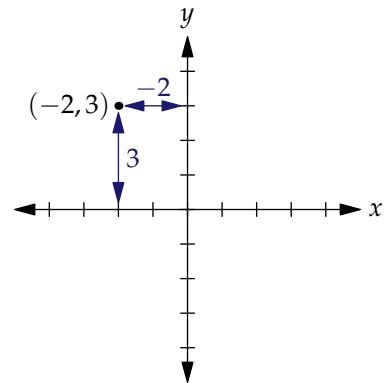


Figure 1.2 Coordinates in \mathbb{R}^2

The set \mathbb{R}^3 of ordered triples of real numbers (x, y, z) can be visualized as (3D) space, since a point in space can be specified by three real numbers: its signed distances from three planes which meet one another at right angles. These planes are called the *xy-plane*, the *yz-plane*, and the *xz-plane*, and their lines of intersection are called the *x-axis*, the *y-axis*, and the *z-axis** (Figure 1.3).

The set \mathbb{R}^4 is defined to be the set of ordered quadruples of real numbers, and similarly for $\mathbb{R}^5, \mathbb{R}^6$, and so on. If n is a positive integer*, we refer to \mathbb{R}^n as a **Euclidean space**. The superscript n is called the **dimension** of \mathbb{R}^n .

Exercise 1.1.1

Write down a formula for the distance on the number line between two real numbers x and y . Repeat for two points in the plane, and for two points in \mathbb{R}^3 .

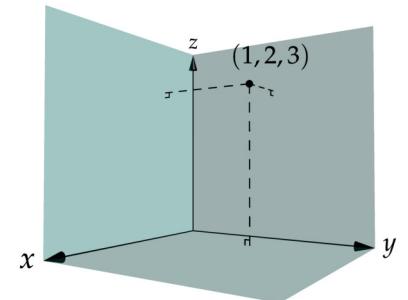


Figure 1.3 Coordinates in \mathbb{R}^3



The **graph** of an equation in \mathbb{R}^n is the set of all points in \mathbb{R}^n whose coordinates satisfy the equation. For example, the graph of the equation $x + y = 1$ in \mathbb{R}^2 is a line. The graph of $x^2 + y^2 = 1$ in \mathbb{R}^3 is a *cylinder*, since a point (x, y, z) satisfies the equation if and only if it is directly above or below a point on the unit circle in the *xy-plane*. Note that the dimension of the Euclidean space under consideration must be specified, since the equation need not involve all of the variables.

Exercise 1.1.2

- Describe the graph of the equation $x(x - 1)(x + 1) = 0$ in \mathbb{R}^1 .
- Interpret $x(x - 1)(x + 1) = 0$ as an equation in \mathbb{R}^2 and describe its graph.
- Interpret $x(x - 1)(x + 1) = 0$ as an equation in \mathbb{R}^3 and describe its graph.

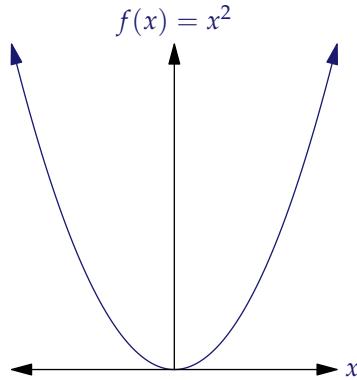


Figure 1.4 The graph of a function from \mathbb{R}^1 to \mathbb{R}^1

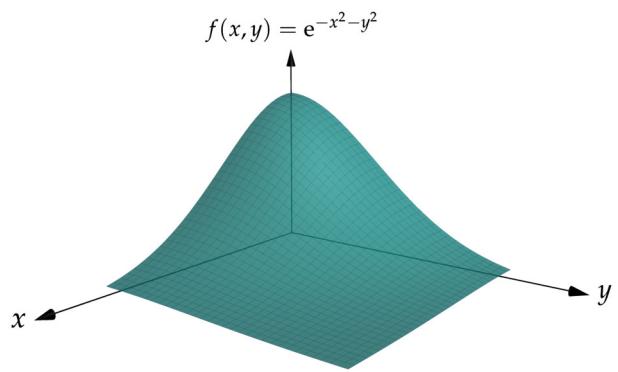


Figure 1.5 A graph of a function from \mathbb{R}^2 to \mathbb{R}^1

1.2 Functions from \mathbb{R}^n to \mathbb{R}^m

1.2.1 VISUALIZING FUNCTIONS

A **function** from \mathbb{R}^1 to \mathbb{R}^1 takes a real number x as input and returns another real number, denoted $f(x)$, as output. We can draw the **graph*** of such a function in $1 + 1 = 2$ dimensions, by associating the horizontal axis with input values and the vertical axis with output values. For example, see Figure 1.4 for a graph of the squaring function.

A function from \mathbb{R}^2 to \mathbb{R}^1 can be visualized in $2 + 1 = 3$ dimensions by using the xy -plane for the input values and the z -axis for the output value. In other words, we plot every triple of the form $(x, y, f(x, y))$ where x and y are real numbers. See Figure 1.5 for a graph of $f(x, y) = e^{-x^2-y^2}$ over a square-shaped region.

The graph of a function from \mathbb{R}^2 to \mathbb{R}^2 would require $2 + 2 = 4$ dimensions to visualize, so we are out of luck there. However, we can visualize a function from \mathbb{R}^2 to \mathbb{R}^2 by drawing a picture of where all the grid lines go*—see Figure 1.6.

* The graph of a function f from \mathbb{R}^1 to \mathbb{R}^1 is defined to be the graph of the equation $y = f(x)$. In other words, (x, y) is on the graph if and only if the equation $y = f(x)$ is satisfied

* One way this method of visualization is different from graphing is that we separate the input values on the left side of the figure from the output values on the right side

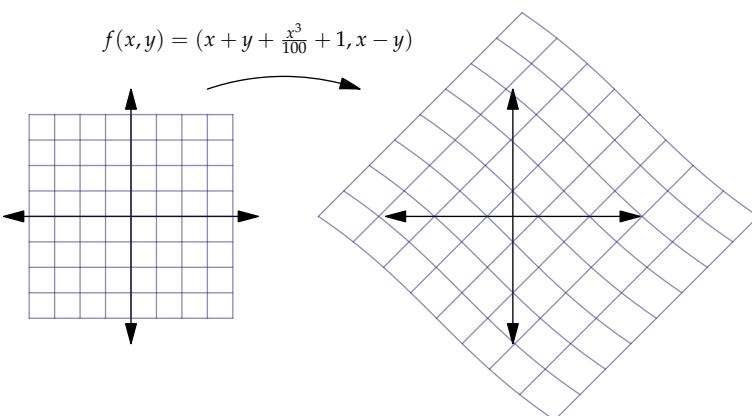


Figure 1.6 A transformation from \mathbb{R}^2 to \mathbb{R}^2

We often refer to a function from \mathbb{R}^2 to \mathbb{R}^2 as a *transformation*, which is a synonym of *function* but is meant to evoke this particular method of visualization. Functions from \mathbb{R}^3 to \mathbb{R}^3 are also called transformations and can be visualized in the same way, but for now we will focus on transformations from \mathbb{R}^2 to \mathbb{R}^2 .

A **level set** of a function is the set of all values in the function's domain which map to a particular value in the codomain. For example, the level sets of the function $f(x, y) = x^2 + y^2$ are concentric circles, as shown in Figure 1.7. The level sets of $f(x, y, z) = x^2 + y^2 - z^2$ are surfaces, as shown in Figure 1.8. Level sets are often called level curves or level surfaces, as appropriate. A

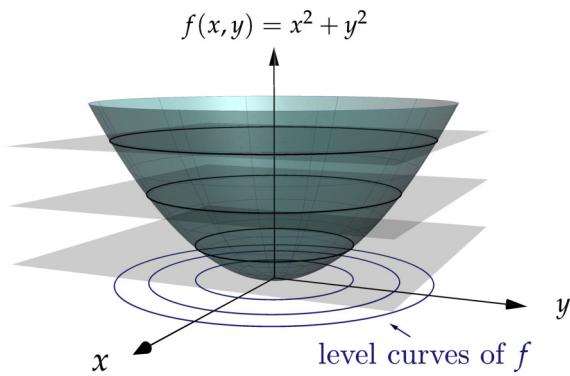


Figure 1.7 The level sets of $f(x, y) = x^2 + y^2$ are concentric circles in \mathbb{R}^2 . They are the projections onto the xy -plane of intersections of the graph of f with “ $z = \text{constant}$ ” planes, as shown

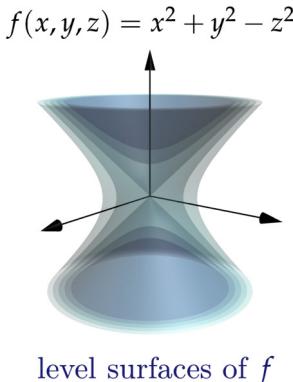


Figure 1.8 The level sets of $f(x, y, z) = x^2 + y^2 - z^2$ are surfaces in \mathbb{R}^3 . The graph of f can't be visualized, but the level surfaces contain some geometric information about f

* This means that Figure 1.7 involves an abuse of notation: the level sets should be points in \mathbb{R}^2 rather than points in the xy -plane in \mathbb{R}^3 . However, it is common practice to make such implicit use of the natural association $(x, y) \leftrightarrow (x, y, 0)$ between \mathbb{R}^2 and the xy -plane

picture showing the level sets corresponding to several equally spaced codomain values is called a **contour plot**.

A function from \mathbb{R}^2 to \mathbb{R}^1 can be visualized using its graph or using a contour plot, whereas the graph of a function from \mathbb{R}^3 to \mathbb{R}^1 can't be visualized spatially. Thus contour plotting provides an important visualization tool for functions of three variables.

Contour plots have some shortcomings: unless the output value corresponding to each level set is identified using a labeling or coloring scheme, information about the values of the function is missing from the picture. In other words, graphing a function combines input and output into a single figure, while level sets are drawn in the function's *domain* only*.

1.2.2 LINEAR TRANSFORMATIONS

One of the key ideas of differential calculus is to uses *linear* functions to approximate functions with curvy graphs. This approximation is useful because (i) all differentiable functions look increasingly linear as you zoom in around a point, and (ii) linear functions are very simple. We will apply the same principle in higher dimensions and use linear transformations to approximate non-linear ones. Just as you learned about linear functions before learning single-variable calculus, we will explore linear transformations before discussing how differentiation generalizes to the multivariable setting.

So what is a linear function from \mathbb{R}^2 to \mathbb{R}^2 ? Single-variable functions of the form $f(x) = mx + b$, where m and b are constants, are often called linear. However, we will take a slightly different view by requiring $b = 0$, so that only functions of the form $f(x) = mx$ are considered linear.* Our definition of linearity in higher dimensions is analogous: only terms of the form “constant times variable” are allowed:

Definition 1.2.1: Linearity

A function f from \mathbb{R}^2 to \mathbb{R}^2 is **linear** if there exist* $a, b, c, d \in \mathbb{R}$ so that

$$f(x, y) = (ax + by, cx + dy).$$

* Actually, this view is ubiquitous in mathematics, because the essence of linearity is the pair of properties $f(x + y) = f(x) + f(y)$ and $f(ax) = af(x)$, and we only get these if we insist $b = 0$

* This notation means that a, b, c, d are in the set of real numbers

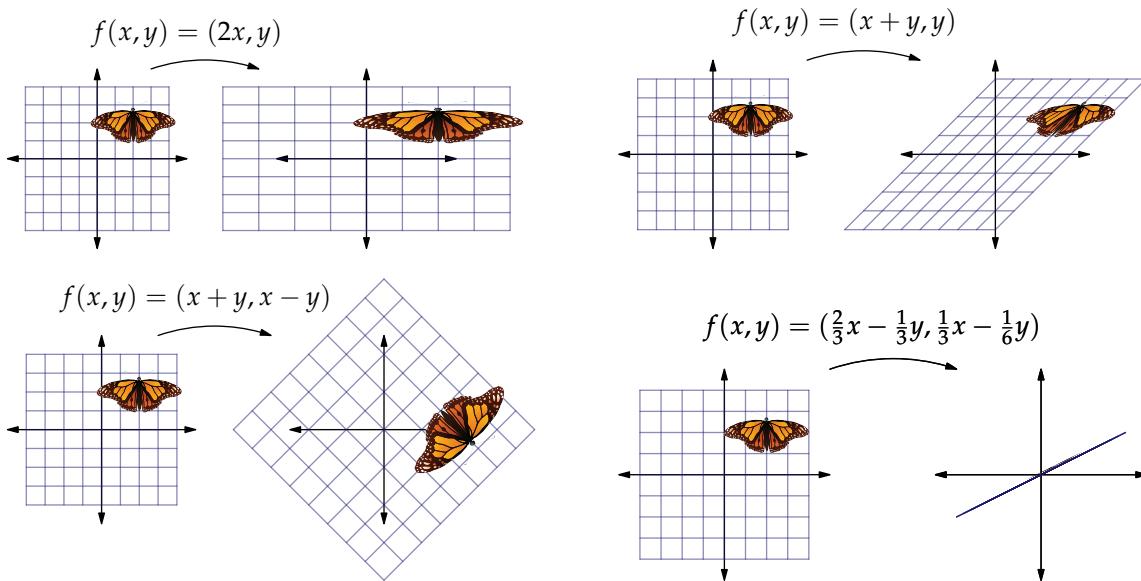


Figure 1.9 Four linear transformations from \mathbb{R}^2 to \mathbb{R}^2

Figure 1.9 shows four examples of linear transformations from \mathbb{R}^2 to \mathbb{R}^2 . The transformations in Figure 1.9 *scale*, *shear*, *rotate/scale*, and *project*, respectively. Not shown is a *reflection*, such as $f(x,y) = (x,-y)$. This list is complete in the sense that every linear transformation from \mathbb{R}^2 to \mathbb{R}^2 can be written as a composition of transformations of these basic types.

We might guess, based on Figure 1.9, that linear transformations map equally spaced lines to equally spaced lines, where coincident lines count as equally spaced with a spacing of zero (as in the last example). This is almost accurate: sometimes equally spaced lines can map to equally spaced *points* (Exercise 1.2.2). The following theorem gives a geometric characterization of linearity.

Theorem 1.2.1

A function from \mathbb{R}^2 to \mathbb{R}^2 is linear if and only if it maps the origin to the origin and equally spaced lines* to equally spaced lines or points.

* Equally spaced lines are lines which are parallel and for which the distances between consecutive lines are the same

Exercise 1.2.1

Use Theorem 1.2.1 to show that if $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ rotates every point counterclockwise about the origin by 30° , there necessarily exist $a, b, c, d \in \mathbb{R}$ such that $f(x,y) = (ax + by, cx + dy)$ for all $(x,y) \in \mathbb{R}^2$.

Exercise 1.2.2

Show that the linear function $f(x,y) = (2x,0)$ maps any collection of equally spaced vertical lines to a collection of equally spaced points.

1.3 The determinant

The *slope* of a linear function from \mathbb{R}^1 to \mathbb{R}^1 measures how it distorts length. For example, the function $f(x) = 3x$ maps any interval $[a, b]$ to the interval $[3a, 3b]$ which is three times as long. The function $g(x) = -\frac{1}{2}x$ maps any interval to an interval which is half as long, and it also flips the interval around. We can say that the absolute value of the slope of a linear function is the *factor by which lengths are transformed*, and the sign of the slope tells us whether the function reverses the real number line.

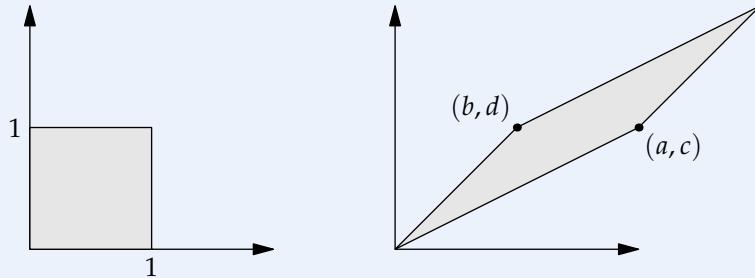
So, what's the corresponding idea for transformations from \mathbb{R}^2 to \mathbb{R}^2 ? Can we look at a linear transformation and conveniently calculate the factor by which that transformation distorts *areas*? Yes!

For each linear transformation in Figure 1.9, the quadrilaterals on the image side of the picture are all congruent. This suggests that the linear transformation does indeed transform every area by the same factor. Taking this fact as given, it suffices for us to consider the image of a single square, which we will take to be $[0, 1]^2$, the set of points both of whose coordinates are between 0 and 1.

Example 1.3.1

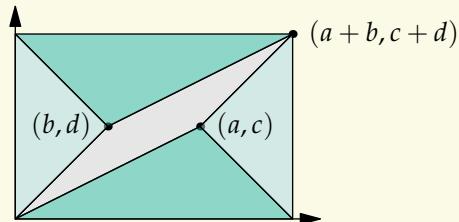
Find the area of the image of the unit square $[0, 1]^2$ under the transformation

$$f(x, y) = (ax + by, cx + dy).$$



Solution

The area of the unit square can be calculated by filling in some triangles to get a complete rectangle, as follows:



The area of the larger rectangle is $(a + b)(c + d)$, and the total area of the triangles we added is $2 \cdot \frac{1}{2}(a + b)(c) + 2 \cdot \frac{1}{2}(c + d)(b)$. Subtracting, we get that the area of the parallelogram is $ad - bc$.

We are not quite finished, however. Note that we assumed in our diagram that the line segment from the origin to (a, c) is clockwise from the line segment from the origin to (b, d) . If we switched these line segments around, the same reasoning would have given us the formula $bc - ad$. We can put this

* Reversing the orientation of three points A , B , and C means that if these points are given in counterclockwise order around the triangle ABC , then their images are in clockwise order around the triangle they form

all together by saying that the factor by which areas are transformed is $|ad - bc|$.

We can interpret the result of Example 1.3.1 as follows: $ad - bc$ tells us how $f(x, y) = (ax + by, cx + dy)$ transforms areas (via its absolute value) and whether applying f to the three points $(0, 0)$, $(1, 0)$, and $(0, 1)$ reverses their orientation* (via its sign). This idea is important enough to deserve its own name. To simplify the definition, we refer to length as 1-dimensional volume and area as 2-dimensional volume.

Definition 1.3.1: Determinant

The **determinant** of a linear transformation from \mathbb{R}^n to \mathbb{R}^n is the signed factor by which it transforms n -dimensional volumes.

!!!

We have already figured out that the determinant of $f(x) = mx$ from \mathbb{R}^1 to \mathbb{R}^1 is the slope m , and the determinant of a function $f(x, y) = (ax + by, cx + dy)$ from \mathbb{R}^2 to \mathbb{R}^2 is given by the formula $ad - bc$.

For convenience, we sometimes represent a linear function by arranging its coefficients into a grid of numbers called a *matrix*. By convention, rows correspond to coordinates of the output of the function, and columns correspond to the variables. So, for example, $f(x, y) = (ax + by, cx + dy)$ is represented by the matrix $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$. So we have

$$\det[m] = m, \text{ and } \det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc.$$

Exercise 1.3.1

Find the determinant of each of the following matrices, and draw the image of the unit square under the corresponding linear transformations to see that the value of the determinant you computed makes sense.

$$(a) \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$

$$(b) \begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix}$$

$$(c) \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$$

$$(d) \begin{bmatrix} 2 & 1 \\ 4 & 2 \end{bmatrix}$$

We define the determinant of a linear transformation

$$f(x, y, z) = (ax + by + cz, dx + ey + fz, gx + hy + iz)$$

to be the product of two quantities:

- (i) the volume of the three-dimensional shape, called a *parallelepiped*, whose vertices are the images under f of the vertices of the unit cube $[0, 1]^3$ (Figure 1.10).
- (ii) a factor of ± 1 which is equal to -1 if and only if the orientation of a small loop drawn on a face is reversed (see Figure 1.10), from the point of view of a small person standing on the solid with their head pointing toward the outside.

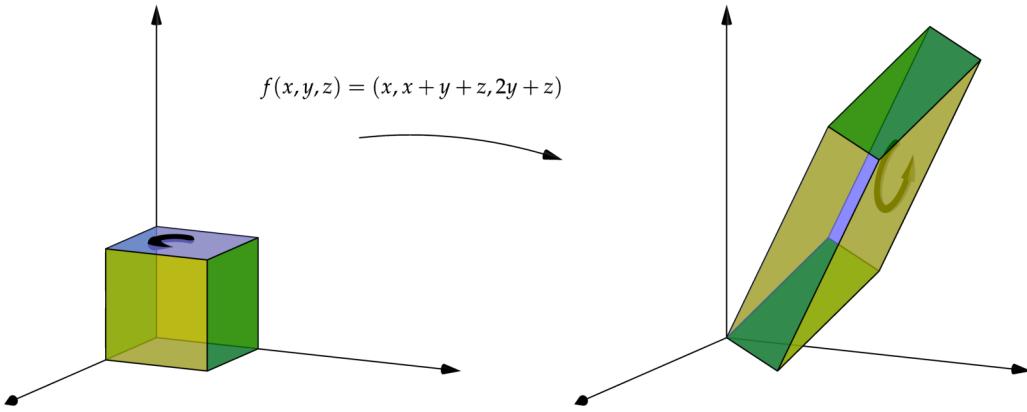


Figure 1.10 A linear transformation from \mathbb{R}^3 to \mathbb{R}^3

* The Geometry of Determinants section of the free online book *Linear Algebra* by Jim Hefferon includes a lengthy discussion of $n \times n$ determinants

The formula for the determinant turns out to be*

$$\det \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} = aei - afh - bdi + bfg + cdh - ceg.$$

Unlike the formula $ad - bc$ for the 2×2 matrix, this formula is not easy to memorize. Let's abbreviate $\det[\]$ to $| |$ and write this formula as

$$\left| \begin{array}{ccc} a & b & c \\ d & e & f \\ g & h & i \end{array} \right| = +a \left| \begin{array}{cc} e & f \\ h & i \end{array} \right| - b \left| \begin{array}{cc} d & f \\ g & i \end{array} \right| + c \left| \begin{array}{cc} d & e \\ g & h \end{array} \right|.$$

This formula is *still* not easy to memorize, so let's break it down: each term on the right-hand side includes three factors: (i) a $+1$ or -1 which alternates starting with $+1$, (ii) an entry from the top row (going from left to right), and (iii) the determinant of the matrix you get when you remove the row and column of that entry from the original matrix. These smaller matrices are called *minors*, and this method of calculating the determinant is called **expansion by minors** along the first row. You can also expand by minors along any row or column (see Exercise 1.3.3 below), but if it's an even-numbered row or column, then the signs start with -1 instead of $+1$.

Exercise 1.3.2

Calculate each determinant.

(a)
$$\begin{vmatrix} 1 & 2 & 3 \\ 0 & 4 & 5 \\ 0 & 0 & 6 \end{vmatrix}$$

(b)
$$\begin{vmatrix} -4 & 2 & 1 \\ 5 & 0 & 3 \\ -2 & 1 & 3 \end{vmatrix}$$

Exercise 1.3.3

Expand by minors along the first *column* of this matrix, and show that you get the same result as when you expand by minors along the first row.

$$\begin{vmatrix} -2 & 1 & 4 \\ 1 & 1 & 2 \\ 2 & 0 & -1 \end{vmatrix}$$

Exercise 1.3.4

Find the values of t for which the determinant of the following matrix is zero.

$$\begin{vmatrix} -2 & t^2 & 4 \\ 3 & 1 & 0 \\ 2 & 0 & -1 \end{vmatrix}$$

Exercise 1.3.5

The **transpose** A^T of a matrix A is obtained by swapping rows and columns. In other words,

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}^T = \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix}.$$

Show that $\det A = \det A^T$.

Exercise 1.3.6

Show that if two rows of a 3×3 matrix A are the same, then $\det A = 0$.

2.1 Introduction to vectors

A **vector** in \mathbb{R}^n is an arrow from one point in \mathbb{R}^n (the *tail*) to another (the *head*). See Figure 2.1). The **length*** of a vector is the distance from the head to the tail. Two vectors are considered equivalent if they have the same length and the same direction.

The **components** of a vector are the coordinates of its head when it is translated so that its tail is at the origin. In other words, to find the components of a vector, we subtract each coordinate of its tail from the corresponding coordinate of the head. The components of the vector in Figure 2.1 are $\langle \frac{3}{2}, 1 \rangle$ —note that we use the pointy brackets to distinguish components of a vector from coordinates of a point.* We can calculate the components by subtracting the coordinates of the tail from the coordinates of the head. Two vectors are equivalent if and only if their components are the same.

The main things we will do with vectors are (i) add two of them together and (ii) multiply a vector by a real number (which is called a **scalar** in this context). These are defined componentwise:

$$\langle u_1, u_2 \rangle + \langle v_1, v_2 \rangle = \langle u_1 + v_1, u_2 + v_2 \rangle, \text{ and}$$

$$c\langle u_1, u_2 \rangle = \langle cu_1, cu_2 \rangle.$$

These definitions of vector addition and scalar multiplication lead to natural geometric interpretations, as shown in Figures 2.2 and 2.3.

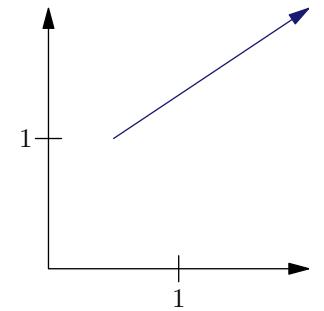


Figure 2.1 A vector in \mathbb{R}^2

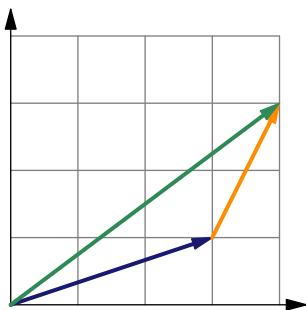


Figure 2.2 Vector addition: $\langle 3, 1 \rangle + \langle 1, 2 \rangle = \langle 4, 3 \rangle$

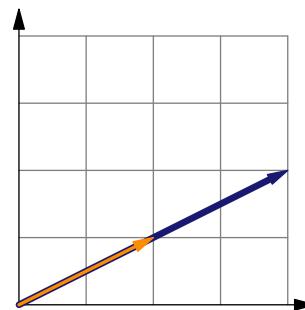


Figure 2.3 Scalar multiplication: $2\langle 2, 1 \rangle = \langle 4, 2 \rangle$

We typically assign names for vectors which are lowercase boldface letters, like \mathbf{u} or \mathbf{v} . Looking at Figure 2.3, we make the following observation.

Observation 2.1.1: Parallel vectors

Two vectors \mathbf{u} and \mathbf{v} are parallel if $\mathbf{u} = c\mathbf{v}$ for some scalar c .



* That is, we use the same notation as we use for multiplication/addition of real numbers because these operations satisfy many of the same properties

Vector operations satisfy several properties suggested by the notation,* such as commutativity ($\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$), associativity ($(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$), the distributive property of scalar multiplication across vector addition (Exercise 2.1.1) and so on. One simple strategy for proving such property-verification exercises is to write out what each side of the equation means in terms of components and then simplify both sides until it's clear that they are equal.

Exercise 2.1.1

Show that scalar multiplication distributes across vector addition. In other words, show that for all $c \in \mathbb{R}$ and vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n , we have

$$c(\mathbf{u} + \mathbf{v}) = c\mathbf{u} + c\mathbf{v}.$$

Exercise 2.1.2

Show that scalar multiplication distributes across scalar addition. In other words, show that for all $c \in \mathbb{R}$, $d \in \mathbb{R}$, and vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n , we have

$$(c + d)\mathbf{u} = c\mathbf{u} + d\mathbf{u}.$$

Exercise 2.1.3

Choose two vectors \mathbf{u} and \mathbf{v} with small integer coordinates and draw a figure to show how \mathbf{u} , \mathbf{v} , and $\mathbf{u} - \mathbf{v}$ fit together to form a triangle.

Exercise 2.1.4

Suppose that \mathbf{u} and \mathbf{v} are vectors in \mathbb{R}^2 or \mathbb{R}^3 , and suppose that $\mathbf{w} = c\mathbf{u} + d\mathbf{v}$, where c and d are both in $[0, 1]$. If the tail of \mathbf{w} is at the origin, then what is the set of possible locations for the head of \mathbf{w} ?

The following example shows how vector ideas can be applied to geometry problems.

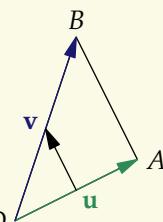
Example 2.1.1

Use vectors to prove that the line segment joining the midpoints of two sides of a triangle is parallel to the third side and half its length.

Solution

Define \mathbf{u} and \mathbf{v} to be two vectors with a common tail at one vertex O of the triangle and heads at the other two vertices A and B as shown. Then the vectors from O to the midpoints of OA and OB are $\frac{1}{2}\mathbf{u}$ and $\frac{1}{2}\mathbf{v}$, since the midpoint of a line segment is defined to be the point which is halfway between the endpoints.

Therefore, the vector \mathbf{w} from one midpoint to another is $\frac{1}{2}\mathbf{v} - \frac{1}{2}\mathbf{u}$. By the distributive property, this is equal to $\frac{1}{2}(\mathbf{v} - \mathbf{u})$. The vector from A to B is $\mathbf{v} - \mathbf{u}$. Therefore, \mathbf{w} has the same direction as the vector from A to B (by Observation 2.1.1) and is half as long.



Exercise 2.1.5

Use vectors to show that the diagonals of a parallelogram bisect one another.

Exercise 2.1.6

A median of a triangle is a line segment from a vertex of the triangle to the midpoint of the opposite side. Use vectors to show that for any triangle, there is a point on all three medians. (Hint: This point will split each median into two segments, one of which is twice as long as the other.)

2.2 The dot product

The fundamental vector operations of scalar multiplication and vector addition are not sufficient to capture information about a really important geometric concept: *angle*. So we introduce a new vector operation.

Definition 2.2.1: Dot product

The **dot product** of two three-dimensional vectors $\mathbf{u} = \langle u_1, u_2, u_3 \rangle$ and $\mathbf{v} = \langle v_1, v_2, v_3 \rangle$ is defined by

$$\mathbf{u} \cdot \mathbf{v} = \langle u_1, u_2, u_3 \rangle \cdot \langle v_1, v_2, v_3 \rangle = u_1 v_1 + u_2 v_2 + u_3 v_3.$$

The dot product distributes across vector addition, and it is closely related to length, as shown in the following exercise. We denote by $|\mathbf{u}|$ the length of \mathbf{u} .

Exercise 2.2.1

Verify that $\mathbf{u} \cdot (\mathbf{v} + \mathbf{w}) = \mathbf{u} \cdot \mathbf{v} + \mathbf{u} \cdot \mathbf{w}$ and that $\mathbf{u} \cdot \mathbf{u} = |\mathbf{u}|^2$.

Now we establish the relationship between the dot product and angle.

Example 2.2.1

Use the law of cosines to show that $\mathbf{u} \cdot \mathbf{v} = |\mathbf{u}||\mathbf{v}| \cos \theta$, where θ is the angle between \mathbf{u} and \mathbf{v} .

Solution

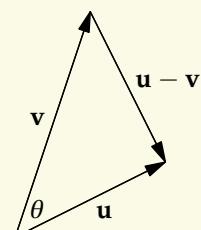
We apply the law of cosines to the triangle with sides \mathbf{u} , \mathbf{v} , and $\mathbf{u} - \mathbf{v}$. We get

$$|\mathbf{u} - \mathbf{v}|^2 = |\mathbf{u}|^2 + |\mathbf{v}|^2 - 2|\mathbf{u}||\mathbf{v}| \cos \theta$$

The left-hand side works out to

$$|\mathbf{u} - \mathbf{v}|^2 = (\mathbf{u} - \mathbf{v}) \cdot (\mathbf{u} - \mathbf{v}) = |\mathbf{u}|^2 + |\mathbf{v}|^2 - 2 \mathbf{u} \cdot \mathbf{v}.$$

Subtracting these equations yields $\mathbf{u} \cdot \mathbf{v} = |\mathbf{u}||\mathbf{v}| \cos \theta$.



Particularly noteworthy is the case where θ is a right angle: We say that two vectors are **perpendicular** or **orthogonal** or **normal** if they meet at a right angle.

Observation 2.2.1: Perpendicular vectors

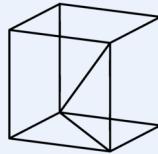
Two vectors \mathbf{u} and \mathbf{v} are perpendicular if and only if $\mathbf{u} \cdot \mathbf{v} = 0$.

!!!

The following example shows how we can use the relation $\mathbf{u} \cdot \mathbf{v} = |\mathbf{u}||\mathbf{v}| \cos \theta$ to find an angle when information about coordinates is available.

Example 2.2.2

Find the angle between the diagonal of a cube and a diagonal of one of its faces.



Solution

The vector from the origin to the opposite corner of the cube is $\langle 1, 1, 1 \rangle$. The vector from the origin to the opposite corner of the bottom face is $\langle 1, 1, 0 \rangle$. Therefore, the angle is given by

$$\theta = \cos^{-1} \left(\frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}||\mathbf{v}|} \right) = \cos^{-1} \left(\frac{1+1+0}{\sqrt{1^2+1^2+1^2} \sqrt{1^2+1^2+0^2}} \right) = \boxed{\cos^{-1} \left(\frac{2}{\sqrt{6}} \right)}.$$

Exercise 2.2.2

Sketch the vectors $\mathbf{u} = \langle 4, 2 \rangle$ and $\mathbf{v} = \langle -1, 2 \rangle$ and show geometrically that they are perpendicular. Then verify that the coordinate formula for dot product indeed gives $\mathbf{u} \cdot \mathbf{v} = 0$ for these two vectors.

We conclude this section by addressing the out-of-nowhere step where we defined the formula for the dot product. We could have begun with the geometric formula $\mathbf{u} \cdot \mathbf{v} = |\mathbf{u}||\mathbf{v}| \cos \theta$ as the *definition* of the dot product, derived the distributive property of the dot product across vector addition using geometry, and then obtained the formula for the dot product in the following way: define $\mathbf{i} = \langle 1, 0, 0 \rangle$, $\mathbf{j} = \langle 0, 1, 0 \rangle$, and $\mathbf{k} = \langle 0, 0, 1 \rangle$. Then a vector $\mathbf{u} = \langle u_1, u_2, u_3 \rangle$ can be written as

$$\mathbf{u} = u_1 \mathbf{i} + u_2 \mathbf{j} + u_3 \mathbf{k},$$

and similarly for \mathbf{v} . Then

$$\begin{aligned} \mathbf{u} \cdot \mathbf{v} &= (u_1 \mathbf{i} + u_2 \mathbf{j} + u_3 \mathbf{k})(v_1 \mathbf{i} + v_2 \mathbf{j} + v_3 \mathbf{k}) \\ &= u_1 v_1 \mathbf{i} \cdot \mathbf{i} + u_1 v_2 \mathbf{i} \cdot \mathbf{j} + u_1 v_3 \mathbf{i} \cdot \mathbf{k} + \\ &\quad u_2 v_1 \mathbf{j} \cdot \mathbf{i} + u_2 v_2 \mathbf{j} \cdot \mathbf{j} + u_2 v_3 \mathbf{j} \cdot \mathbf{k} + \\ &\quad u_3 v_1 \mathbf{k} \cdot \mathbf{i} + u_3 v_2 \mathbf{k} \cdot \mathbf{j} + u_3 v_3 \mathbf{k} \cdot \mathbf{k}, \end{aligned}$$

by the distributive property. This looks like a mess, but since \mathbf{i} , \mathbf{j} , and \mathbf{k} are perpendicular, six of these nine

terms are zero. Furthermore, since $\mathbf{i} \cdot \mathbf{i} = 1$ and similarly for \mathbf{j} and \mathbf{k} , we end up with $\mathbf{u} \cdot \mathbf{v} = u_1v_1 + u_2v_2 + u_3v_3$, as desired.

Exercise 2.2.3

Show that the vectors $\mathbf{u} = \langle \cos \theta, \sin \theta \rangle$ and $\mathbf{v} = \langle -\sin \theta, \cos \theta \rangle$ are unit vectors (meaning that the length of each is 1). Show that \mathbf{u} and \mathbf{v} are orthogonal.

Exercise 2.2.4: Orthogonal Projection

Suppose vectors \mathbf{u} and \mathbf{v} are given with $|\mathbf{u}| = 1$. If $\mathbf{v} = \mathbf{v}_1 + \mathbf{v}_2$, where \mathbf{v}_1 and \mathbf{v}_2 are perpendicular and \mathbf{v}_1 is parallel to \mathbf{u} , then \mathbf{v}_1 is called the *orthogonal projection* of \mathbf{v} onto \mathbf{u} . The magnitude of \mathbf{v}_1 is called the *component* of \mathbf{v} in the \mathbf{u} direction.

- Draw a figure illustrating the relationship between \mathbf{u} , \mathbf{v} , \mathbf{v}_1 , and \mathbf{v}_2 .
- Use right-triangle trigonometry to find a formula for the length of \mathbf{v}_1 in terms of the angle θ between \mathbf{u} and \mathbf{v} , and then use dot products to find a formula for \mathbf{v}_1 in terms of \mathbf{u} and \mathbf{v} .
- Use your findings to support the statement *dotting with a unit vector \mathbf{u} gives the component in the \mathbf{u} direction*.

2.3 The cross product

In the last section we introduced a vector product which reveals information about *angle*; in this section we'll see a new vector product which gives us information about *area*.

* We put *determinant* in scare quotes because the matrix entries are not numbers

A beautiful explanation of this formula

The **cross product** of $\mathbf{u} = (u_1, u_2, u_3)$ and $\mathbf{v} = (v_1, v_2, v_3)$ is defined by expanding the following 'determinant' by minors along the first row:^{*}

$$\mathbf{u} \times \mathbf{v} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix} = (u_2v_3 - u_3v_2)\mathbf{i} - (u_1v_3 - u_3v_1)\mathbf{j} + (u_1v_2 - u_2v_1)\mathbf{k}.$$

Note that the dot product of two vectors is a scalar, while the cross product of two vectors is another vector. It turns out that this vector is orthogonal to *both* of the first two.

Example 2.3.1

Confirm that $\mathbf{u} \times \mathbf{v}$ is orthogonal to \mathbf{u} and to \mathbf{v} .

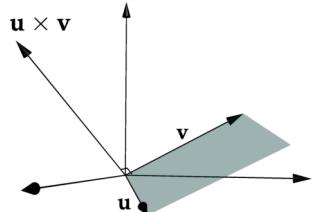


Figure 2.4 The cross product of \mathbf{u} and \mathbf{v} , whose length is equal to the area of the parallelogram shown

Solution

We compute $\mathbf{u} \cdot (\mathbf{u} \times \mathbf{v})$ as follows:

$$\begin{aligned} \langle u_1, u_2, u_3 \rangle \cdot \langle u_2v_3 - u_3v_2, -(u_1v_3 - u_3v_1), u_1v_2 - u_2v_1 \rangle = \\ (u_2v_3 - u_3v_2)u_1 - (u_1v_3 - u_3v_1)u_2 + (u_1v_2 - u_2v_1)u_3 = 0. \end{aligned}$$

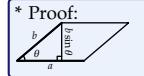
This implies that $\mathbf{u} \times \mathbf{v}$ is orthogonal to \mathbf{u} . Swapping out \mathbf{u} for \mathbf{v} , we see that $\mathbf{u} \times \mathbf{v}$ is orthogonal to \mathbf{v} too.

Alternatively, we could note that

$$\mathbf{u} \cdot (\mathbf{u} \times \mathbf{v}) = \begin{vmatrix} u_1 & u_2 & u_3 \\ u_1 & u_2 & u_3 \\ v_1 & v_2 & v_3 \end{vmatrix}$$

and conclude using Exercise 1.3.6 that this determinant equals zero.

The following exercise provides the advertised connection to area. Recall from geometry that the area of a parallelogram with sides of length a and b meeting at an angle θ is equal to $ab \sin \theta$.*



Exercise 2.3.1

Verify that $|\mathbf{u} \times \mathbf{v}|^2 = |\mathbf{u}|^2 |\mathbf{v}|^2 - (\mathbf{u} \cdot \mathbf{v})^2$. Use this fact to show that

$$|\mathbf{u} \times \mathbf{v}| = |\mathbf{u}| |\mathbf{v}| \sin \theta,$$

where θ is the angle between \mathbf{u} and \mathbf{v} .

So to sum up: $\mathbf{u} \times \mathbf{v}$ is a vector which is orthogonal to both \mathbf{u} and \mathbf{v} and whose length is equal to the area of the parallelogram spanned by \mathbf{u} and \mathbf{v} . Note that there are only two vectors satisfying both of these conditions. To determine which one is $\mathbf{u} \times \mathbf{v}$, we use the *right-hand rule*: imagine orienting your right hand so that you can curl your fingers from \mathbf{u} towards \mathbf{v} . The direction of your thumb (if it's orthogonal to your fingers) is the direction of $\mathbf{u} \times \mathbf{v}$.

Exercise 2.3.2

Find the volume of the parallelepiped spanned by $\langle 3, 4, 1 \rangle$, $\langle -2, 4, 0 \rangle$, and $\langle -5, 5, 2 \rangle$. (Hint: first find the area of the base, then figure out how to use dot products to find the height.)

Exercise 2.3.3: Cross product distributive property

Show that $\mathbf{u} \times (\mathbf{v} + \mathbf{w}) = \mathbf{u} \times \mathbf{v} + \mathbf{u} \times \mathbf{w}$, if \mathbf{u}, \mathbf{v} , and \mathbf{w} are vectors.

3 Three-dimensional Geometry

3.1 Lines and planes

There are various ways to describe a line in 2D space using an equation, including point-slope form and y -intercept form. In this section we will learn the 3D analogue: equation descriptions of lines and planes in space. We begin with an example.

Example 3.1.1

Describe the line L in \mathbb{R}^3 passing through the points $A = (3, -4, 1)$ and $B = (2, -1, 4)$.

Solution

We can tell whether a given point (x, y, z) in \mathbb{R}^3 is on the line L using vectors: (x, y, z) is on L if and only if the vector from $(3, -4, 1)$ to (x, y, z) is a scalar multiple of the vector from $(3, -4, 1)$ to $(2, -1, 4)$ (see Figure 3.1). We can turn this into an equation: a point (x, y, z) is on L if and only if there exists $t \in \mathbb{R}$ such that

$$t\langle 2 - 3, -1 - (-4), 4 - 1 \rangle = \langle x - 3, y - (-4), z - 1 \rangle.$$

Setting components equal, we find that (x, y, z) is on L if and only if there exists t so that

$$\begin{aligned} x &= 3 - t, \\ y &= -4 + 3t, \text{ and} \\ z &= 1 + 3t. \end{aligned} \tag{3.1.1}$$

Note that the solution above involves a new variable t ; this is called a *parameter*, and the form we gave as an answer is called **parametric form**. You can imagine drawing the line by starting with $t = 0$, so that your pen begins at A , and then sweeping t through the values from 0 to 1, changing the location of your pen according to the parametric equations (3.1.1). Your pen will sweep out the line segment from A and B . Then you can let t vary beyond 1 to get the rest of the ray past B , and you can let t vary over the negative numbers to get the part of the line on the other side of A .

If we didn't want to involve t , note that we could solve for t in one equation and substitute into the other two, thereby obtaining *two* equations involving x , y , and z . This makes sense: starting from the plane, imposing one equation on x and y cuts the dimension down by one and gives a line. However, starting from 3D space, we need to reduce the dimension by *two*. So we need two equations.

The procedure developed in Example 3.1.1 works in general: the line through $A = (a, b, c)$ and B has para-

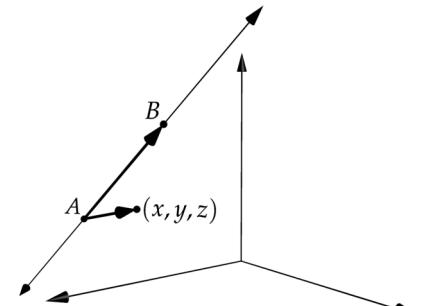


Figure 3.1 Checking whether (x, y, z) is on the line through A and B

metric form*

$$\begin{aligned}x &= a + v_1 t, \\y &= b + v_2 t \\z &= c + v_3 t,\end{aligned}$$

* Note that this representation is not unique, since we could've used B (or any other point on the line) in place of A

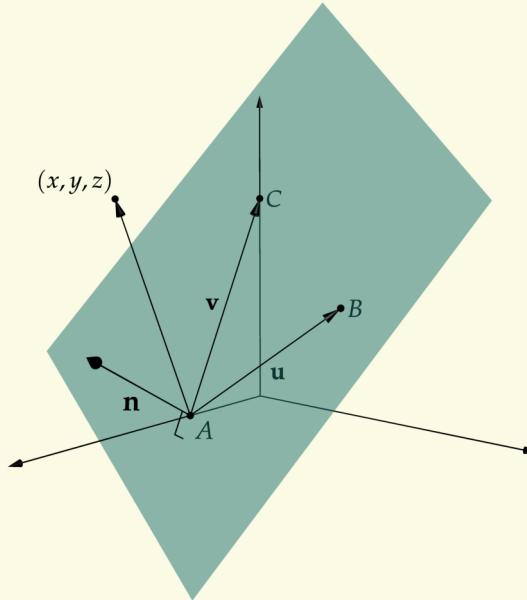
where $\langle v_1, v_2, v_3 \rangle = \overrightarrow{AB}$ is the vector from A to B .

Example 3.1.2

Describe the plane P passing through the points $A = (1, 0, 0)$, $B = (0, 1, 1)$, and $C = (0, 0, 2)$.

Solution

We can tell whether (x, y, z) is on P using vectors. Define \mathbf{u} and \mathbf{v} to be the vectors from A to B and from A to C , respectively.



If we can find a vector \mathbf{n} which is orthogonal to P , then we can say (x, y, z) is on P if and only if the vector from A to (x, y, z) is orthogonal to \mathbf{n} . But we can take $\mathbf{n} = \mathbf{u} \times \mathbf{v}$, since the cross product of two vectors is orthogonal to both of them. So

$$\mathbf{n} = \langle -1, 1, 1 \rangle \times \langle -1, 0, 2 \rangle = \langle 2, 1, 1 \rangle.$$

Now we can say that (x, y, z) is on P if and only if*

$$\mathbf{n} \cdot \langle x - 1, y - 0, z - 0 \rangle = 0,$$

which simplifies to $2x + y + z = 2$.

* Note that we could use B or C instead of $A = (1, 0, 0)$ here and we'd get the same equation for the plane

!!!

Observation 3.1.1: Vector normal to a plane

A vector \mathbf{n} normal to the plane $ax + by + cz = d$ can be read off from the coefficients:

$$\mathbf{n} = \langle a, b, c \rangle.$$

* We can ask about the distance from a point to a point, a point to a line, a point to a plane, a line to a line, a line to a plane, or a plane to a plane

One important 3D geometry problem is to find distances between points, lines, and planes.* We define the distance between two sets to be the *minimum* distance between any pair of points from the respective sets.

Example 3.1.3

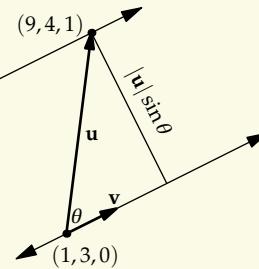
Consider the line ℓ given by the parametric equation $(x, y, z) = (1 - 2t, 3, t)$. Find the distance from ℓ to the line m which is parallel to ℓ and which passes through the point $(9, 4, 1)$.

Solution

The parametric equations give us convenient access to a point on each of the two lines as well as a vector \mathbf{v} which is parallel to both lines. So we make a figure with this information.

If we define \mathbf{u} to be the vector from connecting the two given points, we can see by applying right-triangle trigonometry to the figure that the desired distance d is equal to $|\mathbf{u}| \sin \theta$. Therefore,

$$d = \frac{|\mathbf{u}| |\mathbf{v}| \sin \theta}{|\mathbf{v}|} = \frac{|\mathbf{u} \times \mathbf{v}|}{|\mathbf{v}|} = \frac{\sqrt{105}}{\sqrt{5}} = \boxed{\sqrt{21}}.$$



Note the basic strategy: (i) draw a figure containing the information that the problem gives us (a schematic diagram suffices; there is no need to make it particularly precise), (ii) use right triangle trigonometry to express the desired distance terms of vectors we have, and (iii) use vector formulas to calculate the desired quantity using a dot or cross product.

A second important operation is finding points of intersection.

Example 3.1.4

Find the point where the line $\langle 3 + t, -2t, 3 \rangle$ intersects the plane $x + y + z = 7$.

Solution

A point is on the intersection of two graphs if and only if it's on both graphs. Therefore, we're looking for a point (x, y, z) that satisfies $x + y + z = 7$ and has the property that there exists $t \in \mathbb{R}$ such that $x = 3 + t, y = -2t$, and $z = 3$. Thus we have four equations and four unknowns (x, y, z, t) , and we can substitute the last three equations into the first to find that $6 - t = 7$, which implies $t = -1$. Therefore, the point of intersection is $(2, 2, 3)$.

Exercise 3.1.1

Find the equation of the plane passing through the points $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$. Find the distance from that plane to the origin.

Exercise 3.1.2

Find the distance between the planes $x + y - 2z = 3$ and $x + y - 2z = 0$.

Exercise 3.1.3

Find the distance between the lines $(x, y, z) = (2t, 1 - t, 4)$ and $(x, y, z) = (1 + t, -2t, -1 - t)$. Hint: these lines are *skew*, meaning that they are not parallel but do not intersect. Begin by using a cross product to find a vector which is perpendicular to both lines.

Exercise 3.1.4

Describe parametrically the intersection of the planes $2x + z = 3$ and $x + y - 2z = 4$.

3.2 Vector-valued functions

3.2.1 PATHS IN SPACE

Consider a particle moving along the number line in such a way that its position at time t is given by $r(t)$. Then the velocity of the particle at time t is given by the first derivative $v(t) = r'(t)$. The velocity specifies the *speed* of the particle as well as its *direction* (left if negative, right if positive).

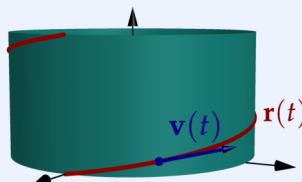
* The derivative of a path r is the vector obtained by taking the derivative of each component

The same is true of a particle moving in 2D or 3D space: its location is specified by a function customarily denoted $\mathbf{r}(t)$ from \mathbb{R} to either \mathbb{R}^2 or \mathbb{R}^3 , and its derivative* $\mathbf{v}(t) = \mathbf{r}'(t)$ at time t tells us the speed of the particle at that time (via its length) as well as the direction. We call \mathbf{r} a *path**

Example 3.2.1

Consider a bug which crawls counterclockwise around a cylinder of radius 1 meter from $(1, 0, 0)$ to $(1, 0, 1)$ as shown:

* For a function from \mathbb{R}^1 to \mathbb{R}^2 or \mathbb{R}^3 to count as a path, we require that each of its components be continuous. For example, $\mathbf{r}(t) = (e^t, t^2 \sin t)$ is a path



Assuming the bug moves at constant speed and makes the whole journey in one second, find a for-

mula for the position and velocity of the bug at time t .

Solution

We can see that the z -coordinate of the bug's position increases at a constant rate from 0 to 1 as t goes from 0 to 1, so the z -coordinate of $\mathbf{r}(t)$ is t .

For the x and y coordinates, we need a pair of functions $(x(t), y(t))$ that traces out the unit circle in one second. Recall that cosine and sine are defined to be the functions that trace out the unit circle according to angle, so we can scale them so they make it around in 1 second instead of 2π seconds:

$$(x(t), y(t)) = (\cos 2\pi t, \sin 2\pi t).$$

So all together we have

$$\mathbf{r}(t) = \langle \cos 2\pi t, \sin 2\pi t, t \rangle,$$

which means that

$$\mathbf{v}(t) = \langle -2\pi \sin 2\pi t, 2\pi \cos 2\pi t, 1 \rangle.$$

Exercise 3.2.1

Find the times $t \geq 0$ when a particle whose location at time t is $\mathbf{r}(t) = \langle t^3 - t^2, t \rangle$ is at its slowest.

Exercise 3.2.2

The acceleration $\mathbf{a}(t)$ of a particle whose position at time t is given by $\mathbf{r}(t)$ is defined to be $\mathbf{v}'(t) = \mathbf{r}''(t)$. Suppose that the acceleration $\mathbf{a}(t)$ of a particle is given by $\mathbf{a}(t) = \langle 2t, t^2 \rangle$. If $\mathbf{r}(0) = \mathbf{0}$ and $\mathbf{v}(0) = \langle 1, 2 \rangle$, then find a formula for $\mathbf{r}(t)$.

Exercise 3.2.3

An astronaut is using a rope to move in space in such a way that his position at time t is given by $\mathbf{r}(t) = (2+t)\mathbf{i} + (2+\ln t)\mathbf{j} + \left(7 - \frac{4}{t^2+1}\right)\mathbf{k}$. The coordinates of the space station doorway are $(5, 4, 9)$. When should the astronaut let go of the rope so as to drift into the doorway?

Given a curve in \mathbb{R}^n , a path \mathbf{r} that traces it out is called a parametric equation for the curve.

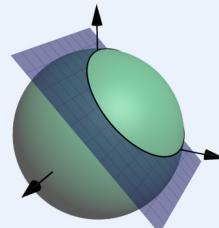
Example 3.2.2

Find a parametric equation describing the intersection of the sphere

$$x^2 + y^2 + z^2 = 1$$

and the plane

$$y + z = 1.$$

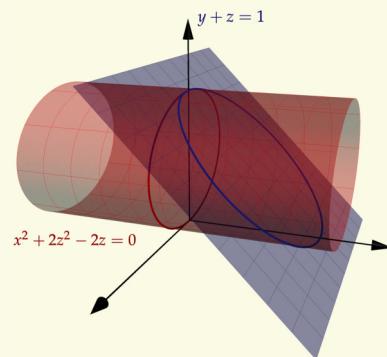


Solution

The main idea is to reduce the 3D problem to a 2D problem by finding a shadow of the curve that we can parameterize. Looking at the figure, it appears that the shadow of this curve in the xz plane is an ellipse. If we substitute $y = 1 - z$ into $x^2 + y^2 + z^2 = 1$, we see that every point on the curve indeed satisfies* $x^2 + 2z^2 - 2z = 0$.

Completing the square to get this equation in standard form, we get

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



* This equation describes an elliptical cylinder, because it is the set of all points whose shadow in the xz -plane is in the red ellipse shown in the figure

The intersection of this surface with the xz -plane is an ellipse which can be parameterized as

$$\left(\frac{1}{\sqrt{2}} \cos t, 0, \frac{1}{2} + \frac{1}{2} \sin t\right),$$

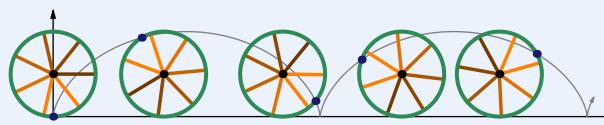
where t ranges over $[0, 2\pi]$. Using the equation $y = 1 - z$, we see that the point on the curve whose shadow is $\left(\frac{1}{\sqrt{2}} \cos t, 0, \frac{1}{2} + \frac{1}{2} \sin t\right)$ is $\left(\frac{1}{\sqrt{2}} \cos t, \frac{1}{2} - \frac{1}{2} \sin t, \frac{1}{2} + \frac{1}{2} \sin t\right)$. As t ranges from 0 to 2π , this point goes all the way around the cylinder, so this formula indeed parameterizes the curve.

Exercise 3.2.4

The intersection of the cylinder of unit radius centered along the x -axis and the cylinder of unit radius centered along the y -axis consists of four curves connecting the points $(0, 0, 1)$ and $(0, 0, -1)$. Choose one of them and parametrize it.

Example 3.2.3

Consider a wheel of unit radius centered at $(0, 1)$ in a coordinate plane. At time $t = 0$, the wheel begins rolling (without slipping) along the x -axis to the right at a rate of one unit per second. Find the location $\mathbf{r}(t)$ of the point on the wheel which was originally located at the origin.



Solution

Because the wheel is moving at a unit rate of speed, the location of its center at time t is $(t, 1)$. After t seconds, the wheel has rotated t units along its perimeter, which corresponds to rotating t radians. Therefore, the angle of the point on the wheel originally at the origin starts at $\frac{3\pi}{2}$ and advances in the clockwise direction at a rate of one radian per second. Thus the vector from the wheel's center to the desired point is $\langle \cos(\frac{3\pi}{2} - t), \sin(\frac{3\pi}{2} - t) \rangle$, which simplifies to $\langle -\sin t, -\cos t \rangle$. So the location of the desired point is the sum of the point $(t, 1)$ and the vector $\langle -\sin t, -\cos t \rangle$:

$$\mathbf{r}(t) = (t - \sin t, 1 - \cos t).$$

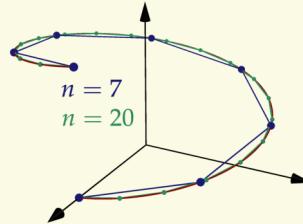
Example 3.2.4

Find the speed of the bug in Example 3.2.1.

Solution

Since the speed of the bug is constant and the whole journey takes a second, the speed in meters per second is equal to the length of the bug's path divided by 1 second.

Consider any time $t \in [0, 1]$ and a very small time Δt . Over the interval $[t, t + \Delta t]$, the bug's speed is approximately $|\mathbf{v}(t)|$, so the distance it moves is approximately equal to $|\mathbf{v}(t)|\Delta t$. Fixing a large integer n and splitting up the interval $[0, 1]$ into n intervals of width $\Delta t = 1/n$, total distance is approximately



$$\overbrace{|\mathbf{v}(0)|\Delta t}^{\text{over } [0, \Delta t]} + \overbrace{|\mathbf{v}(\Delta t)|\Delta t}^{\text{over } [\Delta t, 2\Delta t]} + \overbrace{|\mathbf{v}(2\Delta t)|\Delta t}^{\text{over } [2\Delta t, 3\Delta t]} + \cdots + \overbrace{|\mathbf{v}(1 - \Delta t)|\Delta t}^{\text{over } [1 - \Delta t, 1]}.$$

This expression is a Riemann sum approximating the integral $\int_0^1 |\mathbf{v}(t)| dt$, so as $\Delta t \rightarrow 0$, it converges to

$$\int_0^1 |\mathbf{v}(t)| dt = \int_0^1 \sqrt{4\pi^2 \sin^2 2\pi t + 4\pi^2 \cos^2 2\pi t + 1} dt = \boxed{\sqrt{4\pi^2 + 1}}.$$

Although we arrived at this answer using approximations that we weren't so careful about, this is indeed the same result we would've gotten if we'd cut a slit in the cylinder and unrolled it to find that the bug's path is equal in length to the hypotenuse of a rectangle with side lengths 1 (the height of the cylinder) and 2π (the circumference of the cylinder).

The method developed in Example 3.2.4 leads to the following definition. A function from an interval to \mathbb{R} is **piecewise differentiable** if its domain can be subdivided into finitely many intervals such that the function is differentiable on each of them. A function from an interval to \mathbb{R}^n is defined to be piecewise differentiable if all of its components are.

Definition 3.2.1: Arclength

The **arclength** of a piecewise differentiable path $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ is defined to be

$$\int_a^b |\mathbf{r}'(t)| dt,$$

assuming this integral exists.

One important property of this definition is that it should not depend on the rate at which the curve is traced out.* In other words, suppose person A draws a curve with a consistent pace, while person B rushes through the first part of the curve and slows down towards the end. As long as the curves they draw are the same, the lengths should be the same. Let's try an example.*

* A property of a path which depends only on the curve it traces out is called **parameterization independent**

* See Appendix A.3.2 for an explanation of why this is true in general

Example 3.2.5: Parameterization independence

Use the arclength formula to find the length of the portion of the unit circle in the first quadrant in two ways: using $\mathbf{r}_1(t) = (\cos t, \sin t)$, and using $\mathbf{r}_2(t) = (t, \sqrt{1-t^2})$.

Solution

For \mathbf{r}_1 , we have

$$\int_0^{\pi/2} \sqrt{(-\sin t)^2 + (\cos t)^2} dt = \left[\frac{\pi}{2} \right],$$

and for \mathbf{r}_2 , we have

$$\int_0^1 \sqrt{1^2 + \left(\frac{-2t}{2\sqrt{1-t^2}} \right)^2} dt = \int_0^1 \frac{1}{\sqrt{1-t^2}} dt = \arcsin(1) - \arcsin(0) = \left[\frac{\pi}{2} \right],$$

where we performed the last integral by recalling that the derivative of the inverse sine function is $\frac{1}{\sqrt{1-t^2}}$.

Exercise 3.2.5

Sketch the path $\mathbf{r}(t) = \langle e^{-t} \cos t, e^{-t} \sin t \rangle$. Show that the length of this path over the interval $[0, a]$ converges as $a \rightarrow \infty$. What does this say about the length of the path over the interval $[0, \infty)$?

3.2.3 CURVATURE*

The curvature of a path at a point on the path is a measure of how curvy the path is at that point. For example, the path shaped like the letter “U” has a small curvature (zero, in fact) near its endpoints and a larger curvature at points along the bend.

One natural way to distinguish points on a path where the path has large curvature from points where the path’s curvature is small is to ascertain *how rapidly the direction of the velocity vector changes, per unit of arclength*. This leads to the following definition.

Definition 3.2.2: Curvature

- (i) The **unit tangent vector** of a differentiable path $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^n$ is defined to be

$$\mathbf{T}(t) = \frac{\mathbf{r}'(t)}{|\mathbf{r}'(t)|}.$$

- (ii) The **arclength function** $s(t)$ is defined to be the length of \mathbf{r} over the interval $[a, t]$.

- (iii) The **curvature** $\kappa(t)$ of \mathbf{r} at time t is defined by

$$\kappa(t) = \frac{|\mathbf{dT}/dt|}{ds/dt}.$$

Example 3.2.6

Find the curvature of a circle of radius a using Definition 3.2.2.

Solution

We parameterize the circle as $(a \cos t, a \sin t)$ where t ranges over $[0, 2\pi]$, and we determine that

$$\mathbf{T} = \frac{\langle -a \sin t, a \cos t \rangle}{|\langle -a \sin t, a \cos t \rangle|} = \langle -\sin t, \cos t \rangle.$$

The arclength function is* $s(t) = \int_0^t |\langle -a \sin \tau, a \cos \tau \rangle| d\tau = at$.

$$\text{Therefore, the curvature is } \kappa(t) = \frac{|\langle -\cos t, -\sin t \rangle|}{|d(at)/dt|} = \frac{1}{a}.$$

This formula makes sense, because a very large circle looks quite flat (consider standing on the equator as opposed to standing on a latitude line a few feet from the North Pole).

Exercise 3.2.6

Suppose that \mathbf{r} is a differentiable path and that \mathbf{T} is its unit tangent vector. Use the fact that the length of $\mathbf{T}(t)$ is constant to show that $\mathbf{T}(t)$ is always orthogonal to $\mathbf{T}'(t)$.

Exercise 3.2.7

Find the curvature at each point on the graph of the function $f : (0, \pi) \rightarrow \mathbb{R}$ defined by $f(x) = \ln(\sin x)$.

3.3 Quadric surfaces

The *graph* of an equation involving the variables x, y, z is the set of points (x, y, z) in \mathbb{R}^3 which satisfy the equation. For example, we have seen that the graph of $x + y + z = 1$ is a plane in \mathbb{R}^3 . More generally, the graph of any linear equation in \mathbb{R}^3 is a plane. So let's step it up a notch and consider *quadratic* equations.* A graph of a quadratic equation in the variables x, y, z is called a *quadric surface*.

Perhaps the simplest quadratic equation to reason about is $x^2 + y^2 + z^2 = 1$. The left-hand side has a geometric interpretation as the *squared distance from (x, y, z) to the origin*. Therefore, a point (x, y, z) satisfies this equation if and only if its squared distance to the origin is 1. We have a name for the set of such points: the *sphere* of radius 1, centered at the origin.

The situation is not always so simple. So here's a key idea for tackling 3D geometry problems: **slice it up**. Consider planes of the form $z = \text{constant}$, $y = \text{constant}$, or $z = \text{constant}$ and see what your graph looks like in these planes. Here's an archetypal example.

* We'll use τ as the variable of integration, since we're already using t as a limit of integration

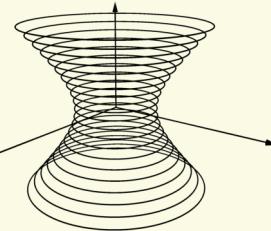
* The 2D analogues are *conic sections*: parabolas, ellipses, and hyperbolas. These are graphs of various quadratic equations in two variables

Example 3.3.1

Figure out what the graph of $x^2 + y^2 - z^2 = 1$ looks like.

Solution

We begin by finding all the points which satisfy this equation and the equation $z = 0$. If (x, y, z) satisfies this equation and $z = 0$, then that means that $x^2 + y^2 = 1$. Furthermore, if $x^2 + y^2 = 1$ and $z = 0$, then (x, y, z) satisfies the equation $x^2 + y^2 - z^2 = 1$. This means that the intersection of the desired graph and the line $z = 0$ is the circle of radius 1 centered at the origin.



Similarly, the intersection of the desired graph and the plane $z = 1$ is a circle which is centered at $(0, 0, 1)$ and has radius $\sqrt{2}$. Drawing in several more of these traces*, we get a picture that looks like the figure above. This is already a pretty clear picture of what the graph looks like: it's rotationally symmetric about the z -axis and “flares out” as you move away from the xy -plane. This graph is called a *one-sheeted hyperboloid*.

* A trace of a figure is an intersection of that figure with a plane

Exercise 3.3.1

Sketch $\frac{z^2}{c} = \frac{x^2}{a^2} + \frac{y^2}{b^2}$, where $a = b = c = 1$. This is called an elliptic paraboloid.

Exercise 3.3.2

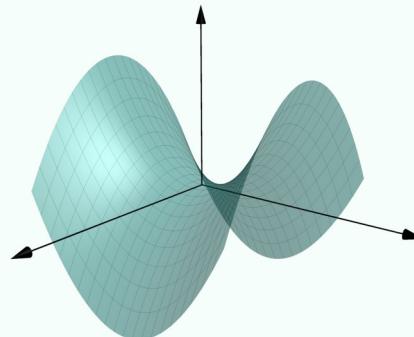
Sketch $\frac{z^2}{c^2} = \frac{x^2}{a^2} + \frac{y^2}{b^2}$, where $a = b = c = 1$. This is called an elliptic cone.

Exercise 3.3.3

Sketch the graph of $x^2 + y^2 - z^2 = -1$. This is called a two-sheeted hyperboloid.

Exercise 3.3.4

Show that the graph of $z = y^2 - x^2$ looks like the figure shown. This is called a hyperbolic paraboloid.



3.4 Polar, cylindrical, and spherical coordinates

A coordinate system is a way of identifying locations using pairs or triples of real numbers. Rectangular coordinates—the ones commonly denoted (x, y) or (x, y, z) —have some nice properties, but some tasks are much more convenient in other coordinate systems.

For example, a captain at sea wishing to communicate the location of a nearby pirate ship would probably describe its location in terms of the distance r between the two ships and an angle θ (which might be given with reference to the ship's orientation or as a cardinal direction). The captain is using *polar coordinates*.

Given a point in the plane, we define r to be its distance from the origin and θ to be the signed angle formed between the positive horizontal axis and the vector from the origin to the point. The word *signed* means that an angle measured clockwise from the positive x -axis counts as negative. The values r and θ are called the radial and angular polar coordinates of the point, respectively.*

Exercise 3.4.1

Show that if the polar coordinates of a point (x, y) are r and θ , then we have

$$x = r \cos \theta, \quad \text{and} \quad y = r \sin \theta.$$

Correspondingly, we can coordinatize three-dimensional space by replacing either one or two spatial coordinates with an angular coordinate. Perhaps the simplest way to do this is leave z the same and replace (x, y) with polar coordinates (r, θ) . In other words, we define* for each point P in \mathbb{R}^3 :

$r(P) = \text{distance from } P \text{ to the } z\text{-axis}$

$\theta(P) = \text{signed angle* of the } P\text{-containing half-plane whose boundary is along the } z\text{-axis}$

$z(P) = \text{signed distance from } P \text{ to the } xy\text{-plane.}$

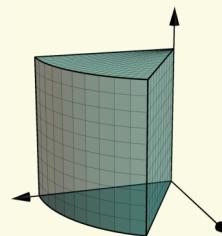
Then we can describe a point by its **cylindrical coordinates** (r, θ, z) rather than its rectangular coordinates. We may also describe a solid in \mathbb{R}^3 by giving inequalities in the variables r, θ , and z such that a point with cylindrical coordinates (r, θ, z) satisfies the inequalities if and only if the point is in the given solid.

Example 3.4.1

Graph* the system of cylindrical coordinate inequalities $r \leq 4$, $0 \leq \theta \leq \pi/3$, $0 \leq z \leq 2$. Find the volume of the resulting region.

Solution

The problem is asking us to find the points whose cylindrical coordinates satisfy all of the given inequalities. Such a point is less than or equal to 4 units from the z -axis, lies between the half-planes $\theta = 0$ and $\theta = \pi/3$, and is above the xy -plane and less than two units away from it. The set of such points is shown to the right.



* So a coordinate is a function from a set of points to \mathbb{R}

* Try coming up with geometric descriptions of the coordinates r, θ , and z in 3D space before looking at the answer below

* ...measured counterclockwise with respect to the positive x -axis

* Familiarity with coordinate slices (Table A.2 in Appendix A.2.2) is helpful for graphing inequalities

This region is one-sixth of a cylinder whose volume is $\pi r^2 h = 32\pi$, so its volume is $\frac{16\pi}{3}$.

Exercise 3.4.2

Graph the system of inequalities $0 \leq r \leq z$, $\pi \leq \theta \leq 2\pi$.

Cylindrical coordinates have two distance coordinates and one angular coordinate. How can we specify a point in space using one distance coordinate and two angular coordinates? The most natural candidate for the distance coordinate is the distance from the origin. In other words, we define $\rho(x, y, z) = \sqrt{x^2 + y^2 + z^2}$. We call this coordinate ρ instead of r to distinguish it from the radial polar coordinate.

As for the angular coordinates, let's use the cylindrical coordinate θ for one of them. For the other, we measure the angle ϕ between the positive z -axis and the vector from the origin to (x, y, z) . This pair of angular coordinates might be familiar: we use them to describe locations on the surface of the earth. In that context, the angle θ is called longitude and the angle ϕ is called latitude.

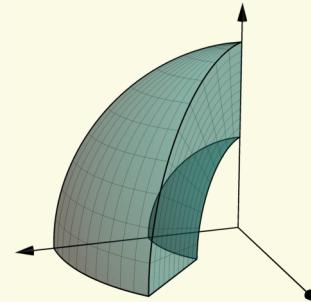
Note that θ varies from 0 to 2π as one loops around the z -axis. However, the angle ϕ varies only from 0 to π as one goes from the north pole to the south pole. Thus the angles θ and ϕ do not play symmetric roles.*

Example 3.4.2

Graph the system of inequalities $\frac{1}{2} \leq \rho \leq 1$, $0 \leq \theta \leq \frac{\pi}{4}$, $0 \leq \phi \leq \frac{\pi}{2}$.

Solution

The set of points with $\rho \leq 1$ is the set of points on or inside of the sphere of radius 1 centered at the origin. Imposing the additional constraint $\rho \geq \frac{1}{2}$ removes the sphere of radius $\frac{1}{2}$ centered at the origin. Then the angular constraints carve out a portion of this hollowed out sphere, as shown.



Exercise 3.4.3

Find a system of inequalities in spherical coordinates to describe the portion of the unit ball* above the plane $z = \frac{1}{2}$.

* This is because ϕ measures the angle required to rotate a vector *freely* so as to align with the positive z -axis, while $-\theta$ measures the signed angle needed to rotate the vector *about the z -axis* to get to the positive half of the xz -plane

* The unit ball is the set of points satisfying $x^2 + y^2 + z^2 \leq 1$

Exercise 3.4.4

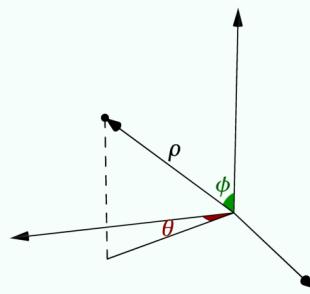
Use the given figure to show that

$$x = \rho \cos \theta \sin \phi$$

$$y = \rho \sin \theta \sin \phi$$

$$z = \rho \cos \phi.$$

Hint: use right-triangle trigonometry to write $\sqrt{x^2 + y^2}$ and z in terms of ρ and ϕ , and then use a different right triangle to write $(x, y, 0)$ in terms of ρ , ϕ , and θ .

**Exercise 3.4.5**

Determine the graph of the spherical-coordinate equation $\rho = 2 \cos \phi$. (Hint: multiply both sides by ρ and then switch to rectangular coordinates.)

Exercise 3.4.6

Determine the graph of $\rho = \sin \phi \sin \theta$.

Exercise 3.4.7

Sketch the set of points satisfying $1 < \rho < 2$ and $\phi < \pi/4$.

4 Multivariable Differentiation

In this chapter, we will be considering functions from \mathbb{R}^n to \mathbb{R}^1 , where $n \geq 2$. The main objectives will be to extend various important notions in single-variable calculus to the higher-dimensional setting.

4.1 Limits

Limits are useful for discussing a function's output values *near*—but not *at*—a point in its domain. For example, if $f : \mathbb{R} \rightarrow \mathbb{R}$ maps 0 to 3 and every nonzero number to 7, then we would say that the limit of $f(x)$ as $x \rightarrow 0$ is equal to 7.

We begin with the notion of a limit at 0 for a *monotone* function from* $(0, 1)$ to \mathbb{R} . We say that a is a **lower bound** for a function f if its range is a subset of $[a, \infty)$. A function is **bounded below** if it has a lower bound. Likewise, b is an **upper bound** of f if the range of f is a subset of $(-\infty, b]$, and a function is **bounded above** if it has an upper bound.*

* The right endpoint of the interval $(0, 1)$ is immaterial here, since we will be interested in the limit at 0

Definition 4.1.1: Limit of a monotone function on $(0, 1)$

Suppose that $f : (0, 1) \rightarrow \mathbb{R}$ is bounded below and has the property that $f(r)$ decreases as r decreases—in other words, we have $f(r) \leq f(s)$ for all $0 < r \leq s < 1$. Then we define L to be the greatest real number such that $L \leq f(r)$ for all $r \in (0, 1)$.

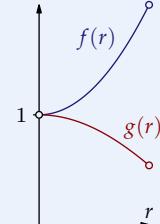
Similarly, if f is bounded above and has the property that $f(r)$ increases as r decreases, then we define L to be the least real number such that $L \geq f(r)$ for all $r \in (0, 1)$.

In either case, we say that the limit as $r \rightarrow 0$ of $f(r)$ exists and is equal to L , that is, $\lim_{r \rightarrow 0} f(r) = L$.

* For example, $f(x) = x$ is neither bounded above nor below, while $g(x) = e^x$ is bounded below but not above, and $h(x) = \sin x$ is bounded above and below (where f , g , and h are from \mathbb{R} to \mathbb{R})

Example 4.1.1

- Find the limit as $r \rightarrow 0$ of the function $f(r) = 1 + r^2$ defined on $(0, 1)$.
- Find the limit as $r \rightarrow 0$ of the function $g(r) = \cos r$ defined on $(0, 1)$.



Solution

- Since $f(r)$ decreases as r decreases and f is bounded below, Definition 4.1.1 says that the desired limit exists and is equal to the greatest lower bound of the range of f . Since the range of f is the interval $(1, 2)$, we see that the greatest lower bound is $\boxed{1}$.

(b) Similarly, since $g(r)$ increases as r decreases and is bounded above, Definition 4.1.1 says that the desired limit exists and is equal to the least upper bound of the range of f . We can see from the geometric definition of cosine (see Section A.1.2) that 1 is an upper bound for the range of \cos and that no smaller number is an upper bound. Therefore, the limit is equal to $\boxed{1}$.

* This means that only values of r to the right of zero are considered. Often the notation $\lim_{r \rightarrow 0^+}$ is used to communicate the exclusion of values to the left, but this notation is not necessary here since the domain of f includes no negative values.

* Put another way, $m(r)$ is the greatest lower bound and $M(r)$ the least upper bound of the range of f restricted to $B^*(\mathbf{a}, r) \cap D$

Definition 4.1.1 has many restrictions: it only applies to monotone, single-variable functions, and it is fundamentally one-sided.* We will drop all these restrictions in one swoop:

Definition 4.1.2: Limit of a function of function of multiple variables

Let $n \geq 1$. Suppose that $D \subset \mathbb{R}^n$, that $\mathbf{a} \in \mathbb{R}^n$ and that $f : D \rightarrow \mathbb{R}$ is a function. For each $r \in (0, 1)$, we consider the **punctured ball** $B^*(\mathbf{a}, r) = \{\mathbf{x} \in \mathbb{R}^n : 0 < |\mathbf{x} - \mathbf{a}| \leq r\}$ of radius r centered at \mathbf{a} . Suppose that $B^*(\mathbf{a}, r) \cap D$ is non-empty for all $r > 0$.

We define* $m(r)$ and $M(r)$ so that $[m(r), M(r)]$ is the smallest closed interval which contains the image of $B^*(\mathbf{a}, r) \cap D$ under f . Then we say that the limit of $f(\mathbf{x})$ as $\mathbf{x} \rightarrow \mathbf{a}$ exists if $m(r)$ and $M(r)$ converge to a common value L as $r \rightarrow 0$. In that case, we say $\lim_{\mathbf{x} \rightarrow \mathbf{a}} f(\mathbf{x}) = L$; otherwise, we say that the limit does not exist.

To illustrate this definition, let's look at a function defined on \mathbb{R}^1 whose limit at the origin doesn't exist.

Example 4.1.2

Show that the limit of $\cos(1/x)$ does not exist as $x \rightarrow 0$.

Solution

For any $r > 0$, the function $\cos(1/x)$ has a maximum value of 1 and a minimum value of -1 in the punctured interval $B^*(0, r)$, because $\cos(1/x) = 1$ whenever $x = \frac{1}{\pi k}$ for some even integer k , and $\cos(1/x) = -1$ whenever $x = \frac{1}{\pi k}$ for some odd integer k .

Therefore, M is the constant function 1 and m is the constant function -1 . These functions converge to different values, so the limit does not exist.

The limit in Example 4.1.2 fails to exist because f is *oscillatory* at the origin. Limit existence failures for functions of multiple variables can look quite different:

Example 4.1.3

Show that* $f(x, y) = -\frac{xy}{x^2+y^2}$ does not have a limit as $(x, y) \rightarrow (0, 0)$.

* See Figure 4.2 for a graph

Solution

Let's represent the point (x, y) as* $(t \cos \theta, t \sin \theta)$. For all $t > 0$, we have

$$f(t \cos \theta, t \sin \theta) = -\frac{t^2 \sin \theta \cos \theta}{t^2(\cos^2 \theta + \sin^2 \theta)} = -\frac{1}{2} \sin 2\theta.$$

* We use t instead of r since we're already using r as the radius of the punctured ball

Therefore, for any $r > 0$, we have $m(r) = -\frac{1}{2}$ and $M(r) = \frac{1}{2}$. Since these functions converge to unequal values as $r \rightarrow 0$, it follows that the limit does not exist.

Example 4.1.4

Use Definition 4.1.2 to show directly that* $\lim_{(x,y) \rightarrow (0,0)} [3 + x^2 - y^2] = 3$.

* See Figure 4.1 for a graph

Solution

Writing $f(x, y) := 3 + x^2 - y^2$ in terms of polar coordinates as $3 + t^2(\cos^2 \theta - \sin^2 \theta) = 3 + t^2 \cos 2\theta$, we see that the least and greatest values of $f(x, y)$ for $(x, y) \in B^*((0, 0), r)$ are $m(r) = 3 - r^2$ and $M(r) = 3 + r^2$, respectively. These functions converge as $r \rightarrow 0$ to a common value of 3, by Definition 4.1.1. Therefore, $\lim_{(x,y) \rightarrow (0,0)} [3 + x^2 - y^2] = 3$.

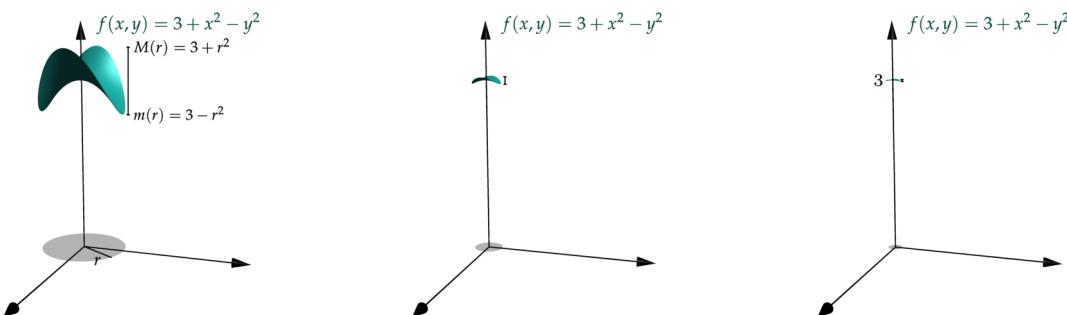


Figure 4.1 The image under f of the ball of radius r is the interval $[3 - r^2, 3 + r^2]$, which shrinks down around 3 as $r \rightarrow 0$. So the limit exists and equals 3

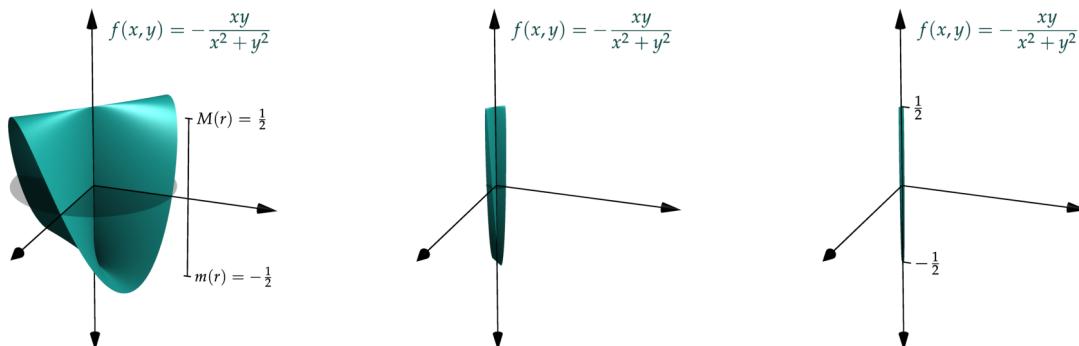


Figure 4.2 The image under f of the ball of radius r is $[-\frac{1}{2}, \frac{1}{2}]$ no matter how small r is. So the limit does not exist

Exercise 4.1.1

Use Definition 4.1.2 to show that $f(x, y) = \sin\left(\frac{1}{x^2+y^2}\right)$ does not have a limit as $(x, y) \rightarrow (0, 0)$.

* Indeed, these values are all equal to $-\frac{1}{2}$

In Example 4.1.3, there are two directions of approach along which f has different limits. Along the ray $\theta = \frac{\pi}{4}$, the values of $f(x, y)$ approach* $-\frac{1}{2}$. Along the ray $\theta = -\frac{\pi}{4}$, $f(x, y)$ converges to $\frac{1}{2}$. This is always an obstruction to the existence of a limit:

Theorem 4.1.1

!!!

Suppose that \mathbf{r}_1 and \mathbf{r}_2 are paths in the plane with the property that $\mathbf{r}_1(0) = (a, b)$ and $\mathbf{r}_2(0) = (a, b)$. If $\lim_{t \rightarrow 0} f(\mathbf{r}_1(t))$ and $\lim_{t \rightarrow 0} f(\mathbf{r}_2(t))$ exist and are unequal, show that $\lim_{(x,y) \rightarrow (a,b)} f(x, y)$ does not exist.

Proof

Let's define $L_i = \lim_{t \rightarrow 0} f(\mathbf{r}_i(t))$ for $i = 1$ and $i = 2$. Note that $m(r) \leq \min(L_1, L_2)$ for all $r > 0$ since any number larger than $\min(L_1, L_2)$ is not a lower bound for the values of the function on $B^*((a, b), r)$. Similarly, $M(r) \geq \max(L_1, L_2)$ for all $r > 0$. Therefore, $\lim_{r \rightarrow 0} m(r) \leq \min(L_1, L_2)$ while $\lim_{r \rightarrow 0} M(r) \geq \max(L_1, L_2)$. So the limit of $f(x, y)$ does not exist as $(x, y) \rightarrow (a, b)$.

With the notion of a multidimensional limit in hand, we can define continuity the same way we did for the one-dimensional case.

Definition 4.1.3

Suppose $n \geq 2$. A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous at a point in \mathbb{R}^n if and only if the limit of f exists at that point and equals the value of the function there.

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be continuous if it is continuous at every point in \mathbb{R}^n .

More generally, a function $f : D \rightarrow \mathbb{R}$ —where $D \subset \mathbb{R}^n$ —is said to be continuous if it is continuous at each point in its domain D . The following theorem gives us some tools for establishing continuity.

Theorem 4.1.2: Continuous functions

$g \circ f$ denotes the composition of g and f . See Appendix A.1.1

1. If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous and $g : \mathbb{R} \rightarrow \mathbb{R}$ is continuous, then* $g \circ f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous.
2. A sum or product of continuous functions is continuous.
3. The “coordinate-extracting” functions $f(x, y, z) = x$, $f(x, y, z) = y$, etc., are continuous.

Example 4.1.5

Show that $\lim_{(x,y,z) \rightarrow (0,0,0)} \left(e^{\sin x} + \frac{xyz}{1+x^2z^2} \right) = 1$.

Solution

We begin by showing that $e^{\sin x} + \frac{xyz}{1+x^2z^2}$ is continuous. Note that $e^{\sin x}$ is a composition of continuous functions:

$$(x, y, z) \mapsto x \mapsto \sin x \mapsto e^{\sin x},$$

Therefore, it's continuous by Theorem 4.1.2. Similarly, $\frac{xyz}{1+x^2z^2}$ is continuous wherever $1+x^2z^2 \neq 0$, which is everywhere since $(xz)^2 \geq 0$. Finally, the sum of two continuous functions is continuous, so $e^{\sin x} + \frac{xyz}{1+x^2z^2}$ is continuous.

Since the function above is continuous, its limit at each point is equal to its value at that point. So we substitute $x = y = z = 0$ and find that the value of the function at the origin is $e^0 + \frac{0}{1+0} = 1$.

Suppose we know that the limits of $f(x, y)$ along every line passing through the origin exist and that they are all equal to some common value L . Does this imply that $\lim_{(x,y) \rightarrow (0,0)} f(x, y) = L$? Perhaps it seems that it should, since we've accounted for every possible angle of approach. Remarkably, this turns out not to be the case:

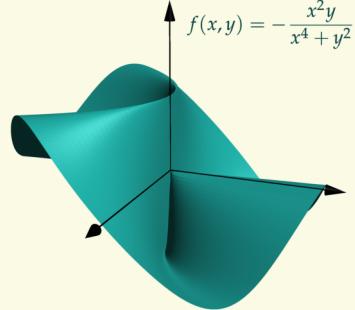
Example 4.1.6

Show that $\lim_{(x,y) \rightarrow (0,0)} \frac{-x^2y}{x^4+y^2}$ does not exist even though the limits along every line through the origin exist and are equal.

Solution

We begin by checking the limit along the line $\mathbf{r}(t) = (t \cos \theta, t \sin \theta)$ (which is the line passing through the origin as well as the point on the unit circle whose angle with respect to the positive x -axis is θ). We find

$$\begin{aligned} f(t \cos \theta, t \sin \theta) &= \frac{-t^3 \cos^2 \theta \sin \theta}{t^4 \cos^4 \theta + t^2 \sin^2 \theta} \\ &= \frac{-t \cos^2 \theta \sin \theta}{t^2 \cos^4 \theta + \sin^2 \theta}. \end{aligned}$$



If we consider the limit of this expression as $t \rightarrow 0$ with θ fixed, then the $\cos \theta$ and $\sin \theta$ factors are constants. So we see that the numerator converges to 0 and the denominator converges to $\sin^2 \theta$. Therefore, as long as $\sin \theta \neq 0$, we have $\lim_{t \rightarrow 0} f(t \cos \theta, t \sin \theta) = 0 / \sin^2 \theta = 0$. However, if $\sin \theta = 0$, then $f(t \cos \theta, t \sin \theta) = 0$ for all t , so $\lim_{t \rightarrow 0} f(t \cos \theta, t \sin \theta) = 0$ in that case too.

However, note that if we consider the limit along the parabolic path $\mathbf{r}(t) = (t, -t^2)$, we get

$$f(t, -t^2) = -\frac{t^2(-t^2)}{t^4 + (-t^2)^2} = \frac{1}{2}.$$

Therefore, the limit along this path is equal to $\frac{1}{2}$. Thus there are two paths (this one, as well as any straight-line path through the origin) along which f has different limits. Therefore, the limit of $f(x, y)$ as $(x, y) \rightarrow (0, 0)$ does not exist.

Note: this makes sense graphically, because this function also has a crease along the z -axis. But now we have to follow a parabolic path to travel along the top "ridge" and realize a limiting value other than zero.

Exercise 4.1.2

Show that $\lim_{(x,y) \rightarrow (0,0)} \frac{x^3 + y^3}{x^2 + y^2} = 0$.

Exercise 4.1.3

Consider the function f defined by $f(x, y) = \frac{x-y}{x^3-y}$ whenever $y \neq x^3$, and $f(x, y) = 1$ when $y = x^3$. Show that f is not continuous at $(1, 1)$. Evaluate the limits along $x = 1$ and along $y = 1$.

4.2 Partial derivatives

 on partial derivatives

Suppose f is a function from \mathbb{R} to \mathbb{R} . The derivative f' of f is the answer to the question “how does $f(x)$ change when x changes just a little?” More precisely, if $a \in \mathbb{R}$, we define

$$f'(a) = \lim_{h \rightarrow 0} \frac{\overbrace{f(a+h) - f(a)}^{\text{how much } f \text{ changes}}}{\underbrace{h}_{\text{how much the input changes}}}$$

This means that if we know $f'(a)$, then we can estimate $f(a+h) - f(a)$ for h very small:

$$f(a+h) - f(a) \approx hf'(a).$$

So the derivative measures **how sensitive $f(x)$ is to small changes in x** .

What is the most natural corresponding idea for the derivative at some point (a, b) of a function f from \mathbb{R}^2 to \mathbb{R} ? We were only able to adjust a value $x \in \mathbb{R}$ by increasing or decreasing it a little. A point in \mathbb{R}^2 , by contrast, can be moved in any direction. Two directions are particularly easy to study: (i) move x a little while holding y fixed, and (ii) move y a little while holding x fixed. Accordingly, we define **partial derivatives***

$$(\partial_x f)(a, b) = \lim_{h \rightarrow 0} \frac{f(a+h, b) - f(a, b)}{h}, \text{ and}$$

$$(\partial_y f)(a, b) = \lim_{h \rightarrow 0} \frac{f(a, b+h) - f(a, b)}{h}.$$

Calculating partial derivatives of elementary functions *isn't actually a new skill*: since one of the two variables is being held constant, we are effectively taking a derivative of a single-variable function.

Example 4.2.1

Find the partial derivatives* f_x and f_y of $f(x, y) = e^x \sin(xy)$ at $(x, y) = (1, 0)$.

* ∂_x is read “partial x ”. Also, the role of x here is purely as a label that means “with respect to the first coordinate”. It does not represent a number, as the symbol x usually does

* f_x is an alternate notation for $\partial_x f$, and similarly for y

Solution

We can find the partial derivative with respect to x at *any* point (x, y) by treating y as constant and applying single-variable differentiation rules:^{*}

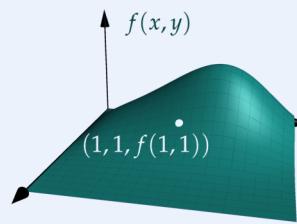
$$\begin{aligned}(\partial_x f)(x, y) &= e^x y \cos(xy) + e^x \sin(xy) \\(\partial_y f)(x, y) &= x \cos(xy) e^x\end{aligned}$$

So the partial derivatives at $(1, 0)$ with respect to x and y are 0 and e , respectively.

^{*} If you have difficulty getting used to holding a variable constant, consider replacing it with some number like 17; then substitute back at the end

Example 4.2.2

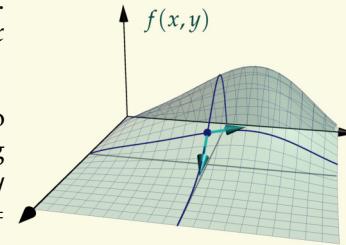
Consider the function f whose graph is shown. Determine the sign of $(\partial_x f)(1, 1)$ and the sign of $(\partial_y f)(1, 1)$.



Solution

If we increase x a little while holding y constant, then f decreases. Therefore, $(\partial_x f)(1, 1) < 0$. If we increase y a little while holding x constant, then f increases. Therefore, $(\partial_y f)(1, 1) > 0$.

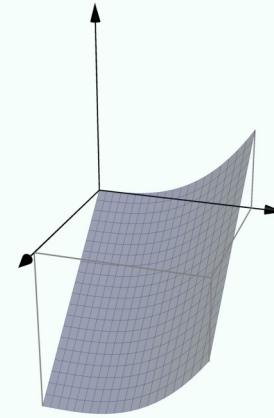
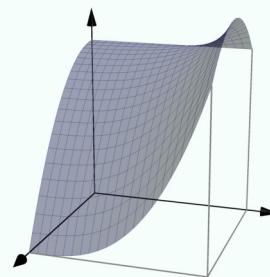
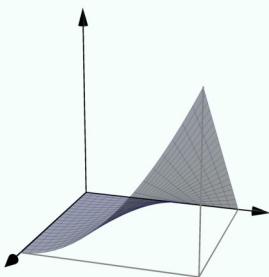
Graphically, the partial derivative with respect to x at a point is equal to the slope of the trace of the graph in the “ $y = \text{constant}$ ” plane passing through that point. Similarly, the partial derivative with respect to y at a point is equal to the slope of the trace of the graph in the “ $x = \text{constant}$ ” plane passing through that point.*



* Thus we can think of partial derivatives as an application of our “slice it up” strategy for understanding three dimensional objects through two dimensional traces

Exercise 4.2.1

The following three graphs represent a function f and its two partial derivatives $\partial_x f$ and $\partial_y f$, in some order. Which is which?



The following theorem says that order doesn't matter when successively taking partial derivatives.

Theorem 4.2.1: Clairaut's theorem

Suppose $f : D \rightarrow \mathbb{R}$, where D is a disk in \mathbb{R}^2 . If $\partial_x \partial_y f$ and $\partial_y \partial_x f$ exist and are continuous, then $\partial_x \partial_y f = \partial_y \partial_x f$ throughout D .

Exercise 4.2.2

Verify the conclusion of Clairaut's theorem for $f(x, y) = e^{xy} \sin y$.

Exercise 4.2.3

Suppose that $f(4.2, 6.8) = 2$, $f(4.3, 6.8) = 3$, $f(4.2, 6.9) = 4$, and $f(4.3, 6.9) = 6$. Justify the approximation $\partial_x \partial_y f(4.2, 6.8) \approx 100$. Apply similar reasoning to obtain the approximation $\partial_y \partial_x f(4.2, 6.8) \approx 100$.

4.3 Linear approximation

The following example shows that partial derivatives don't tell the whole story when it comes to differentiating functions of multiple variables.

Example 4.3.1

Consider the function f for which $f(0, 0) = 0$ and $f(x, y) = -\frac{xy}{x^2+y^2}$ for all $(x, y) \neq (0, 0)$. Show that both partial derivatives of f at the origin are equal to zero.

Solution

If we move x a little from $x = 0$ while holding $y = 0$ fixed, the value of f doesn't change at all. Therefore,

$$(\partial_x f)(0, 0) = \lim_{h \rightarrow 0} \frac{f(h, 0) - f(0, 0)}{h} = \lim_{h \rightarrow 0} \frac{0}{h} = \lim_{h \rightarrow 0} 0 = 0.$$

The same is true for the partial derivative with respect to y .

However, recall from Example 4.1.3 that the function in Example 4.3.1 isn't even continuous at the origin! We haven't said yet what is required for a function of two variables to be considered differentiable, but whatever the definition, we surely cannot allow functions which aren't continuous to be deemed differentiable. This shouldn't be surprising: the partial derivatives only look at the behavior of the function along two slices. A good definition of differentiability at (a, b) should account for how the function behaves all around (a, b) .

Another perspective on differentiability in the single-variable context is that *differentiable functions are the ones which are well-approximated by linear functions*:

Theorem 4.3.1

A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable at $a \in \mathbb{R}$ if and only if there exists a linear function $L(x) = c_0 + c_1(x - a)$ such that*

$$\lim_{x \rightarrow a} \frac{f(x) - L(x)}{|x - a|} = 0.$$

* This equation says that L approximates f so well that the difference between f and L , even after being divided by the tiny number $|x - a|$, still goes to 0 as $x \rightarrow a$

This perspective on differentiability turns out to generalize very nicely to functions of multiple variables. Let's make it a definition.

Definition 4.3.1: Differentiability for a function of two variables

A function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is differentiable at $(a, b) \in \mathbb{R}^2$ if and only if there exists a linear function $L(x, y) = c_0 + c_1(x - a) + c_2(y - b)$ such that

$$\lim_{(x,y) \rightarrow (a,b)} \frac{f(x, y) - L(x, y)}{\sqrt{(x - a)^2 + (y - b)^2}} = 0.$$

Lots of functions are differentiable. The following theorem establishes a handy way to check differentiability.

Theorem 4.3.2: Criterion for differentiability

If the partial derivatives $\partial_x f$ and $\partial_y f$ exist in some disk centered at (a, b) and are continuous at (a, b) , then f is differentiable at (a, b) .

The most common situation is that partial derivatives exist and are continuous everywhere, in which case Theorem 4.3.2 implies that f is differentiable everywhere.

Example 4.3.2

Show that $f(x, y) = e^{xy} \sin(x^2 + y^2)$ is differentiable at every point in \mathbb{R}^2 .

Solution

We can take partial derivatives of f with respect to both x and y and get functions which are built from x and y using addition/multiplication as well as the continuous functions $x \mapsto e^x$, $x \mapsto \sin x$, and $x \mapsto \cos x$. Therefore, the partial derivatives exist and are continuous everywhere. Thus Theorem 4.3.2 implies that f is differentiable everywhere.

In the denominator we replaced $|x - a|$, whose geometric meaning is the distance from x to a on the number line, with the formula for the distance from (x, y) to (a, b) in the plane.

* To interpret this statement as a theorem, we would need to first say what it means for the plane to be *tangent* on some basis other than Definition 4.3.1 (to avoid circular reasoning). So take this as intuition only

Graphically, Definition 4.3.1 says that a function is differentiable at (a, b) if we can draw a plane which is tangent* to the graph of f at the point $(a, b, f(a, b))$.

In Theorem 4.3.1, the coefficients of the approximating function L are the value of the function f at a and the derivative of f at a . The coefficients in Definition 4.3.1 are also quantities that we have names for: as suggested by Figure 4.3, c_1 is the value of the function at (a, b) and c_1 and c_2 are the two partial derivatives at (a, b) .

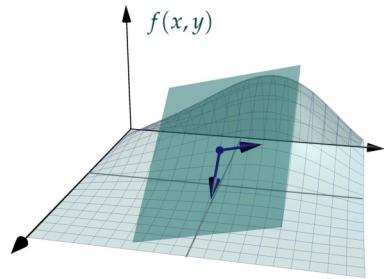


Figure 4.3 A plane tangent to the graph of a function f

Theorem 4.3.3: Linear Approximation

If $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is differentiable at $(a, b) \in \mathbb{R}^2$, then

!!!

$$\lim_{(x,y) \rightarrow (a,b)} \frac{f(x,y) - \overbrace{[f(a,b) + (\partial_x f)(a,b)(x-a) + (\partial_y f)(a,b)(y-b)]}^{L(x,y)}}{\sqrt{(x-a)^2 + (y-b)^2}} = 0.$$

Let's see how this theorem can be used numerically.

Example 4.3.3

Consider the function $f(x, y) = \frac{e^{xy}}{e(1+x^2)}$. Use a tangent plane to approximate $f(0.99, 0.98)$.

Solution

Noticing that $(0.99, 0.98)$ is very close to $(1, 1)$, we differentiate $f(x, y)$ with respect to x and with respect to y and find*

$$\begin{aligned} (\partial_x f)(1, 1) &= \left. \left(\frac{ye^{xy}}{e(x^2+1)} - \frac{2xe^{xy}}{e(x^2+1)^2} \right) \right|_{(x,y)=(1,1)} = 0. \\ (\partial_y f)(1, 1) &= \left. \frac{xe^{xy}}{e(x^2+1)} \right|_{(x,y)=(1,1)} = \frac{1}{2}. \end{aligned}$$

Therefore, $f(0.99, 0.98) \approx f(1, 1) + 0(0.99 - 1) + \frac{1}{2}(0.98 - 1) = \frac{1}{2} + \frac{1}{2} \cdot (-\frac{1}{50}) = 0.49$.

* The bar notation means "substitute"

The actual value is 0.490197...

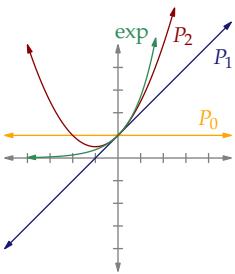


Figure 4.4 The first few Taylor polynomials for the exponential function $\exp(x) = e^x$

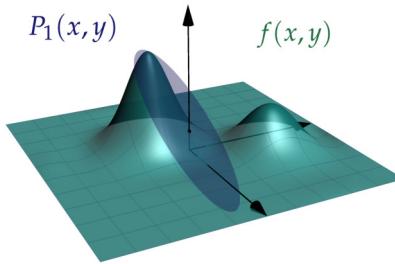


Figure 4.5 The linear Taylor polynomial of a function f centered at the origin

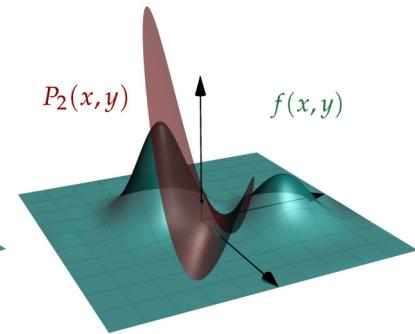


Figure 4.6 The quadratic Taylor polynomial of a function f centered at the origin

4.4 Taylor's theorem*

Let us briefly recall the Taylor* series story for functions of a single variable. The ***k*th order Taylor polynomial** of an n -times-differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$ centered at a point $a \in \mathbb{R}$ is defined to be

$$P_k(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \cdots + \frac{f^{(k)}(a)}{k!}(x - a)^k.$$

The terms on the right-hand side are motivated by the idea that as many derivatives of P_k as possible should match those of f at the point a . Then P_k is an excellent approximation of f for values of x near a :

Theorem 4.4.1: Taylor's theorem, functions of a single variable

If I is an interval in \mathbb{R} and $f : I \rightarrow \mathbb{R}$ is differentiable k times, then

$$\lim_{x \rightarrow a} \frac{f(x) - P_k(x)}{(x - a)^k} = 0.$$

* A Taylor polynomial centered at the origin is called a *Maclaurin* polynomial

The idea in higher dimensions is entirely analogous: given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we define the k th order Taylor polynomial of f at a point $\mathbf{a} \in \mathbb{R}^n$ to be the polynomial p such that all of p 's *mixed* partial derivatives of order k and lower match those of f at \mathbf{a} . For simplicity, we state the theorem for \mathbb{R}^2 -valued functions.

Theorem 4.4.2: Taylor's theorem, functions of multiple variables

If U is an open set* in \mathbb{R}^2 and $f : U \rightarrow \mathbb{R}$ is differentiable n times, then we define

$$\begin{aligned} P_k(x, y) &= f(a, b) + (\partial_x f)(a, b)(x - a) + (\partial_y f)(a, b)(y - b) \\ &\quad + \frac{1}{2!0!}(\partial_x^2 f)(a, b)(x - a)^2 + \frac{1}{1!1!}(\partial_x \partial_y f)(a, b)(x - a)(y - b) + \frac{1}{0!2!}(\partial_y^2 f)(a, b)(y - b)^2 + \cdots, \end{aligned}$$

where we continue until we have all terms of the form $\frac{1}{i!j!}(\partial_x^i \partial_y^j f)(a, b)(x - a)^i(y - b)^j$ where $i + j \leq k$. Then

$$\lim_{(x,y) \rightarrow (a,b)} \frac{f(x, y) - P_k(x, y)}{|\langle x, y \rangle - \langle a, b \rangle|^k} = 0.$$

* An open set is one that contains a small ball around each of its points—in other words, a set that contains none of the points on its boundary

Exercise 4.4.1

How many terms of degree k appear in $P_k(x, y)$? Write out all the terms of degree 3, and show that all the third-order partial derivatives of f and P_3 match at (a, b) .

Exercise 4.4.2

Find the Taylor polynomials P_1 and P_2 for the function $\frac{1}{x^2+y^2+1}$ from \mathbb{R}^2 to \mathbb{R} centered at $(0, 0)$.

4.5 Multivariable optimization

on local extrema

The following problem is a typical example of a single-variable optimization problem.

Example 4.5.1

Find the maximum and minimum of $f(x) = |(x - 1)(3 - x)|$ over the interval $[0, 3]$.

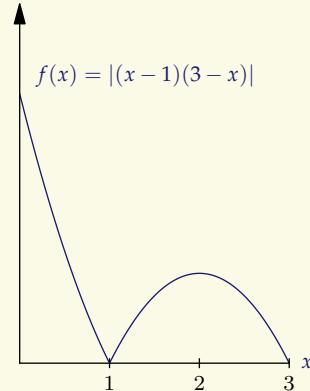
Solution

* The extreme value theorem for single-variable functions says that a continuous function on a closed interval $[a, b]$ achieves a max and a min

Since f is continuous over the closed and bounded interval $[0, 3]$, we know by the extreme value theorem* that it has a maximum and a minimum value over $[0, 3]$. Furthermore, these extrema must be realized at either a *critical point*, at which f is either not differentiable or has derivative zero, or else at an endpoint of the interval.

We check that f is not differentiable at 1 or 3 (see the graph). Also, we can solve $f'(x) = 0$ to find that f has a horizontal tangent line at $x = 2$.

Finally, we can check the values of f at the endpoints 0 and 3, as well as the critical points strictly between them, namely 1 and 2. We find that the maximum value is $f(0) = \boxed{3}$, and the minimum value is $\boxed{0}$, which occurs at $x = 1$ and at $x = 3$.



How does this story change when we consider a function of multiple variables? For concreteness, let's suppose $D = [0, 1]^2$ and that $f : D \rightarrow \mathbb{R}$ is a continuous function. Consider the graph of the function

$$f(x, y) = -x^2 - y^2 + x + \frac{2}{3}y + \frac{23}{36},$$

shown in Example 4.5.2. As in the single-variable case, f does indeed have a maximum value and a minimum value. This is ensured by the following multivariable generalization of the single-variable extreme value theorem.

Theorem 4.5.1: Extreme value theorem

Suppose that $D \subset \mathbb{R}^n$. We say that D is **closed** if it contains all of its boundary points. We say that D is **bounded** if it is contained in a ball of radius R for some $R < \infty$.

If D is closed and bounded and if $f : D \rightarrow \mathbb{R}$ is a continuous function, then f achieves a maximum value and a minimum value on D . In other words, there exists $\mathbf{a} \in D$ such that $f(\mathbf{a}) \geq f(\mathbf{b})$ for all $\mathbf{b} \in D$.

Furthermore, if f is differentiable at a point \mathbf{a} on the inside of D and $(\partial_x f)(\mathbf{a}) > 0$ or $(\partial_x f)(\mathbf{a}) < 0$, then f cannot have a maximum or minimum value at \mathbf{a} since we can increase the function's output value by slightly adjusting the first coordinate of \mathbf{a} in one direction, and we can decrease it by adjusting it in the opposite direction. So $\partial_x f = 0$ at any value where f has an extremum, and similarly for $\partial_y f$.

Theorem 4.5.2: Critical points

If $f : D \rightarrow \mathbb{R}$ achieves its maximum or minimum value at $\mathbf{a} \in D$, then

- (i) $(\partial_x f)(\mathbf{a}) = (\partial_y f)(\mathbf{a}) = 0$, or
- (ii) f is not differentiable at \mathbf{a} , or
- (iii) \mathbf{a} is on the boundary of D .

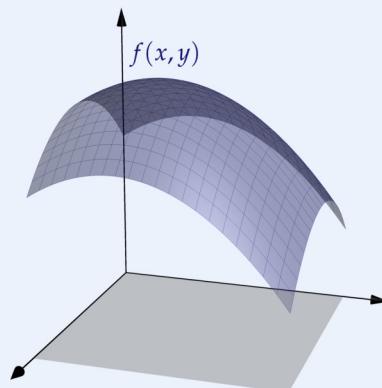
These theorems give rise to a strategy for finding the extrema of a function $f : D \rightarrow \mathbb{R}$, where $D \subset \mathbb{R}^2$: (i) set both partial derivatives of f equal to 0 and solve to find critical points inside D (also include any points where f is not differentiable), and (ii) find the extreme values of f on* ∂D . Let's do an example.

Example 4.5.2

Find the extreme values of the function

$$f(x, y) = -x^2 - y^2 + x + \frac{2}{3}y + \frac{23}{36}$$

over the square $[0, 1]^2$.



* The notation ∂D means "the boundary of D ", which is the set of points $p \in \mathbb{R}^2$ such that any small ball centered at p includes points in D and points not in D .

Solution

We begin by finding the critical points inside the square. We find

$$\begin{aligned}(\partial_x f)(x, y) &= -2x + 1 \\(\partial_y f)(x, y) &= -2y + \frac{2}{3}.\end{aligned}$$

These quantities are both equal to zero only when $x = \frac{1}{2}$ and $y = \frac{1}{3}$. So $(1/2, 1/3)$ is the only critical point inside the square.

To optimize f along the $x = 0$ side, we look at

$$f(0, y) = -y^2 + \frac{2}{3}y + \frac{23}{36},$$

which has a critical point at $y = 1/3$. So $(0, 1/3)$ is a **boundary critical point**, and we should also check the two endpoints $(0, 0)$ and $(0, 1)$. Similarly, for the other three sides, we identify the points $(1, 1/3)$, $(1/2, 0)$, and $(1/2, 1)$, as boundary critical points as well as the other two corners $(1, 1)$ and $(1, 0)$. So, all together:

(x, y)	$(0, 0)$	$(1, 0)$	$(0, 1)$	$(1, 1)$	$(0, 1/3)$	$(1, 1/3)$	$(1/2, 0)$	$(1/2, 1)$	$(1/2, 1/3)$
$f(x, y)$	$23/36$	$23/36$	$11/36$	$11/36$	$3/4$	$3/4$	$8/9$	$5/9$	1
$f(x, y)$	0.64	0.64	0.31	0.31	0.75	0.75	0.89	0.56	1

So the maximum value is 1 and the minimum value is $\frac{11}{36}$.

Exercise 4.5.1

Find the maximum value of $f(x, y) = 10x^2y - x$ over the closed triangle with vertices $(0, 0)$, $(1, 0)$, and $(0, 1)$.

4.6 Second derivative test

Recall from single-variable calculus that the second derivative of a twice-differentiable function can—if it is nonzero—be used to ascertain whether a function has a local minimum* or a local maximum at a given critical point. This is because the convexity of a twice-differentiable function indicates whether the graph of the function is shaped like \cup or \cap .

The situation in higher dimensions is more subtle. Archetypal examples are

- (a) $x^2 + y^2$, which has a bowl-shaped graph like , and so a local minimum at the origin,
- (b) $-x^2 - y^2$, which has an umbrella-shaped graph like , and so a local maximum at the origin, and
- (c) $x^2 - y^2$, which has a saddle-shaped graph like , and so neither a local min nor a local max at the origin.

The following theorem gives us a direct way to distinguish these cases.

* A function f has a local minimum at a point a if there is a small neighborhood of a throughout which the function's values are no smaller than $f(a)$.

Theorem 4.6.1: Second derivative test

Suppose that U is an open set in \mathbb{R}^2 and $f : U \rightarrow \mathbb{R}$ is a twice-differentiable function with a critical point at (a, b) . We define $D = (\partial_x^2 f \partial_y^2 f - [\partial_x \partial_y f]^2)(a, b)$. Then*

- (a) if $D > 0$ and $(\partial_x^2 f)(a, b) > 0$, then f has a local minimum at (a, b) ,
- (b) if $D > 0$ and $(\partial_x^2 f)(a, b) < 0$, then f has a local maximum at (a, b) , and
- (c) if $D < 0$, then f has a saddle point at (a, b) .

* When $D = 0$, this theorem doesn't tell us anything

A proof of the second derivative test is given in Appendix A.3.3.

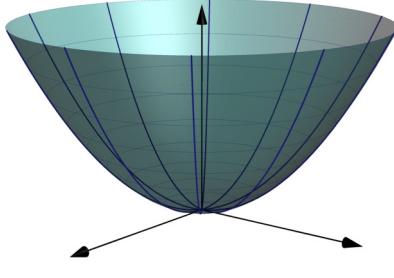


Figure 4.7 If every vertical slice of the graph of f (passing through the origin) is convex, then f has a local minimum at the origin

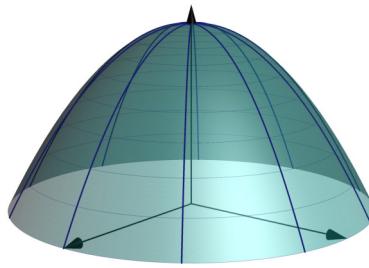


Figure 4.8 If every slice is concave, then f has local maximum

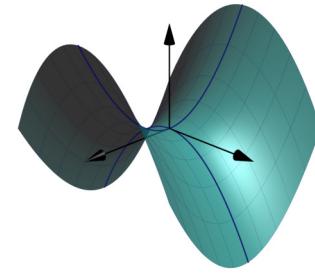


Figure 4.9 If there are both concave and convex slices, then f has a saddle point

Example 4.6.1

Find the critical points of the function $f(x, y) = (2x^2 + 3y^2) e^{-x^2-y^2}$ defined on \mathbb{R}^2 , and classify each critical point as a local minimum, a local maximum, or a saddle point.

Solution

We have

$$\begin{aligned}\partial_x f &= 2x(-2x^2 - 3y^2 + 2)e^{-x^2-y^2}, \text{ and} \\ \partial_y f &= 2y(-2x^2 - 3y^2 + 3)e^{-x^2-y^2}.\end{aligned}$$

To find the critical points of f , we look for all pairs (x, y) for which both of these expressions are equal to zero. We see that $x = 0$ implies $\partial_x f = 0$. In that case, we have $\partial_y f = 0$ if and only if $y \in \{-1, 0, 1\}$. Similarly, $y = 0$ implies $\partial_y f = 0$, and in that case we have $\partial_x f = 0$ only if $x \in \{-1, 0, 1\}$. If neither x nor y is zero, then we may divide the two equations by x and y , respectively, and we arrive at a contradiction since the resulting left-hand sides are unequal and therefore cannot both be equal to zero. Therefore, we have found all the critical points.

To apply the second derivative test, we work out that*

$$\begin{aligned}D = \partial_x^2 f \partial_y^2 f - (\partial_x \partial_y f)^2 &= \left(-32x^6 - 128x^4y^2 + 128x^4 - 168x^2y^4 + 328x^2y^2 - 136x^2 - 72y^6 \right. \\ &\quad \left. + 228y^4 - 156y^2 + 24 \right) e^{-2x^2-2y^2}.\end{aligned}$$

Definitely want to use some computational assistance here

Then we check the sign of D at each critical point:

(x, y)	$(0, 0)$	$(1, 0)$	$(-1, 0)$	$(0, 1)$	$(0, -1)$
$D e^{2x^2+2y^2}$	24	-16	-16	24	24

We see that the points $(1, 0)$ and $(-1, 0)$ are saddle points of f , while f has a local extremum at the origin and at the points $(0, 0)$, $(0, 1)$ and $(0, -1)$. To classify these extrema, we check that $\partial_x^2 f$ at these points is equal to 4 , $-2/e$, and $-2/e$, respectively. Thus f has a local minimum at the origin and a local maximum at $(0, 1)$ and $(0, -1)$.

These classifications accord with the graph of f , shown in Figure 4.10.

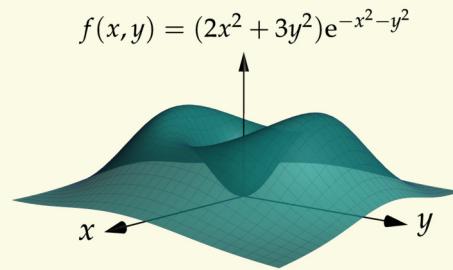


Figure 4.10 The graph of f

Exercise 4.6.1

Find all the critical points of $f(x, y) = x(x+y)(y+y^2)$ and apply the second derivative test to classify as many as possible as a local minimum, a local maximum, or a saddle point.

Exercise 4.6.2

Use the examples $x^4 + y^4$ and $-x^4 - y^4$, which have a critical point at the origin, to show that the second derivative test must be inconclusive when $D = 0$.

Exercise 4.6.3

Show that in parts (a) and (b) of the second derivative test, we have $(\partial_x^2 f)(a, b) > 0$ if and only if $(\partial_y^2 f)(a, b) > 0$.

4.7 Directional derivative and gradient

mat on directional derivatives and the gradient

* That is:
up/down, left-/right

The two partial derivatives of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ tell us how f changes when (x, y) is wiggled a bit, but only in the four cardinal directions.* What about all the other directions? Suppose that \mathbf{u} is a **unit vector** in \mathbb{R}^2 , meaning that its length is 1. (see Figure 4.11).

Definition 4.7.1: Directional derivative

The **directional derivative** of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ in the direction $\mathbf{u} \in \mathbb{R}^2$ is defined by

$$D_{\mathbf{u}}(f)(a, b) = \lim_{h \rightarrow 0} \frac{f((a, b) + h\mathbf{u}) - f(a, b)}{h}.$$

In other words, move (x, y) a small distance h in the \mathbf{u} direction, measure how much f changed, and then divide by h .

If f is differentiable at (a, b) , then it is well-approximated by a linear function $L(x, y) = c_0 + c_1(x - a) + c_2(y - b)$ around (a, b) . This linear function has the property that its slope in the \mathbf{u} direction can be worked out by separating the $h\mathbf{u}$ -step into a $\langle hu_1, 0 \rangle$ step followed by a $\langle 0, hu_2 \rangle$ step. The value of L changes by hc_1u_1 over the $\langle hu_1, 0 \rangle$ step and by hc_2u_2 over the $\langle 0, hu_2 \rangle$ step. Altogether, the change of L is equal to $(c_1u_1 + c_2u_2)h$. This leads to the following theorem.

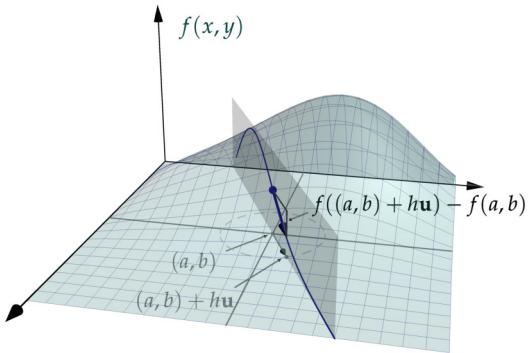


Figure 4.11 The derivative of f in the direction \mathbf{u}

Theorem 4.7.1: Directional derivative formula

If $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is differentiable at (a, b) and $\mathbf{u} \in \mathbb{R}^2$, then

$$D_{\mathbf{u}}f(a, b) = (\partial_x f)(a, b)u_1 + (\partial_y f)(a, b)u_2 = (\nabla f)(a, b) \cdot \mathbf{u},$$

where $(\nabla f)(a, b) = \langle (\partial_x f)(a, b), (\partial_y f)(a, b) \rangle$.

The quantity $(\nabla f)(a, b)$ introduced in Theorem 4.7.1—the vector of partial derivatives of f at (a, b) —is called the **gradient** of f at (a, b) . Observe that the directional derivative of f at (a, b) in the \mathbf{u} direction is equal to $\nabla f \cdot \mathbf{u} = |\nabla f| \cos \theta$, where θ is the angle between ∇f and \mathbf{u} . Since $\cos \theta$ is maximized when $\theta = 0$, we see that the **gradient of f at (a, b) is f 's direction of maximum increase at (a, b)** . Furthermore, the direction opposite to the gradient is the direction of maximum decrease, and f has zero derivative in any direction orthogonal to the gradient.

!!!

The function from \mathbb{R}^2 to \mathbb{R}^2 which maps (a, b) to $(\nabla f)(a, b)$ is called the gradient of f and is denoted ∇f . Thus the **gradient of f at a point** is a vector, while the **gradient of f** is a vector-valued function defined on \mathbb{R}^2 .

Example 4.7.1

Find $(\nabla f)(3, 4)$, where $f(x, y) = x^2 + xy + y^2$. Find all possible values of $D_{\mathbf{u}}f(3, 4)$, where \mathbf{u} is any unit vector.

Solution

We have $(\partial_x f)(x, y) = 2x + y$ and $(\partial_y f)(x, y) = 2y + x$, so

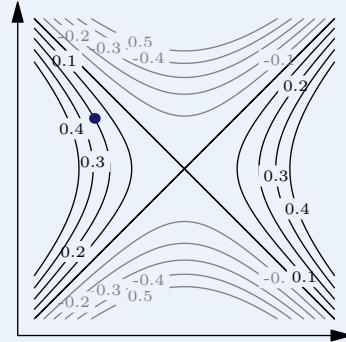
$$\nabla f = \langle 2x + y, 2y + x \rangle.$$

Therefore, the gradient evaluated at $(3, 4)$ is equal to $\langle 10, 11 \rangle$.

Since $D_{\mathbf{u}}f(3, 4) = \langle 10, 11 \rangle \cdot \mathbf{u} = \sqrt{221} \cos \theta$, where θ is the angle between $\langle 10, 11 \rangle$ and \mathbf{u} , we see that as \mathbf{u} ranges over the unit circle, the directional derivative of f in the \mathbf{u} direction ranges over $[-\sqrt{221}, \sqrt{221}]$.

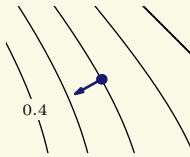
Example 4.7.2

Some of the level curves of a function $g(x, y)$ are shown. Sketch the direction of the gradient at the marked point.



Solution

The key idea is that a function neither increases or decreases along its level curve. Therefore, g has directional derivative equal to 0 in the direction of any line tangent the level curve passing through a given point. This means that the **gradient of g is orthogonal to g 's level curve** at any given point. So the gradient looks like the figure shown (zoomed in).



!!!

The gradient of a function from \mathbb{R}^3 to \mathbb{R} is likewise defined to be the vector of its partial derivatives: $\nabla f = \langle \partial_x f, \partial_y f, \partial_z f \rangle$. The formula $D_{\mathbf{u}} f = \nabla f \cdot \mathbf{u}$ holds for all differentiable functions $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ and $\mathbf{u} \in \mathbb{R}^3$.

Example 4.7.3

Find the equation of a plane tangent to the ellipsoid $x^2 + y^2 + 2z^2 = 4$ at the point $(1, 1, 1)$.

Solution

The ellipsoid is a level set of the function $g(x, y, z) = x^2 + y^2 + 2z^2$. The direction vectors \mathbf{u} contained in the plane of tangency at $(1, 1, 1)$ are characterized by the fact that g is unchanging in the \mathbf{u} direction at $(1, 1, 1)$. Since $D_{\mathbf{u}} g = \nabla g \cdot \mathbf{u}$, this means that* \mathbf{u} is orthogonal to ∇g . So we calculate

$$(\nabla g)(1, 1, 1) = \langle 2x, 2y, 4z \rangle|_{(x,y,z)=(1,1,1)} = \langle 2, 2, 4 \rangle,$$

and see that the desired plane is orthogonal to $\langle 2, 2, 4 \rangle$ and passes through $(1, 1, 1)$. So the equation of the plane is

$$\langle 2, 2, 4 \rangle \cdot \langle x - 1, y - 1, z - 1 \rangle = 0 \implies [x + y + 2z = 4].$$

Exercise 4.7.1

- (a) Confirm that the gradient of $x^2 + y^2$ is orthogonal to the level curves of $x^2 + y^2$ at each point.
- (b) Confirm that the gradient of $x^2 + y^2 + z^2$ is orthogonal to the level surfaces of $x^2 + y^2 + z^2$ at each point.

* All together: the gradient of a function at a point is orthogonal to the function's level set through that point

4.8 The multivariable chain rule

TMJ on the
chain rule

The basic idea of the chain rule is that when considering how $f(g(t))$ changes when we increase t by some small amount h , we can note that $g(t)$ changes by approximately $hg'(t)$, and that change in the input to f induces a change of

$$f(g(t+h)) - f(g(t)) \approx \begin{pmatrix} \text{change in input to } f \\ \overbrace{hg'(t)} \end{pmatrix} \begin{pmatrix} \text{sensitivity of } f \text{ to change in input} \\ \overbrace{f'(g(t))} \end{pmatrix}$$

in the value of $f(g(t))$. Thus $\frac{d}{dt}f(g(t)) = g'(t)f'(g(t))$.

The simplest multivariable generalization of this idea is to define a function from \mathbb{R} to \mathbb{R} by composing a differentiable function $\mathbf{r} : \mathbb{R} \rightarrow \mathbb{R}^2$ with a differentiable function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Let's look at an example.

Example 4.8.1

Suppose $f(x, y) = \sin xy \cos y$ and $\mathbf{r}(t) = (e^t, t^2)$. Find the derivative of $f \circ \mathbf{r}$.

Solution

We can calculate directly

$$(f \circ \mathbf{r})(t) = f(\mathbf{r}(t)) = \sin(t^2 e^t) \cos t^2.$$

So the desired derivative is

$$\begin{aligned} \cos(t^2 e^t) [t^2 e^t + 2te^t] \cos t^2 - \sin(t^2 e^t) 2t \sin t^2 \\ = t^2 e^t \cos(t^2 e^t) \cos t^2 + 2te^t \cos(t^2 e^t) \cos t^2 - 2t \sin(t^2 e^t) \sin t^2. \end{aligned}$$

In our solution to Example 4.8.1, we composed the given functions before differentiating. The **multivariable chain rule** is an alternate approach which describes the general relationship between $(f \circ \mathbf{r})'(t)$ and the derivatives of f and \mathbf{r} . Let's write $\mathbf{r}(t) = \langle r_1(t), r_2(t) \rangle$. When we change t by h , the value of $f(\mathbf{r}(t))$ changes as follows:

$$f \left(\begin{matrix} \text{changes by } hr'_1(t) & \text{changes by } hr'_2(t) \\ \overbrace{r_1(t)} & , \quad \overbrace{r_2(t)} \end{matrix} \right)$$

The change of $hr'_1(t)$ in the first argument induces a change of $hr'_1(t)(\partial_x f)(\mathbf{r}(t))$ in the value of f , while the change of $hr'_2(t)$ in the second argument induces a change of $hr'_2(t)(\partial_y f)(\mathbf{r}(t))$. By linear approximation, the overall change is the sum of these two changes.

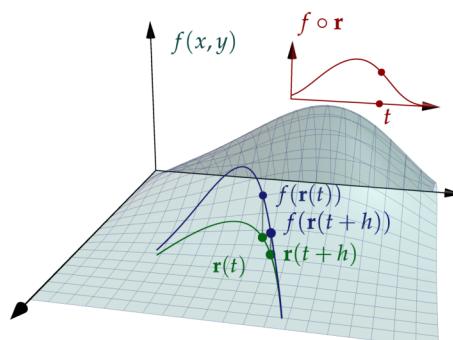


Figure 4.12 The composition $f \circ \mathbf{r}$ can be visualized using a slice of the graph of f along the path \mathbf{r} . The graph of the composition $f \circ \mathbf{r}$ is shown in the red inset figure.

* In the first step we're essentially swapping out f for L , the idea being that they're essentially the same when zoomed way in around $\mathbf{r}(t)$.

More precisely, if we define L to be the linear approximation of f at $\mathbf{r}(t)$, and then we have*

$$\begin{aligned}(f \circ \mathbf{r})'(t) &= \lim_{h \rightarrow 0} \frac{f(\mathbf{r}(t+h)) - f(\mathbf{r}(t))}{h} \\ &= \lim_{h \rightarrow 0} \frac{L(\mathbf{r}(t+h)) - L(\mathbf{r}(t+h)) + f(\mathbf{r}(t+h)) - f(\mathbf{r}(t))}{h} \\ &= \lim_{h \rightarrow 0} \left[\frac{L(\mathbf{r}(t+h)) - f(\mathbf{r}(t))}{h} + \frac{f(\mathbf{r}(t+h)) - L(\mathbf{r}(t+h))}{h} \right].\end{aligned}$$

The second term goes to 0 as $h \rightarrow 0$ since f is differentiable, so we're left with*

$$\lim_{h \rightarrow 0} \frac{L(\mathbf{r}(t+h)) - f(\mathbf{r}(t))}{h} = \lim_{h \rightarrow 0} \frac{(\partial_x f)(a, b)(r_1(t+h) - r_1(t)) + (\partial_y f)(a, b)(r_2(t+h) - r_2(t))}{h}.$$

Taking $h \rightarrow 0$ gives us a factor of $r'_1(t)$ in the first term and $r'_2(t)$ in the second term, yielding the following theorem.

Theorem 4.8.1: Multivariable chain rule

If $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $\mathbf{r} = \langle r_1, r_2 \rangle : \mathbb{R} \rightarrow \mathbb{R}^2$, are differentiable, then

$$(f \circ \mathbf{r})'(t) = (\partial_x f)(\mathbf{r}(t))r'_1(t) + (\partial_y f)(\mathbf{r}(t))r'_2(t) = (\nabla f)(\mathbf{r}(t)) \cdot \mathbf{r}'(t). \quad (4.8.1)$$

* $\frac{\partial f}{\partial x}$ means $\partial_x f$

The chain rule (4.8.1) can be written with the more suggestive notation*

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt},$$

where x and y represent r_1 and r_2 . Although this formula is more memorable, it does involve some abuse of notation: the symbols x and y are being used* as independent variables (in the partial derivative expressions) and as function names (in dx/dt and dy/dt). Also, on the left-hand side f is being treated as a function of a single variable; actually this instance of f is shorthand for the single-variable function $f \circ \mathbf{r}$.

Exercise 4.8.1

Verify that applying the multivariable chain rule to Example 4.8.1 gives the same result we found by calculating that derivative directly.

Exercise 4.8.2

Find the derivative with respect to t of the function $g(t) = t^t$ by writing the function as $f(x(t), y(t))$ where $f(x, y) = x^y$ and $x(t) = t$ and $y(t) = t$.

4.9 Optimization with Lagrange multipliers

Consider the function

$$f(x, y) = -x^2 - y^2 + x + \frac{2}{3}y + \frac{23}{36},$$

which we optimized over the square $[0, 1]^2$ in Example 4.5.2. In that case, we identified possible extreme values on the boundary of the square by doing a single-variable optimization along each edge of the square. But suppose that we want to find the maximum and minimum values of f over a disk D (see Figure 4.13)? Let's consider the disk D of radius $\frac{1}{2}$ centered at $(\frac{1}{2}, \frac{1}{2})$.

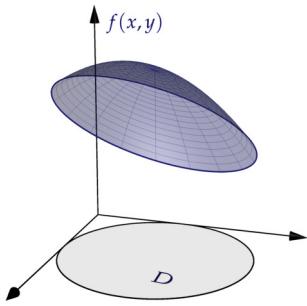


Figure 4.13 The graph of a function f defined on a disk D

We could take a similar approach to this problem. We can parameterize* the boundary of the disk as

$$\mathbf{r}(t) = \left\langle \frac{1}{2} + \frac{1}{2} \cos t, \frac{1}{2} + \frac{1}{2} \sin t \right\rangle, \quad 0 \leq t < 2\pi.$$

* To parameterize a curve means to find a path which traces it out

Then the single-variable function $t \mapsto f(\mathbf{r}(t))$ can be optimized over $[0, 2\pi]$ using the standard single-variable technique (as in Example 4.5.1).

The minimum is 0.5

However, this approach is limited because it requires a parameterization of the boundary of D , which is not always convenient. Suppose that ∂D is specified as a level set of some function $g : \mathbb{R}^2 \rightarrow \mathbb{R}$. For example, the circle in Figure 4.13 is a level set $\{(x, y) : g(x, y) = \frac{1}{2}\}$ of the function

$$g(x, y) = \left(x - \frac{1}{2} \right)^2 + \left(y - \frac{1}{2} \right)^2.$$

Let's derive an approach to finding the extreme values on the boundary which begins with the functions f (the *objective* function) and g (the *constraint* function).

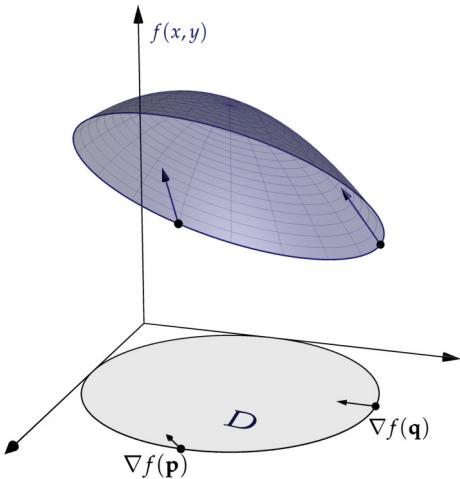


Figure 4.14 q is a boundary critical point and p is not

Imagine a bug moving around on the edge of the graph in Figure 4.13. How can it tell that it is *not* at a maximum or minimum? One approach is to calculate the gradient of f at its location. If the gradient of f is not orthogonal to ∂D , then the value of the function can be increased by sliding a bit in one direction* and can be decreased by sliding a bit in the opposite direction. So, for example, in Figure 4.14, a bug at the point p could increase the value of f at its location by moving slightly clockwise and decrease the value of f by moving slightly counterclockwise around ∂D .

* Specifically, the direction where ∇f is leaning (that is, the direction whose dot product with ∇f is positive)

Therefore, if the gradient of f at a point is *not* orthogonal to ∂D , then f does not have an extreme value there. So to find points where f might have an extreme value on ∂D , we can restrict our attention to the points where ∇f is orthogonal to ∂D .

We can simplify this idea further: recall that the gradient of g at each point is orthogonal to the level set of g passing through that point. It follows that if ∂D is a level set of g and $p \in \partial D$ is a point where f has an extreme value, then ∇g and ∇f are both orthogonal to ∂D . If $\nabla g \neq \mathbf{0}$, this means that they are parallel! By Observation 2.1.1, then, there exists

a scalar λ such that $\nabla f = \lambda \nabla g$.

Theorem 4.9.1: Method of Lagrange multipliers

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ are differentiable functions and $c \in \mathbb{R}$. If the restriction of f to the level set* $\{\mathbf{x} \in \mathbb{R}^n : g(\mathbf{x}) = c\}$ has a local extremum at $\mathbf{x} \in \mathbb{R}^n$, then either $\nabla g(\mathbf{x}) = \mathbf{0}$ or

$$\nabla f(\mathbf{x}) = \lambda \nabla g(\mathbf{x}).$$

* So far we've been considering the restriction of f to the boundary of a region D , but any level set of a differential function g will do

Theorem 4.9.1 implies that we can optimize $f(\mathbf{x})$ subject to the constraint $g(\mathbf{x}) = c$ by solving the system of **Lagrange equations**

!!!

$$\begin{cases} \nabla f(\mathbf{x}) = \lambda \nabla g(\mathbf{x}) \\ g(\mathbf{x}) = c, \end{cases}$$

with the only caveat that we have to watch out for the possibility that $\nabla g = \mathbf{0}$. Let's see how this works out for the example from the beginning of the section.

Example 4.9.1

Find the maximum and minimum values of

$$f(x, y) = -x^2 - y^2 + x + \frac{2}{3}y + \frac{23}{36}$$

over the disk of radius $\frac{1}{2}$ centered at $(\frac{1}{2}, \frac{1}{2})$.

Solution

The only interior critical point is $(\frac{1}{2}, \frac{1}{3})$, as in Example 4.5.2. To find boundary critical points, we use $g(x, y) = (x - \frac{1}{2})^2 + (y - \frac{1}{2})^2$ as our constraint function and set up the Lagrange equations

$$\begin{aligned} \partial_x f &= \lambda \partial_x g \\ \partial_y f &= \lambda \partial_y g \implies \end{aligned} \tag{4.9.1}$$

$$-2x + 1 = \lambda(2x - 1) \tag{4.9.1}$$

$$-2y + \frac{2}{3} = \lambda(2y - 1). \tag{4.9.2}$$

We're looking for pairs (x, y) which satisfy both of these equations **and** the equation

$$g(x, y) = (x - \frac{1}{2})^2 + (y - \frac{1}{2})^2 = \frac{1}{4} \tag{4.9.3}$$

since (x, y) must be on ∂D . So we have three equations and three variables: x , y , and λ . The first equation implies that either $\lambda = -1$ or $x = \frac{1}{2}$.

In the case $x = \frac{1}{2}$, we can use the constraint equation to conclude that $y = 0$ or $y = 1$. In either case, we can substitute into (4.9.2) to get a value for λ so that all three equations are satisfied. So, we have $(x, y) = (1/2, 0)$ and $(1/2, 1)$ as boundary critical points.

If $x \neq \frac{1}{2}$, then $\lambda = -1$. Substituting into (4.9.2) gives a contradiction, which means that we've already found all the boundary critical points.

Finally, we evaluate f at the interior critical point and the two boundary critical points:

(x, y)	$(\frac{1}{2}, 0)$	$(\frac{1}{2}, 1)$	$(\frac{1}{2}, \frac{1}{3})$
$f(x, y)$	$\frac{8}{9}$	$\frac{5}{9}$	1

So the maximum of f over D is 1, and the minimum is $\frac{5}{9}$.

The following example is a 3D application of Lagrange multipliers.

Example 4.9.2

Find the maximum possible volume of a box made with 72 square centimeters of cardboard and having sides and a bottom but no top.

Solution

Denote by x, y , and z the dimensions (in centimeters) of the cardboard. Then the amount of cardboard used is

$$g(x, y, z) = 2yz + 2xz + xy = 72,$$

while the objective function is the volume $f(x, y, z) = xyz$. Setting up the Lagrange equations, we get $yz = \lambda(2z + y)$, $xz = \lambda(2x + y)$, $xy = \lambda(2y + x)$, and $2yz + 2xz + xy = 72$, where the last one is the constraint equation. Multiplying the first two equations by x and y , respectively, and setting the resulting right-hand sides equal implies that either $\lambda = 0$ or $z = 0$ or $x = y$. Since $\lambda = 0$ or $z = 0$ clearly give zero volume (and thus not the maximum volume), it follows that $x = y$. Substituting y for x in the third equation gives*

$$x^2 = 4\lambda x \implies \lambda = \frac{x}{4}.$$

Substituting this into the second equation and simplifying, we get $x = 2z$. Finally, substituting into the constraint equation gives $z = \sqrt{6}$, which in turn implies $x = y = 2\sqrt{6}$.

* Once again, we can divide by x because we know that $x = 0$ wouldn't make sense for the maximum volume

Exercise 4.9.1

Find the set of all critical points of $f(x, y, z) = 3 - x^2 - 2y^2 - z^2$ subject to the constraint $2x + y + z = 2$.

Exercise 4.9.2

Find the points on the ellipse $\left(\frac{x-1}{2}\right)^2 + (y-2)^2 = 1$ which are nearest and farthest from the origin.
Hint: for the objective function, use *squared* distance rather than distance.

Exercise 4.9.3

Find the maximum value of $f(x, y) = y$ for any point (x, y) on the curve $y^2 - 4x^3 + 4x^4 = 0$ in two ways: (i) using Lagrange multipliers, and (ii) writing the upper half of the given curve as the graph of a function and maximizing that function using standard single-variable techniques.

Exercise 4.9.4

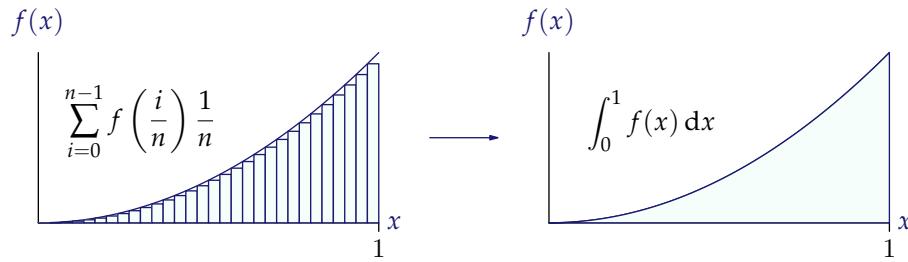
Suppose that we want to maximize a differentiable function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ subject to the constraint $g(x, y) = c$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is differentiable and $c \in \mathbb{R}$. Consider the *Lagrangian* function

$$\mathcal{L}(x, y, \lambda) = f(x, y) - \lambda g(x, y).$$

Show that the equation $\nabla \mathcal{L} = 0$ is equivalent to the Lagrange equations.

5 Multivariable Integration

To find the area under the graph of a continuous function f over the unit interval $[0, 1]$, we first approximate the area by splitting $[0, 1]$ into many short intervals and sum up the areas of rectangles approximating the area under the graph over each short interval:



This approximation converges to the actual area under the graph as $n \rightarrow \infty$.

In this section we will work out how to generalize this concept to integrate of functions of multiple variables over regions in \mathbb{R}^2 or \mathbb{R}^3 .*

5.1 Double integration

* See Appendix A.4.4 for a more in-depth discussion of integration

on double integrals

* Recall our convention that 1D volume is length and 2D volume is area

* It doesn't ultimately matter where we evaluate the function, since the piece is very small and the function is continuous. If we want an overestimate/underestimate of the integral, we can use the maximum/minimum of the function on each piece

We can state the definition of the integral, described above, more informally and generally: split the region of integration into many tiny pieces, multiply the volume* of each piece by the value of the function at some point on that piece*, and add up the results. If we take the number of pieces to ∞ and the piece size to zero, then this sum should converge to a number, and if it does then we declare that number to be the value of the integral.

Stated at this level of generality, the idea of the integral definition applies to a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ over a bounded region $D \subset \mathbb{R}^2$. See Figure 5.1 for an illustration, and see Appendix A.3.5 for a more general definition.

Definition 5.1.1: Integral over a 2D region

Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a continuous function, and that D is a bounded region in \mathbb{R}^2 such that ∂D has zero area. Then the limit*

$$\lim_{n \rightarrow \infty} \sum_{(i,j) : \left(\frac{i}{n}, \frac{j}{n}\right) \in D} f\left(\frac{i}{n}, \frac{j}{n}\right) \overbrace{\frac{1}{n^2}}^{\Delta A} \quad (5.1.1)$$

exists. We define the integral of f over D , denoted $\iint_D f dA$ or $\int_D f dA$, to be the value of this limit.

* The notation means that the sum includes one term for each integer pair (i, j) such that $(i/n, j/n)$ is in the region D

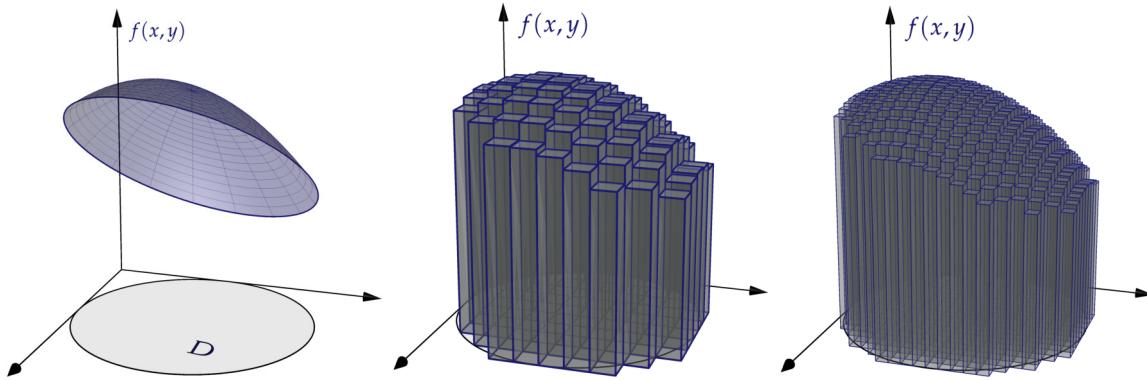


Figure 5.1 The integral of f over a disk D , defined as a limit of sums of volumes of narrow boxes

The sums appearing on the right-hand side of (5.1.1) are called **Riemann sums**. As in the single-variable case, the Riemann-sum definition is not generally practical for exact evaluation of integrals. The fundamental theorem of calculus* is the primary tool for evaluating integrals in single-variable calculus, and fortunately we can bootstrap our way up from 1D integration to 2D integration by applying our primary strategy for tackling higher dimensional problems: slicing. Let's start by considering integrals over rectangular regions D .

Example 5.1.1

Find the integral of $f(x,y) = y \sin(\pi xy)$ over the square $[0,1]^2$.

* The fundamental theorem of calculus says that the integral of f from a to b is equal to $F(b) - F(a)$ where F is an antiderivative of f

Solution

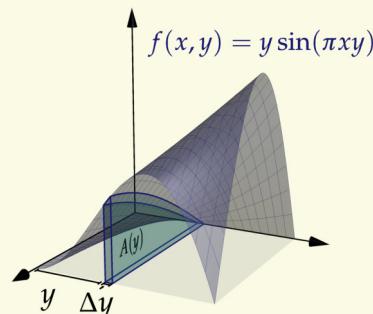
Let's slice up the desired solid using many ' $y = \text{constant}$ ' cuts, producing many thin slices like the one shown. The volume of one of these slices, situated at a particular y -value, is given by* the thickness Δy times the area $A(y)$ under the graph of the single-variable function $x \mapsto f(x,y)$. So we can use the fundamental theorem of calculus to compute

$$\begin{aligned} A(y) &= \int_0^1 y \sin(\pi xy) dx = -\frac{\cos(\pi xy)}{\pi} \Big|_0^1 \\ &= \frac{1 - \cos \pi y}{\pi}. \end{aligned}$$

Once we have each area $A(y)\Delta y$, we can add them all up and take $\Delta y \rightarrow 0$ (as the number of slices goes to ∞) to find that the desired volume is

$$\sum_{\text{all slices}} A(y)\Delta y \rightarrow \int_0^1 A(y) dy.$$

We can again evaluate this integral using the fundamental theorem to get a final answer of $\frac{1}{\pi}$.



* ...ignoring an error, having to do with the top of the slice not being flat—this error tends to zero as the number of slices tends to infinity

For some confidence that our answer is reasonable, we can calculate a Riemann sum for this integrand.

We can express this process more succinctly as

$$\int_0^1 \int_0^1 y \sin(\pi xy) dx dy = \int_0^1 \frac{1 - \cos \pi y}{\pi} dy = \frac{1}{\pi}. \quad (5.1.2)$$

The first expression in (5.1.2) is called an **iterated integral**, since it expresses an integral over a 2D region in terms of two successive single-variable integrals.

Let's see how this works over a non-rectangular region.

Example 5.1.2

Find the integral over the triangle T with vertices $(0, 0)$, $(2, 0)$, and $(0, 3)$ of the function $f(x, y) = x^2y$, by first finding the area under each ' $y = \text{constant}$ ' slice.

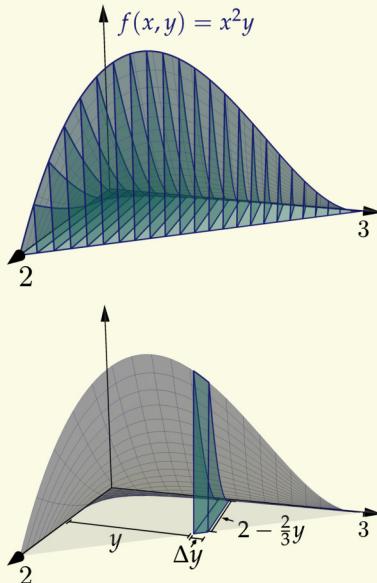
Solution

As in the previous example, we slice up the desired volume making many ' $y = \text{constant}$ ' cuts of thickness Δy , yielding thin slices such that each one has volume (very close to) $A(y)\Delta y$, where y is the slice's signed distance from the xz -plane and $A(y)$ is the area of the cross-section (see figure). Since this cross section is an area under a curve, we can find it by integrating $x \mapsto f(x, y)$ over the set of relevant x -values.*

$$A(y) = \int_0^{2-\frac{2}{3}y} f(x, y) dx.$$

Thus $A(y) = \frac{1}{3} (2 - \frac{2}{3}y)^3 y$. Finally, adding up all these areas and taking $\Delta y \rightarrow 0$ gives the result

$$\int_0^3 A(y) dy = \int_0^3 \left(-\frac{8}{81} y^4 + \frac{8}{9} y^3 - \frac{8}{3} y^2 + \frac{8}{3} y \right) dy = \boxed{\frac{6}{5}}.$$



Let's summarize what we figured out in Example 5.1.2.

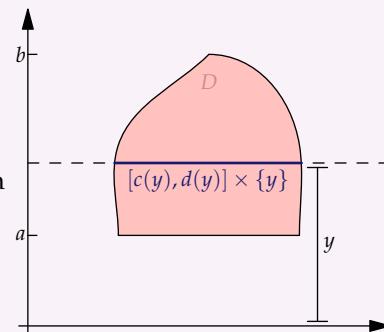
Theorem 5.1.1: Iterated integrals for two-variable functions

Suppose that

- D is a region in \mathbb{R}^2 ,
- $f : D \rightarrow \mathbb{R}$ is a continuous function, and
- for all $y \in \mathbb{R}$, the intersection of D and horizontal line through $(0, y)$ is a segment $[c(y), d(y)] \times \{y\}$.

Then

$$\iint_D f dA = \int_a^b \int_{c(y)}^{d(y)} f(x, y) dx dy.$$



* We can find the formula $2 - \frac{2}{3}y$ by writing an equation for the line connecting $(2, 0)$ to $(0, 3)$ and solving for x

In light of Theorem 5.1.1, we sometimes write the area differential* as $dA = dx dy$. We can describe the procedure in Theorem 5.1.1 less formally:

Observation 5.1.1: Limits of integration over a 2D region

To set up an iterated integral to evaluate $\iint_D f dA$ (where f is continuous and D is a region such that the intersection of every horizontal line with D is a segment):*

1. Find the least and greatest y values for any point in D . These are your **outer limits** of integration.
2. For each fixed horizontal line which intersects D , identify the least and greatest values of x for any point which is in D *and on that line*, expressed in terms of the vertical position y of the line. These are the **inner limits** of integration, and they may depend on y .

* Think of dA merely as a reminder that the positive quantity ΔA involved in the corresponding Riemann sums represents an area

* You should be prepared for these steps to spawn geometric sub-problems that you might need to solve on the side

The role of x and y in Observation 5.1.1 can be reversed (in which case we have vertical rather than horizontal lines in Step 2). The following exercise shows how this can be useful.

Exercise 5.1.1

Find

$$\int_0^{1/2} \int_{2y}^1 4e^{x^2} dx dy$$

by first rewriting it as an integral over a 2D region and then reversing the order of integration.

5.2 Triple integration



We interpret the integral of a single-variable function as an area and the integral of a two-variable function as a volume. So how should we interpret the integral of a function of *three* variables over a region D in \mathbb{R}^3 ? *Four-dimensional volume* is a reasonable answer, but of course this is unsatisfactory from a visualization point of view, since we don't have access to four spatial dimensions with which to visualize.

Therefore, let's consider a physics interpretation of integration which permits a visualization *not* involving the graph of the function being integrated.

Example 5.2.1

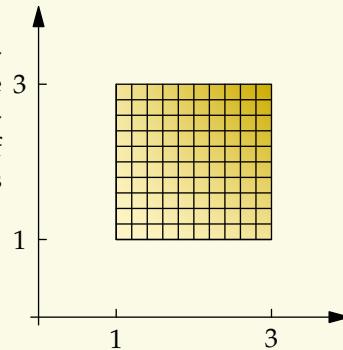
Consider a square plate occupying the square $[1, 3]^2$ whose density at each point is $\sigma(x, y) = xy$ kilograms per square unit.* Find the mass of the plate.

* See the figure in the solution, where darker color indicates a denser portion of the plate

Solution

Let's imagine physically cutting the plate into small squares, computing the mass of each one, and adding up the resulting masses. The mass of a small plate of area $\Delta x \Delta y$ containing the point (x, y) is approximately the area density times the area: $\sigma(x, y) \Delta x \Delta y$. The sum of these masses is a Riemann sum (see Definition 5.1.1) which converges as the number of small squares goes to ∞ to the integral

$$\int_1^3 \int_1^3 xy \, dx \, dy = \boxed{16} \text{ kilograms.}$$



Let's do a three-dimensional example.

Example 5.2.2

Consider a cubical block occupying $D = [1, 2]^3$ whose density at each point is $\delta(x, y, z) = x^2 + y^2 + z^2$ kilograms per cubic unit. Find the mass of the block.

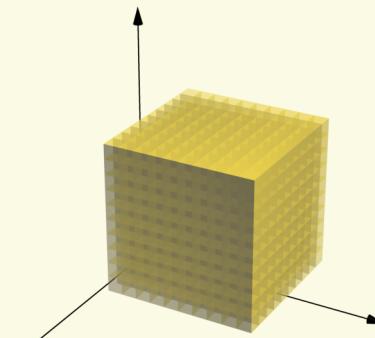
Solution

We cut the cube into n^3 small cubes, where n is a large integer. The mass of one of these cubes with bottom, back* corner (x, y, z) is approximately equal to the product of its volume $\frac{1}{n^3}$ and the approximate density $\delta(x, y, z)$ throughout the small cube. So the approximate volume is

$$\sum_{\text{all cubes}} \delta(x, y, z) \frac{1}{n^3}.$$

Intuitively, this sum should converge to a limit as $n \rightarrow \infty$, and if so, then we should define the limiting value to be the integral of δ over D . Let's state this idea for any continuous function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$: we define the integral of f over D by

$$\iiint_D f(x, y, z) \, dV = \lim_{n \rightarrow \infty} \sum_{(i, j, k) : \left(\frac{i}{n}, \frac{j}{n}, \frac{k}{n}\right) \in D} f\left(\frac{i}{n}, \frac{j}{n}, \frac{k}{n}\right) \frac{1}{n^3}.$$



We can calculate the integral by slicing up the region of integration into thin slabs along 'z = constant' slices, and then performing double integrals to find the area of each slab. This works the same as double iterated integration, but with one extra step. Rather than writing Δz and then taking a limit

* Any corner, or indeed any point in the cube, would give the same result—we choose the bottom, back corner (the point nearest the origin) for concreteness

* This is a general theme: we contract the following two steps into a single step (by writing dz instead of Δz from the outset): (i) reason about sums involving a small but positive quantity Δz , and (ii) replace the sum with an integral over the relevant z values and replace Δz with dz

to turn Δz into dz , we'll skip to the limit and work directly with dz^*

$$\begin{aligned} \text{mass} &= \int_1^2 \overbrace{\int_1^2 \int_1^2 (x^2 + y^2 + z^2) dx dy dz}^{\text{mass of slice from } z \text{ to } z + dz} \\ &= \int_1^2 \int_1^2 \left(y^2 + z^2 + \frac{7}{3} \right) dy dz \\ &= \int_1^2 \left(z^2 + \frac{14}{3} \right) dz \\ &= \boxed{7} \text{ kilograms.} \end{aligned}$$

The following theorem summarizes the idea of integrating in 3D by breaking down the 3D region of integration into 2D slices.

Theorem 5.2.1: Iterated integrals for three-variable functions

Suppose f is a continuous function over a region D which is bounded between the planes $z = a$ and $z = b$. For each $z \in (a, b)$, define $D_z \subset \mathbb{R}^2$ to be the region*

$$D_z = \{(x, y) \in \mathbb{R}^2 : (x, y, z) \in D\}.$$

Then

$$\iiint_D f dV = \int_a^b \left[\iint_{D_z} f(x, y, z) dx dy \right] dz,$$

* D_z is the region obtained by intersecting D with the plane which is z units from the xy -plane and then dropping off the third coordinate

Let's break this theorem down into a simple algorithm (the following observation is the 3D analogue of Observation 5.1.1):

Observation 5.2.1: Limits of integration over a 3D region

To set up an iterated integral to evaluate $\iiint_D f dV$ (where f is continuous and D is a region which intersects every line parallel to the x -axis in an interval):

1. Find the least and greatest z values for any point in D . These are your **outer limits** of integration.
2. For each fixed ' $z = \text{constant}$ ' plane which intersects D , identify the least and greatest values of y for any point which is in D and on that plane, expressed in terms of the vertical position z of the plane. These are the **middle limits** of integration, and they may depend on z .
3. For each line of the form ' $z = \text{constant}$ and $y = \text{constant}$ ' which intersects D , find the least and greatest values of x for any point which is in D and on that line. These are your **inner limits** of integration, and they may depend on both z and y .

Example 5.2.3

Find the volume of the tetrahedron with vertices $(0, 0, 0)$, $(2, 0, 0)$, $(0, 3, 0)$, and $(0, 0, 4)$ using a triple integral.

Solution

The volume of a region is equal to the integral of the constant function 1 over that region:

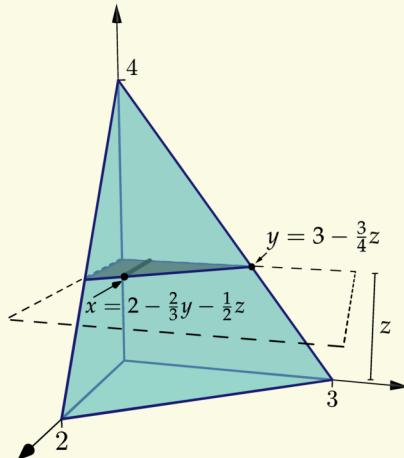
$$\text{volume}(D) = \iiint_D 1 \, dV,$$

because $\iiint_D 1 \, dV$ is equal to the mass of a solid occupying the region D and having density 1 at every point. But if a solid has a constant mass density of 1, then its mass is equal to its volume.*

So we set up our iterated integral: the least and greatest values of z are 0 and 4, so those are our outer limits. For a fixed value of z , the least and greatest values of y for a point in D are 0 and $3 - \frac{3}{4}z$, respectively. Finally, for fixed y and z , the least and greatest values of x for a point in D are 0 and the point on the plane $6x + 4y + 3z = 12$ with the given values of y and z (see figure).

So we get

$$\begin{aligned} \text{volume}(D) &= \int_0^4 \int_0^{3-\frac{3}{4}z} \int_0^{2-\frac{2}{3}y-\frac{1}{2}z} 1 \, dx \, dy \, dz \\ &= \int_0^4 \int_0^{3-\frac{3}{4}z} \left(2 - \frac{2}{3}y - \frac{1}{2}z \right) \, dy \, dz \\ &= \int_0^4 \frac{3}{16} (z-4)^2 \, dz \\ &= \boxed{4}. \end{aligned}$$



There is nothing special about the order $dx \, dy \, dz$ —any way of slicing up the region gives the same result. The region of integration for the following exercise can be sliced up six different ways, and you can check that the integral is the same with respect to all the different orders of integration.

Exercise 5.2.1

Write the iterated integral

$$\int_0^1 \int_{\sqrt{x}}^1 \int_0^{1-y} f(x, y, z) \, dz \, dy \, dx.$$

as an integral over a 3D region. Then sketch that region and use your figure to rewrite the integral in five other ways, using the five other permutations of (x, y, z) .

5.3 Polar, cylindrical, and spherical integration

Some regions in \mathbb{R}^2 are more conveniently described in polar coordinates than rectangular coordinates. If we are integrating a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ over such a region, it is helpful to work directly in polar coordinates. Let's do an example.

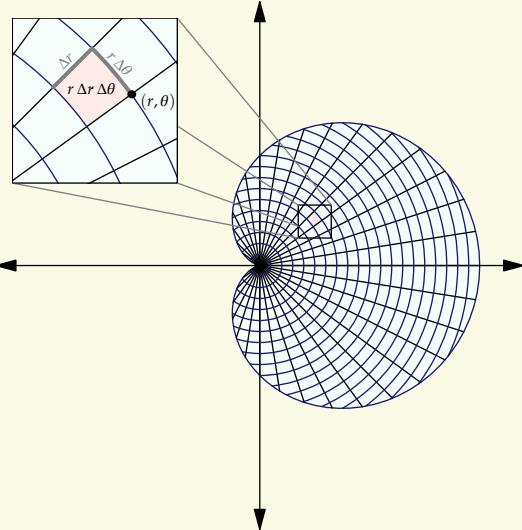
Example 5.3.1

- (i) Find the area of the region D enclosed by the solution set of the polar coordinate equation $r = 1 + \cos \theta$.
- (ii) Integrate the function $f(x, y) = x + y$ over D .

Solution

(i) Let's slice D into small pieces using equally spaced cuts along rays and circles of the form $r = \text{constant}$ and $\theta = \text{constant}$, as shown. This decomposes D into a set of **coordinate patches**. The figure suggests that these pieces farther away from the origin are larger than the ones that are close to the origin, which leads us to investigate the area of each patch.

To find the area of the set of points with radial polar coordinate between r and $r + \Delta r$ and angular polar coordinate between θ and $\theta + \Delta\theta$, we note that this region is approximately a rectangle. The straight side length is Δr , and the curvy side length is $r\Delta\theta$, because the perimeter of the circle of radius r is $2\pi r$, and the angle represents $\frac{\theta}{2\pi}$ of the whole circle. So the area is approximately $r\Delta r\Delta\theta$.



Now, for fixed θ , we can add up all the coordinate patches between θ and $\theta + \Delta\theta$, and this sum of areas is approximately equal to the product of $\Delta\theta$ and the integral

$$\int_0^{1+\cos\theta} r \, dr.$$

Adding up these areas over all θ from 0 to 2π , we get

$$\int_0^{2\pi} \int_0^{1+\cos\theta} r \, dr \, d\theta = \int_0^{2\pi} \frac{1}{2}(1 + \cos\theta)^2 \, d\theta = \boxed{\frac{3\pi}{2}}.$$

(ii) We can find this integral using the same procedure as above, except that at the step where we calculate the area of a patch, we also need to multiply it by the value of the function at some point in the patch. Since our function is defined in terms of x and y , we need to substitute $x = r \cos \theta$ and $y = r \sin \theta$ to discover the value of f at the point whose polar coordinates are (r, θ) . So we get

$$\int_0^{2\pi} \int_0^{1+\cos\theta} (r \cos \theta + r \sin \theta) r \, dr \, d\theta = \boxed{\frac{5\pi}{4}}.$$

This one is tedious if done by hand, so we use computer algebra assistance

We can see from this example that the ideas for setting up an iterated polar integral are similar to those for rectangular integration:

Observation 5.3.1: Iterated polar integration

To write the integral of a function f on a region D in \mathbb{R}^2 as a double iterated integral, we may

- (i) find the least and greatest values of θ for any point in the region of integration,
- (ii) for each fixed value of θ , find the least and greatest values of r for any point which is in D and on the ray of angle θ ,
- (iii) include the area differential* $dA = r dr d\theta$, and
- (iv) plug $x = r \cos \theta$ and $y = r \sin \theta$ into f , that your integrand varies appropriately as r and θ vary

* Don't forget the extra factor of r in the polar area differential!

This same basic idea can be carried out in cylindrical and spherical coordinates. The ingredients we need are (i) the volume differential dV expressed in terms of cylindrical and spherical coordinates, and (ii) the formulas for x, y, z in terms of r, θ, z and in terms of ρ, θ, ϕ . This information is listed in Appendix A.2.2. The only surprising entry in the tables of that appendix is the spherical coordinate volume differential $dV = \rho^2 \sin \phi d\rho d\phi d\theta$.

Example 5.3.2: Spherical coordinate volume differential

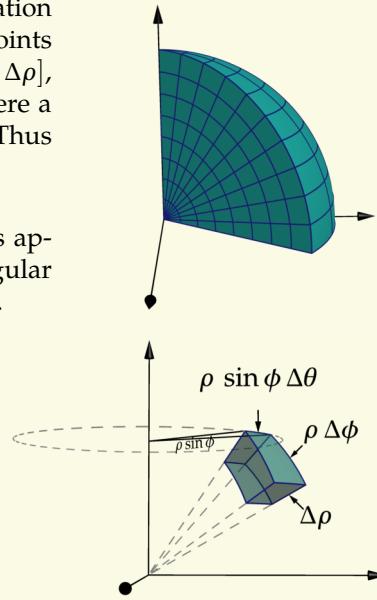
Explain why the volume differential in spherical coordinates is $dV = \rho^2 \sin \phi d\rho d\phi d\theta$.

Solution

The volume differential arises from slicing up the region of integration into small coordinate “patches”, each of which consists of those points whose three spherical coordinates lie in the intervals* $[\rho, \rho + \Delta\rho]$, $[\theta, \theta + \Delta\theta]$, and $[\phi, \phi + \Delta\phi]$, respectively (see the top figure, where a wedge has been decomposed into spherical coordinate patches). Thus we must calculate the approximate volume of one such patch.

When $\Delta\rho$, $\Delta\theta$, and $\Delta\phi$ are all very small, the coordinate patch is approximately a rectangular prism. The dimensions of this rectangular prism, as marked in the lower figure, are $\Delta\rho$, $\rho \Delta\phi$, and $\rho \sin \phi \Delta\theta$.

To see why the top edge length is $\rho \sin \phi \Delta\theta$, note that the dashed circle in the figure has radius $\rho \sin \phi$, since the cylindrical radial coordinate r satisfies the equation $r = \rho \sin \phi$. Thus the volume of the patch is approximately $\rho^2 \sin \phi \Delta\rho \Delta\phi \Delta\theta$.



* Here (ρ, θ, ϕ) are the spherical coordinates of one of the corners of the patch

Example 5.3.3

Consider a solid whose density at each point (x, y, z) is $\delta(x, y, z) = \frac{1}{x^2+y^2+z^2}$ and which occupies the region enclosed by the cone $z = \sqrt{x^2 + y^2}$ and the plane $z = 1$. Find the mass of the solid.

Solution

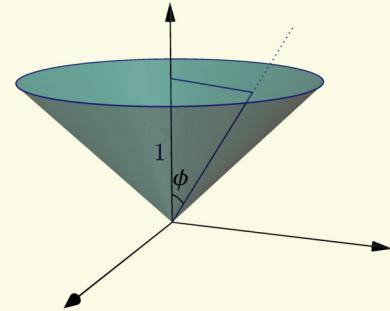
Let's set up the iterated integral with the order $d\rho d\phi d\theta$. The solid has points with θ values as small as 0 and as large as 2π , so the outer limits will be 0 and 2π .

For any given value of θ , there are points with that θ value whose ϕ value is as small as zero (for the points on the positive z -axis) and as large as $\frac{\pi}{4}$ (for the points on the cone $z = \sqrt{x^2 + y^2}$). So the middle limits are 0 and $\frac{\pi}{4}$.

Finally, for any given ϕ and θ , the solid contains points with z as small as 0 and as large as $\frac{1}{\cos \phi}$ (by right-triangle trigonometry; see figure).

For the integrand, we should substitute the spherical coordinate formulas for x , y , and z . However, we know that it will simplify to $\frac{1}{\rho^2}$, since $\rho^2 = x^2 + y^2 + z^2$. So we get

$$\int_0^{2\pi} \int_0^{\pi/4} \int_0^{\sec \phi} \rho^{-2} \rho^2 \sin \phi d\rho d\phi d\theta = \int_0^{2\pi} \int_0^{\pi/4} \sec \phi \sin \phi d\phi d\theta = (2\pi)(\frac{1}{2} \ln 2) = [\pi \ln 2].$$



Help with calculating this integral

Exercise 5.3.1

What proportion of the volume of the unit sphere lies above the plane $z = \frac{1}{2}$?

Exercise 5.3.2

Find

$$\int_{-2}^2 \int_{-\sqrt{4-x^2}}^{\sqrt{4-x^2}} \int_{\sqrt{x^2+y^2}}^2 (x^2 + y^2) dz dy dx$$

by writing it as integral over a 3D region and then rewriting that integral using cylindrical coordinates.

5.4 Integration in custom coordinates

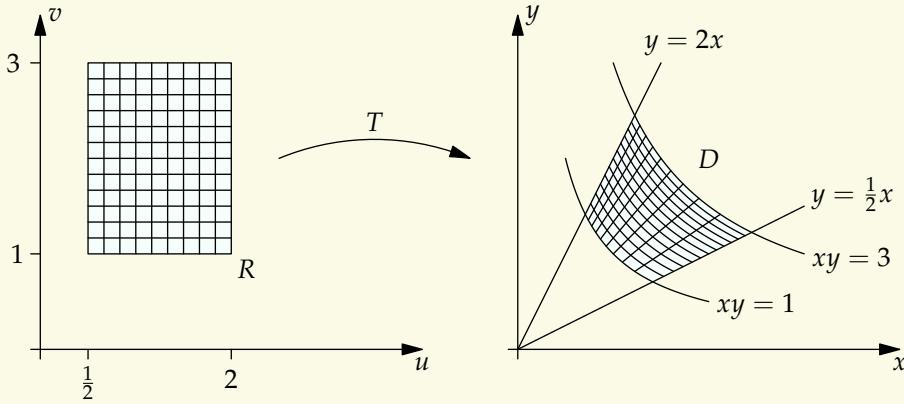
on change of variables

Sometimes we want to integrate a function over a region which is not conveniently described using any of the standard coordinate systems. In this section we will develop a general program for integrating with respect to a custom-tailored coordinate system.

Example 5.4.1

Find $\iint_D y^2 dA$, where D is the region in the first quadrant bounded by the lines through the origin of slope $\frac{1}{2}$ and 2, as well as the hyperbolas $xy = 1$ and $xy = 3$.

Solution



Looking at the shape of the region (in the axes on the right above), it makes sense to slice it up along lines of the form $y = ux$ where u ranges over equally spaced values between $\frac{1}{2}$ and 2, and along hyperbolas of the form $xy = v$ where v ranges over $[1, 3]$.

Note that each point (x, y) in D can be identified by its u and v values.* In other words, u and v provide a coordinate system for the first quadrant. We can represent the relationship between (u, v) and (x, y) as a transformation T that maps each (u, v) pair to its corresponding (x, y) pair, as shown in the figure. If we want to find a formula for this map, we can solve the system $y = ux$ and $xy = v$ for x and y to find that $y = \sqrt{uv}$ and $x = \sqrt{v/u}$.

Now, to integrate $f(x, y) = y^2$ over D , we want to find the area of each of the small pieces we sliced D into, multiply each such area by the value of f somewhere on the piece, and sum the resulting products. Each patch is the image under T of a rectangle of the form $[u, u + \Delta u] \times [v, v + \Delta v]$. The area of this rectangle is $\Delta u \Delta v$, and the transformation distorts its area by an amount that we can approximate by treating the transformation as linear around (u, v) and using the fact that area distortion is measured by the determinant. Writing $T(u, v) = (\sqrt{v/u}, \sqrt{uv}) = (g(u, v), h(u, v))$, we linearly approximate g at a point (u, v) as the function L_g defined by

$$L_g(\tilde{u}, \tilde{v}) = g(u, v) + \partial_u g(u, v)(\tilde{u} - u) + \partial_v g(u, v)(\tilde{v} - v),$$

and similarly for the linear approximation L_h of h . Thus for small Δu and Δv , T behaves like

$$T(u + \Delta u, v + \Delta v) \approx T(u, v) + (\partial_u g(u, v)\Delta u + \partial_v g(u, v)\Delta v, \partial_u h(u, v)\Delta u + \partial_v h(u, v)\Delta v). \quad (5.4.1)$$

So the area of the image of $[u, u + \Delta u] \times [v, v + \Delta v]$ under T is approximately*

$$\left| \det \begin{bmatrix} \partial_u g & \partial_v g \\ \partial_u h & \partial_v h \end{bmatrix} \right| \Delta u \Delta v = \left| \det \begin{bmatrix} -\frac{1}{2}\sqrt{v}u^{-3/2} & \frac{1}{2\sqrt{uv}} \\ \frac{\sqrt{v}}{2\sqrt{u}} & \frac{\sqrt{u}}{2\sqrt{v}} \end{bmatrix} \right| \Delta u \Delta v = \frac{1}{2u} \Delta u \Delta v.$$

So the contribution of each piece to the value of the integral is $uv \frac{1}{2u} \Delta u \Delta v$. Then summing and taking $(\Delta u, \Delta v) \rightarrow (0, 0)$ yields $\int_1^3 \int_{1/2}^2 uv \frac{1}{2u} du dv = [3]$.

* These specify, respectively, which line through the origin and which "xy = constant" hyperbola the point is on

* The four entries of the matrix below come from the coefficients of Δu and Δv in (5.4.1), or equivalently from the coefficients of \tilde{u} and \tilde{v} in L_g and L_h

The matrix $\begin{bmatrix} \frac{\partial_u g}{\partial_u h} & \frac{\partial_v g}{\partial_v h} \end{bmatrix}$ is called the *Jacobian matrix*, and its determinant is called the *Jacobian determinant*.

Often we just say “Jacobian”, relying on context to distinguish. It can be written as $\left| \frac{\partial(x,y)}{\partial(u,v)} \right|$ for short.*

The following theorem summarizes the technique we developed in Example 5.4.1.

Theorem 5.4.1: Multivariable change of variables

Suppose that $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is a differentiable transformation that maps a region R one-to-one onto a region D . Then for any continuous function f , we have

$$\iint_D f(x,y) dx dy = \iint_R f(T^{-1}(x,y)) \left| \frac{\partial(x,y)}{\partial(u,v)} \right| du dv.$$

The biggest challenge in custom coordinate problems typically lies in choosing suitable coordinates. The strategy of Example 5.4.1 is pretty broadly useful: try to express all boundary arcs of D as level sets of at most two functions of x and y . These two functions* are good choices for the coordinates u and v .

Exercise 5.4.1

Find the integral of $\frac{(x-y)^2}{x+y+2}$ over the square whose vertices are the four points of intersection between the axes and the unit circle.

* Note that in this new notation, the symbols g and h are replaced by x and y . Thus we are abusing notation by regarding x and y as functions of u and v , even though they also represent independent variables

* If all the boundary arcs are level sets of a single function u , then you'll have some flexibility in choosing the other; often $v = x$ or $v = y$ works

Exercise 5.4.2

Use the change of variables $x = u^2 - v^2$, $y = 2uv$ to evaluate the integral $\iint_R y dA$, where R is the region above the x -axis bounded by the parabolas $y^2 = 4 - 4x$ and $y^2 = 4 + 4x$.

5.5 Applications of integration*

5.5.1 AVERAGE VALUE

The most basic example of averaging is summing a list of numbers and dividing by the number of entries in the list. While this concept is usually used without reference to functions, one way to express it is to say that the average of a function $f : \{1, 2, \dots, n\} \rightarrow \mathbb{R}$ is defined to be

$$\text{average}(f) = \frac{f(1) + f(2) + \dots + f(n)}{n}. \quad (5.5.1)$$

Now suppose that the domain of f is a region $R \subset \mathbb{R}^n$. The most natural way to adapt (5.5.1) to this setting is to replace the sum in the numerator with an integral, and replace the “size” of the domain n in the denominator with the volume of R . We will state the definition for a function defined on a region in \mathbb{R}^3 :

Definition 5.5.1: Average value of a function over a solid

If $R \subset \mathbb{R}^3$ is a region and $f : R \rightarrow \mathbb{R}$, then the **average value** of f is

$$\text{average}(f) = \frac{\iiint_R f dV}{\iiint_R 1 dV}.$$

Let's explain the connection between Definition 5.5.1 and (5.5.1) more carefully: intuitively, the average value of a continuously varying function $f : R \rightarrow \mathbb{R}$ should be very close to the result of overlaying R with a fine grid, evaluating the function at all the corners of the cells of that grid, and averaging the results. By (5.5.1), this latter average is equal to the quotient* of (1) the sum of all the function's values at the grid points and (2) the number of grid points. If we multiply numerator and denominator of this quotient by the volume of a single small grid cell, then the numerator becomes a Riemann sum approximating $\iiint_R f dV$, and the denominator becomes a Riemann sum approximating the volume $\iiint_R 1 dV$ of R . Therefore, the limit of this quotient as the mesh of the grid tends to zero is equal to the formula given in Definition 5.5.1.

Exercise 5.5.1

- (a) Find the average squared distance from the origin to a point in the unit sphere.
- (b) Find the average distance from the origin to a point in the unit sphere.
- (c) Is the square of the average distance equal to the average of the squared distance?

If we want some of the numbers being averaged to count more than others, we can replace (5.5.1) with a **weighted average**. If $f : \{1, 2, \dots, n\} \rightarrow \mathbb{R}$ and $w : \{1, 2, \dots, n\} \rightarrow \mathbb{R}$, then the w -weighted average of f is defined to be

$$\text{average}_w(f) = \frac{f(1)w(1) + f(2)w(2) + \dots + f(n)w(n)}{w(1) + w(2) + \dots + w(n)}. \quad (5.5.2)$$

For example, if your scores on 3 exams are $f(1) = 90, f(2) = 90, f(3) = 100$ with weights $w(1) = 30, w(2) = 30, w(3) = 40$, then your weighted average is $\frac{90 \cdot 30 + 90 \cdot 30 + 100 \cdot 40}{30 + 30 + 40} = 94$. We can adapt (5.5.2) for continuously varying functions over a domain in \mathbb{R}^3 as we did in Definition 5.5.1.

Definition 5.5.2: Weighted average of a function over a solid

If $R \subset \mathbb{R}^3$ is a region, and that $f : R \rightarrow \mathbb{R}$ and $w : R \rightarrow [0, \infty)$ are functions. Then the w -**weighted average** of f is

$$\text{average}_w(f) = \frac{\iiint_R fw dV}{\iiint_R w dV}.$$

Exercise 5.5.2

Define $w : [0, 1]^3 \rightarrow \mathbb{R}$ via $w(x, y, z) = x + y + z$. (a) Find the w -weighted average of the constant function $f(x, y, z) = 1$ on $[0, 1]^3$. (b) Find the w -weighted average of $f(x, y, z) = x$ on $[0, 1]^3$.

* We're sticking with verbal descriptions of the quantities involved here to avoid the somewhat unwieldy notation associated with Riemann sums

Exercise 5.5.3

Identify the changes to the paragraph following Definition 5.5.1 necessary to obtain the corresponding explanation of Definition 5.5.2.

5.5.2 CENTER OF MASS

The center of mass of a physical rod is the point along the length of the rod where it balances. The center of mass of a rod with constant mass density is the point halfway between the rod's ends. But what if the rod has a varying mass density?



Figure 5.2 A rod with mass density indicated by color

Clearly the rod in Figure 5.2 will tip to the right if we support it in the middle. To find its center of mass, recall from physics that *torque** applied by a force to an object rotating about a fixed pivot is equal to the product of the force applied* and the distance to the pivot. The rod is balanced if the torque applied by gravity on the two sides is equal. If the rod is situated along the interval $[a, b]$, has density $\delta : [a, b] \rightarrow \mathbb{R}$, and is supported at the point x_0 , then the clockwise torque is equal to

$$\text{gravitational acceleration } \overbrace{g}^{\text{length of moment arm}} \int_{x_0}^b \overbrace{(x - x_0)}^{\text{density of small piece } [x, x+dx]} \overbrace{\delta(x)}^{\text{length of small piece }} dx \overbrace{\text{mass of small piece}}^{dx}.$$

* See Section 5.5.3 for more discussion on torque

* Actually, we include only the component of the force in the direction perpendicular to the vector from the pivot to the point where the force is being applied

Similarly, the counterclockwise torque is

$$g \int_a^{x_0} (x_0 - x) \delta(x) dx.$$

Setting these two integrals equal and collecting terms we get

$$0 = \int_{x_0}^b (x - x_0) \delta(x) dx + \int_a^{x_0} (x - x_0) \delta(x) dx = \int_a^b (x - x_0) \delta(x) dx.$$

Now we can distribute the $\delta(x)$ factor, apply linearity of the integral, and solve for x_0 to find

$$x_0 = \frac{\int_a^b x \delta(x) dx}{\int_a^b \delta(x) dx}.$$

In other words, the location of the center of mass is equal to the **density-weighted average value of the coordinate function x** .

We can apply the same idea (one axis at a time) to a **lamina**, which is a plate* occupying a region $L \subset \mathbb{R}^2$ and having mass density $\sigma : L \rightarrow [0, \infty)$. And similarly for a solid in \mathbb{R}^3 with a mass density function δ . In all cases, the coordinates of the center of mass are equal to the density-weighted averages of the corresponding coordinate functions:

* By slight abuse of notation, we may use L to refer to either the lamina or the region in \mathbb{R}^2 that it occupies

Definition 5.5.3: Center of mass

The center of mass of a lamina L with mass density $\sigma : L \rightarrow [0, \infty)$ is

$$\left(\frac{\iint_L x \sigma(x, y) dx dy}{\iint_L \sigma(x, y) dx dy}, \frac{\iint_L y \sigma(x, y) dx dy}{\iint_L \sigma(x, y) dx dy} \right).$$

The center of mass of a solid occupying a region R with mass density $\delta : R \rightarrow [0, \infty)$ is

$$\left(\frac{\iiint_R x \delta(x, y, z) dx dy dz}{\iiint_R \delta(x, y, z) dx dy dz}, \frac{\iiint_R y \delta(x, y, z) dx dy dz}{\iiint_R \delta(x, y, z) dx dy dz}, \frac{\iiint_R z \delta(x, y, z) dx dy dz}{\iiint_R \delta(x, y, z) dx dy dz} \right).$$

Example 5.5.1

Find the center of mass of the triangle T with vertices $(0, 0)$, $(1, 0)$, and $(0, 1)$. Assume constant density.*

Solution

If the density is k , then by Definition 5.5.3, the x -coordinate of the center of mass is

$$\frac{\iint_T kx dx dy}{\iint_T k dx dy} = \frac{k/6}{k/2} = \frac{1}{3}.$$

By symmetry, the y -coordinate of the center of mass is also $\frac{1}{3}$. Therefore, the center of mass of the lamina is $\boxed{\left(\frac{1}{3}, \frac{1}{3}\right)}$.

Exercise 5.5.4

Find the center of mass of a tetrahedron with vertices $\mathbf{0}$, \mathbf{i} , \mathbf{j} , and \mathbf{k} .

Exercise 5.5.5

Show that the ρ -value of the center of a mass of an origin-centered sphere is **not** equal to the average ρ value of a point in the sphere.

5.5.3 MOMENTS OF INERTIA

Newton's second law says that

$$\mathbf{F} = m\mathbf{a}, \quad (5.5.3)$$

where \mathbf{F} , m , and \mathbf{a} are the net force acting upon a particle, its mass, and its resulting acceleration, respectively. The rotational analogue of Newton's second law says that, for any line ℓ in space,

$$\tau = I\alpha \quad (5.5.4)$$

where τ is the torque about ℓ being applied to the particle, I is the particle's **moment of inertia** about ℓ , and α is the angular acceleration of the particle about ℓ . Thus the moment of inertia is the rotational analogue of mass.

Let's derive (5.5.4) from (5.5.3) as context for providing mathematical definitions for the quantities in (5.5.4). Denote by \mathbf{r} the vector from the nearest point on ℓ to the location of the particle. Then crossing \mathbf{r} with both sides of (5.5.3), we find that

$$\mathbf{r} \times \mathbf{F} = m\mathbf{r} \times \mathbf{a}.$$

We define the torque about ℓ to be $\tau := \mathbf{r} \times \mathbf{F}$. Letting $r = |\mathbf{r}|$ be the distance from ℓ to the particle, we define the angular acceleration of the particle about ℓ to be $\alpha = \frac{\mathbf{r}}{r} \times \frac{\mathbf{a}}{r}$. Then we obtain (5.5.4) with $I = mr^2$.

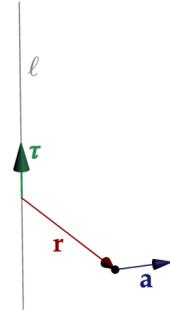


Figure 5.3 The vectors τ , \mathbf{a} , and \mathbf{r}

Exercise 5.5.6: Angular acceleration

Show that $\left| \frac{\mathbf{r}}{r} \times \frac{\mathbf{a}}{r} \right|$ is the tangential component of the acceleration vector \mathbf{a} .

Moment of inertia, like mass, is additive. That is, the moment of inertia of an ensemble of particles about a line ℓ is equal to the sum of the moments of inertia about ℓ of the particles. Therefore, the moment of inertia of a solid with continuously varying mass density may be approximated by subdividing the solid into small pieces and approximating the moment of inertia of each piece as its mass times the squared distance from some point in the small piece to ℓ . Summing the results and taking the size of the pieces to zero leads to the following formula:

Definition 5.5.4: Moment of inertia

The moment of inertia about a line ℓ of a solid occupying a region $R \subset \mathbb{R}^3$ and having density $\delta : \mathbb{R}^3 \rightarrow [0, \infty)$ is

$$I = \iiint_R r(x, y, z)^2 \delta(x, y, z) dx dy dz,$$

where $r(x, y, z)$ is the distance from (x, y, z) to ℓ .

Example 5.5.2

Find the moment of inertia of the unit-density, unit sphere about the z -axis.

Solution

The distance from the z -axis to (x, y, z) is $r = \rho \sin \phi$. Therefore, the moment of inertia is

$$\int_0^{2\pi} \int_0^{\pi} \int_0^1 (\rho \sin \phi)^2 \cdot 1 \cdot \rho^2 \sin \phi \, d\rho \, d\phi \, d\theta = \left(\frac{1}{5}\right) \left(\frac{4}{3}\right) (2\pi) = \boxed{\frac{8\pi}{15}}.$$

Exercise 5.5.7

Find the moment of inertia about the z -axis of the cylinder described by the inequalities $(x - 1)^2 + y^2 \leq 1$ and $0 \leq z \leq 1$. Repeat with the unit cube $[0, 1]^3$. Which is harder to rotate about the z -axis?

Exercise 5.5.8: Parallel axis theorem

Suppose that I_k is the moment of inertia of a solid of mass M around a line ℓ_k , for $k \in \{1, 2\}$. Show that if ℓ_1 passes through the solid's center of mass and ℓ_1 and ℓ_2 are parallel, then

$$I_2 = I_1 + md^2,$$

where d is the distance between the lines ℓ_1 and ℓ_2 .

5.5.4 PROBABILITY

A dart thrown at a dartboard $D \subset \mathbb{R}^2$ strikes a *random* point P in D . We model this state of affairs by describing a **probability density function*** $f : D \rightarrow [0, \infty)$ with the property that the probability that P lies in any region A given by integrating f over A .

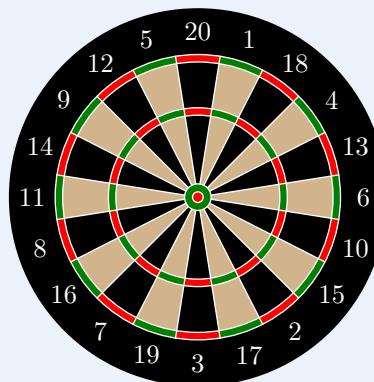
Example 5.5.3

Suppose that the probability density function for the random point where your dart hits the dartboard* $D = \mathbb{R}^2$ is given by

$$f(x, y) = \frac{1}{\pi} e^{-x^2-y^2},$$

where the origin is situated at the dartboard's bull's eye, and where x and y are measured in inches. Find the probability of scoring triple 20 on your next throw.

Note: the triple 20 region is the smaller of the two thin red strips in the sector labeled "20". The inner and outer radii of this thin strip are 3.85 inches and 4.2 inches, respectively.



* This function is positive everywhere in \mathbb{R}^2 , so the "dartboard" includes the disk shown as well as the (infinite) wall it is mounted on—this is realistic insofar as one can indeed hit the wall with a dart throw.

Solution

The region in question is described most easily in polar coordinates: it is the set of points whose polar coordinates (r, θ) satisfy $r_i \leq r \leq r_o$ and^{*} $81^\circ \leq \theta \leq 99^\circ$, where $r_i = 3.85$ and $r_o = 4.2$.

Therefore, we can obtain the probability of hitting the triple 20 by expressing the density function in polar coordinates and integrating

$$\int_{9\pi/20}^{11\pi/20} \int_{r_i}^{r_o} \frac{1}{\pi} e^{-r^2} r dr d\theta = \left(\frac{\pi}{10}\right) \left(\frac{1}{\pi}\right) \left(-\frac{1}{2}e^{-r_o^2} - \left(-\frac{1}{2}e^{-r_i^2}\right)\right).$$

Substituting the given values of r_o and r_i yields a probability of approximately 1.717×10^{-8} of scoring 60 on a single throw.

* Note that the width of each sector is $360^\circ / 20 = 18^\circ$, so the angles of the rays bounding the sector labeled 20 are $90^\circ \pm \frac{18^\circ}{2}$.

The integral of a probability density function over its whole domain must be equal to 1, since that is the probability that the random point P is *somewhere* in D .

Exercise 5.5.9

Verify that $f(x, y) = \frac{1}{\pi} e^{-x^2-y^2}$ defined on \mathbb{R}^2 is a valid probability density function (that is, integrating it over its domain yields the value 1).

Exercise 5.5.10

Suppose that f is a probability density function for a random point P in a domain $D \subset \mathbb{R}^2$. Is it possible that there exists a point $(x, y) \in D$ such that $f(x, y) = 100$? Explain why or why not. (Hint: is it possible for a solid whose total mass is 1 kilogram to have a density of 100 kilograms per cubic centimeter at a particular point in the solid?)

The analogy between probability density and mass density suggested in Exercise 5.5.10 is a broadly useful conceptual tool: you can imagine probability as some kind of abstract dust lying around on D , and regions which contain more of it are more likely to include P . To figure out exactly how much is in a given region, we integrate over that region.

Exercise 5.5.11

Explain why the answer found in Example 5.5.3 was unrealistically low for a typical dart thrower. If you were actually good enough that the pdf given in that problem describes your accuracy well, where should you consider aiming?

Note that the green bullseye is worth 25 points and the red bullseye is worth 50 points.

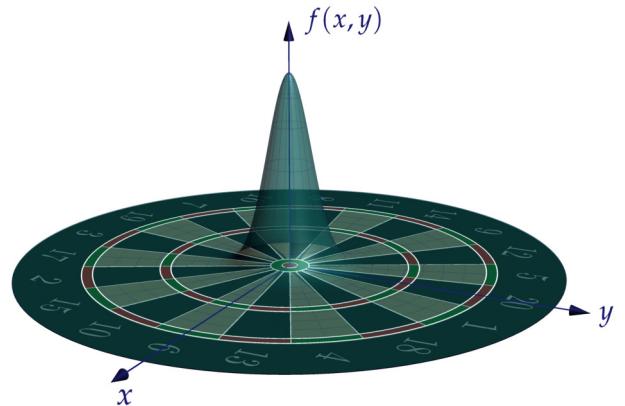


Figure 5.4 The pdf $\frac{1}{\pi} e^{-x^2-y^2}$

Given a function $g : D \rightarrow \mathbb{R}$, the value of $g(P)$ depends on the value of P and is therefore random. We call $g(P)$ a **random variable** and define its **expected value** to be the probability-weighted integral of g .

Example 5.5.4

Find the average distance from the origin to a point P selected uniformly at random from the triangle T with vertices at $(0, 0)$, $(1, 0)$, and $(1, 1)$.

Solution

The word “uniform” means that P ’s pdf is a constant function. To determine the value of this constant, we use the fact that the pdf integrates to 1. If $f(x, y) = k$ for all $(x, y) \in T$, then

$$\iint_T k \, dA = k \operatorname{area}(T) = 1,$$

which implies that $k = 2$. We are calculating the expected value of $g(P)$, where $g(x, y) = \sqrt{x^2 + y^2}$. So the expected value of $g(P)$ is given by $\iint_T g \, dA$. Because the integrand is conveniently expressed in polar coordinates, we represent T in polar coordinates as the set of points where θ is between 0 and $\pi/4$ and where r is between 0 and $\sec \theta$. Then we have

$$\iint_T g \, dA = \int_0^{\pi/4} \int_0^{\sec \theta} r^2 r \, dr \, d\theta = \frac{2}{3} \int_0^{\pi/4} \sec^3 \theta \, d\theta.$$

We can use integration by parts to work out $\int \sec^3 \theta \, d\theta = \frac{1}{2} (\sec \theta \tan \theta + \ln |\sec \theta + \tan \theta|)$. Substitution yields an expected value of

$$\boxed{\frac{1}{3} (\sqrt{2} + \ln(\sqrt{2} + 1)) \approx 0.765}.$$

Exercise 5.5.12

Explain why the answer to Example 5.5.4 is also the average distance from the origin to a point selected uniformly at random from the unit square $[0, 1]^2$.

Example 5.5.5

Find the expected value of the latitude (as measured from the north pole) of a point chosen uniformly at random from the upper half of the earth (the part above the plane passing through the equator). (Model the earth as a sphere.)

Solution

The answer doesn’t change if we change the radius of the sphere, so we might as well choose unit radius. The pdf is some constant k , and the probability of the randomly selected point lying somewhere in the upper half ball H is

$$\int_H k \, dV = k \operatorname{vol}(H) = \frac{2\pi k}{3}.$$

Since the pdf must integrate to 1, we have $k = \frac{3}{2\pi}$.

Now the function on H that we want to calculate the expected value of is the coordinate function ϕ .

So we have

$$\int_H k\phi \, dV = \frac{3}{2\pi} \int_0^{2\pi} \int_0^{\pi/2} \int_0^1 \phi \rho^2 \sin \phi \, d\rho \, d\phi \, d\theta = 1,$$

where we've used the integration-by-parts side-problem $\int_0^{\pi/2} \phi \sin \phi \, d\phi = 1$. Therefore, the expected value of the latitude is 1 radian.

Exercise 5.5.13

Find the expected value of $g(P)$ where $g(x, y) = xy$ and P is a random point in $[0, 1]^2$ whose pdf is given by $f(x, y) = xy(1 - x)(1 - y)$.

Exercise 5.5.14

Suppose that X and Y represent your score and your friend's score on an upcoming exam. If the pair of scores $P = (X, Y)$ has pdf given by $f(x, y) = \frac{363}{310} - \frac{30}{31} \left(x - y - \frac{1}{10} \right)^2$ on $[0, 1]^2$, then what is the probability that your friend scores higher than you on the exam?*

 help with calculating this integral

6 Vector Calculus

on vector fields

6.1 Vector fields and line integrals

So far we have been considering functions from \mathbb{R}^n to \mathbb{R}^1 (where n is 2 or 3). In this chapter we work with functions from \mathbb{R}^n to \mathbb{R}^2 or \mathbb{R}^3 . How should such functions be visualized? Let's begin by considering how they arise in applications.

To describe the gravitational force* felt by a particle, we would use a function with a three-dimensional input (to specify the particle's location) as well as a three-dimensional output, to specify the direction and magnitude of the force. It is natural to represent this function by drawing a small arrow indicating the output vector at several points in space, because this makes it easy to imagine how the force changes as the particle moves around (see Figure 6.1).

This picture suggests the term **vector field** for a function from \mathbb{R}^n to \mathbb{R}^n , where $n > 1$. The gravitational vector field plotted in Figure 6.1 is

$$\mathbf{F}(x, y, z) = -\frac{GMm}{(x^2 + y^2 + z^2)^{3/2}} \langle x - 1, y - 1, z - 1 \rangle. \quad (6.1.1)$$

where G is the gravitational constant and Mm is the product of the masses of particle and the attracting body at $(1, 1, 1)$.

Exercise 6.1.1

The vector plot shown represents the velocity of water on the surface of a river. The water is flowing due east, and it is flowing faster near the south end of the river than the north. Come up with a vector field \mathbf{F} whose vector plot looks approximately like the one shown.

on line integrals

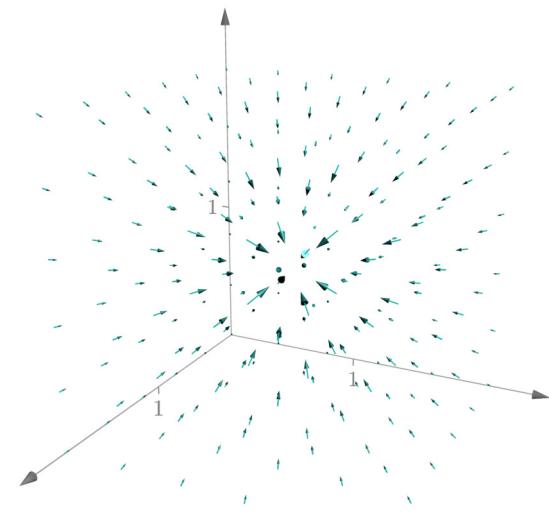
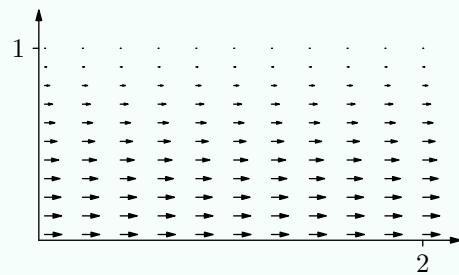


Figure 6.1 A gravitational vector field



Now suppose that rather than remaining stationary, our particle moves along a path in the presence of a force field (see Figure 6.2).* Sometimes the particle is moving with the force field and getting a boost from it, whereas other times it's working against the force field. How much net work does it take to move along the path?

If the force field were constant and the path were straight, then physics tells us that work is equal to the product of the magnitude F of the force, the distance d traveled, and the cosine of the angle θ between the force and the path. Alternatively, we may interpret the force and distance as vectors \mathbf{F} and \mathbf{d} and write the

* A vector field in which the vectors represent a physical force

work as a dot product:

$$W = Fd \cos \theta = \mathbf{F} \cdot \mathbf{d}.$$

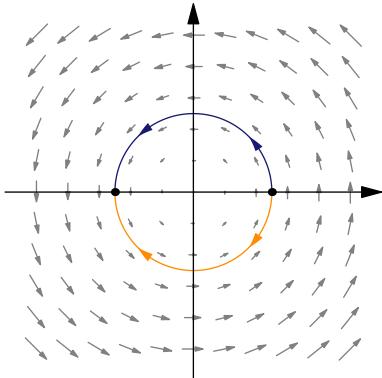


Figure 6.2 The path of a particle moving through a vector field

So how do we bootstrap our way from constant force and a straight path to varying force and a curvy path? We can cut up the path into small pieces, handle each small piece by treating the force as approximately constant and the path as approximately straight, and then add up the amount of work for each small piece. We will assume that our path $\mathbf{r}(t)$ is differentiable.*

Suppose we have a path C parameterized as $\mathbf{r}(t)$ where t ranges from a to b . Over the time interval $[t, t + \Delta t]$, the particle is displaced by* the vector $\mathbf{r}'(t) \Delta t$, and the force it feels over that time is* $\mathbf{F}(\mathbf{r}(t))$. Therefore, the contribution from the time period $[t, t + \Delta t]$ is equal to

$$\mathbf{F}(\mathbf{r}(t)) \cdot \mathbf{r}'(t) \Delta t.$$

Summing all these contributions and taking $\Delta t \rightarrow 0$, we arrive at the formula

$$W = \int_a^b \mathbf{F}(\mathbf{r}(t)) \cdot \mathbf{r}'(t) dt = \int_C \mathbf{F} \cdot d\mathbf{r},$$

where the last expression is an abbreviation for the middle expression.

We call an integral of the form $\int_C \mathbf{F} \cdot d\mathbf{r}$ a **line integral**.

Example 6.1.1

Find the line integral of $\mathbf{F}(x, y) = \langle xy, y \rangle$ along the parabola $y = x^2$ from $(0, 0)$ to $(2, 4)$.

Solution

Let's parameterize the parabola using the x -coordinate as the parameter:

$$\mathbf{r}(t) = \langle t, t^2 \rangle.$$

Note that the point $(2, 4)$ is visited at time $t = 2$, while the origin is visited at time $t = 0$. Therefore,

$$W = \int_0^2 \langle t(t^2), t^2 \rangle \cdot \langle 1, 2t \rangle dt = \int_0^2 (t^3 + 2t^3) dt = \boxed{12}.$$

The following theorem states that the choice of parameterization of a curve doesn't matter when computing a line integral. This makes sense physically, since the formula $W = Fd$ does not involve time, and our derivation of the line integral formula was based on $W = Fd$. The role of the parameterization was merely to provide a convenient way to split up the path into short pieces. Exercise 6.1.2 below gives an example.

* The curve could alternatively be *piecewise* differentiable, meaning that the curve is non-differentiable at only finitely many points

** Approximately, with an error that vanishes as $\Delta t \rightarrow 0$

Theorem 6.1.1: Independence of parameterization

If C is a curve parameterized by \mathbf{r}_1 over $[a, b]$ and also by \mathbf{r}_2 over $[c, d]$, then

$$\int_a^b \mathbf{F}(\mathbf{r}_1(t)) \cdot \mathbf{r}'_1(t) dt = \int_c^d \mathbf{F}(\mathbf{r}_2(t)) \cdot \mathbf{r}'_2(t) dt.$$

In other words, $\int_C \mathbf{F} \cdot d\mathbf{r}$ depends only on the curve C , not the choice of parameterization.

Exercise 6.1.2

- (i) Compute the line integral of $\mathbf{F} = \langle x^2, -xy \rangle$ over the portion of the unit circle in the first quadrant, using the parameterization $\mathbf{r}(t) = \langle \sin t, \cos t \rangle$.
- (ii) Perform the same line integral using the parameterization $\mathbf{r}(t) = \langle t, \sqrt{1-t^2} \rangle$.

Exercise 6.1.3

Consider the vector field \mathbf{F} and path C shown in Figure 6.2. Is $\int_C \mathbf{F} \cdot d\mathbf{r}$ positive or negative?

6.2 The fundamental theorem of vector calculus

mei on the gradient theorem for line integrals

In general, the line integral of \mathbf{F} over a path between two points depends on the path, not just the starting and ending points. For example, in Figure 6.3, the line integral along the blue (top) path is positive, while the line integral along the orange (bottom) path is negative.

However, there is an important class of vector fields which are path-independent, meaning that the value of $\int_C \mathbf{F} \cdot d\mathbf{r}$ depends only on the starting and ending points of C . These are the vector fields which can be written as the gradient of a function from \mathbb{R}^n to \mathbb{R}^1 . For example, $\mathbf{F}(x, y, z) = \langle -2x, -2y, z \rangle$ is the gradient of the function

$$f(x, y, z) = -x^2 - y^2 + \frac{1}{2}z^2.$$

Such vector fields are called **conservative**.

If we calculate the line integral of ∇f along a curve C parameterized by $\mathbf{r}(t) = \langle r_1(t), r_2(t), r_3(t) \rangle$, then the contribution from the portion of the curve from $\mathbf{r}(t)$ to $\mathbf{r}(t + \Delta t)$ is*

$$\langle \partial_x f, \partial_y f, \partial_z f \rangle \cdot \langle r'_1(t), r'_2(t), r'_3(t) \rangle \Delta t,$$

which by the chain rule is* the change in $f(\mathbf{r}(t))$ over that interval. Therefore, the line integral of ∇f along a path is equal to the change in f from the beginning to the end of the path.

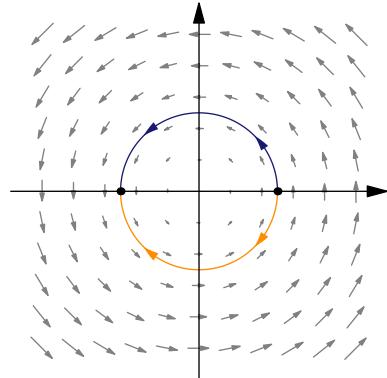


Figure 6.3 The vector field $\mathbf{F}(x, y) = \langle -y, x \rangle$ and two semicircular paths

* approximately, with an error that vanishes as $\Delta t \rightarrow 0$

Theorem 6.2.1: Fundamental theorem for line integrals

If C is a path from \mathbf{a} to \mathbf{b} and f is a differentiable function, then

$$\int_C \nabla f \cdot d\mathbf{r} = f(\mathbf{b}) - f(\mathbf{a}).$$

Example 6.2.1

Suppose $\mathbf{F}(x, y, z) = \langle 2xy^3z, 3x^2y^2z + y, x^2y^3 \rangle$ and that C is the circular arc from the origin to the point $(1, 1, 1)$ and passing through the point $(1/2, 1/2, 1)$. Find $\int_C \mathbf{F} \cdot d\mathbf{r}$.

Solution

Finding a parameterization for C seems computationally messy. However, if \mathbf{F} is conservative, then we can use Theorem 6.2.1. Integrating $2xy^3z$ with respect to x , we see that if $\mathbf{F} = \nabla f$ for some function f , then we would have

$$f(x, y, z) = x^2y^3z + C_1(y, z),$$

where $C_1(y, z)$ denotes a function not depending on x . Similarly, we can integrate the second and third components with respect to y and z to find that

$$\begin{aligned} f(x, y, z) &= x^2y^3z + \frac{1}{2}y^2 + C_2(x, z) \\ f(x, y, z) &= x^2y^3z + C_3(x, y). \end{aligned}$$

We see that these three conditions are simultaneously satisfied by the function $f(x, y, z) = x^2y^3z + \frac{1}{2}y^2$. So the desired line integral is equal to

$$f(1, 1, 1) - f(0, 0, 0) = \frac{3}{2} - 0 = \left[\frac{3}{2} \right].$$

The following theorem provides a convenient way to check whether a two-dimensional vector field is conservative.

Theorem 6.2.2

A vector field $\mathbf{F}(x, y) = \langle M(x, y), N(x, y) \rangle$ which is differentiable on* \mathbb{R}^2 is conservative if and only if

$$\partial_x N = \partial_y M. \quad (6.2.1)$$

* Here it is important that \mathbf{F} is differentiable on all of \mathbb{R}^2 , rather than an arbitrary subset thereof—see Exercise 6.2.2

To see where (6.2.1) comes from, note that this equation follows directly from Clairaut's theorem for conservative fields \mathbf{F} . So the more interesting aspect of Theorem 6.2.2 is the converse direction: merely checking $\partial_x N = \partial_y M$ establishes existence or nonexistence of a gradient function.

Exercise 6.2.1

Show that the gravitational force in (6.1.1) is conservative.

Exercise 6.2.2

(i) Try to apply Theorem 6.2.2 to the vector field

$$\mathbf{F}(x, y) = \left\langle -\frac{y}{x^2 + y^2}, \frac{x}{x^2 + y^2} \right\rangle.$$

(ii) Show by plotting this vector field that it is not conservative. How does this square with Theorem 6.2.2?

on Green's theorem

6.3 Green's theorem

Is it possible to engineer a simple mechanical device that measures the area bounded by whatever curve it traces out on paper? This seems surprising, since computing the area would seem to require some inspection of the region inside the curve. However, the *planimeter** can record the area of a region based on the motion of its wheels as its tip traverses the boundary of the region. The design of the planimeter takes advantage of the following beautiful relationship between line integrals along the boundary of a curve and double integrals over the enclosed region.

* Invented in 1854



Figure 6.4 A planimeter

Theorem 6.3.1: Green's theorem

If $\mathbf{F} = \langle M, N \rangle$ is a vector field* on \mathbb{R}^2 with continuous partial derivatives, and if D is a region bounded by a simple, counterclockwise-oriented, piecewise smooth curve C , then

$$\int_C \mathbf{F} \cdot d\mathbf{r} = \iint_D (\partial_x N - \partial_y M) \, dA.$$

Example 6.3.1

Verify Green's theorem in the case where D is the unit disk and $\mathbf{F}(x, y) = \langle 0, x \rangle$.

Solution

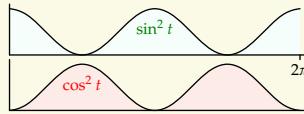
We parameterize the unit disk trigonometrically as $(\cos t, \sin t)$, and we calculate the line integral

$$\int_0^{2\pi} \langle 0, \cos t \rangle \cdot \langle -\sin t, \cos t \rangle \, dt = \int_0^{2\pi} \cos^2 t \, dt = \pi.$$

This last integral can be done with a trick: note that

$$\int_0^{2\pi} (\cos^2 t + \sin^2 t) \, dt = \int_0^{2\pi} 1 \, dt = 2\pi.$$

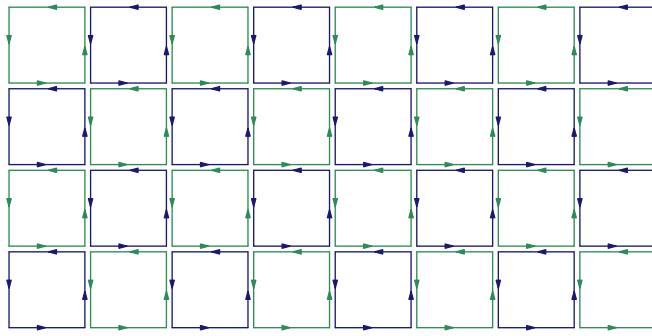
However, the contributions of $\int_0^{2\pi} \cos^2 t \, dt$ and $\int_0^{2\pi} \sin^2 t \, dt$ are equal, since their graphs over the region of integration are the same up to a shift. So each is equal to π .



The integrand for the double integral is* $\partial_x N - \partial_y M = 1 - 0 = 1$, so the value of the double integral is the area of the unit disk, which is equal to π . Thus the conclusion of Green's theorem is satisfied.

Proving Green's theorem

The idea of the proof of Green's theorem is to cut D into small rectangles along grid lines (shown with small gaps for visual clarity). Green's theorem holds approximately on each small rectangle $R = [x - \frac{\Delta x}{2}, x + \frac{\Delta x}{2}] \times [y - \frac{\Delta y}{2}, y + \frac{\Delta y}{2}]$, because the left and right sides of R contribute to $\int_R \mathbf{F} \cdot d\mathbf{r}$ approximately



$$\underbrace{N\left(x + \frac{\Delta x}{2}, y\right) \Delta y}_{\mathbf{F} \text{ at midpoint dotted with } \langle 0, \Delta y \rangle \text{ step}} - \underbrace{N\left(x - \frac{\Delta x}{2}, y\right) \Delta y}_{\mathbf{F} \text{ at midpoint dotted with } \langle 0, -\Delta y \rangle \text{ step}} \approx (\partial_x N)(x, y) \Delta x \Delta y.$$

Similarly, the contribution of the top and bottom sides is $-(\partial_y M)(x, y) \Delta x \Delta y$. So all together, the circulation* of \mathbf{F} around R is $(\partial_x N - \partial_y M) \Delta x \Delta y$. This is also approximately equal to the integral of $\partial_x N - \partial_y M$ over R , since $\partial_x N - \partial_y M$ is approximately constant over R .

The line integrals of \mathbf{F} over the small rectangles sum to the line integral of \mathbf{F} around the boundary* of D , because each interior segment is integrated along twice (once for each adjoining rectangle) and in opposite directions. These contributions sum to zero, leaving only the integrals along the outer edges. Since these outer edges fit together to form ∂D , the line integrals along them sum to the line integral along ∂D .

Since the integral of $\partial_x N - \partial_y M$ over D is also equal to the sum of the integrals of $\partial_x N - \partial_y M$ over the small rectangles, the (approximate) Green's theorem for the small rectangles implies (approximate) Green's theorem for D . Letting the sizes of the rectangles tend to 0, this approximation becomes exact and yields Green's theorem.

* Using $\mathbf{F} = \langle 0, x \rangle$ gives an integrand of 1 on the right-hand side of Green's theorem, but: (i) $\langle 0, x \rangle$ isn't the only vector field with this property, and (ii) Green's theorem is also useful when the integrand is non-constant

* Circulation is a synonym of line integral which is used only when curve starts and ends at the same point.

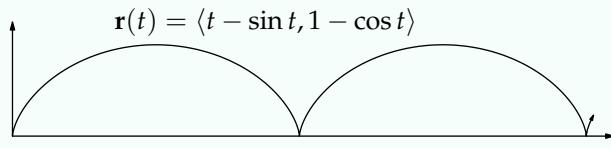
* In other words, the circulation around the boundary of a region is additive

Exercise 6.3.1

Use Green's theorem to find the line integral of $\mathbf{F} = \langle \sqrt{x^2 + 1}, \arctan x \rangle$ along a counterclockwise traversal of the triangle with vertices $(0, 0)$, $(1, 0)$, and $(0, 1)$.

Exercise 6.3.2

Use Green's theorem to find the area under each arch of the cycloid shown below.



6.4 Surface integrals and flow

on surface integrals

6.4.1 SURFACE INTEGRALS

What is the average temperature on the surface of the earth? Let's look past the scientific challenges of this problem and imagine that the earth is a sphere S and that we have a reading of the temperature T at every point on its surface at a particular point in time. The average temperature at that moment should then be the integral of T over S divided by the surface area of S . But what does it mean to integrate over a surface?

We can use the same approach we use throughout calculus when studying quantities which are additive and continuously varying: split S into tiny patches over which T may be treated as constant, multiply the area of each patch by the value of T somewhere on that patch, and sum the resulting products. As the size of the patches tends to zero, we expect this sum to converge to some limiting value, and we can declare that limit to be the value of the **surface integral*** of f over S , denoted $\iint_S f \, dA$.

Example 6.4.1

Find the surface integral of $f(x, y, z) = 2x^2 + 2y^2 + 2z^2$ over the unit sphere S . You may assume that the surface area of a sphere of radius R is $4\pi R^2$.

Solution

This function is equal to 2 everywhere on the unit sphere, so if we split S into many small pieces and consider one of them, then the product of the value of f somewhere on the piece and the area of the piece is just twice the area of the piece. When we sum these contributions over all the pieces, we will end up with twice the total surface area of the sphere. In other words, we have

$$\iint_S 2 \, dA = 2 \iint_S 1 \, dA = 2 \times \text{surface area}(S) = 8\pi.$$

Example 6.4.1 was special because the function happened to be constant over the surface. Another special situation which simplifies the process of finding a surface integral is when we're integrating over a surface which is contained in a plane.

* Or scalar surface integral, to distinguish from vector surface integral introduced in the next subsection

Example 6.4.2

Find the surface integral of $f(x, y, z) = xyz$ over the rectangular prism $[0, 1] \times [0, 2] \times [0, 3]$.

Solution

The value of the function f at every point in a coordinate plane is zero. So the contribution from these faces is zero.

For the top face, the value of the function at each point $(x, y, 3)$ is $3xy$. Therefore, the integral of f over this face is what you get when you split the rectangle $[0, 1] \times [0, 2]$ into many small rectangles, multiply the value of $3xy$ by the area of each one (where x and y are the first and second coordinates of some point in the small rectangle), sum the results, and take the rectangle size to zero. In other words, the contribution from the top face is

$$\int_0^1 \int_0^2 3xy \, dy \, dx = 3.$$

Likewise, the contributions from the other two faces are

$$\begin{aligned} \int_0^3 \int_0^2 1yz \, dy \, dz &= 9, \text{ and} \\ \int_0^3 \int_0^1 2xz \, dx \, dz &= \frac{9}{2} \end{aligned}$$

So altogether the surface integral is equal to $3 + 9 + \frac{9}{2} = \boxed{\frac{33}{2}}$.

Let's develop a general-purpose method for evaluating surface integrals. Recall that a parameterization of a curve C in \mathbb{R}^3 is a function $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^3$ for which $\mathbf{r}(t)$ traces out the points of C as t goes from a to b . Likewise, a **parameterization** of a surface S in \mathbb{R}^3 is a function \mathbf{r} from some planar domain $D \subset \mathbb{R}^2$ to \mathbb{R}^3 such that $\mathbf{r}(u, v)$ sweeps out* the surface S as (u, v) varies over D .

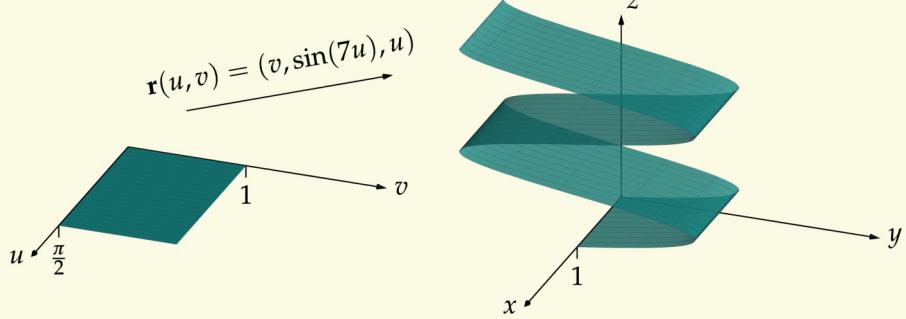
* More precisely, a parameterization \mathbf{r} is a bijective map from D to S

Example 6.4.3

Find a parameterization of the surface consisting of the points $(x, y, z) \in \mathbb{R}^3$ such that $y = \sin 7z$, $0 \leq x \leq 1$, and $0 \leq z \leq \frac{\pi}{2}$.

Solution

One way to map a pair of points (u, v) to a triple of points (x, y, z) satisfying the given equations is to use $u = z$ and $v = x$ as our two parameters. Then $\mathbf{r}(u, v) = (v, \sin(7u), u)$ maps the rectangle $[0, \frac{\pi}{2}] \times [0, 1]$ to the desired surface, as shown.



Example 6.4.4

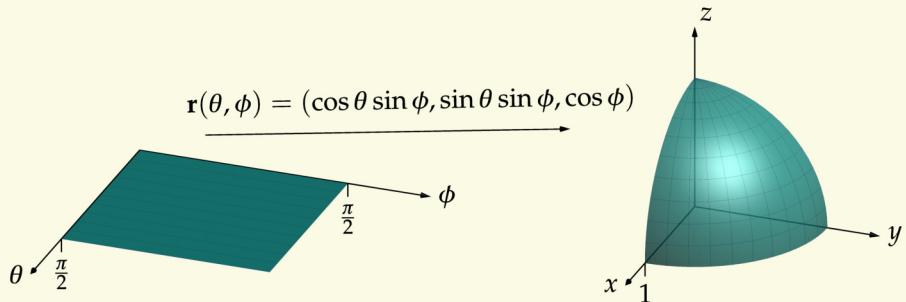
Find a parameterization of the portion of the unit sphere in the first octant.

Solution

A parameterization is a way of using two numbers to identify points on the surface. We have a standard way of doing this on the sphere we live on: latitude and longitude. So we can parametrize a sphere or a portion thereof by using θ and ϕ as parameters.* As we discovered in Exercise 3.4.4, the point on the sphere with spherical coordinates $1, \theta$, and ϕ is

$$\mathbf{r}(\theta, \phi) = (\cos \theta \sin \phi, \sin \theta \sin \phi, \cos \phi).$$

This function \mathbf{r} maps $[0, \frac{\pi}{2}] \times [0, \frac{\pi}{2}]$ onto the portion of the sphere in the first octant:



The lesson of Examples 6.4.3 and 6.4.4 is that **coordinates make good parameters**. It might be convenient to use Cartesian coordinates, cylindrical or spherical coordinates, or even custom coordinates. But if we find two coordinate functions whose values conveniently specify a location (x, y, z) on the surface, then we can define \mathbf{r} to be the map which sends each pair of coordinate values to the triple (x, y, z) .

Exercise 6.4.1

- (i) Parametrize the portion of the plane $x + 2y + 3z = 6$ in the first octant.
- (ii) Sketch and parametrize the set of points (x, y, z) satisfying $0 \leq z \leq x$ and $x^2 + y^2 = 1$.

If we can parametrize a surface, then we can split S into many small pieces by splitting D into many small

rectangles. Each small rectangle $[u, u + \Delta u] \times [v, v + \Delta v]$ maps under the parameterization to an approximate parallelogram spanned by

$$\mathbf{r}(u + \Delta u, v) - \mathbf{r}(u, v) \approx \partial_u \mathbf{r}(u, v) \Delta u \quad \text{and} \quad \mathbf{r}(u, v + \Delta v) - \mathbf{r}(u, v) \approx \partial_v \mathbf{r}(u, v) \Delta v.$$

The area spanned by this small parallelogram is the cross product of these two vectors. Multiplying this area by the value of f somewhere in the parallelogram and summing over all the tiny rectangles, we get the following theorem.

Theorem 6.4.1: (Surface integral formula)

 on surface parameterization

If $\mathbf{r} : D \rightarrow \mathbb{R}^3$ is a parameterization of a surface S , then

$$\iint_S f \, dA = \iint_D f(\mathbf{r}(u, v)) |\partial_u \mathbf{r} \times \partial_v \mathbf{r}| \, du \, dv. \quad (6.4.1)$$

Example 6.4.5

Find the average value of the function $f(x, y, z) = z$ on the upper unit half-sphere.

Solution

Let's use the parametrization

$$\mathbf{r}(\theta, \phi) = (\cos \theta \sin \phi, \sin \theta \sin \phi, \cos \phi),$$

as θ ranges from 0 to 2π and ϕ ranges from 0 to $\frac{\pi}{2}$. Then

$$|\partial_\theta \mathbf{r} \times \partial_\phi \mathbf{r}|^2 = \left(-\sin \phi \sin^2 \theta \cos \phi - \sin \phi \cos \phi \cos^2 \theta \right)^2 + \sin^4 \phi \sin^2 \theta + \sin^4 \phi \cos^2 \theta = \sin^2 \phi.$$

Theorem 6.4.1 gives

$$\int_S z \, dA = \int_0^{2\pi} \int_0^{\frac{\pi}{2}} \overbrace{(\cos \phi)}^z \overbrace{\sin \phi}^{|\partial_\theta \mathbf{r} \times \partial_\phi \mathbf{r}|} \, d\phi \, d\theta = \pi.$$

The surface area of the upper half-sphere is*

$$\int_S 1 \, dA = \int_0^{2\pi} \int_0^{\frac{\pi}{2}} \overbrace{1}^1 \overbrace{\sin \phi}^{|\partial_\theta \mathbf{r} \times \partial_\phi \mathbf{r}|} \, d\phi \, d\theta = 2\pi.$$

Therefore, the average value of z is $\frac{\pi}{2\pi} = \boxed{\frac{1}{2}}$.

* We could also use the formula for the surface area of the sphere, but let's do it from scratch

Exercise 6.4.2

Suppose that S is the surface consisting of the points $(x, y, z) \in \mathbb{R}^3$ satisfying $y^2 + z^2 = 4$ and $-1 \leq x \leq 2$. Equip S with the density function $\sigma(x, y, z) = y + 2$, and find the resulting center of mass of S .

Exercise 6.4.3

- Find the surface area of the portion of the origin-centered unit sphere lying above the line $z = \frac{1}{2}$.
- Find the value z_0 such that half of the sphere lies above the plane $z = z_0$.

Exercise 6.4.4

Suppose that S is the portion of the graph of the function $f(x, y) = 3 - x - y$ lying above the disk $(x - 2)^2 + (y - 1)^2 \leq 1$. Suppose that a point P is selected uniformly at random from S . Find the expected value of $g(P)$, where $g(x, y, z) = 2x + z$.

Exercise 6.4.5

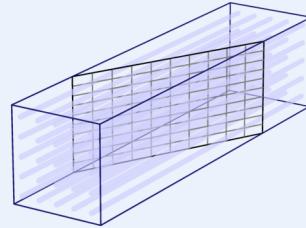
Show that if S is a planar domain, then the formula (6.4.2) agrees with the custom-coordinate integration formula (Theorem 5.4.1).

on surface
integrals of a
vector field

6.4.2 FLOW

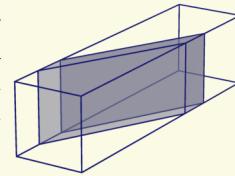
Example 6.4.6

Consider a constant-velocity river flowing through a net as shown. Find the volume of water flowing through the net per unit time, in terms of the area A of the net, the velocity v , and the angle θ between the direction of the river's flow and a vector normal to the plane of the net.



Solution

Imagine letting the water flow for one time unit and then taking a snapshot. The locations of all the water molecules which have flowed through the net during this period occupy a parallelepiped, as shown. The base area of this parallelepiped is A , while its height is equal to* $v \cos \theta$. Therefore, the volume of water passing through the net per unit of time is $Av \cos \theta$.



Let's define the vector \mathbf{A} whose length is equal to A and whose direction is orthogonal to the net. Then the **flow** $Av \cos \theta$ can be written as

$$\text{flow} = \mathbf{A} \cdot \mathbf{v},$$

where \mathbf{v} is the river's velocity vector.

* This is a right-triangle trigonometry exercise

Example 6.4.7

Suppose that the velocity field of a body of water is given by $\mathbf{F} = \langle -2y, 4z, x \rangle$ (in meters per second) and that a rectangular net is positioned in the water with corners at $(1, 1, 3)$, $(1, 4, 3)$, $(1, 4, 5)$, and $(1, 1, 5)$ (in meters). Find the volume of water flowing through the frame of the net per second.

Solution

Since the velocity field isn't constant, we divide the net into small patches and treat the velocity as constant on each one. Since the rectangle is contained in the plane $x = 1$, the vector $\langle 1, 0, 0 \rangle$ is normal to the rectangle. Therefore, the flow through a small patch of area ΔA located at (x, y, z) is approximately equal to

$$\mathbf{A} \cdot \mathbf{v} = \langle \Delta A, 0, 0 \rangle \cdot \langle -2y, 4z, x \rangle = -2y \Delta A.$$

If we sum the flow through each patch across the whole rectangular region occupied by the net, we get a Riemann sum that converges as $\Delta A \rightarrow 0$ to

$$\int_1^4 \int_3^5 (-2y) dz dy = -30.$$

Therefore, the volume of water is 30 cubic meters per second, and the net flow is in the direction towards the side facing the yz -plane.

The ideas in Example 6.4.7 yield the following definition, which develops the notion of vector field integration over surfaces in terms of the scalar surface integral.

Definition 6.4.1

The **flow** of a vector field \mathbf{F} through a surface S from one side s to the other side t is defined by

$$\iint_S \mathbf{F} \cdot d\mathbf{A} = \iint_S \mathbf{F} \cdot \mathbf{n} dA,$$

where $\mathbf{n} = \mathbf{n}(x, y, z)$ is a unit vector which is orthogonal to S at each point (x, y, z) and points in the direction from s to t .

A surface with a distinguished side is called an *oriented surface*, and not every surface is orientable. For example, the Möbius strip:



Exercise 6.4.6

Find the flow of the vector field $\mathbf{F} = \langle x^2, y^2, z^2 \rangle$ from the inside to the outside of the unit sphere.

We can write down an integral formula for calculating flow through an arbitrary parametrized surface S . Since both $\partial_u \mathbf{r}$ and $\partial_v \mathbf{r}$ are tangent to the surface S , the cross product of $\partial_u \mathbf{r}$ and $\partial_v \mathbf{r}$ is orthogonal to S . When we divide this cross product by its length and substitute into Theorem 6.4.1, we get a cancellation of the $|\partial_u \mathbf{r} \times \partial_v \mathbf{r}|$ factors. In other words, $d\mathbf{A} = \partial_u \mathbf{r} \times \partial_v \mathbf{r} du dv$, so to find the flow we can dot \mathbf{F} with $\partial_u \mathbf{r} \times \partial_v \mathbf{r}$ and integrate over D :

Theorem 6.4.2: (Flow integral formula)

If $\mathbf{r} : D \rightarrow \mathbb{R}^3$ is a parameterization of a surface S such that at each point of S , the vector $\partial_u \mathbf{r} \times \partial_v \mathbf{r}$ points from one side s of S to the other side t , then the flow of \mathbf{F} through S from s to t is

$$\iint_S \mathbf{F} \cdot d\mathbf{A} = \iint_D \mathbf{F}(\mathbf{r}(u, v)) \cdot (\partial_u \mathbf{r} \times \partial_v \mathbf{r}) du dv. \quad (6.4.2)$$

Example 6.4.8

Find the flow of the vector field $\mathbf{F} = \langle y^2, 0, z \rangle$ up through the part of the paraboloid $z = 4 - x^2 - y^2$ above the xy -plane.

Solution

Let's parametrize the surface using the cylindrical coordinates $u = r$ and $v = \theta$, so that

$$\mathbf{r}(u, v) = \langle u \cos v, u \sin v, 4 - u^2 \rangle,$$

where (u, v) ranges over $[0, 2] \times [0, 2\pi]$. Then

$$\partial_u \mathbf{r}(u, v) = \langle \cos v, \sin v, -2u \rangle \quad \text{and} \quad \partial_v \mathbf{r}(u, v) = \langle -u \sin v, u \cos v, 0 \rangle,$$

which yields

$$\partial_u \mathbf{r} \times \partial_v \mathbf{r} = \langle 2u^2 \cos v, 2u^2 \sin v, u \rangle.$$

This vector indeed points from the bottom side of the surface to the top, so Theorem 6.4.2 tells us that the desired flow is

$$\int_0^{2\pi} \int_0^2 \langle u^2 \sin^2 v, 0, 4 - u^2 \rangle \cdot \langle 2u^2 \cos v, 2u^2 \sin v, u \rangle du dv,$$

which simplifies to

$$\int_0^{2\pi} \int_0^2 (2u^4 \sin^2 v \cos v + 4u - u^3) du dv = \boxed{8\pi}.$$

Exercise 6.4.7

Find the flow of the vector field $\mathbf{F} = \langle -y, x, z^{3 \sin xyz} \rangle$ from the inside to the outside of the cylinder $\{(x, y, z) \in \mathbb{R}^3 : 0 \leq z \leq 3 \text{ and } x^2 + y^2 = 1\}$. Explain why your answer makes sense physically.

6.5 Divergence and curl

As suggested by the title of the classic book *div, grad, curl, and all that* by H.M. Schey*, the gradient is one of the main characters in the vector calculus story. In this section we will develop the two other fundamental vector calculus derivative operators: *divergence* and *curl*. Like Schey, we will emphasize the underlying physical intuition.

* Highly rec-
ommended

The **divergence** of a vector field $\mathbf{F} = \langle M, N, P \rangle$ is the function whose value at a point (x, y, z) is equal to the **net flow density** of \mathbf{F} out of a small region located at (x, y, z) . In other words, if we put a small box of dimensions $\Delta x \times \Delta y \times \Delta z$ centered at (x, y, z) , then the ratio of the net flow of \mathbf{F} out of the box to the volume of the box converges to $\operatorname{div} \mathbf{F}(x, y, z)$ as $(\Delta x, \Delta y, \Delta z) \rightarrow 0$.

Let's work out a formula for $\operatorname{div} \mathbf{F}$. The net flow through the top of the box of dimensions $\Delta x \times \Delta y \times \Delta z$ centered at (x, y, z) is approximately* $P(x, y, z + \Delta z/2) \Delta x \Delta y$, and the net flow through the bottom of the box is approximately $P(x, y, z - \Delta z/2) \Delta x \Delta y$. Thus the difference is approximately $(\partial_z P)(x, y, z) \Delta x \Delta y \Delta z$, and the difference per unit volume is $\partial_z P(x, y, z)$. Similarly, the front/back and left/right sides contribute $\partial_x M(x, y, z)$ and $\partial_y N(x, y, z)$ to the net flow density out of the box. Therefore, $\operatorname{div} \mathbf{F} = \partial_x M + \partial_y N + \partial_z P$. This formula suggests the alternate notation $\nabla \cdot \mathbf{F}$ for the divergence of \mathbf{F} .

* The normal vector for the top is $\langle 0, 0, 1 \rangle$, and that dots with \mathbf{F} to give P .

Definition 6.5.1: Divergence

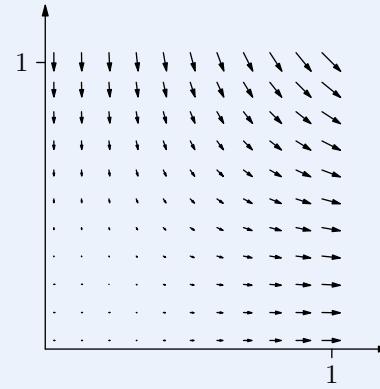
The **divergence** of a vector field $\mathbf{F} = \langle M, N, P \rangle$ is a scalar function defined by

$$\nabla \cdot \mathbf{F} = \partial_x M + \partial_y N + \partial_z P.$$

For example, the divergence of $\langle x^2, xy, z \rangle$ is $2x + x + 1 = 3x + 1$.

Example 6.5.1

Figure out where $\nabla \cdot \mathbf{F}$ is positive for the vector field \mathbf{F} shown.



Solution

We can see that in the top left of the diagram that there is more flow into each small region than out of it, since the vectors are downward-pointing and longer than the vectors below them. Therefore, the divergence is negative in the top left.

By similar reasoning, we see that the divergence is positive in the bottom-right part of the figure. The dividing line between regions of positive and negative divergence is $y = x$, since points along that line have vectors of equal length pointing towards and away from them.

Exercise 6.5.1

Find the divergence of the gravitational vector field in (6.1.1).

Exercise 6.5.2

Look at a vector plot* to figure out where $\nabla \cdot \mathbf{F} > 0$, using the approach of Example 6.5.1, for the vector field $\mathbf{F} = \langle xy, y^2 \rangle$. Then evaluate $\nabla \cdot \mathbf{F}$ and find where $\nabla \cdot \mathbf{F} > 0$ algebraically.

for drawing a vector plot

on curl

6.5.2 CURL

The **curl** of a vector field $\mathbf{F} = \langle M, N, P \rangle$ is the vector field whose value at a point (x, y, z) describes the **circulation density** of \mathbf{F} around (x, y, z) . Specifically, consider a curve ℓ centered at a point (x, y, z) looping counterclockwise around the vector $\langle 0, 0, 1 \rangle$. We define the *third* component of $\text{curl } \mathbf{F}(x, y, z)$ to be the limit as the size of the loop tends to zero of the ratio of $\int_{\ell} \mathbf{F} \cdot d\mathbf{r}$ to the area enclosed by ℓ . We define the first two components of $\text{curl } \mathbf{F}(x, y, z)$ similarly, with ℓ running counterclockwise around $\langle 1, 0, 0 \rangle$ or $\langle 0, 1, 0 \rangle$, respectively, instead of $\langle 0, 0, 1 \rangle$.

* Since Green's theorem says that circulation equals $\partial_x N - \partial_y M$ integrated, it can be interpreted as the assertion "circulation density is given by $\partial_x N - \partial_y M$ ".

By Green's theorem, third component of the curl is equal to* $\partial_x N - \partial_y M$. If we swap out (x, y) for (y, z) or (z, x) to get similar formulas for the first two coordinates, we find that $\text{curl } \mathbf{F} = \langle \partial_y P - \partial_z N, -\partial_x P + \partial_z M, \partial_x N - \partial_y M \rangle$. This formula suggests the notation $\nabla \times \mathbf{F}$ for the curl of \mathbf{F} .

Definition 6.5.2: Curl

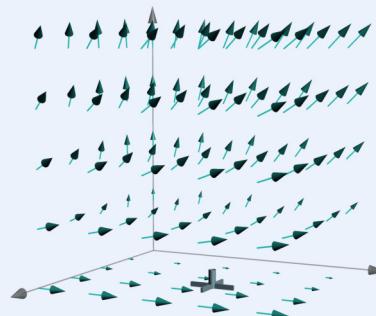
The **curl** of a vector field $\mathbf{F} = \langle M, N, P \rangle$ is a vector field on \mathbb{R}^3 defined by

$$\nabla \times \mathbf{F} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ \partial_x & \partial_y & \partial_z \\ M & N & P \end{vmatrix} = \langle \partial_y P - \partial_z N, -\partial_x P + \partial_z M, \partial_x N - \partial_y M \rangle.$$

Another way to visualize curl physically is to interpret the vector field as a fluid velocity field and place a small paddle wheel (with an axis of rotation in a coordinate direction) at a particular point in this field. The corresponding component of the curl measures how rapidly and in which direction this paddle wheel turns.

Example 6.5.2

Consider the vector field \mathbf{F} shown. Find the sign of the z -component of the curl of \mathbf{F} at any point in the xy -plane.



Solution

We can see that if we place a small paddle wheel at a point of interest (as shown in the figure above) it will rotate in the counterclockwise direction. This is because the vectors on the right (meaning the side where x is larger) push harder than the vectors on the left. Therefore, the z -component of the curl is positive.

When we studied gradients, we learned that directional derivatives of a function in the coordinate directions actually determine its directional derivatives in all directions: the derivative in the $\langle u_1, u_2 \rangle$ direction is equal to a linear combination with weights u_1 and u_2 of the derivatives in the coordinate directions. The same idea holds for the curl: if \mathbf{u} is a unit vector, let's define $\text{curl}_{\mathbf{u}} \mathbf{F}$ to be the circulation density of \mathbf{F} in the counterclockwise direction around \mathbf{u} . More precisely, if ℓ is a small loop which is perpendicular to \mathbf{u} and oriented counterclockwise as viewed from the head of \mathbf{u} , then the ratio of $\int_{\ell} \mathbf{F} \cdot d\mathbf{r}$ to the area enclosed by ℓ converges to $\text{curl}_{\mathbf{u}} \mathbf{F}$ as the size of ℓ tends to zero.

Theorem 6.5.1

The circulation density $\text{curl}_{\mathbf{u}} \mathbf{F}$ of a vector field \mathbf{F} around a unit vector \mathbf{u} is equal to $(\nabla \times \mathbf{F}) \cdot \mathbf{u}$.

!!!

Thus $\nabla \times \mathbf{F}$ is a vector field whose **main purpose is to be dotted with unit vectors to compute circulation**. We will revisit this idea in Section 6.7 when we discuss Stokes' theorem.

Example 6.5.3

Find the orientation for a paddle wheel at the point $(1, 1, 1)$ in the velocity field $\langle xyz, x^2 - y, z \rangle$ which will maximize how fast it spins.

Solution

We calculate the curl: $\nabla \times \mathbf{F} = \langle 0, xy, -xz + 2x \rangle$, which at the point $(1, 1, 1)$ is equal to $\langle 0, 1, 1 \rangle$. Since the dot product of a fixed vector \mathbf{v} with a unit vector is maximized when the unit vector is aligned with \mathbf{v} , we see that the paddle wheel should be oriented so that its axis is in the direction

$$\left\langle 0, \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right\rangle.$$

Exercise 6.5.3

Calculate $\nabla \times \mathbf{F}$, where $\mathbf{F} = \langle e^{\sin \log x} + y^2, -2z, y^3 + \cos z \rangle$.

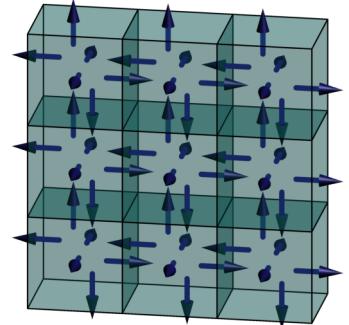
Exercise 6.5.4

Show that the curl of a conservative vector field is zero.

6.6 Divergence theorem

Suppose that $\mathbf{F} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is a vector field. Just as integrating mass density yields total mass, or integrating charge density yields total charge, integrating flow density yields total flow. This fact is called the *divergence theorem*.

The property of physical mass which makes it compatible with the idea of computing a total by integrating a density function is **additivity**: any way you divide up a solid, its mass is equal to the sum of the masses of its parts. Flow density works similarly: the net flow out of a region D is equal to the sum of the net flows out any set of regions in a subdivision of D . This is because any surface connecting adjoining regions contributes two opposite terms to the sum (see Figure 6.5).



Theorem 6.6.1: Divergence theorem

If $\mathbf{F} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is a vector field with continuous partial derivatives and D is a region in \mathbb{R}^3 bounded by a piecewise smooth surface $S = \partial D$, then

$$\overbrace{\iiint_D \nabla \cdot \mathbf{F} dV}^{\text{net flow density integrated over } D} = \overbrace{\iint_{\partial D} \mathbf{F} \cdot d\mathbf{A}}^{\text{total flow out through } \partial D}.$$

Example 6.6.1

Verify that the divergence theorem holds in the case where $\mathbf{F} = \langle x^2, 3z^2, 2z^2 + y^2 \rangle$ and $D = [0, 1]^3$.

Solution

The divergence of \mathbf{F} is $2x + 4z$, so the divergence theorem asserts that

$$\iiint_D (2x + 4z) dV = \iint_{\partial D} \langle x^2, 3z^2, 2z^2 + y^2 \rangle \cdot d\mathbf{A}.$$

The left-hand side equals

$$\int_0^1 \int_0^1 \int_0^1 (2x + 4z) dx dy dz = 3.$$

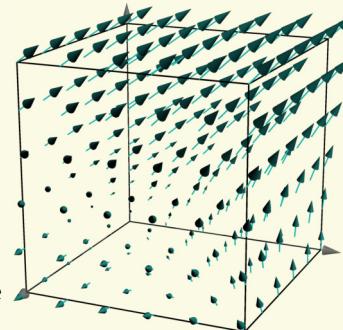
To evaluate the right-hand side directly, we split the boundary of the cube into its six square faces. For the top face, we get

$$\iint_{\text{top face}} \mathbf{F} \cdot d\mathbf{A} = \iint_{\text{top face}} \langle x^2, 3z^2, 2z^2 + y^2 \rangle \cdot \langle 0, 0, 1 \rangle dA = \int_0^1 \int_0^1 (2 + y^2) dx dy = \frac{7}{3},$$

where we've substituted 1 for z since $z = 1$ for every point in the top face. Likewise, the integral over the bottom face is

$$-\int_0^1 \int_0^1 (0 + y^2) dx dy = -\frac{1}{3},$$

where the negative sign comes from the fact that the outward-pointing normal on the bottom face is



$\langle 0, 0, -1 \rangle$.

Similarly, the integral over the $x = 1$ face is

$$\int_0^1 \int_0^1 1 \, dy \, dz = 1,$$

while the $x = 0$ face contributes 0. The $y = 1$ face yields

$$\int_0^1 \int_0^1 3z^2 \, dx \, dz = 1,$$

and the $y = 0$ face gives $-\int_0^1 \int_0^1 3z^2 \, dx \, dz = -1$. Indeed,

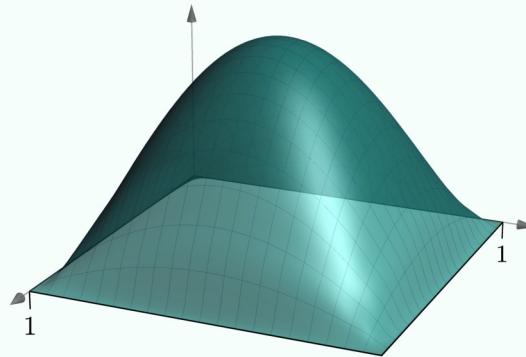
$$3 \stackrel{?}{=} \frac{7}{3} + \left(-\frac{1}{3}\right) + 1 + 0 + 1 + (-1).$$

Exercise 6.6.1

Verify the divergence theorem in the case where $\mathbf{F} = \langle x^2, y^2, z^2 \rangle$ and S is the unit sphere.

Exercise 6.6.2

Consider the vector field $\mathbf{F} = \langle x^3, xz, 1 - 3zx^2 \rangle$. Verify that the divergence of \mathbf{F} is zero everywhere. Then use the divergence theorem to calculate the flow of \mathbf{F} through the surface S shown. Note that this is not a closed surface: it excludes the square $[0, 1]^2 \times \{0\}$ on the bottom.



6.7 Stokes' theorem

The planar domain D in Green's theorem can be thought of as a surface S , in which case the conclusion of Green's theorem can be written as

$$\iint_S \nabla \times \mathbf{F} \cdot d\mathbf{A} = \int_{\partial S} \mathbf{F} \cdot d\mathbf{r}.$$

The argument for Green's theorem then applies even if S doesn't lie in a plane, because Theorem 6.5.1 tells us that $\nabla \times \mathbf{F} \cdot \mathbf{n} dA = \nabla \times \mathbf{F} \cdot dA$ yields the circulation around a small patch of the surface S . As we discussed for Green's theorem, circulation is additive: if you divide a surface into many small patches and sum the circulations around all of them, you get the circulation around the boundary of the surface. Thus we arrive at a generalization of Green's theorem known as *Stokes' theorem*.*

* If you take a differential geometry course, you will learn a far more general result of the same name which implies the divergence theorem and Theorem 6.7.1 as special cases

* The orientation used for the surface integral must be compatible with the orientation on ∂S used for the line integral, as in Figure 6.6. See also Example 6.7.2

Theorem 6.7.1: Stokes' theorem

If $\mathbf{F} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is a vector field with continuous partial derivatives and S is a surface in \mathbb{R}^3 , then*

$$\underbrace{\iint_S \nabla \times \mathbf{F} \cdot d\mathbf{A}}_{\text{circulation density integrated over } S} = \underbrace{\int_{\partial S} \mathbf{F} \cdot d\mathbf{r}}_{\text{circulation around } \partial S}.$$

Example 6.7.1

Let $\mathbf{F} = \langle x \sin(\pi y), e^x, -\cos(\pi z) \rangle$. Find the flow of $\nabla \times \mathbf{F}$ up through the surface shown in Exercise 6.6.2.

Solution

Stokes' theorem tells us that the flow of $\nabla \times \mathbf{F}$ upwards through S is equal to the line integral of \mathbf{F} around ∂S in the counterclockwise direction (see Figure 6.6). Since ∂S consists of four line segments, we calculate $\int \mathbf{F} \cdot d\mathbf{r}$ along each edge and sum the results. Integrating from $(0, 0, 0)$ to $(1, 0, 0)$, we get

$$\int_0^1 \mathbf{F}(x, 0, 0) \cdot \langle 1, 0, 0 \rangle dx = 0.$$

From $(1, 0, 0)$ to $(1, 1, 0)$, we get

$$\int_0^1 \mathbf{F}(1, y, 0) \cdot \langle 0, 1, 0 \rangle dy = e.$$

From $(1, 1, 0)$ to $(0, 1, 0)$, we get

$$\int_0^1 \mathbf{F}(x, 1, 0) \cdot \langle -1, 0, 0 \rangle dx = 0.$$

And finally from $(0, 1, 0)$ back to the origin, we get

$$\int_0^1 \mathbf{F}(0, y, 0) \cdot \langle 0, -1, 0 \rangle dy = -1.$$

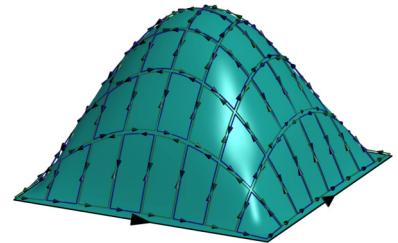


Figure 6.6 The circulation of \mathbf{F} (not shown) around each patch sums to the circulation around the boundary of the surface

So altogether the circulation of \mathbf{F} around the boundary of S is $e - 1$.

Theorem 6.7.1 implies that for any vector field \mathbf{F} , a surface can be deformed without changing the flow of $\nabla \times \mathbf{F}$ through it, as long as it is deformed in such a way that its boundary is preserved:

Observation 6.7.1

In the context of Stokes' theorem, if S_1 and S_2 are surfaces whose boundaries are the same, then

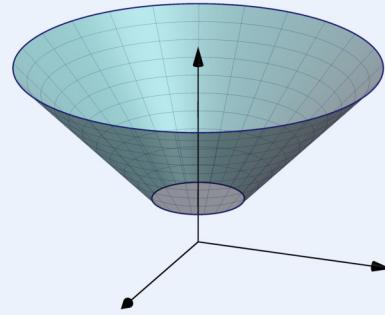
$$\iint_{S_1} (\nabla \times \mathbf{F}) \cdot d\mathbf{A} = \iint_{S_2} (\nabla \times \mathbf{F}) \cdot d\mathbf{A}.$$

Exercise 6.7.1

Redo Example 6.7.1 using Observation 6.7.1 in place of the calculation of the line integral around ∂S .

Example 6.7.2

Verify Stokes' theorem for the portion of the cone $z = \sqrt{x^2 + y^2}$ between the planes $z = 1$ and $z = 4$ and the vector field $\mathbf{F} = \langle x^2y, z, y \rangle$.



Solution

We parametrize the cone as $\mathbf{r}(u, v) = \langle u \cos v, u \sin v, u \rangle$ as (u, v) ranges over $[1, 4] \times [0, 2\pi]$. We also calculate $\nabla \times \mathbf{F} = \langle 0, 0, -x^2 \rangle$. Then applying the parametrization formula for flow (Theorem 6.4.2), we get

$$\iint_S \nabla \times \mathbf{F} \cdot d\mathbf{A} = \int_0^{2\pi} \int_1^4 \langle 0, 0, -u^2 \cos^2 v \rangle \cdot \langle -u \cos v, -u \sin v, u \rangle du dv = -\frac{255\pi}{4}.$$

Since the normal vector $\langle -u \cos v, -u \sin v, u \rangle$ points in/up for all $(u, v) \in [1, 4] \times [0, 2\pi]$, the value $-\frac{255\pi}{4}$ represents the flow of the curl \mathbf{F} from the *outside* of the cone to the *inside*. We should be careful to orient the boundary arcs in a manner which is compatible with this orientation. We can do this by following the derivation of Stokes' theorem: we subdivide the surface into many small patches, and we orient each little boundary in the counterclockwise direction around the normal vector used in the surface integral (in this case, in/up). Then the boundary of S is obtained by putting together the oriented segments whose contributions aren't canceled (see the small oriented red loop in Figure 6.7). We can see from this procedure that the top edge of the surface is oriented counterclockwise when viewed from above.

Since the bottom edge of each patch is oriented *clockwise* when viewed from above, the bottom edge of the surface should be oriented clockwise. So we can parametrize the top edge as $\langle 4 \cos t, 4 \sin t, 4 \rangle$ and the bottom edge as* $\langle \cos t, -\sin t, 1 \rangle$, in both cases as t ranges from 0 to 2π . We get

$$\int_{\partial S} \mathbf{F} \cdot d\mathbf{r} = \int_0^{2\pi} \mathbf{F}(4 \cos t, 4 \sin t, 4) \cdot \langle -4 \sin t, 4 \cos t, 0 \rangle dt + \int_0^{2\pi} \mathbf{F}(\cos t, -\sin t, 1) \cdot \langle -\sin t, -\cos t, 0 \rangle dt$$

which works out to

$$\begin{aligned} \int_0^{2\pi} [(4 \cos t)^2 (4 \sin t) (-4 \sin t) + (4)(4 \cos t) + 0] dt + \int_0^{2\pi} [\cos^2 t (-\sin t)^2 + (1)(-\cos t) + 0] dt \\ = -64\pi + \frac{\pi}{4} = -\frac{255\pi}{4}, \end{aligned}$$

as desired.

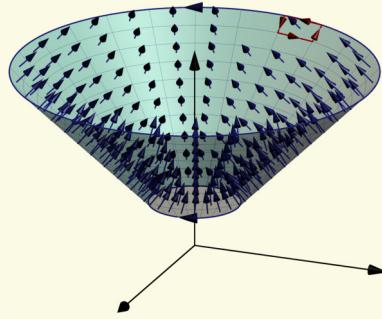


Figure 6.7 The normal vectors for the surface integral and the corresponding orientations of the boundary arcs in Stokes' theorem

We reverse the orientation by replacing t with $-t$, and then we can use the identities $\cos(-t) = \cos t$ and $\sin(-t) = -\sin t$

The following consequence of Stokes' theorem tells us that the flow of the curl of a vector field through a closed surface (such as a sphere, a donut, or a rectangular prism) is zero:

Observation 6.7.2

If \mathbf{G} is a vector field which is equal to the curl of some vector field with continuous partial derivatives and if S is a piecewise smooth closed surface (in other words, a surface with no boundary), then the flow of \mathbf{G} through S is zero.

Exercise 6.7.2

Suppose that S is the surface consisting of the points on the sphere $x^2 + y^2 + z^2 = 1$ which are not inside the sphere $x^2 + y^2 + (z+1)^2 = 1$. Find $\iint_S \nabla \times \mathbf{F} \cdot d\mathbf{A}$, where $\mathbf{F} = \langle yz, x, e^{xyz} \rangle$.

Exercise 6.7.3

Suppose $\mathbf{F} = \langle xy, y, xz \rangle$. Find $\iint_S \nabla \times \mathbf{F} \cdot d\mathbf{A}$ where S is the portion of the unit sphere $x^2 + y^2 + z^2 = 1$ in the first octant.

Colophon

This text was typeset with Lua \LaTeX , using `tcolorbox` and a version of the `mathpazo` package's Palatino fonts which was modified to borrow Greek symbols from Utopia and blackboard bold symbols from Computer Modern. The cover art was rendered using TikZ.

The figures are all produced in Asymptote and are included using the `asymptote` \LaTeX package. All the files necessary to produce this document are available at github.com/sswatson.

A.1 Review

A.1.1 SETS AND FUNCTIONS

A **set** is a collection of elements. These elements can be numbers, points, shapes, vectors, other sets, whatever. For example,

$$A = \{1, 4, 9\}$$

is the set consisting of the positive, single-digit perfect squares. The main thing you can do with a set is check whether a particular element is in it. For example, we say that $1 \in A$ (read “1 is an element of A ”), while $2 \notin A$ (“2 is not an element of A ”).

Some sets with standard and specially typeset names include

- \mathbb{R} , the set of real numbers,
- \mathbb{Q} , the set of rational numbers,
- \mathbb{Z} , the set of integers, and
- \mathbb{N} , the set of natural numbers.

Subsets and set equality

We say that $A \subset B$ (read “ A is a **subset** of B ”) if every element of A is an element of B . For example,

$$\{1, 4, 9\} \subset \{1, 2, 3, 4, 9, 10\}.$$

We say that two sets A and B are **equal** if $A \subset B$ and $B \subset A$. Note that

$$\{1, 1, 2\} = \{1, 2\} = \{2, 1\}.$$

since each element of each set is in the others. Thus we can see that sets “don’t care” about repeated elements or order. All that matters is what is in and what is not. It is customary to writesets with repeats omitted, for clarity.

Intersections and unions

We write $A \cap B$, the **intersection** of A and B , for the set of all the elements that are in both A and B . So, for example,

$$\{1, 4, 9\} \cap \{x \in \mathbb{R} : x^2 > 15\} = \{4, 9\}.$$

That second set on the left-hand side, which is written in *set-builder* notation, is read as “the set of all real numbers x such that the square of x is greater than 15”.

We write $A \cup B$, the **union** of A and B , for the set of all the elements that are in either A or B . So, for example,

$$\{1, 4, 9\} \cup \{1, 9, 25\} = \{1, 4, 9, 25\}.$$

Functions

If A and B are sets, then a function $f : A \rightarrow B$ is a rule that assigns a single element of B to each element of A . The set A is called the **domain** of f and B is called the **codomain** of f . Given a subset A' of A , we define the **image** $f(A')$ to be

$$f(A') = \{b \in B : \text{there exists } a \in A' \text{ so that } f(a) = b\}. \quad (\text{A.1.1})$$

This is the set of all elements of B that get mapped to from some element of A' . The **range** of f is defined to be the set $f(A)$, which contains all the elements of B that get mapped to at least once.

Similarly, if $B' \subset B$, then the **preimage** $f^{-1}(B')$ of B' is defined by

$$f^{-1}(B') = \{a \in A : f(a) \in B'\}.$$

This is the subset of A consisting of every element of A that maps to some element of B' .

A function f is **injective** if no two elements in the domain map to the same element in the codomain; in other words if $f(a) = f(a')$ implies $a = a'$.

A function f is **surjective** if the range of f is equal to the codomain of f ; in other words, if $b \in B$ implies that there exists $a \in A$ with $f(a) = b$.

A function f is **bijective** if it is both injective and surjective. This means that for every $b \in B$, there is exactly one $a \in A$ such that $f(a) \in b$. If f is bijective, then the function from B to A that maps $b \in B$ to the element $a \in A$ that satisfies $f(a) = b$ is called the **inverse** of f .

If $f : A \rightarrow B$ and $A' \subset A$, then the **restriction** of f to A' is the function $f|_{A'} : A' \rightarrow B$ defined by $f|_{A'}(x) = f(x)$ for all $x \in A'$.

If $f : A \rightarrow B$ and $g : B \rightarrow C$, then the function $g \circ f$ which maps $x \in A$ to $g(f(x)) \in C$ is called the **composition** of g and f .

If the rule defining a function is sufficiently simple, we can describe the function using **anonymous function notation**. For example, $x \in \mathbb{R} \mapsto x^2 \in \mathbb{R}$, or $x \mapsto x^2$ for short, is the squaring function from \mathbb{R} to \mathbb{R} . Note that bar on the left edge of the arrow, which distinguishes the arrow in anonymous function notation from the arrow between the domain and codomain of a named function.

Cartesian product

The **Cartesian product** of two sets A and B , denoted $A \times B$, is the set of all pairs (a, b) where $a \in A$ and $b \in B$. For example, $[0, 3] \times [0, 2]$ is a rectangle in the plane. We sometimes use exponents for a Cartesian product of a set with itself. Thus $[0, 1]^2$ is a unit square in \mathbb{R}^2 , and $[0, 1]^3$ is a unit cube in \mathbb{R}^3 .

This appendix provides a streamlined presentation of trig which is intended to provide enough starting off points to recover everything else you need.

Trig Review

- Cosine and sine.** The basic trig functions are $\cos \theta$ and $\sin \theta$. The most important definition of these functions is the following: the cosine of an angle θ is equal to the **x-coordinate** of the point obtained by rotating $(1, 0)$ an angle of θ about the origin. Sine is the same, but with the y -coordinate instead of x .

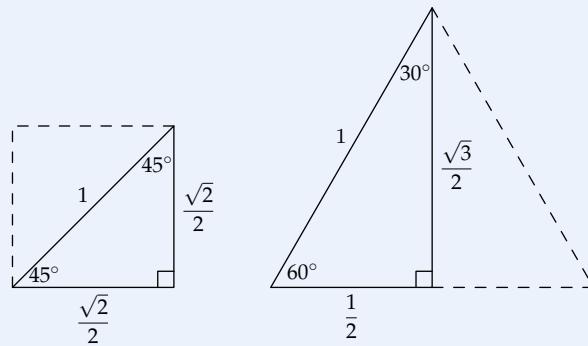
!!!

This idea bears repeating: **the point on the unit circle obtained by rotating $(1, 0)$ an angle θ about the origin is equal to $(\cos \theta, \sin \theta)$, by definition of cosine and sine.**

- The other ones.** The other four trig functions are simply abbreviations for various combinations of sine and cosine:

$$\begin{aligned}\sin \theta &= \sin \theta & \sec \theta &= \frac{1}{\cos \theta} \\ \cos \theta &= \cos \theta & \csc \theta &= \frac{1}{\sin \theta} \\ \tan \theta &= \frac{\sin \theta}{\cos \theta} & \cot \theta &= \frac{\cos \theta}{\sin \theta}\end{aligned}$$

- Special right triangles.** The following two triangles, each half of a regular polygon, can be handy for evaluating trig functions at special angles.



- Pythagorean identities.** The famous identity $\sin^2 \theta + \cos^2 \theta = 1$ follows from the definition of sine and cosine combined with the Pythagorean theorem. Dividing both sides of this equation by $\sin^2 \theta$ or $\cos^2 \theta$, we get

$$\tan^2 \theta + 1 = \sec^2 \theta \quad \text{and} \quad 1 + \cot^2 \theta = \csc^2 \theta.$$

- Sum-angle formulas.** The sine sum-angle formula is worth memorizing: for all α and β , we have

$$\boxed{\sin(\alpha + \beta) = \sin \alpha \cos \beta + \sin \beta \cos \alpha}.$$

The cosine sum-angle formula is worth memorizing too, although it can be derived fairly easily from the sine formula by substituting $\frac{\pi}{2} - \alpha$ for α and $-\beta$ for β . We get

$$\cos(\alpha + \beta) = \cos \alpha \cos \beta - \sin \alpha \sin \beta.$$

From the above identities, we can derive many others. For example, setting $\alpha = \beta$ in the cosine sum-angle formula, we get

$$\cos 2\alpha = \cos^2 \alpha - \sin^2 \alpha.$$

Substituting $\cos^2 \alpha = 1 - \sin^2 \alpha$, we find that

$$\cos 2\alpha = 1 - 2\sin^2 \alpha.$$

which can be solved to express $\sin^2 \alpha$ in terms of $\cos 2\alpha$.

A.1.3 SUMMATION NOTATION

Some shorthand is essential for writing sums with many terms. Perhaps the most common approach is to use ellipses:

$$1 + 2 + 3 + \cdots + 99 + 100 = 5050.$$

However, this approach is not ideal because the reader is left to infer the pattern.

When more precision is required, we would like to specify a formula for the k th term as well as a starting and ending value. For example, the sum

$$1 + 4 + 9 + 16 + \cdots + 100$$

can be written as “the sum of k^2 as k ranges from 1 to 10”. The math notation that has been adopted for abbreviating this English phrase is the following:

$$\sum_{k=1}^{10} k^2$$

The variable k is called a *dummy variable*, since it is only there as a way to specify a formula for generating the terms. We could change each k to a different symbol without changing the essential meaning, which is “sum the first 10 positive perfect squares”.

Exercise A.1.1

Find $\sum_{k=1}^5 \frac{1}{k(k+1)}$.

Exercise A.1.2

Express $\frac{1}{1} + \frac{3}{4} + \frac{5}{9} + \frac{7}{16} + \frac{9}{25} + \cdots$ using summation notation. Hint: your upper limit will be ∞ .

Exercise A.1.3

Find $\sum_{k=1}^5 \sum_{j=1}^k j$.

A.2 Reference

A.2.1 VISUALIZING FUNCTIONS

Graphical visualization is an important conceptual tool for reasoning about the behavior of functions. There are a variety of different methods for visualizing functions (see the table on the next page for pictures):

Function Visualization Methods

1. **Graphs.** The graph of a function $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$ is the set of points of the form $(x, f(x))$, where x is in the domain of f . The graph of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}^1$ is the set of points of the form $(x, y, f(x, y))$, where (x, y) is in the domain of f . The graph uses one or two dimensions for the input and one dimension for the output, so it only works (as a visualization tool) for $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ if $m + n \leq 3$.

The graph involves no loss of information; in principle, you can read off anything you want to know about a function from its graph. It depicts the domain and the codomain in the same picture.

2. **Level sets.** A level set of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is the solution set of an equation of the form $f(x, y) = c$, where c is some constant. For example, the $c = 1$ level set of the function $x^2 + y^2$ is the unit circle. The level set of a function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ is typically a *surface*. For example, the level sets of $x^2 + y^2 + z^2$ are spheres.

We can visualize a function by drawing its level sets. However, there are a couple drawbacks: we have to choose a discrete number of level sets to draw, and the picture doesn't tell us which c value corresponds to each level set, unless we draw that information in with colors or labels. When we visualize a function in this way, we are looking only at the *domain* of the function.

3. **Grid lines.** For a function T from \mathbb{R}^2 to \mathbb{R}^2 , we can understand T as a transformation which moves points in the plane to other points in the plane, and we can visualize this geometric action by drawing the images of various grid lines under T . This visualization is drawn on the codomain side and loses some information about which grid lines match to which images.

4. **Traces.** We can visualize a function \mathbf{r} from \mathbb{R}^1 to \mathbb{R}^2 or \mathbb{R}^3 by highlighting every point in \mathbb{R}^2 or \mathbb{R}^3 which is equal to $\mathbf{r}(t)$ for some $t \in \mathbb{R}$. This set of points is called the *trace* of \mathbf{r} .

The trace is drawn entirely on the codomain side, which means that this visualization lacks information about which t value or values mapped to each highlighted point.

5. **Vector fields.** For functions from \mathbb{R}^2 to \mathbb{R}^2 or \mathbb{R}^3 to \mathbb{R}^3 , we can visualize them by interpreting the output value as a vector and depicting a discrete set of these vectors as small arrows drawn in place at the corresponding input values. Doing this requires scaling the vectors down so the picture doesn't get chaotic. This visualization incorporates inputs and outputs in the same picture, and information is lost about the absolute size of each vector and about what happens between the discrete set of input values shown.

Table A.1 below shows examples of each type of visualization, with the input (domain) dimension varying by row and the output (codomain) dimension by column. An example of a common type of visualization is shown for each input-output pair of dimensions.

In a couple cases, the method shown isn't the only one in common use: a function from \mathbb{R}^2 to \mathbb{R}^2 can also

be drawn as a vector field (particularly if one is thinking of the outputs as vectors rather than points in \mathbb{R}^2), and the level set method shown for a function from \mathbb{R}^3 to \mathbb{R}^1 can also be applied to a function from \mathbb{R}^2 to \mathbb{R}^1 .

Conventional names are used for each function; note that these vary by input and output dimension. Functions from \mathbb{R}^3 to \mathbb{R}^2 don't have a dedicated visualization method, although one could visualize each component of such a function separately, or identify the codomain \mathbb{R}^2 with the xy -plane in \mathbb{R}^3 to make a vector field representation.

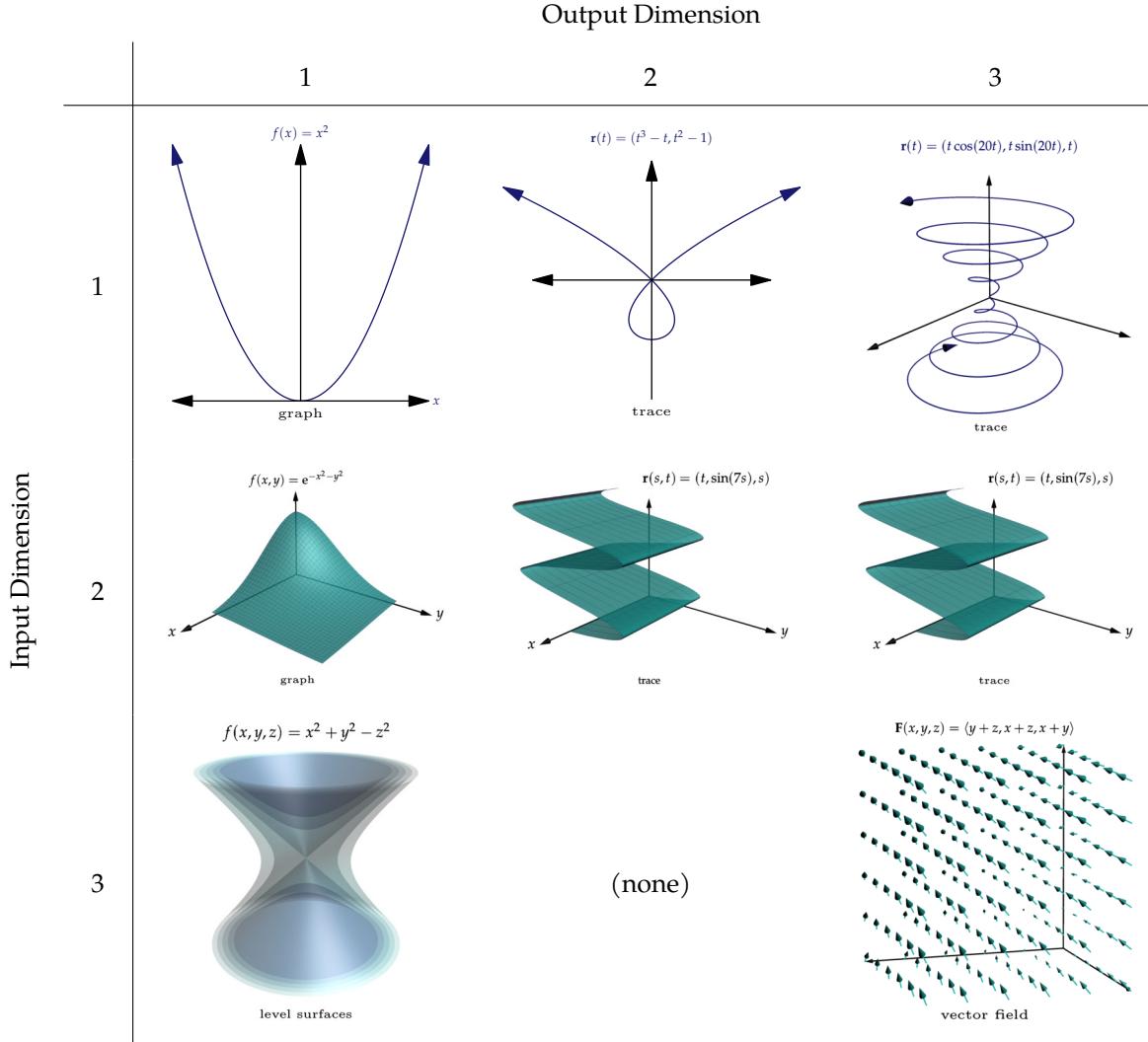


Table A.1 Different methods of visualizing functions from \mathbb{R}^n to \mathbb{R}^m , arranged by (n, m) pairs

Polar to Cartesian

$$\begin{cases} x = r \cos \theta \\ y = r \sin \theta \end{cases}$$

Cylindrical to Cartesian

$$\begin{cases} x = r \cos \theta \\ y = r \sin \theta \\ z = z \end{cases}$$

Spherical to Cartesian

$$\begin{cases} x = \rho \cos \theta \sin \phi \\ y = \rho \sin \theta \sin \phi \\ z = \rho \cos \phi \end{cases}$$

Area/Volume differentials

$$\begin{cases} dA = r dr d\theta \\ dV = r dr d\theta dz \\ dV = \rho^2 \sin \phi d\rho d\phi d\theta \end{cases}$$

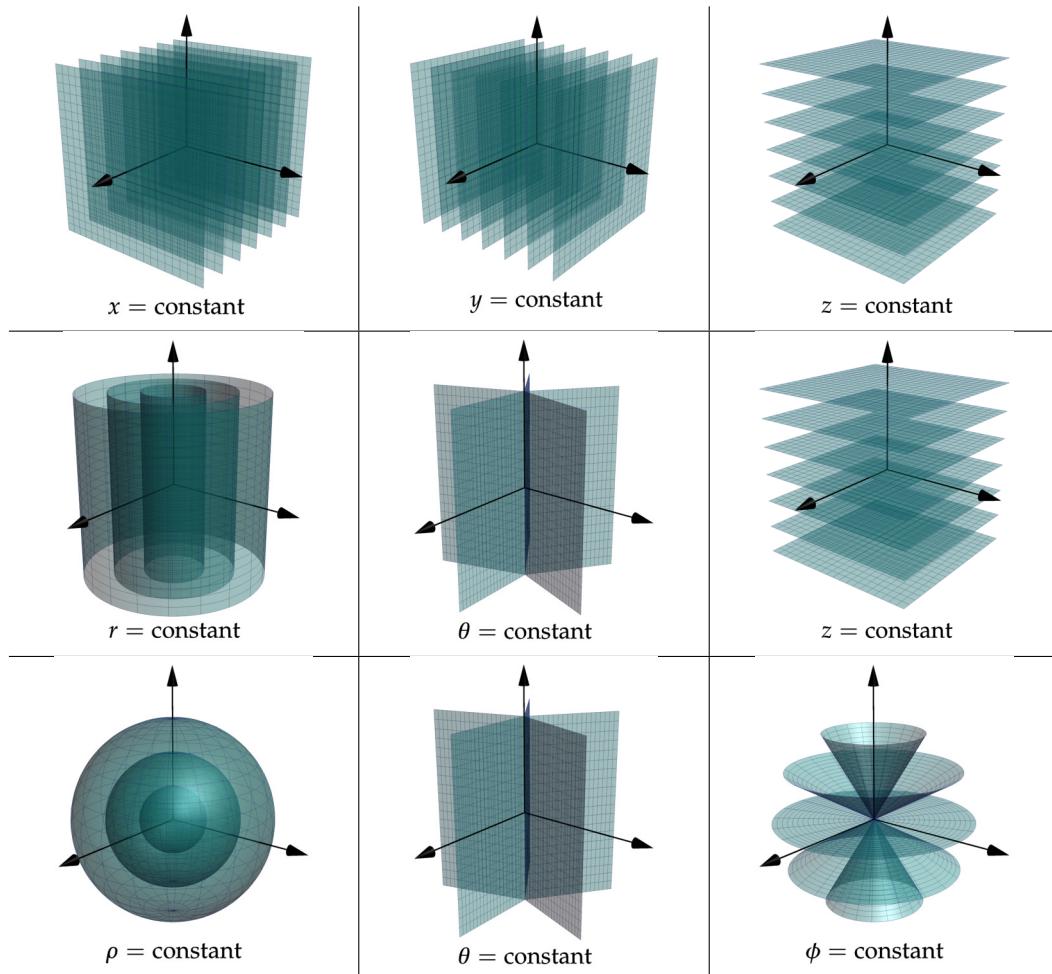


Table A.2 Level surfaces for each coordinate in the rectangular, cylindrical, and spherical systems

A.3 Technical Appendix

A.3.1 THE CONVENTIONAL DEFINITION OF A LIMIT

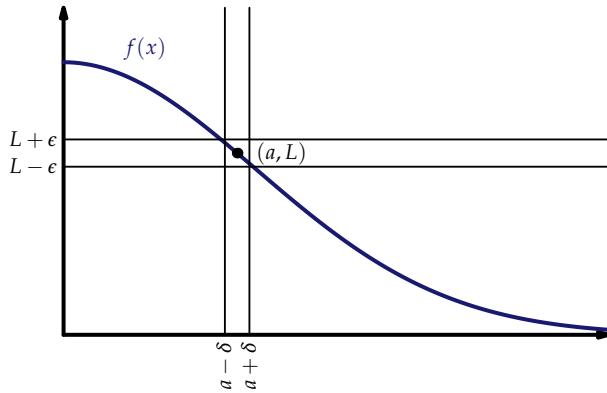


Figure A.1 The ϵ - δ definition of a limit.

Recall that $f(x)$ converges to L as $x \rightarrow a$ if $f(x)$ can be made as close to L as desired by restricting x to a sufficiently small neighborhood of a . More precisely, $f(x)$ converges to L as $x \rightarrow a$ if and only if for every $\epsilon > 0$, there is $\delta > 0$ so that $|f(x) - L| < \epsilon$ for all x satisfying $0 < |x - a| < \delta$ (see Figure A.1).

It can be helpful to think of this definition as a game against an adversary: the adversary chooses a positive real number ϵ which can be as small as they like. Then, after seeing the ϵ value, you get to choose a number $\delta > 0$, as small as you like. Finally, the adversary chooses an x value in the interval $(a - \delta, a + \delta)$ which is not equal to a . If it turns out that $|f(x) - L| \geq \epsilon$, then the adversary wins. Otherwise, you win. We call this the **limit game**.

* We abbreviate " $f(x)$ converges to L as $x \rightarrow a$ " to $\lim_{x \rightarrow a} f(x) = L$

If you have a strategy for winning this game no matter how the adversary plays, then the limit of $f(x)$ as x approaches a exists and equals L . If the adversary has a strategy for winning, then it is not true that $\lim_{x \rightarrow a} f(x) = L$ (either because the limit does not exist or because the limit exists and equals a number other than L).

Exercise A.3.1

Show that $\lim_{x \rightarrow 3} \frac{x^2 - 9}{x - 3} = 6$ by explaining the winning strategy in the limit game.

Exercise A.3.2

Suppose that f is a function from \mathbb{R} to \mathbb{R} and that $a \in \mathbb{R}$. Show that if $f(x)$ converges to L as $x \rightarrow a$ and $f(x)$ converges to L' as $x \rightarrow a$, then $L = L'$. This fact is called *uniqueness of limits*.

We reviewed the definition of a limit for a single-variable function so we could think about how to generalize the definition for a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. For simplicity, let's take $n = 2$. What should it mean to say $\lim_{(x,y) \rightarrow (a,b)} f(x, y) = L$? The only aspect of the definition that requires revision is the part about x being within δ of a , and we can use standard Euclidean distance to compare (x, y) to (a, b) . This leads to the following definition.

Definition A.3.1: Limit of a function of two variables

We say $\lim_{(x,y) \rightarrow (a,b)} f(x, y) = L$ if and only if for every $\epsilon > 0$, there is $\delta > 0$ so that $|f(x, y) - L| < \epsilon$ for all (x, y) satisfying $0 < \sqrt{(x - a)^2 + (y - b)^2} < \delta$.

* $B((a, b), r)$
means the ball
of radius r centered
at (a, b)

[†]Shadow here
means the
set of points
 $(x, y, f(x, y))$
where $(x, y) \in$
 $B((a, b), \delta)$

One way to think about this definition is to consider the shadow[†] of the disk* $B((a, b), \delta)$ on the graph (see Figure A.2). The limit exists and equals L if for every ϵ , there exists δ small enough that this shadow lies entirely in the slab $L - \epsilon < z < L + \epsilon$.

Exercise A.3.3

Show that Definitions 4.1.2 and A.3.1 are equivalent.

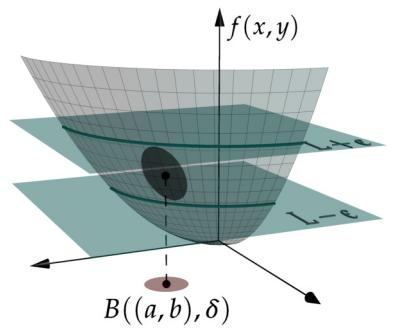


Figure A.2 The definition of a limit
for a two-variable function

* To appreciate the difference, note that the path contains information about how quickly the trajectory is traced out, while knowing the trajectory only tells us which points were passed through.

In Section 3.2.2, we defined arclength only for piecewise differentiable paths. In this appendix, we provide a more natural and general definition which also resolves the question of parameterization independence. Recall that a **path** in \mathbb{R}^n is a continuous function from some interval $[a, b]$ to \mathbb{R}^n . We distinguish a path from its **trajectory**, which is the set of points that the path passes through.*

Consider an arbitrary path $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^3$. Consider any finite sequence S of points P_1, P_2, \dots, P_n in order along the path, and connect them with line segments. The length of this path of line segments is

$$|\overrightarrow{P_0P_1}| + |\overrightarrow{P_1P_2}| + \cdots + |\overrightarrow{P_{n-1}P_n}|,$$

and we will call this quantity $\ell(S)$.

Exercise A.3.4

Although we haven't defined the length of \mathbf{r} yet, explain intuitively why $\ell(S)$ *should* be less than the length of \mathbf{r} .

Also, explain qualitatively how one might choose the points in S so as to make $\ell(S)$ very close to the length of \mathbf{r} (again assuming an intuitive notion of the length of \mathbf{r}).

Motivated by Exercise A.3.4, we make the following definition.

Definition A.3.2: Arclength, for a general path

The **length** of a path is defined to be the smallest real number which is greater than or equal to $\ell(S)$ for every finite sequence S of points along the path.

If no such number exists, then we define the length of the path to be $+\infty$.

Since S is just a sequence of points along the path, the set of sequences under consideration in Definition A.3.2 is the same regardless of the parameterization.

Theorem A.3.1: Parameterization Independence

If \mathbf{r}_1 and \mathbf{r}_2 are two paths which visit the same points in the same order*, then their lengths are the same.

* Rigorously,
this means that
there exists an
increasing, con-
tinuous func-
tion σ from the
domain of \mathbf{r}_1
to the domain
of \mathbf{r}_2 such that
 $\mathbf{r}_1 = \mathbf{r}_2 \circ \sigma$

The following theorem shows that the formula given in Section 3.2.2 indeed agrees with Definition A.3.2.* For simplicity, we assume that $\mathbf{r}'(t)$ exists for all t and that \mathbf{r}' is continuous on $[a, b]$ —in other words, that \mathbf{r} is *continuously differentiable*.

* Although we took it as a definition in Section 3.2.2, the formula in Definition 3.2.1 must be treated as a *theorem* in the present context, since we should not have two different definitions for the same term.

Theorem A.3.2

Suppose $\mathbf{r} : [a, b] \rightarrow \mathbb{R}^3$ is continuously differentiable. Then the length of \mathbf{r} is equal to $\int_a^b |\mathbf{r}'(t)| dt$.

Proof (idea)

Recall from single-variable calculus that $\int_a^b |\mathbf{r}'(t)| dt$ is defined to be a limit of Riemann sums of the form

$$|\mathbf{r}'(s_1)|(t_1 - t_0) + |\mathbf{r}'(s_2)|(t_2 - t_1) + \cdots + |\mathbf{r}'(s_n)|(t_n - t_{n-1}) \quad (\text{A.3.1})$$

where (i) $a = t_0 < t_1 < \dots < t_n = b$ is a partition and (ii) s_k is any number in the interval $[t_{k-1}, t_k]$, for each k from 1 to n . The limit is taken as the *mesh* $\{t_0, \dots, t_n\}$ (that is, the largest gap between consecutive numbers in the partition) tends to zero.

The main idea of the proof is to show that the segment path lengths $\ell(S)$ in Definition A.3.2 and the Riemann sums (A.3.1) are very close. Roughly speaking, the reason is that $|\mathbf{r}'(s_1)|(t_1 - t_0)$ is very close to $|\mathbf{r}(t_1) - \mathbf{r}(t_0)|$ (and similarly for the other terms) by the **mean value theorem**.

For each k from 1 to n , we let t_k be the time when the path visits P_k . Define x and y to be the components of \mathbf{r} ; then the mean-value theorem implies that there is some value s_k between t_{k-1} and t_k such that

$$x(t_k) - x(t_{k-1}) = x'(s_k)(t_k - t_{k-1}).$$

Similarly, there are values \tilde{s}_k and \hat{s}_k between t_{k-1} and t_k such that

$$\begin{aligned} y(t_k) - y(t_{k-1}) &= y'(\tilde{s}_k)(t_k - t_{k-1}), \text{ and} \\ z(t_k) - z(t_{k-1}) &= z'(\hat{s}_k)(t_k - t_{k-1}). \end{aligned}$$

Squaring these three equations, summing them, and taking the square root of both sides, we get

$$|\mathbf{r}(t_k) - \mathbf{r}(t_{k-1})| = \sqrt{x'(s_k)^2 + y'(\tilde{s}_k)^2 + z'(\hat{s}_k)^2} (t_k - t_{k-1}).$$

The expression $\sqrt{x'(s_k)^2 + y'(\tilde{s}_k)^2 + z'(\hat{s}_k)^2}$ on the right-hand side can be made as close as desired to $\sqrt{x'(s_k)^2 + y'(s_k)^2 + z'(s_k)^2}$ for all partitions with suitably small mesh. Thus for any $\epsilon > 0$, we have

$$\begin{aligned} \left(\sqrt{x'(s_k)^2 + y'(s_k)^2 + z'(s_k)^2} - \epsilon \right) (t_k - t_{k-1}) &\leq |\mathbf{r}(t_k) - \mathbf{r}(t_{k-1})| \\ &\leq \left(\sqrt{x'(s_k)^2 + y'(s_k)^2 + z'(s_k)^2} + \epsilon \right) (t_k - t_{k-1}) \end{aligned} \quad (\text{A.3.2})$$

for sufficiently fine-mesh partitions.

By definition of the Riemann integral, the quantity $\int_a^b |\mathbf{r}'(t)| dt + \epsilon$ exceeds the Riemann sum (A.3.1) for all partitions with sufficiently small mesh. Summing (A.3.2) over k , we see that $\int_a^b |\mathbf{r}'(t)| dt + \epsilon$ is greater than or equal to $\ell(S)$ for all sufficiently granular S , and thus* for all S . By Definition A.3.2, this means that the length of \mathbf{r} is no greater than $\int_a^b |\mathbf{r}'(t)| dt + \epsilon$. Since ϵ is arbitrary, this means that

$$\text{length of } \mathbf{r} \leq \int_a^b |\mathbf{r}'(t)| dt.$$

* Why are lengths for coarse-mesh partitions less than those for fine-mesh partitions?

Similarly, (A.3.2) tells us that for all $\epsilon > 0$, there is a sequence S of points on the path such that $\ell(S)$ exceeds $\int_a^b |\mathbf{r}'(t)| dt - \epsilon$. Thus

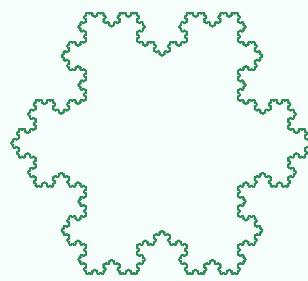
$$\text{length of } \mathbf{r} \geq \int_a^b |\mathbf{r}'(t)| dt.$$

These two inequalities imply the theorem.

Definition A.3.2, unlike Definition 3.2.1, applies to *fractal* curves:

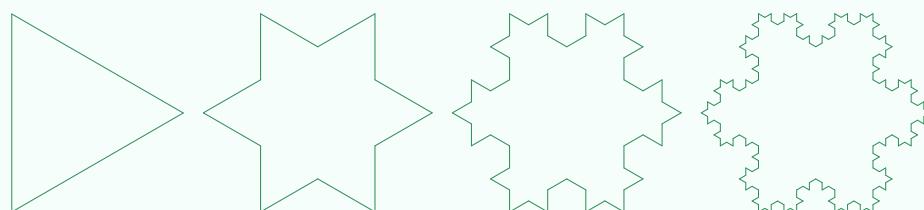
Exercise A.3.5: The Koch snowflake

The Koch snowflake is a fractal, meaning that it continues to look wiggly no matter how far you zoom in on it:



It is created by starting with an equilateral triangle with side length 1 and by repeating the following process: (1) Divide each line segment in the figure into three equal-length pieces. (2) For each line segment from Step 1, draw an equilateral triangle whose base is the middle segment from Step 1. (3) For each triangle you draw in Step 2, erase its base.

The result of the n th iteration of this procedure is shown below for $n = 0, 1, 2, 3$. The Koch snowflake is defined to be the set of all points which are eventually drawn and then never erased as this procedure is repeated for all $n = 0, 1, 2, \dots$



Find the length of the Koch snowflake.

In this section we prove Theorem 4.6.1, which we restate here for the reader's convenience.

Theorem A.3.3

Suppose that U is an open set in \mathbb{R}^2 and $f : U \rightarrow \mathbb{R}$ is a twice-differentiable function with a critical point at (a, b) . We define $D = (\partial_x^2 f \partial_y^2 f - [\partial_x \partial_y f]^2)(a, b)$. Then*

- (a) if $D > 0$ and $(\partial_x^2 f)(a, b) > 0$, then f has a local minimum at (a, b) ,
- (b) if $D > 0$ and $(\partial_x^2 f)(a, b) < 0$, then f has a local maximum at (a, b) , and
- (c) if $D < 0$, then f has a saddle point at (a, b) .

* When $D = 0$, this theorem doesn't tell us anything

Proof

We begin by investigating the simpler case where f is a quadratic polynomial. We can then leverage Taylor's theorem to apply that knowledge to the much more general class of twice-differentiable functions. It suffices to assume $(a, b) = (0, 0)$ and that $f(0, 0) = 0$, since $f(x, y)$ has the same local behavior at (a, b) as the function $f(x + a, y + b) - f(a, b)$ has at the origin.

Consider the quadratic polynomial

$$p(x, y) = ax^2 + bxy + cy^2 + dx + ey,$$

where a, b, c, d, e, f are real numbers. Our goal is to determine whether p has a local minimum, a local maximum, or neither, at each of its critical points. The main idea is to slice the graph of p along a line passing through the critical point in the direction of an arbitrary vector $\mathbf{v} \in \mathbb{R}^2$. If p has a local minimum along each of these lines, then it has a local minimum at the critical point. If p has a local maximum along each of these lines, then it has a local maximum at the critical point. If it has a local minimum in some directions and a local maximum in other directions, then it has neither a local min nor a local max at the critical point.

Note that $d = e = 0$, since d and e are the first-order partial derivatives of p at the origin. Fix a vector $\mathbf{v} = \langle v_1, v_2 \rangle$, and consider the function from \mathbb{R} to \mathbb{R} that maps t to $p(v_1 t, v_2 t)$. The shape of this function's graph is the same shape you get by slicing the graph of p along a vertical plane containing \mathbf{v} and passing through the critical point. We get

$$p(tv_1, tv_2) = t^2(av_1^2 + bv_1v_2 + cv_2^2). \quad (\text{A.3.3})$$

This single-variable function has a local min or max at $t = 0$ according to whether $av_1^2 + bv_1v_2 + cv_2^2$ is positive or negative. So,

- (a) if $av_1^2 + bv_1v_2 + cv_2^2$ is positive for all \mathbf{v} , then f has a local minimum at the origin
- (b) if $av_1^2 + bv_1v_2 + cv_2^2$ is negative for all \mathbf{v} , then f has a local maximum at the origin, and
- (c) if $av_1^2 + bv_1v_2 + cv_2^2$ is positive for some vectors \mathbf{v} and negative for others, then f has a saddle point.

We can determine whether $av_1^2 + bv_1v_2 + cv_2^2$ is always positive, always negative, or neither by assuming* that $v_1 \neq 0$, in which case we can define $m = v_2/v_1$ and rewrite the expression in the form $v_1^2(a + bm + cm^2)$. This expression, viewed as a quadratic function of m , is

- (a) positive for all m if and only if the discriminant $b^2 - 4ac$ is negative and $a > 0$,
- (b) negative for all m if and only if the discriminant $b^2 - 4ac$ is negative, and $a < 0$, and

* To be complete, the case $v_1 = 0$ must be checked separately

(c) positive for some values of m and negative for others if and only if $b^2 - 4ac$ is positive.
 By Taylor's theorem, f can be written as

$$f(x, y) = ax^2 + bxy + cy^2 + R(x, y),$$

where R is a function satisfying $\lim_{(x,y) \rightarrow (0,0)} \frac{R(x,y)}{|\langle x,y \rangle|^2} = 0$, and $a = \frac{1}{2}(\partial_x^2 f)(0,0)$, $b = (\partial_x \partial_y f)(0,0)$, and $c = \frac{1}{2}(\partial_y^2 f)(0,0)$. As we just figured out, $ax^2 + bxy + cy^2$ has a local min or max or saddle according to signs of a and $b^2 - 4ac = [\partial_x \partial_y f(0,0)]^2 - \partial_x^2 f(0,0)\partial_y^2 f(0,0)$. So the proof is done once we can show that the error term can be safely ignored. We will outline how to do this only for the case where $p(x, y) = ax^2 + bxy + cy^2$ has a local minimum at the origin. The basic idea is that value of $p(x, y)$ is increasing quadratically in all directions from its value at the origin, and $R(x, y)$ isn't nearly large enough to cancel this increase.

If we write an arbitrary point (x, y) as $(r \cos \theta, r \sin \theta)$, then (A.3.3) shows us that

$$p(x, y) = (x^2 + y^2)(a \cos^2 \theta + b \cos \theta \sin \theta + c \sin^2 \theta).$$

We've assumed we're in the case where $a \cos^2 \theta + b \cos \theta \sin \theta + c \sin^2 \theta$ is positive for all $\theta \in [0, 2\pi]$. By the extreme value theorem, this quantity has some positive minimum value which we'll call α . Since $\lim_{(x,y) \rightarrow (0,0)} \frac{R(x,y)}{|\langle x,y \rangle|^2} = 0$, we know that $R(x, y) \geq -\frac{1}{2}\alpha(x^2 + y^2)$ for (x, y) sufficiently close to the origin. Therefore, for such (x, y) , we have

$$f(x, y) \geq \alpha(x^2 + y^2) + R(x, y) \geq \alpha(x^2 + y^2) - \frac{\alpha}{2}(x^2 + y^2) > 0.$$

Thus f has a local minimum at the origin.

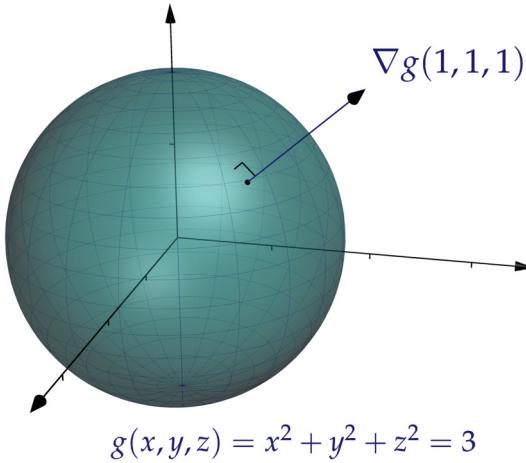


Figure A.3 The gradient of a function at a point is orthogonal to the level set of that function through that point

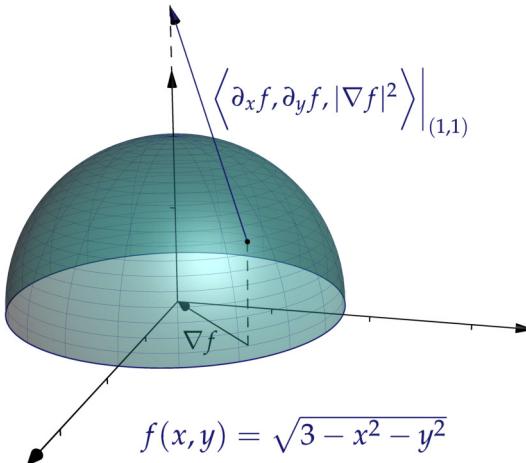


Figure A.4 The gradient of a function specifies the direction of maximum increase, which corresponds to moving along its graph in the steepest direction

You might have a couple different images come to mind when you think about gradients, one with a vector which is tangent to a surface, and the other with a vector normal to a surface. So is the gradient normal or tangent? Both!* To specify the picture, we have to specify *what function* we're taking the gradient of, and *what surface* our vector is tangent or normal to.

* Sort of: the tangent vector is only closely related to the gradient

The most important idea is that the gradient of a function at a point specifies its direction of maximum increase. This has two implications:

(1) (Figure A.3) If g is a differentiable function from \mathbb{R}^3 to \mathbb{R}^1 , then the gradient of g is orthogonal to the **level surface** of g at each point (x_0, y_0, z_0) .

(2) (Figure A.4) If f is a differentiable function from \mathbb{R}^2 to \mathbb{R}^1 , then the gradient of f at (x_0, y_0) tells us the direction *in the xy-plane* in which we should move so that the point $(x, y, f(x, y))$ on the **graph** of f moves in the direction of maximum z increase.*

In the second case, if we want a vector which is tangent to the graph of f , we need to include the third component (since the gradient only has two). This component should be chosen so that the vector is contained within the plane tangent to the graph at that point, which means we should substitute into the tangent plane equation

$$z - z_0 = (\partial_x f)(x_0, y_0)(x - x_0) + (\partial_y f)(x_0, y_0)(y - y_0).$$

Substituting $x - x_0 = (\partial_x f)(x_0, y_0)$ and $y - y_0 = (\partial_y f)(x_0, y_0)$, we find that the third component should be the squared gradient of f at (x_0, y_0) .

* If the graph of f corresponds to a hilly landscape, the gradient at each point only specifies the cardinal direction of maximum increase. The hiker moving in this direction naturally moves tangent to the hill, because gravity keeps them pinned to the graph

When we defined the integral of a function f defined on a two-dimensional region D (Definition 5.1.1), some features seemed a bit arbitrary: why did we choose $\frac{1}{n}$ for the box width, as opposed to $\frac{1}{2^n}$ or some other expression which converges to 0 as $n \rightarrow \infty$? Also, why did we evaluate the function at the corner points $(\frac{i}{n}, \frac{j}{n})$? In this section we present a more canonical definition of the Riemann integral. Basically, the idea is to encompass all ways of subdividing the region and choosing points in each subdivision.

Suppose that D is contained in the box $[a, b] \times [c, d]$, and that $f : D \rightarrow \mathbb{R}$ is a bounded function. We define a **partition** P of $[a, b]$ to be a finite set of points between a and b , denoted $a = p_1 < p_2 < \dots < p_m = b$. Similarly, let Q be a partition of $[c, d]$, denoted $c = q_1 < q_2 < \dots < q_n = d$.

For each pair (i, j) , we define $M_{i,j}$ to be the maximum value and $m_{i,j}$ the minimum value of f on $[p_i, p_{i+1}] \times [q_j, q_{j+1}]$ (unless $[p_i, p_{i+1}] \times [q_j, q_{j+1}]$ does not intersect D , in which case we set $M_{i,j}$ and $m_{i,j}$ both to 0). We define the upper Riemann sum $\bar{I}(f, P, Q)$ associated with the partitions P and Q by

$$\bar{I}(f, P, Q) = \sum_{i=1}^m \sum_{j=1}^n M_{i,j} (p_{i+1} - p_i)(q_{i+1} - q_i),$$

and the lower Riemann sum by

$$\underline{I}(f, P, Q) = \sum_{i=1}^m \sum_{j=1}^n m_{i,j} (p_{i+1} - p_i)(q_{i+1} - q_i).$$

Clearly we will have $\underline{I}(f, P, Q) \leq \bar{I}(f, P, Q)$ for all partitions P and Q , since $m_{i,j} \leq M_{i,j}$. Furthermore, as we refine the partitions P and Q (that is, put more points in them), $\bar{I}(f, P, Q)$ decreases and $\underline{I}(f, P, Q)$ increases. Therefore, we are motivated to define $\int_D f \, dA$ to be the number that $\bar{I}(f, P, Q)$ decreases toward as P and Q get increasingly fine—more precisely, we define $\int_D f \, dA$ to be the greatest number which is less than or equal to $\bar{I}(f, P, Q)$ for all partitions P and Q . Likewise, we define $\underline{\int}_D f \, dA$ to be the least number which is greater than or equal to $\underline{I}(f, P, Q)$ for all partitions P and Q .

If $\underline{\int}_D f \, dA$ and $\int_D f \, dA$ are equal, we say that f is **Riemann integrable** and define its integral $\int_D f \, dA$ to be the common value of $\underline{\int}_D f \, dA$ and $\int_D f \, dA$. Otherwise, we say that the function is not Riemann integrable.

Theorem A.3.4

If $f : D \rightarrow \mathbb{R}$ is continuous and if the boundary of D has area zero*, then f is Riemann integrable.

Exercise A.3.6

Assuming that f is Riemann integrable, explain why the sum in Definition 5.1.1 indeed converges to the value $\int_D f \, dA$ defined in this section.

Exercise A.3.7

Consider the function $f : [0, 1]^2 \rightarrow \mathbb{R}$ for which $f(x, y)$ is equal to 1 when x and y are both rational and 0 otherwise. Find $\underline{\int}_{[0,1]^2} f \, dA$ and $\int_{[0,1]^2} f \, dA$, and explain why f is not Riemann integrable.

* Yes, this is pretty much always the case in practice: a region has to be pretty wild for its boundary to have positive area

A.4 Big picture

A.4.1 HOW TO SOLVE MATH PROBLEMS

Many of us have experienced a mode of solving math problems that goes something like this: (1) see several similar examples solved out using some sort of formula or sequence of steps, (2) pick up on cues to recognize which numbers in the problem statement go where in the solution frame, and (3) try our best to carry out these steps for the problem at hand. Hopefully you've also experienced problem solving that involves reasoning about the underlying ideas and tightly integrates the manipulation of numbers and symbolic expressions with the logical or physical or geometric concepts they represent. This second mode is much more versatile, useful, and fun. Think of math ideas as **tools, not templates**.

Becoming a composer requires developing facility with how individual notes come together to make chords, and how notes and chords work together to make music. Achieving such agency in the context of mathematical problem solving, likewise, involves the incorporation of individual techniques into your skill set and practice with weaving them together to compose solutions.

Some practical tips:

1. **Work forwards and backwards.** Your goal in solving a math problem is to arrive at a conclusion based on some given data. It is often possible to work forwards from the given information *and* to work backwards to simplify the desired conclusion. If you try one of these and get stuck, switch to the other, repeatedly if necessary. When you write up your final solution, however, it will flow better if you start at the beginning and reason your way to the end.
2. Ask “**what does this mean?**” often. A major part of the process of doing mathematics is packaging ideas into definitions and unpacking them. For example, if we’re asked to show that $(0, 1]$ is equal to the range of the function $x \mapsto \frac{1}{1+x^2}$ from \mathbb{R} to \mathbb{R} , then the first question we need to answer is “what does it mean to say that the range of a function is equal to a particular set?” If we apply the definition of range to the function at hand, we get a concrete statement which is easier to approach.
3. Engage your **physical intuition** early and often throughout the problem solving process. Ask questions like *can I tell if the final answer will be positive or negative?* and *this expression increases as t increases; does that make sense in the context of the figure?* to establish frequent checks on your calculations and to develop your conceptual comprehension.
4. Make sure you **know the type and role** of each variable you use. For example, you might ask *is x a function or a vector or a number?* (In other words, what is its type?) *Is x an unknown that we are trying to solve for, or is it a dummy variable being used in the definition of a function or in the statement of an identity?* (In other words, what is its role?) To follow a plot line, it is essential to bear in mind basic information about the characters and how they relate to the story as a whole. It’s the same with math.
5. **Don’t be afraid** of dead ends or missteps. Interesting problems often require us to proceed without seeing how the subsequent steps of our solution are going to unfold. When you solve a jigsaw puzzle, you try pieces in lots of different places before your fragments begin to coalesce into a coherent whole. Trying things that don’t work is part of the process.

For further reading, see <http://www.math.ucla.edu/%7Etao/preprints/problem.ps>.

Whenever you write mathematics, including when you submit solutions for homework or exams, your primary goal should be **communication**. Specifically, you should bear in mind a reader who is a fellow student in the course. This imagined reader will be well acquainted with the ideas developed elsewhere in the course, but they have not yet managed to solve the question at hand. Your writing should contain enough reasoning and details to bring them comfortably around to an understanding of your solution.

Furthermore, unless explicitly stated otherwise, your solutions should be **complete**. This means that you have convincingly explained why your assertions are true; there are no remaining holes that your reader could poke in your argument. If you assumed at some point that $x \neq -1$, for example, you should devote the necessary space to addressing the possibility that $x = -1$. All problem statements should be read with an implicit "...and explain your reasoning".

- **Use complete sentences.** Writing in complete sentences helps your reader understand what you mean. This is as true for mathematical exposition as it is for prose. Find a way to incorporate your equations and expressions into sentences which properly contextualize them. A good rule of thumb is that all math notation (equations and expressions) should be incorporated into an English sentence.
- **Focus on what's important.** Reflect on which steps in your solution are important to the overall flow of the solution and which steps are routine, and allocate space accordingly. For example, it's fine to say "...by the quadratic formula, $x^2 - x - 1 = 0$ implies $x = \frac{1 \pm \sqrt{5}}{2}$."
- **Watch out for unquantified variables.** Quantifiers are phrases like "there exists" or "for all" which tell us how to interpret a variable in an equation. For example, consider the statement "for all x , there exists y such that $x^2 = 4y$ ". The same equation might appear in a completely different statement like "for all $y \geq 0$, there exists at least one value of x such that $x^2 = 4y$." Writing the equation $x^2 = 4y$ without any quantifiers, therefore, does not communicate the meaning of the equation.
- **Don't write too much.** If you're able to reach the desired conclusion in two steps, write those two steps clearly and leave it at that.

Example A.4.1

Show that the range of the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ defined by $f(x, y) = (2x + y, -4x - 2y)$ is a line in \mathbb{R}^2 .

Good solution

The range of f is the line $\ell = \{(s, t) \in \mathbb{R}^2 : t = -2s\}$. To establish this, we need to show that (i) for every pair $(s, t) \in \ell$, there is some point $(x, y) \in \mathbb{R}^2$ that f maps to (s, t) , and (ii), the image under f of every pair $(x, y) \in \mathbb{R}^2$ is in ℓ .

For (i), we note that for all points (s, t) , the point $(s/2, 0)$ maps to (s, t) under f . For (ii), we note that $-2(2x + y) = -4x - 2y$ for all x and y , so if (s, t) is a point in \mathbb{R}^2 such that $(s, t) = f(x, y)$, then we have $t = -2s$.

Solution that needs improvement

line: $y = -2x$

$$-2(2x + y) \stackrel{?}{=} -4x + 2y$$

$$(x/2, 0) \mapsto (x, y)$$

The second solution is much less effective than the first. It requires some inference to figure out what the equations are supposed to mean, since they don't appear in the context of sentences. We're missing the step where we unpack the definition of range and arrive at the two statements that need to be checked. Furthermore, the variables that appear there aren't quantified at all: are we talking about all x and y , or some x and y , or what? The biggest problem is that the reasoning which binds together the equations is left up to the reader to guess.

Example A.4.2

Find the minimum and maximum of the function $f : [0, 2] \rightarrow \mathbb{R}$ defined by $f(x) = x^2 - x$.

Good solution

Since $[0, 2]$ is closed and bounded, the extreme value theorem tells us that f achieves a maximum value and a minimum value. Furthermore, by the first derivative test, each extremum occurs at either a critical point or an endpoint of the interval.

To find the set of critical points of f , we set $f'(x) = 2x - 1$ equal to zero and solve. The only solution to this equation is $x = \frac{1}{2}$. Therefore, the maximum and minimum values occur at some value in the set $\{0, \frac{1}{2}, 2\}$. Evaluating the function at each of these points gives 0, $-\frac{1}{4}$, and 2, respectively.

So we conclude that the maximum value of the function is $\boxed{2}$, realized at $x = 2$. The minimum value is $\boxed{-\frac{1}{4}}$, realized at $x = \frac{1}{2}$.

Solution that needs improvement

We differentiate $f(x)$ using the power rule: $2x - 1$

$$2x - 1 = 0$$

$$2x = 1$$

$$\frac{2x}{2} = \frac{1}{2}$$

$x = \frac{1}{2}$. Then $f(0) = 0$ and $f(\frac{1}{2}) = -\frac{1}{4}$ and $f(2) = 2$. Min $-\frac{1}{4}$ Max 2.

This solution is more readable than the needs-improvement solution to Example A.4.1, but it still focuses on the least important ideas. We don't need to see that linear equation written out step by step. What we *do* want to understand is why we're setting the derivative equal to zero in the first place, how we know it suffices to check the points 0, $\frac{1}{2}$ and 2, etc. There is also an issue with the way the equations are presented: there is no logical connector from one equation to the next, and the sequence of equations collectively is not incorporated into an English sentence. This leaves the reader to guess the role played by these equations.

The following documents present more detailed advice on writing good solutions.

<https://sites.math.washington.edu/~lee/Writing/writing-proofs.pdf>, or

<http://www.ohiouniversityfaculty.com/mohlenka/goodproblems/goodproblem.pdf>

A.4.3 THINKING ABOUT FUNCTIONS

If asked to write down some examples of functions, you might produce a list like

$$\sin(x^2) \quad \frac{1}{1-x-x^2} \quad e^x \ln \sin 2x \quad 8x - x^2 + x^{-1}$$

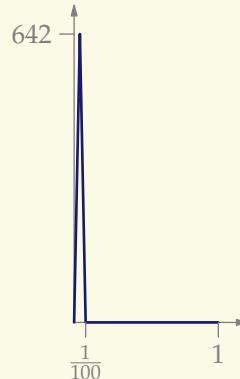
By $\sin(x^2)$, for example, we mean “the function that squares its input and takes the sine of the result”. These are *elementary* functions, meaning that they can be written as a composition of a finite number of arithmetic operations ($\times, \div, +, -$), exponentials, logarithms, trigonometric functions and their inverses, constants, and solutions of polynomial equations (like square roots). Functions that can be defined purely using expressions in this way make up a *tiny* sliver of the set of all functions: if we’re defining a function on a domain D , we can define $f(x)$ however we want for each $x \in D$. Thinking of functions in terms of expressions obscures this fact, because changing the expression defining an elementary function changes many values of the function simultaneously.

Example A.4.3

Show that there exists a continuous function from $[0, 1]$ to \mathbb{R} which is equal to zero on the interval $[0.01, 1]$ and satisfies the integral equation $\int_0^1 f(x) dx = 3.21$.

Solution

Consider the path which begins at $(0, 0)$, proceeds in a straight line to $(0.005, 642)$, then in a straight line to $(0.01, 0)$, and from there directly across to $(1, 0)$. Every vertical line between $x = 0$ and $x = 1$ intersects this path exactly once; therefore, this path is the graph of a function. The area under the graph of this function is the area of a triangle with base 0.01 and height 642, so its area is indeed 3.21.



Note that in our solution to Example A.4.3, we specified our function by drawing its graph rather than by giving a formula. This flexibility is a powerful problem solving tool when it comes to proving the existence of functions with desired properties.

Exercise A.4.1

Show that there exists a *differentiable* function $[0, 1] \rightarrow \mathbb{R}$ which is equal to zero on the interval $[0.01, 1]$ and satisfies the integral equation $\int_0^1 f(x) dx = 3.21$.

Exercise A.4.2

The intermediate value theorem states that if f is a continuous function from \mathbb{R} to \mathbb{R} , then for every interval $[a, b] \subset \mathbb{R}$ and for every value y which is between $f(a)$ and $f(b)$, there exists $c \in [a, b]$ such that $f(c) = y$. Show that the continuity hypothesis is essential. In other words, show that there exists a discontinuous function which does not satisfy the conclusion of the intermediate value theorem.

Example A.4.4: (requires notion of multivariable limits)

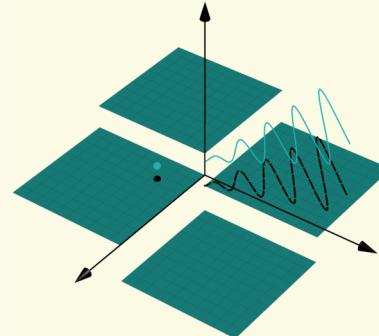
Show that there is a function from \mathbb{R}^2 to \mathbb{R} which is discontinuous at a point (x, y) if and only if either (i) x or y is zero, (ii) $(x, y) = (\frac{1}{4}, -\frac{1}{4})$, or (iii) (x, y) lies on the path $\mathbf{r}(t) = \left\langle -t \left(1 + \frac{1}{2} \sin 40t\right), t \left(1 - \frac{1}{2} \sin 40t\right) \right\rangle$.

Solution

We can just directly define the function to do what we want. There are lots of options, so let's get a little creative. For each $t \geq 0$, we define*

$$f \left(-t \left(1 + \frac{1}{2} \sin 40t\right), t \left(1 - \frac{1}{2} \sin 40t\right) \right) = \frac{1}{8} + \frac{t}{3}.$$

For every point (x, y) in the closed second quadrant not on the curve, we define $f(x, y) = -\frac{1}{10}$. Let's also define $f(x, y)$ to be $-\frac{1}{3}$ if (x, y) is in the open first quadrant, $\frac{1}{5}$ if (x, y) is in the open third quadrant, and 0 if (x, y) is in the closed fourth quadrant, *unless* $(x, y) = (0, 0)$ or $(x, y) = (\frac{1}{4}, -\frac{1}{4})$, in which case we set $f(x, y) = \frac{1}{8}$. The figure shown illustrates the graph of this function.



* The only concern here would be that \mathbf{r} maps two t values to the same (x, y) value, in which case we'd be trying to assign two different values to the same input pair (x, y) .

Exercise A.4.3

Suppose that A is a finite set with m elements, and B is a finite set with n elements. How many functions $f : A \rightarrow B$ exist?

Exercise A.4.4

Show that there exists a continuous function from \mathbb{R}^2 to \mathbb{R} such that $f(x, y) = 0$ if (x, y) lies on the upper half of the unit circle, and $f(x, y) > 0$ for all other points $(x, y) \in \mathbb{R}^2$. (Hint: provide a geometric definition of f that directly involves the desired half circle.)

Mass is an **additive** property of an object. This means that if we subdivide the object into smaller pieces and sum their masses, we get the same result as if we had measured the mass of the original object. Many other quantities of interest in calculus are additive, such as the length of a path, the area of a surface, the volume of a 3D region, the flow of a vector field across a surface, the integral of a function over an interval or region, the moment of inertia, torque, etc.

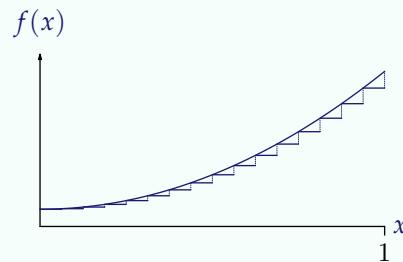
The central idea of integral calculus is that an additive property of an object can be calculated by **subdividing the object into many tiny pieces**, approximating the quantity of interest for each piece, summing the results, and letting the number of pieces go to infinity to obtain an expression which is amenable to formal integration techniques. There are two pretty amazing things that happen here: (1) we end up with an *exact* answer even though our approach is based on approximations, and (2) performing the integration ends up being far more tractable than working with the approximating sums.

The force behind (2) is the fundamental theorem of calculus, which provides a link between limits of approximating sums and the often-quite-easy task of finding an antiderivative. But (1) typically remains mysterious until one takes a course or reads a book in which the hand-waving justifications of an introductory calculus course are replaced with rigorous proofs. Nevertheless, the basic idea behind (1) can be appreciated at an elementary level. It is helpful to begin by considering what can go wrong.

Exercise A.4.5: A cautionary tale

For each positive integer n , consider the greatest function f_n which is piecewise constant over each interval $\left[\frac{i}{n}, \frac{i+1}{n}\right)$, where i ranges from 0 to $n - 1$, and which is less than or equal to the function $f(x) = x^2 + 1$ at each point $x \in [0, 1]$.

Find the arc length of the graph of f_n over $[0, 1]$ (either (i) count only the horizontal steps or (ii) count the horizontal and vertical steps). Does this approximate the arc length of f over $[0, 1]$ increasingly well as $n \rightarrow \infty$?



Exercise A.4.5 shows that some properties of f_n can approximate those of f quite poorly even if f_n approximates f very well. What goes wrong here with this approximation is that the **relative error** is large. In other words, the ratio of the arc length of f over a short interval and the arc length of the approximating “stairstep” over that interval does not converge to 1 as the width of the interval is decreased to zero.

Contrast Exercise A.4.5 with Figure A.5, which illustrates approximating the area under the graph of f with the area under f_n . In this case the relative error of the approximation for each slice decreases to zero. In other words, quotient of the green area and the red+green area converges to 1 as the number of slices converges to infinity. By Exercise A.4.6 below, this implies that the relative error of the approximating sum converges to zero as the number of steps goes to infinity. Therefore, the limit of the sequence of approximating sums is exactly equal to the area under the curve.

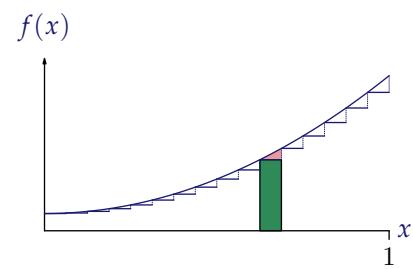


Figure A.5 The area under the step function approximates the area under the curve

Exercise A.4.6: Small relative error for the part implies small relative error for the whole

Let $\epsilon > 0$. Suppose that A_i and A'_i , where i ranges from 1 to n , are positive numbers satisfying $1 - \epsilon \leq \frac{A'_i}{A_i} \leq 1 + \epsilon$ for all i from 1 to n . Show that $1 - \epsilon \leq \frac{A'_1 + \dots + A'_n}{A_1 + \dots + A_n} \leq 1 + \epsilon$.

Exercise A.4.7

Verify that if $f : [a, b] \rightarrow (0, \infty)$ is a continuous function and $x_0 \in [a, b]$, then the quotient of the area under the graph of f over $[x_0, x_0 + h]$ and the area under the graph of the constant function $f(x_0)$ over $[a, b]$ converges to 1 as $h \rightarrow 0^+$.

The above discussion does omit some important details. For example, Exercise A.4.6 stipulates that we achieve the same small relative error ϵ across the whole interval, while in Exercise A.4.7 we focus on each point x_0 one at a time. Also, what if f isn't assumed to be positive, and the area under the curve over a particular short interval happens to be zero?

Despite these shortcomings, the relative error perspective can clarify what's going on when we do our "split it into many tiny pieces" derivations throughout the course.

A.5 SageMath

Math is more fun when you learn how to take advantage of computational resources to assist your learning. Some problem solving tasks are done much better by computers than people, and while in some cases it is important to gain facility with performing such calculations by hand, at some point you want to delegate tedious tasks to the computer and spend your time and attention on the more creative aspects of problem solving. The open source project which has arguably made the most concerted and successful effort to be broadly useful to math students with minimal fuss is SageMath (or just Sage). You can use Sage in your browser without having to install anything. The webpage sagecell.sagemath.org doesn't require sign in, while cocalc.com gives you some hard drive space and saves your work. You can use Sage essentially as a calculator for many tasks. However, it is quite convenient if you need to do something more involved, since Sage uses the beginner-friendly language Python. This is where Sage really shines in comparison to a query tool like Wolfram Alpha, which is great for one-liners but is not well-suited to several-step computations.

Here are some examples of calculations you can do with Sage. You have to tell it that you're going to use `x` as a symbol using the `var` function; you can declare several symbols in this way using spaces. The text following the hashtag is a comment and is ignored by Sage.

```
> var("x y") # tell sage that you want to use x and y as symbols
> integrate(sqrt(x^2+1),x) + integrate(e^y * sin(y), y)
-1/2*(cos(y) - sin(y))*e^y + 1/2*sqrt(x^2 + 1)*x + 1/2*arcsinh(x)

> limit(sin(x^2)/x^2,x=0) # find a limit
1

> cos(3*x).trig_expand() # work with trig functions
cos(x)^3 - 3*cos(x)*sin(x)^2

> diff(x^x,x) # differentiate the function x^x
x^x*(log(x) + 1)

> find_local_maximum(sin(x) + cos(x), 0, 2*pi) # find the maximum of sin+cos over [0,2pi]
(1.414213562373095, 0.78539814681742492)

> plot(sin(x) + cos(x), 0, 2*pi) # plot a function

> [factor(x^n - 1) for n in [1..5]] # factor the first five polynomials of the form x^n - 1
[x - 1,
 (x + 1)*(x - 1),
 (x^2 + x + 1)*(x - 1),
 (x^2 + 1)*(x + 1)*(x - 1),
 (x^4 + x^3 + x^2 + x + 1)*(x - 1)]
```

Throughout the text, some computations are performed using Sage.* They are indicated with the CoCalc icon  which can be clicked to open a page at cocalc.com showing the result as well as the code used to generate it.

If you want to learn more about Sage, I recommend the (freely available) book *Sage for Undergraduates* by Gregory Bard.

* Actually, some of these linked code snippets are in Julia, which is a newer language more suited to numerical work

A.6 Additional exercises

SECTION 1.1

Exercise A.1.1.1

Find a formula for the distance from the point (x, y, z) to the z -axis.

Exercise A.1.1.2

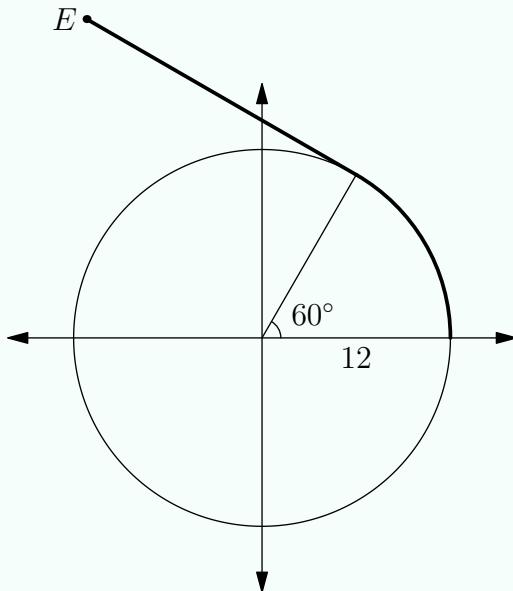
Find a formula for the distance from the point (x, y, z) to the plane $x = -1$.

Exercise A.1.1.3

Find an equation whose solution set is the set of points whose distance to $(-3, 2, 5)$ is equal to 4. What is the radius of the intersection of this sphere with the yz -plane?

Exercise A.1.1.4

A rope of length 12π units is partially wrapped around a tree of radius 12 units, as shown in the figure below. The part of the rope not touching the tree is pulled tight. Find the coordinates of the end of the rope, labeled E .



SECTION 1.2

Exercise A.1.2.1

Find linear transformations with the following geometric descriptions:

- reflect across the x -axis
- rotate 180 degrees and double the distance from the origin
- halve the distance from the origin while preserving the angle between (x, y) , the origin, and the positive x -axis
- rotate 90 degrees counterclockwise
- project a point in \mathbb{R}^3 onto the xy -plane

SECTION 1.3

Exercise A.1.3.1

A 2×2 matrix is said to be *symmetric* if its top right and lower left entries are equal. The *diagonal* entries are the top left and lower right entries. Show that the diagonal entries of a 2×2 symmetric matrix with positive determinant must have the same sign.

Exercise A.1.3.2

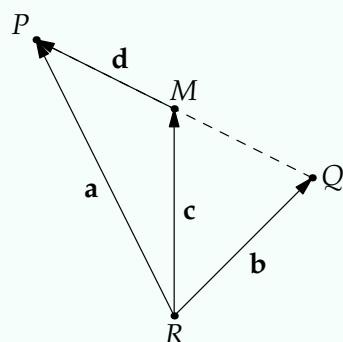
Find all values of λ for which

$$\det \begin{bmatrix} 3 - \lambda & 0 & 0 \\ 0 & 2 - \lambda & 1 \\ 1 & 1 & 1 - \lambda \end{bmatrix} = 0.$$

SECTION 2.1

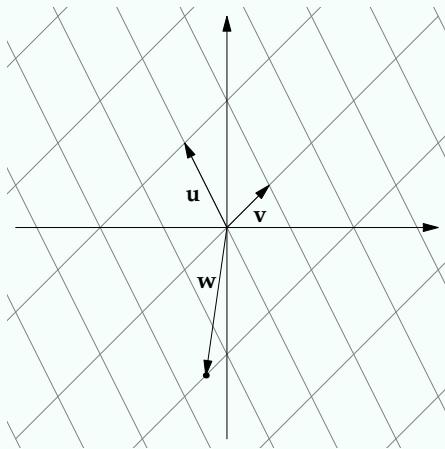
Exercise A.2.1.1

The point M is the midpoint of P and Q . Express \mathbf{c} in terms of \mathbf{a} and \mathbf{b} , and express \mathbf{d} in terms of \mathbf{a} and \mathbf{b} .



Exercise A.2.1.2

Consider the vectors \mathbf{u} , \mathbf{v} and \mathbf{w} in the figure below.



- Find scalars a and b such that $\mathbf{w} = a\mathbf{u} + b\mathbf{v}$. You may assume that both a and b are integer multiples of 0.5. Express your answer as an ordered pair (a, b) , written in the box.
- Find scalars a and b such that $\mathbf{w} - \mathbf{u} + \mathbf{v} = a\mathbf{u} + b\mathbf{v}$. You may assume that both a and b are integer multiples of 0.5. Express your answer as an ordered pair (a, b) , written in the box.

Exercise A.2.1.3

Suppose that all three coordinate planes are outfitted with a mirror surface, and that a laser beam is directed from a point in the first octant toward the origin. Suppose that the beam is slightly misaimed, so that it strikes one of the mirror planes first. Then the beam strikes a second plane, and then a third, before being reflected back out into the first octant.

Show that the reflected beam returns along a path which is parallel to its incoming path. The physics idea you need here is that when a light beam hits a reflective plane, its reflection angle is the same as its incident angle, and the plane containing the incoming and outgoing beams meets the reflective plane at a 90-degree angle.

SECTION 2.2

Exercise A.2.2.1

Use dot products to find a vector which is orthogonal to both $\langle 1, -2, 4 \rangle$ and $\langle -3, 1, 1 \rangle$. Hint: represent your vector as $\langle a, b, c \rangle$ and solve for a , b , and c .

Exercise A.2.2.2

Suppose that a force $\mathbf{F} = (1, -2)$ is acting on an object moving parallel to the vector $(4, 1)$. Decompose \mathbf{F} into a sum of vectors \mathbf{F}_1 and \mathbf{F}_2 , where \mathbf{F}_1 points along the direction of motion and \mathbf{F}_2 is perpendicular to the direction of motion.

Exercise A.2.2.3

Let A, B, C , and D be four points in \mathbb{R}^3 . Use vectors to show that

$$AB^2 + BC^2 + CD^2 + DA^2 \geq AC^2 + BD^2.$$

(This generalizes the fact that the sum of the squares of the sides of a quadrilateral is at least the sum of the squares of its diagonals.) Make a statement about when equality holds.

SECTION 2.3

Exercise A.2.3.1

Consider a triangle T with vertices $A = (1, 0)$, $B = (0, 1)$, and C on the line $y = -x$. Use a cross product to find the area of T , and show that it does not depend on the location of C .

Exercise A.2.3.2

Is the cross-product associative? Is the dot product associative? Prove or give a counterexample for each.

Exercise A.2.3.3

Suppose that the four vectors $\mathbf{a}, \mathbf{b}, \mathbf{c}$, and \mathbf{d} in \mathbb{R}^3 are coplanar. Show that
 $(\mathbf{a} \times \mathbf{b}) \times (\mathbf{c} \times \mathbf{d}) = \mathbf{0}$.

Exercise A.2.3.4

The volume of a parallelepiped is the product of the area of its base and its height. Consider the parallelepiped spanned by \mathbf{a}, \mathbf{b} , and \mathbf{c} . You may suppose for simplicity that \mathbf{b}, \mathbf{c} , and \mathbf{a} form a right-handed triple of vectors, which means that a right-handed screw rotated an angle less than 180° from \mathbf{b} to \mathbf{c} advances in the direction of \mathbf{a} .

- Let us think of the parallelogram spanned by \mathbf{b} and \mathbf{c} as the base of the parallelepiped. What is the (signed) area of this parallelogram?
- What is the height of the parallelogram, in terms of \mathbf{a} and the unit vector pointing in the direction of $\mathbf{b} \times \mathbf{c}$?
- Put together parts (a) and (b) to derive the triple scalar product formula for the volume of a parallelepiped.

Exercise A.2.3.5

Suppose that $\ell_1(t) = t\mathbf{a} + \mathbf{b}_1$ and $\ell_2(t) = t\mathbf{a} + \mathbf{b}_2$ are parallel lines in \mathbb{R}^2 or \mathbb{R}^3 . Show that the distance D between them is given by

$$D = \frac{|\mathbf{a} \times (\mathbf{b}_2 - \mathbf{b}_1)|}{|\mathbf{a}|}.$$

Exercise A.2.3.6

Find the point P on the line $(3 - t, 2 + 2t, -4t)$ which is closest to the point $Q = (3, 7, 1)$.

Exercise A.2.3.7

It is possible to prove the vector triple product formula

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \cdot \mathbf{c})\mathbf{b} - (\mathbf{a} \cdot \mathbf{b})\mathbf{c}$$

in a tedious way using coordinates. This problem outlines a more conceptual proof, taken from a note written by William C. Schulz.

- (a) Use the parallelepiped interpretation of the triple scalar product to show that

$$\mathbf{b} \cdot (\mathbf{c} \times \mathbf{n}) = \mathbf{c} \cdot (\mathbf{n} \times \mathbf{b}) = \mathbf{n} \cdot (\mathbf{b} \times \mathbf{c})$$

- (b) Use the right-hand rule to observe that if \mathbf{c} is perpendicular to \mathbf{n} , then

$$\mathbf{n} \times (\mathbf{c} \times \mathbf{n}) = |\mathbf{n}|^2 \mathbf{c}.$$

- (c) Show that it suffices to consider the case where \mathbf{a} , \mathbf{b} , and \mathbf{c} form a basis for \mathbb{R}^3 .

- (d) Write $\mathbf{a} \times (\mathbf{b} \times \mathbf{c})$ as a linear combination of \mathbf{a} , \mathbf{b} , and \mathbf{c} , so that

$$\mathbf{a} \times (\mathbf{b} \times \mathbf{c}) = \kappa \mathbf{a} + \lambda \mathbf{b} + \mu \mathbf{c}, \quad (\text{A.6.1})$$

- (e) The easiest coefficient to determine is κ . What is it?

- (f) To find λ , dot both sides of (A.6.1) with $\mathbf{c} \times \mathbf{n}$, where $\mathbf{n} = \mathbf{b} \times \mathbf{c}$.

- (g) To find μ , dot both sides of (A.6.1) with $\mathbf{b} \times \mathbf{n}$, where $\mathbf{n} = \mathbf{b} \times \mathbf{c}$.

SECTION 3.1**Exercise A.3.1.1**

Find an equation for the plane containing the points $(3, -1, 2)$, $(2, 0, 5)$, and $(1, -2, 4)$.

Exercise A.3.1.2

Find an equation for the plane that contains the lines

$$\begin{cases} x(t) = 5 + t \\ y(t) = 1 - t \\ z(t) = 4 - 3t \end{cases}$$

and

$$\begin{cases} x(t) = 5 - 4t \\ y(t) = 1 + t \\ z(t) = 4 - 3t \end{cases}$$

Exercise A.3.1.3

Let $O = (0, 0, 0)$ be the origin in \mathbb{R}^3 . If the vector \overrightarrow{OP} has length 3, what is the greatest possible distance from P to the line $\ell = \{-1 - t, 1 + 2t, t\} : t \in \mathbb{R}\}$?

SECTION 3.2**Exercise A.3.2.1**

A chickadee starts at the point $(2, -4, 1)$ and flies in the direction of the vector $\left\langle \frac{3}{\sqrt{10}}, 0, \frac{1}{\sqrt{10}} \right\rangle$ at a rate of $\sqrt{10}$ units per second. A hummingbird starts at the point $(8, 20, 7)$ and flies in the direction of the vector $\left\langle \frac{1}{\sqrt{5}}, -\frac{2}{\sqrt{5}}, 0 \right\rangle$ at a rate of $3\sqrt{5}$ units per second.

- Do the paths of the chickadee and the hummingbird intersect?
- Do the hummingbird and the chickadee collide?

Exercise A.3.2.2

Suppose that a particle is revolving clockwise around the point $(4, 2)$ at a rate of 3 revolutions per second. Write parametric equations describing the location of the particle at time t , assuming that it starts at the point $(6, 2)$.

Exercise A.3.2.3

Sketch the image of the path $\mathbf{x}(t) = (\cos t, e^t)$.

Exercise A.3.2.4

Derive an integral formula for the arc length of a differentiable path $\mathbf{r}(t)$ where t ranges from a to b .

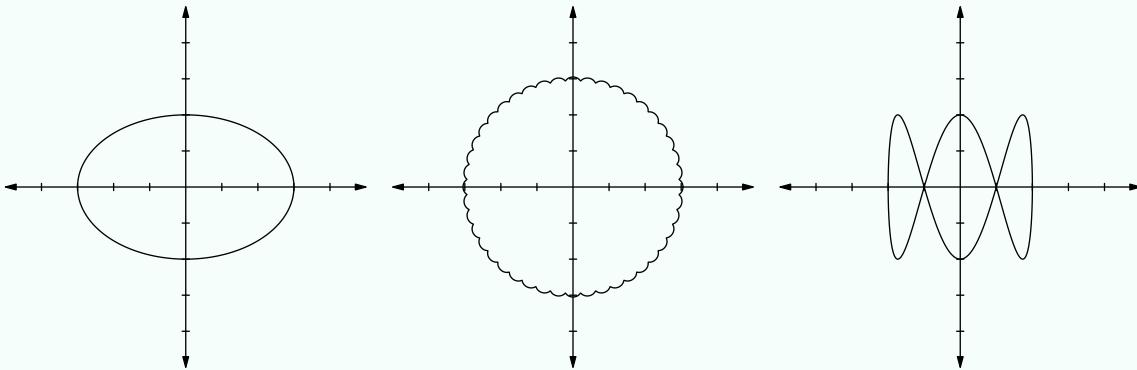
Hint: use the Pythagorean theorem to find the distance from $\mathbf{r}(t)$ to $\mathbf{r}(t + \Delta t)$. Divide the curve into many such segments and approximate its length as the sum of the length of those segments. Then let $\Delta t \rightarrow 0$.

Exercise A.3.2.5

Calculate the total length of the curve represented by the parametric equation $\mathbf{r}(t) = (a \cos^3 t, a \sin^3 t)$, where a is a positive constant.

Exercise A.3.2.6

All the figures below were made by graphing $\{(x(t), y(t)) : 0 \leq t \leq 2\pi\}$ for simple expressions $(x(t), y(t))$. For each graph, find a curve whose image looks (at least roughly) like the figure shown.

**SECTION 3.3****Exercise A.3.3.1**

Prove that the projection onto the xy -plane of the intersection of the plane $x + y + z = 1$ and the ellipsoid $x^2 + 4y^2 + 4z^2 = 4$ is an ellipse.

Exercise A.3.3.2

Find an equation for the set of points whose distance to the plane $z = -1$ and to the point $(0, 0, 1)$ are equal. What kind of shape is this surface?

Exercise A.3.3.3

Show that if the point (a, b, c) lies on the hyperbolic paraboloid $z = y^2 - x^2$, then the line with parametric equation $(a + t, b + t, c + 2(b - a)t)$ lies on the hyperbolic paraboloid. (Thus, even though the paraboloid is curvy, it contains lots of lines!)

SECTION 3.4

Exercise A.3.4.1

Write inequalities describing the unit sphere in cylindrical coordinates.

Exercise A.3.4.2

Consider a sphere of radius 10 centered at the origin. Suppose that the portion of the sphere above the plane $z = 8$ is removed. Furthermore, a sphere of radius 2 is removed from the center of the solid. Write inequalities in spherical coordinates to describe the resulting shape.

SECTION 4.1

Exercise A.4.1.1

Evaluate $\lim_{(x,y) \rightarrow (0,0)} \frac{2xy^2 + xy}{x^2 + y^2}$, or explain why the limit fails to exist.

Exercise A.4.1.2

Define $f(x) = (x^2 - 4)/(x - 2)$ when $x \neq 2$ and $f(2) = c_1$. Determine the value of the constant c_1 for which f is continuous. Do the same for

$$g(x,y) = \begin{cases} \frac{3|x|^3 + 3|y|^3 - x^{10} \arctan(y)}{|x|^3 + |y|^3} & \text{if } (x,y) \neq (0,0) \\ c_2 & \text{if } (x,y) = (0,0). \end{cases}$$

Exercise A.4.1.3

Suppose that f, g , and h are real-valued functions on \mathbb{R}^n , and suppose that $f(\mathbf{x}) = g(\mathbf{x})h(\mathbf{x})$. Prove that if g is bounded and $\lim_{\mathbf{x} \rightarrow \mathbf{0}} h(\mathbf{x}) = 0$, then $\lim_{\mathbf{x} \rightarrow \mathbf{0}} f(\mathbf{x}) = 0$.

Exercise A.4.1.4

In this problem, we will see how to discover the nonexistence of the limit in Example 4.1.6 without needing the hint to look at parabolic paths through the origin.

Find $m(r)$ and $M(r)$ for the function $f(x, y) = -x^2y/(x^4 + y^2)$ by showing that $f(s \cos \theta, s \sin \theta) = \frac{s \cos^2 \theta \sin \theta}{s^2 \cos^4 \theta + \sin^2 \theta}$. Then use single-variable calculus to find the maximum and minimum value of this expression as s ranges over $[0, r]$, for each fixed value of θ . Finally, find the minimum of those minima as θ ranges over $[0, 2\pi]$, and similarly for the maxima.

SECTION 4.2

Exercise A.4.2.1

Find $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$, and $\frac{\partial f}{\partial z}$ for $f(x, y, z) = \frac{1}{xy} + \log(xyz) + e^x \sin(yz)$.

Exercise A.4.2.2

Find the partial derivatives of

$$F(a, b) = \int_a^b \sqrt{t^3 + 1} dt$$

with respect to a and with respect to b .

Exercise A.4.2.3

Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Is it possible for $\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$ to exist at $(0, 0)$ while f is not differentiable at $(0, 0)$? Prove that it isn't possible, or provide an example to show that it is possible.

SECTION 4.3

Exercise A.4.3.1

Verify the linear approximation

$$\frac{2x+3}{4y+1} \approx 3 + 2x - 12y$$

at $(x, y) = (0, 0)$.

Exercise A.4.3.2

Find the linear approximation of the function $f(x, y) = 1 - xy \cos(\pi y)$ at $(1, 1)$ and use it to approximate $f(1.02, 0.97)$.

SECTION 4.5

Exercise A.4.5.1

For functions of one variable, it is impossible for a continuous function to have two local maxima and no local minimum: between every two peaks there must be a trough. Show, however, that the function

$$f(x, y) = -(x^2 - 1)^2 - (x^2y - x - 1)^2$$

has only two critical points and has a local maximum at each (!!).

Exercise A.4.5.2

The Hessian of a twice-differentiable function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is defined by

$$H(x, y) = \begin{bmatrix} \partial_x^2 f & \partial_x \partial_y f \\ \partial_x \partial_y f & \partial_y^2 f \end{bmatrix}.$$

The *second derivative test* says that if $\det H(x, y) < 0$ at a critical point (x, y) then f has neither local minimum nor a local maximum at (x, y) , and if $\det H(x, y) > 0$ then f has a local minimum or a local maximum at (x, y) . The minimum and maximum cases are distinguished by whether the diagonal entries of H are positive or negative, respectively.

Verify that the second derivative test gives the correct results for the functions $x^2 + y^2$, $-x^2 - y^2$, and $x^2 - y^2$ at the origin.

SECTION 4.7

Exercise A.4.7.1

Find the equation for the plane tangent to $z = x^2 - 6x + y^3$ at $(x, y) = (4, 3)$.

Exercise A.4.7.2

Suppose that the temperature in a room $[0, 5]^3$ is given as a function of position by $T(x, y, z) = 50 + x^2 + (y - 3)^2 + 2z$. You are a bug starting at position $(3, 2, 2)$, and you are cold. You decide to move in the direction of greatest temperature increase at all times.

- What vector describes the direction in which you initially move?
- Do you first hit the ceiling, the floor, or a wall?

Exercise A.4.7.3

Show that the sum of the x -, y - and z -intercepts of any tangent plane to the surface $\sqrt{x} + \sqrt{y} + \sqrt{z} = \sqrt{c}$ is a constant.

Exercise A.4.7.4

The second directional derivative of $f(x, y)$ in the direction \mathbf{u} is defined to be

$$D_{\mathbf{u}}^2 f(x, y) = D_{\mathbf{u}}[D_{\mathbf{u}} f(x, y)].$$

Find $D_{\mathbf{u}}^2 f(x, y)$ if $f(x, y) = x^3 + 5x^2y + y^3$ and $\mathbf{u} = \frac{1}{5}\langle 3, 4 \rangle$.

SECTION 4.8**Exercise A.4.8.1**

Let $f(x, y) = e^{3x+y}$, and suppose that $x = s^2 + t^2$ and $y = 2 + t$. Find $\partial f / \partial s$ and $\partial f / \partial t$ by substitution and by means of the chain rule. Verify that the results are the same for the two methods.

Exercise A.4.8.2

Use the chain rule to find $\partial z / \partial s$, $\partial z / \partial t$, and $\partial z / \partial u$ when $s = 4$, $t = 2$, $u = 1$, given

$$z = x^4 + x^2y, \quad x = s + 2t - u, \quad y = stu^2.$$

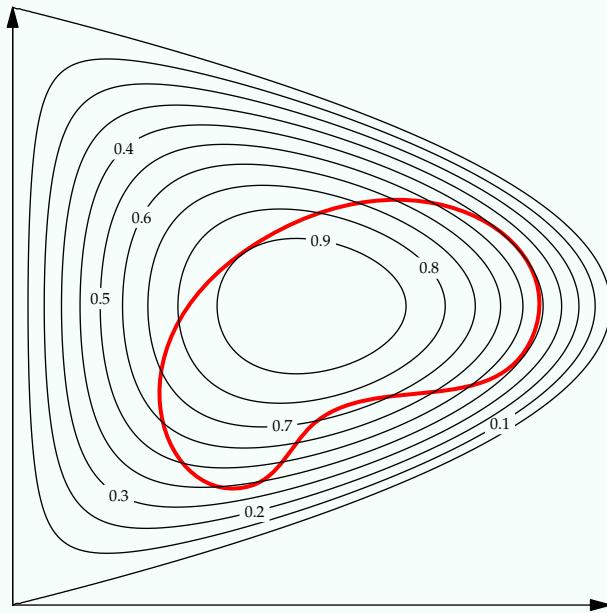
Exercise A.4.8.3

A conical ice sculpture melts in such a way that its height decreases at a rate of 0.001 meters per second and its radius decreases at a rate of 0.002 meters per second. At what rate is the volume of the sculpture decreasing when its height reaches 3 meters, assuming that its radius is 2 meters at that time? Express your answer in terms of π and in units of cubic meters per second.

SECTION 4.9

Exercise A.4.9.1

Some of the level curves of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ are shown. Find the minimum and maximum values of $f(x, y)$ subject to the constraint that (x, y) lies on the red curve.



Exercise A.4.9.2

Minimize the function $f(x, y) = (x - y)^2$ subject to the constraint $xy = 1$ without using Lagrange multipliers. Verify that the method of Lagrange multipliers gives the same result.

Exercise A.4.9.3

Verify the hypotheses of the extreme value theorem and find the absolute maximum and minimum values of $f(x, y) = x^2 + y^2 - 2x$ on the closed triangle D with vertices $(2, 0)$, $(0, 2)$, and $(4, 0)$.

Exercise A.4.9.4

Use Lagrange multipliers to find the maximum and minimum values of $f(x, y) = x^2 + y^2$ subject to $xy = 1$.

Exercise A.4.9.5

Use Lagrange multipliers to find the maximum and minimum values of $x^2 + y^2 + z^2$ subject to the constraint $x + y + z = 12$. Confirm your answer by solving the same problem using vector methods.

Exercise A.4.9.6

Use Lagrange multipliers to show that the rectangle with maximum area that has a given perimeter p is a square.

SECTION 5.1**Exercise A.5.1.1**

Evaluate $\int_0^1 \int_0^{y^2} x^2 y \, dx \, dy$ and sketch the region of integration in \mathbb{R}^2 indicated by the limits of integration.

Exercise A.5.1.2

Evaluate $\int_0^\pi \int_y^\pi \frac{\sin x}{x} \, dx \, dy$.

Exercise A.5.1.3

Evaluate $\int_{-\infty}^{\infty} e^{-x^2} \, dx$. (Hints: write down the product of this integral with itself, the second time using y as the dummy variable. Rewrite this as an area integral and switch to polar coordinates.)

Exercise A.5.1.4

Define a function $f(x, y)$ on $[0, 1] \times [0, 2]$ by

$$f(x, y) = \begin{cases} 1 & \text{if } x \text{ is rational} \\ 0 & \text{if } x \text{ is irrational and } y \leq 1 \\ 2 & \text{if } x \text{ is rational and } y > 1. \end{cases}$$

Show that the iterated Riemann integral $\int_0^1 \int_0^2 f(x, y) \, dy \, dx$ exists, and find its value. Show that the iterated Riemann integral $\int_0^2 \int_0^1 f(x, y) \, dx \, dy$ does not exist.

Exercise A.5.1.5

For $(x, y) \neq (0, 0)$, we define

$$f(x, y) = \frac{xy(x^2 - y^2)}{(x^2 + y^2)^3}.$$

Calculate the iterated integrals of f over $[0, 2] \times [0, 1]$.

SECTION 5.2

Exercise A.5.2.1

Explain why the integral of the function $f(x, y, z) = \frac{1}{x+y+z+1}$ over the cube $[0, 1] \times [0, 1] \times [0, 1]$ is equal to $\lim_{n \rightarrow \infty} S_n$, where

$$S_n = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \frac{f(i/n, j/n, k/n)}{n^3}.$$

(b) At sagecell.sagemath.org, use the code

```
n = 20
var("x", "y", "z", "i", "j", "k")
f(x,y,z) = 1/(x+y+z+1)
assume(0<x<1);assume(0<y<1);assume(0<z<1)
I = integrate(integrate(integrate(f(x,y,z),x,0,1),y,0,1),z,0,1)
S = sum(sum(sum(f(i/n,j/n,k/n)/n^3,k,1,n),j,1,n),i,1,n)
(I,S,N(I),N(S))
```

which gives the exact values of the integral and S_{20} followed by their decimal representations, to show that S_{20} is close to the value of the integral. Increase n (just change the first line of code to assign a different value to n) by multiples of 5 to **find the least value of n such that n is a multiple of 5 and S_n differs from I by less than 0.01.**

Exercise A.5.2.2

Find the region E in \mathbb{R}^3 for which

$$\iiint_E (1 - x^2 - 2y^2 - 3z^2) dV$$

is as large as possible. (Hint: this problem is not as difficult as it may seem.)

Exercise A.5.2.3

Sketch the solid whose volume is given by the integral $\int_0^1 \int_0^{1-x} \int_0^{2-2z} dy dz dx$.

Exercise A.5.2.4

Integrate $f(x, y, z) = 1 - z^2$ over the tetrahedron W with vertices at the origin, $(1, 0, 0)$, $(0, 2, 0)$, and $(0, 0, 3)$.

SECTION 5.3

Exercise A.5.3.1

Find the volume of the solid that is bounded by the paraboloid $z = 9 - x^2 - y^2$, the xy -plane, and the cylinder $x^2 + y^2 = 4$.

Exercise A.5.3.2

Express as a triple integral the volume of the wedge in the first octant that is cut from the cylinder $y^2 + z^2 = 1$ by the planes $y = x$ and $x = 1$.

Exercise A.5.3.3

Sketch the solid whose volume is given by the integral $\int_{-\pi/2}^{\pi/2} \int_0^2 \int_0^{r^2} r \, dz \, dr \, d\theta$, and evaluate the integral.

Exercise A.5.3.4

Evaluate $\int_0^1 \int_0^{\sqrt{1-x^2}} \int_{\sqrt{x^2+y^2}}^{\sqrt{2-x^2-y^2}} xy \, dz \, dy \, dx$ by changing to spherical coordinates.

Exercise A.5.3.5

Find the mass of a ball B given by $x^2 + y^2 + z^2 \leq a^2$ if the density at any point is proportional to its distance from the z -axis. (Hint: use cylindrical coordinates.)

Exercise A.5.3.6

Find the volume of the region W that represents the intersection of the solid cylinder $x^2 + y^2 \leq 1$ and the solid ellipsoid $2(x^2 + y^2) + z^2 \leq 10$.

Exercise A.5.3.7

Use polar coordinates to combine the sum

$$\int_{1/\sqrt{2}}^1 \int_{\sqrt{1-x^2}}^x xy \, dy \, dx + \int_1^{\sqrt{2}} \int_0^x xy \, dy \, dx + \int_{\sqrt{2}}^2 \int_0^{\sqrt{4-x^2}} xy \, dy \, dx$$

into one double integral. Then evaluate the double integral.

SECTION 5.4

Exercise A.5.4.1

Let D be a parallelogram with vertices $(0, 0)$, $(1, 0)$, $(1, 1)$, and $(2, 1)$. Calculate $\iint_D 1 \, dA$ in two ways:

- Find $\iint_D 1 \, dA$ without using calculus.
- Find $\iint_D 1 \, dA$ using the change of variables $u = 2x - 2y$ and $v = 2y$.

Exercise A.5.4.2

Find $\iint_R \cos\left(\frac{y-x}{y+x}\right) \, dA$, where R is the trapezoidal region with vertices $(1, 0)$, $(2, 0)$, $(0, 2)$, and $(0, 1)$.

Exercise A.5.4.3

Evaluate the integral $\iint_R e^{(x+y)/(x-y)} \, dA$, where R is the trapezoidal region with vertices $(1, 0)$, $(2, 0)$, $(0, -2)$, and $(0, -1)$.

Exercise A.5.4.4

Find $\iint_R y^2 \, dA$, where R is the region bounded by the curves $xy = 1$, $xy = 2$, $xy^2 = 1$, and $xy^2 = 2$.

Exercise A.5.4.5

Find $\iint_R e^{x+y} \, dA$ where R is given by the inequality $|x| + |y| \leq 1$.

Exercise A.5.4.6

Evaluate the integral $\int_0^2 \int_{x/2}^{x/2+1} x^5(2y-x)e^{(2y-x)^2} \, dy \, dx$ by making the substitution $u = x$ and $v = 2y - x$.

SECTION 5.5

Exercise A.5.5.1

Consider the functions f , w_1 , and w_2 from $[1, 2]^2 \rightarrow \mathbb{R}$ defined by

$$\begin{aligned}f(x, y) &= xy, \\w_1(x, y) &= x + y, \\w_2(x, y) &= 4 - x - y.\end{aligned}$$

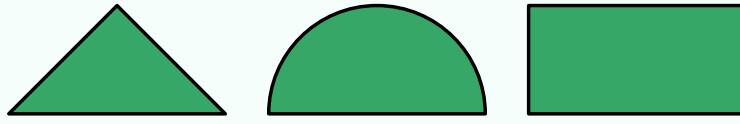
- Describe qualitatively how you can guess which of the following is greater, before doing any calculation (i) the w_1 -weighted average of f , or (ii) the w_2 -weighted average of f .
- Calculate the w_1 -weighted average of f and the w_2 -weighted average of f .
- Confirm that both weighted averages of f computed above lie between the maximum and minimum values of f .

Exercise A.5.5.2

Find the average distance from a point in a ball of radius a to the center.

Exercise A.5.5.3

Which of the following shapes has the highest center of mass? Which has the lowest?



Exercise A.5.5.4

Find the center of mass of the cone $0 \leq z \leq a - b\sqrt{x^2 + y^2}$, where a and b are positive real numbers.

Exercise A.5.5.5

Consider an L-shaped lamina consisting of the points $(x, y, z) \in \mathbb{R}^3$ such that (i) $z = 0$, (ii) $(x, y) \in [0, 1]^2$, and (iii) either x or y is at least $\frac{1}{2}$. Assuming that this lamina has constant density δ , find its moment of inertia about the z -axis.

Exercise A.5.5.6

The dart thrower in Example 5.5.3 is terribly unlikely to hit the triple-20. Let's see how they can increase their chances by becoming less accurate.

- (a) Suppose that the probability density function of the dart's location is given by

$$f_\alpha(x, y) = \frac{1}{\pi\alpha} e^{-\frac{x^2+y^2}{\alpha}},$$

where $\alpha > 0$ is an accuracy parameter. If a player becomes more accurate, does their α value increase or decrease?

- (b) Explain in intuitive terms why a thrower with accuracy α is extremely unlikely to hit the triple-20 either when α is very small or when α is very large.
- (c) Find the value of α that maximizes the probability of hitting the triple-20.

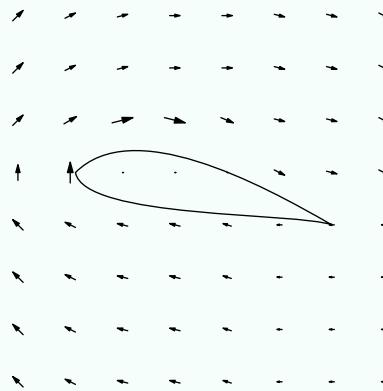
SECTION 6.1

Exercise A.6.1.1

Sketch the vector field $\mathbf{F} = \langle xy, x \rangle$ over $[-1, 1]^2$.

Exercise A.6.1.2

A velocity field \mathbf{v} showing air flow around an airfoil is illustrated below. Sketch a continuous path C along which $\int_C \mathbf{v} \cdot d\mathbf{r}$ is as large as possible.



Exercise A.6.1.3

Denote by $-C$ the reversal of a path C . Show that

$$\int_{-C} \mathbf{F} \cdot d\mathbf{r} = - \int_C \mathbf{F} \cdot d\mathbf{r}.$$

SECTION 6.2

Exercise A.6.2.1

A thin wire is bent into the shape of a semicircle $x^2 + y^2 = 4$, $x \geq 0$. If the linear density of the wire at (x, y) is given by $x + y + 2$, find the mass of the wire.

Exercise A.6.2.2

Let $\mathbf{F} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be the vector field given by $\mathbf{F}(x, y, z) = ay^2\mathbf{i} + 2y(x+z)\mathbf{j} + (by^2 + z^2)\mathbf{k}$.

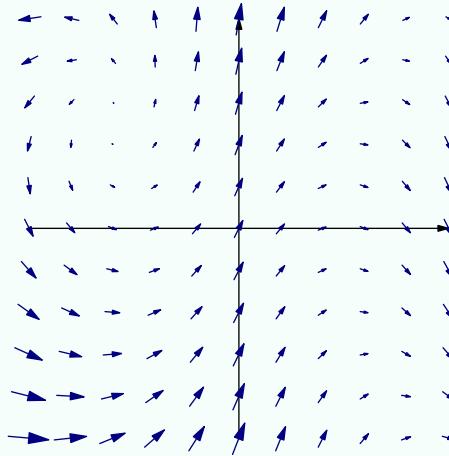
- For which values of a and b is the vector field \mathbf{F} conservative?
- Find a function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ such that $\mathbf{F} = \nabla f$ for these values.
- Find an equation describing a surface S with the property that for every smooth oriented curve C lying on S ,

$$\int_C \mathbf{F} \cdot d\mathbf{r} = 0,$$

for these values.

Exercise A.6.2.3

The vector field \mathbf{F} plotted below is not conservative. Pick two points a and b and sketch two paths from a to b along which the line integrals of \mathbf{F} are clearly different, and explain why the integrals are clearly different.



Exercise A.6.2.4

According to Coulomb's law, the force between a particle of charge q_1 at the origin and a particle of charge q_2 at the point $\mathbf{r} = (x, y, z) \in \mathbb{R}^3$ is given by

$$\mathbf{F} = \frac{q_1 q_2}{4\pi\epsilon_0} \frac{\mathbf{r}}{|\mathbf{r}|^3},$$

where ϵ_0 is a physical constant.

- Is \mathbf{F} a conservative vector field? If so, find a function $\phi : \mathbb{R}^3 \rightarrow \mathbb{R}$ such that $\nabla\phi = \mathbf{F}$.
- If the distance between two charges is tripled, by what factor is the force between them reduced?
- How much work is required to move the second particle along the path

$$\gamma(t) = (1 + (1-t)\cos(t^2), \sqrt{\sin \pi t}, 4t - t^2) \quad 0 \leq t \leq 1?$$

Express your answer in terms of q_1 , q_2 , and ϵ_0 .

Exercise A.6.2.5

Let C be a level set of the function $f(x, y)$. Show that $\int_C \nabla f \cdot d\mathbf{s} = 0$.

SECTION 6.3**Exercise A.6.3.1**

Find the area of the rectangle $D = [0, a] \times [0, b]$ using Green's theorem.

Exercise A.6.3.2

Use Green's theorem to prove the *shoelace formula*, which says that twice the area of a polygon with vertices $(x_1, y_1), \dots, (x_n, y_n)$ listed in counterclockwise order is equal to $x_1y_2 + x_2y_3 + \dots + x_{n-1}y_n + x_ny_1 - (x_2y_1 + x_3y_2 + \dots + x_ny_{n-1} + x_1y_n)$.

Exercise A.6.3.3

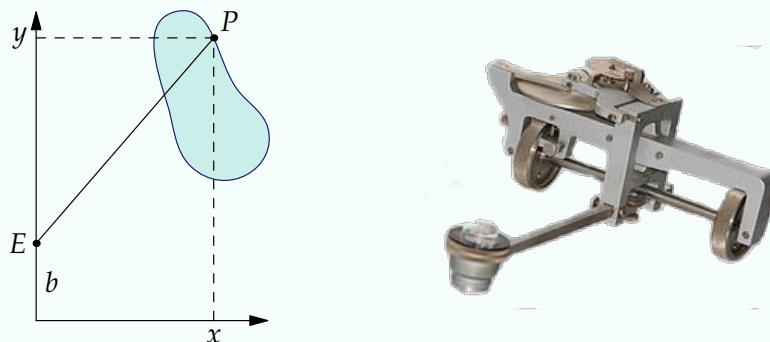
Consider the polygon P whose vertices are listed below, in counterclockwise order. Describe a simple

algorithm for approximating the area of P .

x	y	x	y
1.	0.	-0.999301	0.000100651
0.990436	0.000510364	-0.994893	-0.000198912
0.962169	0.00404386	-0.971584	-0.00262612
0.916455	0.0134308	-0.930411	-0.0101725
0.855307	0.031127	-0.873191	-0.0254117
0.781385	0.059054	-0.802423	-0.0504237
0.697852	0.0984698	-0.721151	-0.0866563
0.608195	0.149879	-0.632792	-0.134824
0.516034	0.212985	-0.540951	-0.194848
0.424993	0.286693	-0.449221	-0.265846
0.338182	0.369151	-0.360994	-0.346165
0.258659	0.457849	-0.279269	-0.433459
0.188645	0.549739	-0.206502	-0.524811
0.129723	0.641405	-0.144474	-0.616889
0.082709	0.729244	-0.0942056	-0.706119
0.0475787	0.809659	-0.0559169	-0.788884
0.0235696	0.879262	-0.0290287	-0.86172
0.00916064	0.935053	-0.0122151	-0.921505
0.00221855	0.974593	-0.00349795	-0.965635
0.000130899	0.996135	0.00000440875	-0.999913
-0.0000248959	0.998721	0.000669899	-0.988538
-0.00129456	0.982236	0.00464303	-0.958545
-0.00665282	0.947413	0.0147221	-0.911264
-0.0188134	0.895794	0.0333185	-0.848777
-0.0400542	0.829645	0.0622956	-0.773799
-0.0720676	0.751827	0.102843	-0.689529
-0.115845	0.665643	0.15539	-0.599477
-0.171601	0.574651	0.219565	-0.507267
-0.238745	0.482473	0.294196	-0.41644
-0.315899	0.392595	0.377366	-0.330264
-0.400957	0.30818	0.466505	-0.251558
-0.491197	0.231899	0.558825	-0.182543
-0.583415	0.165791	0.649984	-0.124732
-0.674106	0.111162	0.73727	-0.0788538
-0.759645	0.0685255	0.816797	-0.0448348
-0.836489	0.0375914	0.885207	-0.0218148
-0.901371	0.0172975	0.939547	-0.00821674
-0.951479	0.00589015	0.977439	-0.00185515
-0.984611	0.00104315	0.997206	-0.0000804643

Exercise A.6.3.4

A planimeter is a device used to calculate the area of a two-dimensional region. In this problem, we explore the mathematics behind how the planimeter works.



The pointer P at one end of the planimeter follows the contour C of the surface S to be measured. For the linear planimeter the movement of the “elbow” E is restricted to the y -axis. Connected to the arm PE is the measuring wheel with its axis of rotation parallel to PE . A movement of the arm PE can be decomposed into a movement perpendicular to PE , causing the measuring wheel to rotate, and a movement parallel to PE , causing the measuring wheel to skid, with no contribution to its reading. You may assume that the length of PE is 1 unit.

Use Green’s theorem to explain why the final reading on the measuring wheel is equal to the area of the surface S .

Hints: (i) define $b(x, y)$ to be the y -coordinate of the point E when the needle is at $P = (x, y)$. Then (ii) find the component of $\langle \Delta x, \Delta y \rangle$ which is perpendicular to \vec{EP} and use your result to set up a line integral whose value equals to final reading on the meter. Then (iii) show that that the value of line

integral is equal to the desired area.

SECTION 6.4

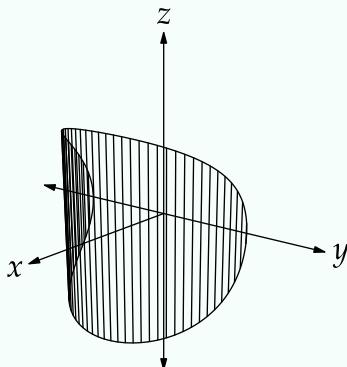
Exercise A.6.4.1

Calculate $\int_{\partial D} xy \, dS$, where $D = [0, 1]^3$.

Exercise A.6.4.2

Consider the surface $S = \left\{ (x, y, z) \in \mathbb{R}^3 : x > 0 \text{ and } r = 1 \text{ and } -\sqrt{\frac{\pi^2}{4} - \theta^2} \leq z \leq \sqrt{\frac{\pi^2}{4} - \theta^2} \right\}$, shown below. (Note that r and θ refer to cylindrical coordinates.)

- (a) Find the surface area of S by splitting it into vertical strips as shown and performing an appropriate integral.
(b) Check your answer by finding a non-calculus method of calculating the area of S .



SECTION 6.5

Exercise A.6.5.1

Find the divergence and curl of $\mathbf{F} = (2x^2, xe^z, -4y)$.

Exercise A.6.5.2

Let $f(x, y) = \log(x^2 + y^2)$ for $(x, y) \in \mathbb{R}^2 \setminus \{(0, 0)\}$. Show that $\nabla \cdot (\nabla f) = 0$.

Exercise A.6.5.3

Find the flow of the vector field $\mathbf{F} = x^2\mathbf{i} + xy\mathbf{j}$ across the surface

$$z = 1 - x^2 - y^2, \quad z \geq 0.$$

SECTION 6.6**Exercise A.6.6.1**

Verify the divergence theorem for the cube bounded by the planes $x = 0, x = 1, y = 0, y = 1, z = 0, z = 1$ and the vector field $\langle 3x, xy, 2xz \rangle$.

Exercise A.6.6.2

Confirm that the divergence theorem holds in the case $\mathbf{F} = \langle y, z, x \rangle$ and $D = \{(x, y, z) : x^2 + y^2 + z^2 \leq 16\}$.

Exercise A.6.6.3

Use the divergence theorem to calculate the flow of $\mathbf{F} = \langle x^4, -x^3z^2, 4xy^2z \rangle$ through the boundary of the region where $x^2 + y^2 \leq 1$ and $0 \leq z \leq x + 2$.

Exercise A.6.6.4

Show that the volume of a three-dimensional region D with piecewise smooth boundary is equal to

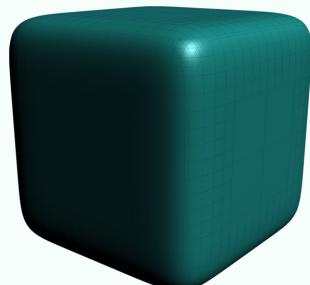
$$\iint_{\partial D} \mathbf{F} \cdot d\mathbf{A},$$

where $\mathbf{F} = \langle x, y, z \rangle$.

Exercise A.6.6.5

Use the divergence theorem to calculate the flow of the vector field $(x^2 \sin y, x \cos y, -xz \sin y)$ through the “rounded cube”

$$x^8 + y^8 + z^8 = 8.$$



SECTION 6.7

Exercise A.6.7.1

Use Stokes' theorem to evaluate $\int_C \mathbf{F} \cdot d\mathbf{r}$, where C is the triangle with vertices $(1, 0, 0)$, $(0, 1, 0)$, and $(0, 0, 1)$, oriented counterclockwise when viewed from above, and $\mathbf{F} = (x + y^2, y + z^2, z + x^2)$.

Exercise A.6.7.2

Verify Stokes' theorem for a closed sphere with the vector field $\mathbf{F} = (-y^3, x, z)$.

Exercise A.6.7.3

Find the work done on a particle by a force $\mathbf{F} = \langle xyz - e^x, -xyz, x^2yz + \sin z \rangle$ on a particle that moves along the line segments from $(0, 0, 0)$, then to $(1, 1, 1)$, then to $(0, 0, 2)$, then back to $(0, 0, 0)$.

Exercise A.6.7.4

Show that if C is a simple smooth curve contained in the plane $x + y + z = 1$, then the line integral

$$\int_C \langle z, 2x, 3z \rangle \cdot d\mathbf{r}$$

depends only on the area of the region enclosed by C . This means that if C_1 and C_2 are any two such curves enclosing the same area, then the integrals around C_1 and C_2 are equal.