

Hubbard & Hubbard's Vector Calculus, Linear Algebra, and Differential
Forms: A Unified Approach - A partial solutions manual

George Lydakis

June 1, 2024

Contents

0	Preliminaries	5
0.5	Real Numbers	5
0.6	Infinite Sets	6
0.7	Complex Numbers	12
1	Vectors, matrices and derivatives	15
1.1	Introducing the actors: Points and vectors	15
1.3	Matrix multiplication as a linear transformation	16
1.4	The geometry of \mathbb{R}^n	16
1.5	Limits and continuity	20
1.6	Five big theorems	36
1.7	Derivatives in several variables as linear transformations	38
1.8	Rules for computing derivatives	42

Chapter 0

Preliminaries

0.5 Real Numbers

Exercise 1

Show that if p is a polynomial of odd degree with real coefficients, then there is a real number c such that $p(c) = 0$.

Solution.

Let $p(x) = a_0 + a_1x + \dots + a_{2k+1}x^{2k+1}$, $k \geq 0$, $a_{2k+1} \neq 0$. For any $x \neq 0$, we can write:

$$p(x) = x^{2k+1} \left(\frac{a_0}{x^{2k+1}} + \frac{a_1}{x^{2k}} + \dots + a_{2k+1} \right)$$

Consider any term of the form $\frac{a_i}{x^m}$, $m > 0$. Consider the corresponding sequence $n \rightarrow \frac{a_i}{n^m}$. If $a_i > 0$, the infimum of this sequence is zero. To see why this is the case, observe that the fraction is always nonnegative for positive n . If we assume that the infimum is instead some positive number c , then for every $n > 0$ it must hold that:

$$\frac{a_1}{n^m} \geq c \implies \frac{a_1}{c} \geq n^m \implies \left(\frac{a_1}{c}\right)^{\frac{1}{m}} \geq n$$

, which implies that integers are bounded above, a clear contradiction. Therefore, the infimum is indeed 0. Furthermore, the sequence is non-increasing since n^m grows larger as $n > 1$ grows larger. Since the sequence is bounded below and non-increasing, it converges to its infimum, that is, zero.

If instead $a_i < 0$, arguments completely symmetrical to the above can show that the corresponding sequence again converges to zero. Clearly this is also the case for $a_i = 0$.

Now, if we consider the sequence $n \rightarrow \frac{a_0}{n^{2k+1}} + \frac{a_1}{n^{2k}} + \dots + a_{2k+1}$, and because each term converges (either to zero or a_{2k+1}), the entire sequence converges to a_{2k+1} . As n increases, n^{2k+1} always remains positive, therefore for sufficiently large values of n the sign of $p(n)$ has to equal the sign of a_{2k+1} .

A completely symmetrical argument shows that for n being the negative integers instead, for sufficiently small (by absolute value, large) values of n the sign of $p(n)$ has to equal the negative sign of a_{2k+1} .

Therefore, p takes at least one positive and one negative value, and since as a polynomial it is continuous, by the intermediate value theorem it also takes the value 0 at some real number c .

Exercise 2

a. Show that the function $f(x) =$

$$f(x) = \begin{cases} \sin(\frac{1}{x}) & , x \neq 0 \\ 0 & , x = 0 \end{cases}$$

is not continuous.

b. Show that f satisfies the conclusion of the intermediate value theorem: if $f(x_1) = a_1$, $f(x_2) = a_2$, then for any number a between a_1 and a_2 , there exists a number x between x_1 and x_2 such that $f(x) = a$.

Solution.

a. Suppose that f is continuous, more specifically that it is continuous at 0. This means that for any $\epsilon > 0$ there exists $\delta > 0$ such that whenever $|x - 0| < \delta$ it holds that $|f(x) - f(0)| < \epsilon$. Pick $\epsilon = \frac{1}{2}$. Then there must exist $\delta > 0$ such that for all x with $|x| < \delta$, it holds that $|\sin(\frac{1}{x})| < \epsilon = \frac{1}{2}$.

Recall that for $y = 2k\pi + \frac{\pi}{2}$, k integer, it is the case that $\sin(y) = 1$. Let us then examine when is it the case that $\frac{1}{x} = 2k\pi + \frac{\pi}{2}$:

$$\frac{1}{x} = 2k\pi + \frac{\pi}{2} \implies \frac{1}{2k\pi + \frac{\pi}{2}} = x \implies x = \frac{2}{4k\pi + \pi}$$

Clearly, for positive k this is always positive. Let us examine if any such numbers are in the range $[0, \delta]$. If there are, then the value of the sine function on them will be $1 > \frac{1}{2} = \epsilon$, which is a contradiction. We have that:

$$\frac{2}{4k\pi + \pi} < \delta \implies \frac{2}{\delta} < 4k\pi + \pi \implies \frac{1}{4\pi}(\frac{2}{\delta} - \pi) < k$$

Because δ must be fixed for this ϵ , and because the integers are not bounded above, we can always find a k that satisfies this inequality, leading to the contradiction described above.

Therefore, f cannot be continuous.

b. For all $x > 0$, it can be shown that f is continuous (it is the composition of two continuous functions). Therefore, for $a_1, a_2 > 0$, the conclusion of the intermediate value theorem follows from the continuity of f . The same holds for $a_1, a_2 < 0$. If $a_1 \leq 0, a_2 \geq 0$, then for any $a < 0, a \in [a_1, a_2]$, the continuity of f in $(-\infty, 0)$ and the intermediate value theorem guarantee that there exists $x \in [a_1, 0) \subset [a_1, a_2]$ such that $f(x) = a$. The same is true for $a > 0, a \in [a_1, a_2]$. For $a = 0$, we know that $f(0) = 0$, therefore we have but to pick $x = 0$. We conclude that f satisfies the conclusion of the intermediate value theorem.

Exercise 3

Suppose $a \leq b$. Show that if $f : [a, b] \rightarrow [a, b]$ is continuous, there exists $c \in [a, b]$ with $f(c) = c$.

Solution.

For every $x \in [a, b]$, it is the case that $f(x) \in [a, b]$, that is, $a \leq f(x) \leq b$. Let $g : [a, b] \rightarrow [a, b]$ with $g(x) = f(x) - x$. Because f is continuous, and the function $h(x) = x$ is —rather trivially— continuous, g is also continuous. Observe also that $g(a) = f(a) - a \geq 0, g(b) = f(b) - b \leq 0$. Therefore, $0 \in [g(b), g(a)]$, and by the intermediate value theorem (slightly modified, since here $g(b) \leq g(a)$) we have that there exists a $c \in [a, b]$ such that:

$$g(c) = 0 \implies f(c) - c = 0 \implies f(c) = c$$

0.6 Infinite Sets

Exercise 1

- Show that the set of rational numbers is countable (i.e., that you can list all rational numbers).
- Show that the set \mathbb{D} of finite decimals is countable.

Solution.

a. Consider a rational number $q \in \mathbb{Q}$. By definition, q can be written as a fraction of two positive integers, multiplied by 1 or -1, that is, $q = \text{sign}(q)\frac{k}{l}, k, l \in \mathbb{Z}$. We know that k, l can be factorized into prime factors uniquely, that is, $k = p_1^{a_1} \dots p_n^{a_n}, l = p_1^{b_1} \dots p_n^{b_n}$, where we list the first n primes that are necessary so that both numbers are factorized, and if one of them needs fewer, we “pad” the rest with zero exponents, and the signs are either. Additionally, q has a sign of either 1 or -1. By combining these observations, we have that q can be written as:

$$q = \text{sign}(q)p_1^{a_1-b_1} \dots p_n^{a_n-b_n}$$

Note that, by convention, we consider the prime factorization of 1 to be simply 1^0 and of 0 to be 0^0 . Define, now, the function $f : \mathbb{Z} \rightarrow \mathbb{N}$ such that $f(z) = 2z$ if $z \geq 0$ and $f(z) = -2z - 1$ if $z < 0$. If $f(z_1) = f(z_2)$ for $z_1, z_2 \in \mathbb{Z}$ we observe that either $f(z_1), f(z_2)$ must both be even or they must both be odd. In both cases, their equality implies $z_1 = z_2$. Therefore, f is injective. Consider now any $n \in \mathbb{Z}$. If n is even, observe

that $f(\frac{n}{2}) = 2\frac{n}{2} = n$ since $\frac{n}{2} \in \mathbb{Z}$, $\frac{n}{2} \geq 0$. If n is odd, $n+1$ is even and $-\frac{n+1}{2} \in \mathbb{Z}$, $-\frac{n+1}{2} < 0$. Therefore, $f(-\frac{n+1}{2}) = -2(-\frac{n+1}{2}) - 1 = n$. In both cases we can find a $z \in \mathbb{Z}$ such that $f(z) = n$, therefore, f is also surjective, and thus bijective.

Define also the function $g : \mathbb{Q} \rightarrow \mathbb{Z}$ such that $g(\text{sign}(q)p_1^{a_1-b_1} \dots p_n^{a_n-b_n}) = \text{sign}(q)p_1^{f(a_1-b_1)} \dots p_n^{f(a_n-b_n)}$, where we have written $q \in \mathbb{Q}$ in the manner we describe above. Observe that this is indeed meaningful because each $a_i - b_i \in \mathbb{Z}$.

If $g(q_1) = g(q_2)$, then we observe that $\text{sign}(q_1) = \text{sign}(q_2)$, since the product of primes raised to integers is always non-negative. Observe, also, that the prime factorization of any number is unique. This means that $g(q_1), g(q_2)$ feature the exact same primes with non-zero exponents, and each corresponding pair of non-zero exponents must be equal. In other words, $f(a_{1i} - b_{1i}) = f(a_{2i} - b_{2i})$, where $a_{1i} - b_{1i}, a_{2i} - b_{2i}$ are the exponents of the i -th prime in the factorization of q_1, q_2 . The injectivity of f now implies $a_{1i} - b_{1i} = a_{2i} - b_{2i}$, and because this holds for every i and $\text{sign}(q_1) = \text{sign}(q_2)$, we have that $q_1 = q_2$.

For any $z \in \mathbb{Z}$, z can be written uniquely as $z = \text{sign}(z)p_1^{e_1} \dots p_n^{e_n}$, where p_i are primes. Because $e_i \in \mathbb{N}$ and f is surjective, we have that there exist $x_i \in \mathbb{Z}$ such that $f(x_i) = e_i$. Let then $q = \text{sign}(z)p_1^{x_1} \dots p_n^{x_n}$. Clearly $q \in \mathbb{Q}$ and $g(q) = z$. Therefore, g is surjective and thus bijective.

Finally, if we define $h : \mathbb{Q} \rightarrow \mathbb{N}$, $h(q) = f(g(q))$ we see that h is a mapping from \mathbb{Q} to \mathbb{N} that is a composition of two bijective mappings, and thus is itself bijective. This therefore proves that \mathbb{Q} is countable, since it has the same cardinality as the natural numbers.

b. Consider a number $d \in \mathbb{D}$. Since d has a finite number of decimals, it can be written as $d = \text{sign}(d)10^{-x}z$, $x \in \mathbb{N}, z \in \mathbb{Z}$. Let $z = 2^{a_1}3^{a_2}5^{a_3} \dots p_n^{a_n}$ be the prime factorization of z , where we explicitly include the first three primes, and if any of them is not present in the factorization the corresponding exponent is 0. Then:

$$d = \text{sign}(d)2^{-x}5^{-x}2^{a_1}3^{a_2}5^{a_3} \dots p_n^{a_n} = \text{sign}(d)2^{a_1-x}3^{a_2}5^{a_3-x} \dots p_n^{a_n}$$

Consider the function $p : \mathbb{D} \rightarrow \mathbb{Z}$, $p(d) = \text{sign}(d)2^{f(a_1-x)}3^{a_2}5^{f(a_3-x)} \dots p_n^{f(a_n)}$. By the same reasoning as in part (a), this function is a bijection from \mathbb{D} to \mathbb{Z} .

Let, then, $r : \mathbb{D} \rightarrow \mathbb{N}$, $r(d) = f(p(d))$. Again, by the same reasoning as in part (a) this is a composition of bijections, thus is itself a bijection from \mathbb{D} to \mathbb{N} , therefore \mathbb{D} is countable.

Exercise 2

- Show that if E is finite and has n elements, then the power set $\mathcal{P}(E)$ has 2^n elements.
- Choose a map $f : \{a, b, c\} \rightarrow \mathcal{P}(\{a, b, c\})$, and compute for that map the set $\{x | x \notin f(x)\}$. Verify that this set is not in the image of f .

Solution.

a. We can do this by induction on n . For $n = 0$, E is the empty set, and its only subset is the empty set. Thus $\mathcal{P} = \{\emptyset\}$, therefore it does indeed have $2^0 = 1$ element. Suppose now that this holds for sets with $n = k \geq 0$ elements. Let E be a set of $k + 1$ elements. We can write E as $E = E' \cup \{e_n\}$, where E' is a set of $n - 1$ elements, and e_n is any of the elements of E . We know then that $\mathcal{P}(E')$ has 2^{n-1} elements, and these are all of the possible subsets of E' .

Now, any subset of E is either a subset of E' or the union of a subset of E' with $\{e_n\}$ (note that because E is a set we can have no duplicate elements). A short proof by contradiction shows why this is true: if a subset S of E could not be formed in this way, it contains at least one element not in any of the subsets of E' and not in $\{e_n\}$. But then this element is not in E' nor in $\{e_n\}$, which means that it is not in E , contradiction.

Therefore, each element of $\mathcal{P}(E')$ yields two subsets of E , one that includes e_n and one that does not. Therefore, the total number of subsets of E is $2 * 2^{n-1} = 2^n$.

b. Consider the function $f(x) = \{x\}$. Then the set $\{x | x \notin f(x)\}$ is the empty set, because by definition this function maps every element of $\{a, b, c\}$ to a set containing just that element. But then, by computing $f(a), f(b), f(c)$ we can see that the image of f is $\{\{a\}, \{b\}, \{c\}\}$, which clearly does not contain the empty set.

Exercise 5

Let $f : A \rightarrow B$ and $g : B \rightarrow A$ be one to one. We will sketch how to construct a mapping $h : A \rightarrow B$ that is one to one and onto.

Let an (f, g) -chain be a sequence consisting alternately of elements of A and B , with the element following an element $a \in A$ being $f(a) \in B$ and the element following an element $b \in B$ being $g(b) \in A$.

- a. Show that every (f, g) -chain can be uniquely continued forever to the right, and to the left can
 1. either be uniquely continued forever, or
 2. can be continued to an element of A that is not in the image of g , or
 3. can be continued to an element of B that is not in the image of f .
- b. Show that every element of A and every element of B is an element of a unique such maximal (f, g) -chain.
- c. Construct $h : A \rightarrow B$ by setting

$$h(a) = \begin{cases} f(a) & , \text{if } a \text{ belongs to a maximal chain of type 1 or 2 above} \\ g^{-1}(a) & , \text{if } a \text{ belongs to a maximal chain of type 3} \end{cases}$$

Show that $h : A \rightarrow B$ is well defined, one to one, and onto.

- d. Take $A = [-1, 1]$ and $B = (-1, 1)$. It is surprisingly difficult to write a mapping $h : A \rightarrow B$. Take $f : A \rightarrow B$ defined by $f(x) = x/2$ and $g : B \rightarrow A$ defined by $g(x) = x$.

What elements of $[-1, 1]$ belong to chains of type 1, 2, 3?

What map $h : [-1, 1] \rightarrow (-1, 1)$ does the construction in part (c) give?

Solution.

a. Let us consider an (f, g) -chain that we are interested in extending to the right. This sequence must be of the form $a, f(a), g(f(a)), \dots$ for some $a \in A$ or of the form $b, g(b), f(g(b)), \dots$ for some $b \in B$. Firstly, it is clear that successive applications of f, g can indeed continue the sequence forever, every application of one of these functions yields a result that is within the domain of the other. If there were multiple ways to continue such a sequence, this would mean that after some index i there are at least two possible continuations. However, these would arise from the application of either f or g to the element at index i , and two or more possible results would imply that either f or g is not a function, contradiction.

Suppose now that we want to continue the sequence to the left. Suppose that the sequence is of the form $a, f(a), g(f(a)), \dots$. We can do this iteratively in the following manner:

- If the first element a is not in $g(B)$, we cannot continue the sequence. The reason is that if we could, there would be an element b that could be prepended before a , and then due to the rules of our sequence creation it would have to be that $g(b) = a \implies a \in g(B)$, contradiction.
- Otherwise, there exists *unique* $b \in g(B)$ such that $g(b) = a$. Note that uniqueness is guaranteed by the fact that g is one-to-one. Prepend b to the sequence. Repeat the previous step, but this time examine whether $b \in f(A)$. The rest of the argument described in these two steps remains the same (since f is also one-to-one).

Clearly if the sequence was initially of the form $b, g(b), f(g(b)) \dots$ we could apply the exact same procedure (it is, in essence, the same case “shifted” by one “prepending step”). It is also clear that the “prepending” described above will terminate in two cases: either we reach an element of A not in the image of g , or an element of B not in the image of f . Otherwise, the procedure will continue forever, and the resulting sequence will be unique due to the fact that at each step there is a unique element that can be prepended (due to f, g being one-to-one).

b. Pick any element a of A or any element b of B and consider it a sequence of length 1 where the element is at index 0. We proved in (a) that the sequence a or b can be uniquely extended infinitely to the right. Therefore, we only need to examine whether the continuation to the left is also unique. Note that every time we took a step in the way described in (a) there was a unique way to do this, due to f, g being one

to one. Therefore, two different finite continuations of the same length are impossible. Furthermore, an infinite continuation is unique, as we observed in (a). Two finite continuations of different lengths or one finite continuation and one infinite one are the only remaining case.

But this would imply, due to step uniqueness, that these are equal for all elements of the shorter of the two. However, this would also imply that the shorter one ends in either $x \in A$ or $y \in B$, and that the longer one has a preceding element z such that $g(z) = x, z \in B$ or $f(z) = y, z \in A$, respectively. But this would imply that the left-most element of the shorter sequence was either in the image of g or f respectively, and then the procedure described in (a) would not have terminated, contradiction.

To complete the proof that a maximal (f, g) -chain containing a or b is unique, note that if some other sequence contained it in *some* index i , then the fact that f, g are one-to-one implies that index $i + k, k \in \mathbb{Z}$ of that sequence would equal index k of the maximal chain which we started with a or b at index 0. In other words the sequence equals the original one up to a translation of the indices by a fixed integer.

c. First we check whether h is well defined as a function.

Consider any $a \in A$. Then in part (b) we proved that there is a unique maximal chain containing a . Since this chain is unique, its type is also unique. If it is of type 1 or 2, $h(a) = f(a)$ which is unique due to f being a function. If it is of type 3, then the chain containing a has to have as leftmost element some element $b \in B$, therefore a cannot be the leftmost element, therefore there is at least one element x before a . Then the rules of how the maximal chain is formed impose that $g(x) = a$. In other words, $g^{-1}(a) = x$ (where $g^{-1}(a)$ denotes the element of B such that $g(g^{-1}(a)) = a$ and is well defined whenever $a \in g(B)$). Therefore, $h(a) = g^{-1}(a)$ is again uniquely defined.

Suppose now that $h(a_1) = h(a_2)$ for $a_1, a_2 \in A$. If a_1, a_2 's maximal chains are of type 1 or 2 (not necessarily the same) then $h(a_1) = f(a_1) = h(a_2) = f(a_2)$. Then, the injectivity of f implies that $a_1 = a_2$. If a_1 's maximal chain is of type 1 or 2 and a_2 's is of type 3, we have the following:

- The maximal chain of a_1 is of the form $\dots, a_1, f(a_1), \dots$ or of the form $a_x, \dots, a_1, f(a_1), \dots$, with $a_x \in A$ (where possibly $a_1 = a_x$).
- The maximal chain of a_2 is of the form $b_x, \dots, b_1, a_2, \dots$, with $b_1, b_x \in B$ (where possibly $b_1 = b_x$).
- It holds that $h(a_1) = f(a_1) = h(a_2) = g^{-1}(a_2) = b_1$.
- Because the maximal chain of a_2 is unique, this implies that the maximal chain of a_1 appears in the left of a_2 , so that $f(a_1) = b_1$. If the maximal chain of a_1 extends infinitely to the left (type 1) then the maximal chain of a_2 cannot be of finite length, contradiction. If it is of type 2, then either the maximal chain of a_2 has to end in a_x or the maximal chain of a_1 has to end in b_x , both of which are contradictions.

Therefore, in this case we get a contradiction.

Finally, if $h(a_1) = h(a_2)$ and a_1, a_2 belong to maximal chains of type 3, then $g^{-1}(a_1) = g^{-1}(a_2)$. This means that there exist $b_1 = g^{-1}(a_1), b_2 = g^{-1}(a_2) \in A$ such that $g(b_1) = a_1, g(b_2) = a_2$. But $b_1 = b_2$, thus $a_1 = a_2$. Therefore, in all cases $h(a_1) = h(a_2)$ implies $a_1 = a_2$, which means h is injective.

Now let $b \in B$. There exists a unique maximal chain containing b .

If this is of type 1 or 2, there is at least one element a immediately to the left of b , such that $f(a) = b$. Then it is also the case that this is the maximal chain of that a , otherwise we could obtain two maximal chains for b . Then, by definition, $h(a) = f(a) = b$.

If this is of type 3, then the element $a = g(b)$ immediately to the right of b also has a maximal chain of type 3, and additionally $g^{-1}(a) = b$. Therefore, $h(a) = g^{-1}(a) = b$.

In both cases, we can find an $a \in A$ such that $h(a) = b$, which means that h is onto.

d. Consider any $a \in A = [-1, 1]$. If its maximal chain is of type 1, then it extends infinitely to the left. Observe that, due to the nature of f, g , this chain is of the form $\dots, 4a, 4a, 2a, 2a, a, a, \dots$. This is because $g(x) = x$ and f divides its argument by 2. Therefore, we observe that $2^n a$ must be an element of A for every $n \in \mathbb{N}$. In other words:

$$-1 \leq 2^n a \leq 1 \implies \frac{-1}{2^n} \leq a \leq \frac{1}{2^n}$$

The only way this can be true for any $n \in \mathbb{N}$ is if $a = 0$. Therefore, only 0 has a maximal chain of type 1. For this element, $h(0) = f(0) = 0$.

If the maximal chain of a is of type 2, then it ends to the left at some $a' \notin g(B)$. The image of g is clearly $(-1, 1)$, whereas $A = [-1, 1]$. Therefore the leftmost element a' can either be 1 or -1. Successive applications of f and g lead to the conclusion that a has to be of the form $a = \frac{-1}{2^n}$ or $a = \frac{1}{2^n}, n \in \mathbb{N}$. For those elements, $h(a) = f(a) = \frac{a}{2}$.

Recall now that every element of A must be an element of a unique maximal chain. Therefore, all elements that have not been accounted for so far must have a maximal chain of type 3. These elements are precisely $x \in [-1, 1], x \neq 0, x \neq \frac{-1}{2^n}, x \neq \frac{1}{2^n}$, for every $n \in \mathbb{N}$. For those elements, $h(x) = g^{-1}(x)$, which must be the immediately preceding element in the chain, for which $g(g^{-1}(x)) = a$, and because $g(x) = x$ we have in the end that $h(x) = x$.

Therefore, we can write the resulting map from $[-1, 1]$ to $(-1, 1)$ as:

$$h(a) = \begin{cases} 0 & , \text{if } a = 0 \\ \frac{a}{2} & , \text{if } a \text{ is of the form } -\frac{1}{2^n} \text{ or } \frac{1}{2^n} \text{ for some } n \in \mathbb{N} \\ a & , \text{otherwise} \end{cases}$$

Exercise 6

Show that the points of the circle $\{(x, y) \in \mathbb{R}^2 | x^2 + y^2 = 1\}$ have the same infinity of elements as \mathbb{R} .

Hint: This is easy if you use Bernstein's theorem (Exercise 0.6.5).

Solution.

Let C denote the set of points on the given circle. By exercise 5, it suffices to show that there exist two injective functions $f : C \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow C$. First, let us consider f . In particular, consider the following function:

$$f((x, y)) = \begin{cases} x & , \text{if } y \geq 0 \\ x + 2 & , \text{if } y < 0 \end{cases}$$

Essentially, this function "squashes" the upper semicircle (intersection points with x' included) to $[-1, 1]$ and the lower semicircle to $(1, 3)$. The first observation is obvious, and the second comes from the fact that for $y < 0$, it must be the case that $x \in (-1, 1)$, therefore $1 < x + 2 < 3$. The mapping is injective, because if $f((x_1, y_1)) = f((x_2, y_2))$ we have the following cases.

- $y_1 \geq 0, y_2 \geq 0$, in which case $x_1 = x_2$, but then the circle equation and the positivity of y_1, y_2 implies $y_1 = y_2$, therefore $(x_1, y_1) = (x_2, y_2)$.
- $y_1 < 0, y_2 < 0$, then $x_1 + 2 = x_2 + 2$, and the circle equation plus the negativity of y_1, y_2 imply that $y_1 = y_2$, therefore again $(x_1, y_1) = (x_2, y_2)$.
- $y_1 \geq 0, y_2 < 0$, in which case it must be true that $x_1 = x_2 + 2$. But $x_2 \in [-1, 1]$ thus $x_2 + 2 \in [1, 3]$ and $x_1 \in [-1, 1]$ which yields that $x_1 = 1$. Then, however, $x_2 = 1$, which would mean that $y_2 = 0 = y_1$, which is a contradiction to $y_2 < 0$.

Therefore f is injective. Now observe that $(-1, 1)$ can be mapped one-to-one to the upper semicircle by mapping each $x \in (-1, 1)$ to $(x, \sqrt{1 - x^2})$. Let this mapping be called h_1 . Now consider the function $h_2 : \mathbb{R} \rightarrow (-1, 1), h_2(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$. For this function we have that:

$$\begin{aligned} h_2(x_1) = h_2(x_2) &\implies \frac{e^{x_1} - e^{-x_1}}{e^{x_1} + e^{-x_1}} = \frac{e^{x_2} - e^{-x_2}}{e^{x_2} + e^{-x_2}} \implies \\ e^{x_1+x_2} + e^{x_1-x_2} - e^{-x_1+x_2} - e^{-x_1-x_2} &= e^{x_2+x_1} + e^{x_2-x_1} - e^{-x_2+x_1} - e^{-x_2-x_1} \\ &\implies 2e^{x_1-x_2} = 2e^{x_2-x_1} \implies x_1 = x_2 \end{aligned}$$

, which means that is injective. Therefore the mapping $g : \mathbb{R} \rightarrow C$ such that $g(x) = h_1(h_2(x))$ is also injective.

By exercise 5 (Schröder - Bernstein theorem) we have that this guarantees the existence of a bijection between C and \mathbb{R} .

Exercise 7

- a. Show that $[0, 1) \times [0, 1)$ has the same cardinality as $[0, 1)$.
- b. Show that \mathbb{R}^2 has the same infinity of elements as \mathbb{R} .
- c. Show that \mathbb{R}^n has the same infinity of elements as \mathbb{R} .

Solution.

a. Consider first the mapping $f : [0, 1) \rightarrow [0, 1) \times [0, 1)$, $f(x) = (x, 0)$. Clearly, this is an injective mapping. Consider also the mapping $g : [0, 1) \times [0, 1) \rightarrow [0, 1)$, $g((0.a_1a_2a_3\dots, 0.b_1b_2b_3\dots)) = 0.a_1b_1a_2b_2a_3b_3\dots$, where we have written the arguments in decimal form. Note that we make the following choice: if a number ends in infinite 9's, i.e. can be written as $0.a_1a_2\dots a_n999\dots$, we choose to write it instead as $0.a_1a_2\dots a_n + 0.00\dots 1$, where the second term in this addition has 1 in the n -th decimal. In essence, we choose to write the number "rounded up", and we know that in the case of infinite 9's the two numbers are equal.

With that in mind, let $(x, y) = (0.a_1a_2\dots, 0.b_1b_2\dots)$, $(z, w) = (0.c_1c_2\dots, 0.d_1d_2\dots)$ such that $g((x, y)) = g((z, w))$. This would mean that the numbers $0.a_1b_1a_2b_2\dots, 0.c_1d_1c_2d_2\dots$ are equal. This could be the case in one of two ways. One, if all of their —potentially infinite— decimals are equal. If this is true, then we directly have that $x = z, y = w$, thus that $(x, y) = (z, w)$. Two, if one of them ends in infinite 9's after decimal n , the other in infinite 0's, and the decimals of the second equal the result of adding $0.0\dots 01$ (1 in the n -th decimal) to the number formed by the first n digits of the first.

However, if this was the case it would mean that after some n either all of a_i, b_i or all of c_i, d_i are 9's, which does not follow the convention we mentioned previously. In other words, if this was the case we would already have written those numbers x, y or z, w such that they end in 0's and not 9's.

Therefore, whenever $g((x, y)) = g((z, w))$ we have that $(x, y) = (z, w)$, therefore g is injective. The existence of f, g and the Schröder-Bernstein theorem imply that there exists a bijective mapping between $[0, 1) \times [0, 1)$ and $[0, 1)$, which means that these have the same cardinality.

b. We have seen, as part of exercise 6, that the real numbers have the same infinity of elements as $[0, 1)$. Using the corresponding bijective mapping $h_1(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$, we can also bijectively map \mathbb{R}^2 to $[0, 1) \times [0, 1)$ by defining $h_2((x, y)) = (h_1(x), h_1(y))$. The existence of these bijective mappings as well as the existence of the bijective mapping $h : [0, 1) \times [0, 1) \rightarrow [0, 1)$ that we proved in part (a) imply that $h \circ h_2 : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a bijective mapping from \mathbb{R}^2 to \mathbb{R} , showing that these have the same infinity of elements.

c. If we now define $h_n : \mathbb{R}^n \rightarrow [0, 1) \times \dots \times [0, 1)$, $h_n((x_1, x_2, \dots, x_n)) = (h_1(x_1), h_1(x_2), \dots, h_1(x_n))$, we can see that this is a bijection from \mathbb{R}^n to $[0, 1) \times \dots \times [0, 1)$ (where the Cartesian product both here and above has n terms). If we follow the same convention as in part (a) and define a mapping from \mathbb{R}^n to $[0, 1) \times \dots \times [0, 1)$ to $[0, 1)$ as $g((a_{11}a_{12}a_{13}\dots, a_{21}a_{22}a_{23}\dots, \dots, a_{n1}a_{n2}a_{n3}\dots)) = 0.a_{11}a_{21}\dots a_{n1}a_{12}a_{22}\dots a_{n2}\dots$, as well as $f : [0, 1) \rightarrow [0, 1) \times \dots \times [0, 1)$, $f(x) = (x, 0, \dots, 0)$, the Schröder-Bernstein theorem tells us that $[0, 1)$ has the same cardinality as $[0, 1) \times \dots \times [0, 1)$, and since we can bijectively map $[0, 1)$ to \mathbb{R} and \mathbb{R}^n to $[0, 1) \times \dots \times [0, 1)$, we can also bijectively map \mathbb{R}^n to \mathbb{R} , thus proving that these have the same infinity of elements.

Exercise 8

Show that the power set $\mathcal{P}(\mathbb{N})$ has the same cardinality as \mathbb{R} .

Solution.

Consider the following mapping $f : \mathcal{P}(\mathbb{N}) \rightarrow \mathbb{R}$: Given a set of natural numbers $\{n_1, n_2, n_3, \dots\}$, $f(\{n_1, n_2, n_3, \dots\}) = 0.a_00a_10a_20\dots$, where a_i is 1 if $i \in \{n_1, n_2, n_3, \dots\}$ and 0 otherwise. Observe that if two sets $s_1, s_2 \in \mathcal{P}(\mathbb{N})$ are such that $f(s_1) = f(s_2)$, then it must be the case that all corresponding decimals must be equal. Importantly for this argument, f by its construction can never yield numbers that end in infinite 1s, and as such $f(s_1) = f(s_2)$ cannot happen for e.g. $0.1 = 0.01111\dots$. Also by construction, all digits can only be 0 or 1. For any natural number n , $n \in s_1$ iff the i -th decimal of $f(s_1)$ is 1, and the same holds for s_2 . These two facts imply that $n \in s_1$ iff $n \in s_2$, therefore the sets contain the exact same elements, therefore they are equal, therefore f is injective.

We also know from previous exercises that \mathbb{R} can be mapped bijectively to $[0, 1]$. Therefore, if we can find an injective mapping from $[0, 1]$ to $\mathcal{P}(\mathbb{N})$, we will also have found an injective mapping from \mathbb{R} to $\mathcal{P}(\mathbb{N})$. Consider then the function $g : [0, 1] \rightarrow \mathcal{P}(\mathbb{N})$, $g(0.a_1a_2a_3\dots) = \{p_1^{a_1}, p_2^{a_2}, \dots\}$, where, again, p_i is the i -th prime number. Note that if one or more a_i are zero, the number $p_i^0 = 1$ is included only once in the set,

since sets cannot contain duplicate elements. Furthermore, if a number ends in infinite 9's, we first write it as ending in zeros by “rounding up” as described in 0.5. Suppose now that $g(0.a_1a_2a_3\dots) = g(0.b_1b_2b_3\dots)$. Then the sets $\{p_1^{a_1}, p_2^{a_2}, p_3^{a_3}, \dots\}, \{p_1^{b_1}, p_2^{b_2}, p_3^{b_3}, \dots\}$ are equal. Because a power of the i -th prime number can appear in the set only originating from the i -th decimal—since it cannot be factorized as a power of any other prime number—we have that for every p_i the corresponding exponents must be equal, that is, $a_i = b_i$. One possible exception is the number 1. However, if the number 1 is included in the set, it can only have originated from one or more a_i or b_i being 0. In particular, every a_i or b_i such that a power of p_i does *not* appear in the set must be zero, because otherwise the sets would contain $p_i^{a_i}$ or $p_i^{b_i}$ respectively. Therefore, because the two sets contain the exact same elements, $a_i = 0$ iff $b_i = 0$. Thus, $a_i = b_i$ for all i , which means that g is injective, which means that there exists also an injective mapping g' from \mathbb{R} to $\mathcal{P}(\mathbb{N})$

By the Bernstein-Schröder theorem, f and g' guarantee the existence of a bijective mapping between \mathbb{R} and $\mathcal{P}(\mathbb{N})$, which means that these have the same cardinality.

0.7 Complex Numbers

Exercise 11

- Describe the loci in \mathbb{C} given by the following equations:
 - $\operatorname{Re}(z) = 1$
 - $|z| = 3$
- What is the image of each locus under the mapping $z \rightarrow z^2$?
- What is the inverse image of each locus under the mapping $z \rightarrow z^2$?

Solution.

a. We first consider case (i). Any number $z = a + bi$ in the locus of $\operatorname{Re}(z) = 1$ must have a real part of 1, and all numbers $z = 1 + bi$ belong in this locus. Therefore, the locus is all complex numbers of the form $z = 1 + bi, b \in \mathbb{R}$, which geometrically is the vertical line $x = 1$.

For case (ii), we seek all complex numbers $z = a + bi$ such that their modulus is 3. It is easy to see that these are precisely the complex numbers on the circle with center $(0, 0)$ and radius 3, i.e. complex numbers of the form $z = 3(\cos(\theta) + i\sin(\theta)), \theta \in \mathbb{R}$.

b. For (i), consider z^2 whenever $z = 1 + bi$. We have that $z^2 = (1 + bi)^2 = 1 - b^2 + 2bi, b \in \mathbb{R}$. To visualize this on the plane, set $x = 1 - b^2, y = 2b$, which implies $x = 1 - \frac{y^2}{4}$. This corresponds to parabola that intersects the y axis at $(0, 2), (0, -2)$ and the x axis at $(1, 0)$.

For (ii), we have that $z^2 = 3^2(\cos(2\theta) + i\sin(2\theta)), \theta \in \mathbb{R}$. Clearly, this corresponds to a circle with radius 9 and center $(0, 0)$.

c. For (i), we need to solve $z^2 = 1 + bi$. If $z = x + yi$, we have that:

$$(x + yi)(x + yi) = 1 + bi \implies x^2 - y^2 + 2xyi = 1 + bi \implies x^2 - y^2 = 1, 2xy = b$$

Because $b \in \mathbb{R}$, x, y can be chosen freely and $2xy = b$ will always be satisfied for some b . Therefore, it suffices to hold that $x^2 - y^2 = 1$, which corresponds to a hyperbola intersecting the x axis at $(-1, 0), (1, 0)$ (we could also work out the foci of this hyperbola). Geometrically, squaring any complex number lying on this hyperbola yields a complex number with real part 1, which thus lies on the original locus.

For (ii), we have to solve $z^2 = 3(\cos(\theta) + i\sin(\theta)), \theta \in \mathbb{R}$. Clearly, any number satisfying this equation has modulus $\sqrt{3}$, and conversely, if a number has modulus $\sqrt{3}$, its square has modulus 3, and therefore belongs in the locus described in (a). This means that the inverse image of $z \rightarrow z^2$ is, in this case, the circle with center $(0, 0)$ and radius $\sqrt{3}$.

Exercise 13

- Find all the cubic roots of 1.
- Find all the 4th roots of 1.
- Find all the 6th roots of 1.

Solution.

The n -th roots of 1 have the form: $z_k = \cos(\frac{0+2k\pi}{n}) + i\sin(\frac{0+2k\pi}{n})$, $k = 0, 1, \dots, n-1$. Therefore, we have but to replace n with the appropriate number to solve the subtasks.

a. For $n = 3$, $z_1 = 1$, $z_2 = \cos(\frac{2\pi}{3}) + i\sin(\frac{2\pi}{3}) = -\frac{1}{2} + \frac{\sqrt{3}}{2}i$, $z_3 = \cos(\frac{4\pi}{3}) + i\sin(\frac{4\pi}{3}) = -\frac{1}{2} - \frac{\sqrt{3}}{2}i$.

Tasks (b) and (c) are solved similarly, so we omit the calculations. More interesting is the fact that these constitute the vertices of a regular n -gon inscribed in the unit circle of the complex plane, such that one vertex is always on $(1, 0)$.

Chapter 1

Vectors, matrices and derivatives

1.1 Introducing the actors: Points and vectors

Exercise 8

Suppose that water is flowing through a pipe of radius r , with speed $r^2 - a^2$, where a is the distance to the axis of the pipe.

- Write the vector field describing the flow if the pipe is in the direction of the z -axis.
- Write the vector field describing the flow if the axis of the pipe is the unit circle in the (x, y) -plane.

Solution.

a. We assume that the axis of the pipe is the z -axis. Therefore, for any given point (x, y, z) inside the vertical cylinder of radius r , its distance from the pipe axis is $\sqrt{x^2 + y^2}$. Water flows downwards, and in parallel to the z -axis. Therefore, any flow vector will have to be of the form $(0, 0, z)$, $z \leq 0$. Furthermore, the speed of “falling” will have to equal the magnitude of the flow vector, which means that z will equal $-r^2 + a^2$. Denote the set of points inside the pipe as $C = \{(x, y, z) \in \mathbb{R}^3 \mid \sqrt{x^2 + y^2} < r\}$. Then, the vector field describing the flow is:

$$F : C \rightarrow \mathbb{R}^3, F(x, y, z) = (0, 0, -r^2 + x^2 + y^2)$$

b. The only things that change here are the domain of the vector field, and the formula for the distance. In particular, a cross section of the pipe with a horizontal plane now looks like a ring with inner radius 1 and outer radius r , where the pipe is the space between the two circles. Also, the distance of any point (x, y, z) of the pipe from the axis is its distance from the circle $\{(w_1, w_2, z) \in \mathbb{R}^3 \mid w_1^2 + w_2^2 = 1\}$, i.e. the unit circle on the horizontal plane with defined by the point's z -coordinate. It is simple enough to compute this distance intuitively, although proving it rigorously would be a bit harder: measure the distance of the point from the circle's center, $(0, 0, z)$, and subtract the radius of the circle. This equals $\sqrt{x^2 + y^2} - 1$. Denote, now, the set of points inside the pipe as $C = \{(x, y, z) \in \mathbb{R}^3 \mid 1 < x^2 + y^2 < r^2\}$. Then the flow vector field will be:

$$F : C \rightarrow \mathbb{R}^3, F(x, y, z) = (0, 0, -r^2 + (\sqrt{x^2 + y^2} - 1)^2) = (0, 0, -r^2 + x^2 + y^2 + 1 - 2\sqrt{x^2 + y^2})$$

1.3 Matrix multiplication as a linear transformation

Exercise 20

Identify \mathbb{R}^2 to \mathbb{C} by identifying $\begin{pmatrix} a \\ b \end{pmatrix}$ to $z = a + ib$.

Show that the following mappings $\mathbb{C} \rightarrow \mathbb{C}$ are linear transformations, and give their matrices:

- $\text{Re} : z \rightarrow \text{Re}(z)$ (the real part of z)
- $\text{Im} : z \rightarrow \text{Im}(z)$ (the imaginary part of z)
- $c : z \rightarrow \bar{z}$ (the complex conjugate of z , i.e. $\bar{z} = a - ib$ if $z = a + ib$)
- $m_w : z \rightarrow wz$, where $w = u + iv$ is a fixed complex number

Solution.

a. Let $z_1 = a_1 + b_1i, z_2 = a_2 + b_2i$. Then $z_1 + z_2 = (a_1 + a_2) + (b_1 + b_2)i$, which means that $\text{Re}(z_1 + z_2) = \text{Re}(z_1) + \text{Re}(z_2)$. Additionally, if $\lambda \in \mathbb{R}$ (the underlying field is assumed to be \mathbb{R}), we have that $\lambda z_1 = \lambda(a_1 + b_1i) = \lambda a_1 + \lambda b_1i$, thus $\text{Re}(\lambda z_1) = \lambda \text{Re}(z_1)$. Thus, the function Re satisfies both homogeneity and additivity, therefore it is linear. To give its matrix, it suffices to note its behavior on the standard basis, that is, $\text{Re}(1 + 0i) = 1, \text{Re}(0 + 1i) = 0$. Therefore, the matrix with respect to this basis is:

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

Note that due to the identification we make, we can write $f(a + bi)$ for a linear map operating in \mathbb{R}^2 without confusion.

b. This can be proved in exactly the same way as (a), and the only thing that changes is that $\text{Im}(1 + 0i) = 0, \text{Im}(0 + 1i) = 1$, thus the matrix will be:

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

c. Again, for $z_1 = a_1 + b_1i, z_2 = a_2 + b_2i$ we have that $\overline{z_1 + z_2} = \overline{(a_1 + a_2) + (b_1 + b_2)i} = (a_1 + a_2) - (b_1 + b_2)i = \bar{z}_1 + \bar{z}_2$. Also, for $\lambda \in \mathbb{R}$, we have $\overline{\lambda z_1} = \overline{\lambda a_1 + \lambda b_1i} = \lambda a_1 - \lambda b_1i = \lambda \bar{z}_1$. This means that the function mapping z to its conjugate is both additive and homogeneous, therefore it is linear. We additionally have that $c(1 + 0i) = 1, c(0 + 1i) = -i$, thus the matrix must be:

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

d. For $z_1 = a_1 + b_1i, z_2 = a_2 + b_2i$ we have that $m_w(z_1 + z_2) = w(z_1 + z_2) = wz_1 + wz_2 = m_w(z_1) + m_w(z_2)$, by the distributive property of complex multiplication on addition. Also, for $\lambda \in \mathbb{R}$ we have $m_w(\lambda z) = w(\lambda z) = \lambda(wz) = \lambda m_w$, by the commutative and associative properties of complex multiplication. Therefore m_w is indeed linear. We also have that $m_w(1 + 0i) = w(1 + 0i) = w = u + iv = u \begin{pmatrix} 1 \\ 0 \end{pmatrix} + v \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ and $m_w(0 + 1i) = w(0 + 1i) = (u + iv)i = -v + ui = -v \begin{pmatrix} 1 \\ 0 \end{pmatrix} + u \begin{pmatrix} 0 \\ 1 \end{pmatrix}$. This means that the corresponding matrix must be:

$$\begin{pmatrix} u & -v \\ v & u \end{pmatrix}$$

1.4 The geometry of \mathbb{R}^n

Exercise 12

- Show that $\|v + w\| \geq \|v\| - \|w\|$
- True or false? $|\det([a, b, c])| \leq \|a\| \cdot \|b \times c\|$. Explain your answer. What does it mean geometrically?

Solution.

a. Suppose v, w are vectors. Let $A = \|v + w\|^2 - (\|v\| - \|w\|)^2$. We then have that:

$$A = \langle v + w, v + w \rangle - \langle v, v \rangle + 2\|v\| \cdot \|w\| - \langle w, w \rangle = \langle v, v \rangle + 2\langle v, w \rangle + \langle w, w \rangle - \langle v, v \rangle + 2\|v\| \cdot \|w\| - \langle w, w \rangle \\ \implies A = 2\langle v, w \rangle + 2\|v\| \cdot \|w\|$$

By the Cauchy-Schwarz inequality we have that $|\langle v, w \rangle| \leq \|v\| \cdot \|w\|$. This implies that $-\|v\| \cdot \|w\| \leq \langle v, w \rangle$, which directly gives us also that $A \geq 0$. But then:

$$\|v + w\|^2 - (\|v\| - \|w\|)^2 \geq 0 \implies \|v + w\| \geq \left| \|v\| - \|w\| \right|$$

This implies, of course, that $\|v + w\| \geq \|v\| - \|w\|$.

b. We know that $|\det([a, b, c])| = |a \cdot (b \times c)|$. By the Cauchy-Schwarz inequality, $|a \cdot (b \times c)| \leq \|a\| \cdot \|b \times c\|$. This then directly yields that $|\det([a, b, c])| \leq \|a\| \cdot \|b \times c\|$ for any three vectors a, b, c .

Geometrically, we know that the left-hand side equals the volume of the parallelepiped formed by a, b, c . The right-hand side equals the volume of the parallelepiped formed by the base defined by b, c and by a third vector orthogonal to their plane, with magnitude equal to the magnitude of a . Therefore, the inequality tells us that the volume of this second parallelepiped is always at least equal to the volume of the first one. Essentially, you can never “lose” volume by making one of the edges orthogonal to the base formed by the other two.

Exercise 13

Show that the cross product of two vectors pointing in the same direction is zero.

Solution.

Recall that the cross product of a, b is defined so that its length is the area of the parallelogram spanned by a, b . If these are colinear, said parallelogram is degenerate and its area is zero. The only vector that can have zero length is the zero vector, thus $a \times b = 0$. An alternative proof would involve writing out the determinant-based definition of the cross product and observing that one of a, b can be written as a scalar multiple of the other.

Exercise 15

Given two vectors $v, w \in \mathbb{R}^3$, show that $(v \times w) = -(w \times v)$.

Solution.

Two proofs can be given here.

One, by going through the properties of the cross product we see that both $w \times v, v \times w$ must have the same magnitude, because this equals the area of the parallelogram spanned by v, w (the order doesn't change anything). We also see that they must be orthogonal to the plane spanned by them, so they are also colinear. Therefore, either they are equal or one equals the additive inverse of the other. However, $v \times w$ must be such that $\det([v, w, v \times w]) > 0$ if v, w are not colinear (if they are, we showed in 13 that $v \times w = 0 = w \times v = -w \times v$). It must also be the case that $\det([w, v, w \times v]) > 0$. If $v \times w = w \times v$, then observe that this second matrix could be obtained by the first by exchanging the first two columns, and thus the two determinants would have opposite signs, which is a contradiction. Therefore $v \times w = -w \times v$. The second proof involves writing out the determinant-based definition of $v \times w$, and then swapping the columns of each 2×2 matrix whose determinant we are computing, which incurs a multiplication with -1 for each coordinate. But then this new determinant-based vector is nothing but $w \times v$, and thus $v \times w = -w \times v$.

Exercise 24

Let $v \in \mathbb{R}^n$ be a non-zero vector and denote by $v^\perp \subset \mathbb{R}^n$ the set of vectors $w \in \mathbb{R}^n$ such that $v \cdot w = 0$.

- Show that v^\perp is a subspace of \mathbb{R}^n .
- Given any vector $a \in \mathbb{R}^n$, show that $a - \frac{\langle a, v \rangle}{\|v\|^2}v$ is an element of v^\perp .
- Define the projection of a onto v^\perp by the formula

$$P_{v^\perp}(a) = a - \frac{\langle a, v \rangle}{\|v\|^2}v$$

Show that there is a unique number $t(a)$ such that $(a + t(a)v) \in v^\perp$. Show that

$$a + t(a)v = P_{v^\perp}(a)$$

Solution.

a. We have that:

- $0 \in v^\perp$, because the zero vector is orthogonal to all vectors.
- If $w_1, w_2 \in v^\perp$, we have that $\langle v, w_1 + w_2 \rangle = \langle v, w_1 \rangle + \langle v, w_2 \rangle = 0 + 0 = 0$. Therefore, $w_1 + w_2 \in v^\perp$.
- If $w \in v^\perp$ and $\lambda \in \mathbb{R}$, we have that $\langle v, \lambda w \rangle = \lambda \langle v, w \rangle = 0$, therefore $\lambda w \in v^\perp$.

Therefore, v^\perp is indeed a subspace of \mathbb{R}^n .

b. We have that:

$$\langle v, a - \frac{\langle a, v \rangle}{\|v\|^2}v \rangle = \langle v, a \rangle - \langle v, \frac{\langle a, v \rangle}{\|v\|^2}v \rangle = \langle v, a \rangle - \frac{\langle a, v \rangle}{\|v\|^2} \langle v, v \rangle = \langle v, a \rangle - \langle a, v \rangle = 0$$

, which means of course that $a - \frac{\langle a, v \rangle}{\|v\|^2}v \in v^\perp$.

c. We must examine when is it the case that $\langle v, a + xv \rangle = 0, x \in \mathbb{R}$. We have that:

$$\langle v, a + xv \rangle = 0 \iff \langle v, a \rangle + \langle v, xv \rangle = 0 \iff \langle v, a \rangle + x \langle v, v \rangle = 0 \iff \langle v, a \rangle + x\|v\|^2 = 0$$

Recall that v is non-zero, thus $\|v\|^2 \neq 0$. The equation above is linear in x , and x 's coefficient is non-zero. Thus, it has a unique solution, namely:

$$x = -\frac{\langle v, a \rangle}{\|v\|^2} = -\frac{\langle a, v \rangle}{\|v\|^2}$$

We can thus say that for every $a \in \mathbb{R}^3$ there exists a unique number $t(a) = -\frac{\langle a, v \rangle}{\|v\|^2}$ such that $(a + t(a)v) \in v^\perp$. Clearly, from the definition of P_{v^\perp} , it holds that $a + t(a)v = P_{v^\perp}(a)$.

Exercise 27

Let $v = \begin{pmatrix} a \\ b \\ c \end{pmatrix}$ be a unit vector in \mathbb{R}^3 , so that $a^2 + b^2 + c^2 = 1$.

- Show that the transformation T_v defined by $T_v(u) = u - 2(v \cdot u)v$ is a linear transformation $\mathbb{R}^3 \rightarrow \mathbb{R}^3$.
- What is $T_v(v)$? If u is orthogonal to v , what is $T_v(u)$? Can you give a name to T_v ?
- Write the matrix M of T_v (in terms of a, b, c of course). What can you say of M^2 ?

Solution.

a. Clearly, $T_v(u)$ is a vector in \mathbb{R}^3 . We have that:

- For $u_1, u_2 \in \mathbb{R}^3$, $T_v(u_1 + u_2) = (u_1 + u_2) - 2(v \cdot (u_1 + u_2))v = u_1 + u_2 - 2(v \cdot u_1 + v \cdot u_2)v = u_1 - 2(v \cdot u_1)v + u_2 - 2(v \cdot u_2)v = T_v(u_1) + T_v(u_2)$, therefore T_v is additive.
- For $u \in \mathbb{R}^3, \lambda \in \mathbb{R}$, $T_v(\lambda u) = \lambda u - 2(v \cdot \lambda u)v = \lambda u - 2\lambda(v \cdot u)v = \lambda(u - 2(v \cdot u)v) = \lambda T_v(u)$, therefore T_v is homogeneous.

We conclude that T_v is linear.

b. We have that:

$$T_v(v) = v - 2(v \cdot v)v = v - 2||v||^2v = v(1 - 2||v||^2) = -v$$

, where we used the fact that v is a unit vector. For u orthogonal to v , we have that:

$$T_v(u) = u - 2(v \cdot u)v = u$$

T_v is a transformation that first projects its argument on the vector v , computes two times this projection, inverts the result and adds the resulting vector to its argument. By this observation, and by its behavior on v and the vectors orthogonal to v , we can see that T_v reflects its argument using as plane of reflection the plane orthogonal to v .

c. We need only examine the values $T_v\left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}\right), T_v\left(\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}\right), T_v\left(\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}\right)$. We have that:

$$T_v\left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}\right) = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - 2\left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \end{pmatrix}\right) \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1 - 2a^2 \\ -2ab \\ -2ac \end{pmatrix}$$

Similarly, we obtain:

$$T_v\left(\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}\right) = \begin{pmatrix} -2ab \\ 1 - 2b^2 \\ -2bc \end{pmatrix}, T_v\left(\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}\right) = \begin{pmatrix} -2ac \\ -2bc \\ 1 - 2c^2 \end{pmatrix}$$

Thus, the matrix of T_v is:

$$M = \begin{pmatrix} 1 - 2a^2 & -2ab & -2ac \\ -2ab & 1 - 2b^2 & -2bc \\ -2ac & -2bc & 1 - 2c^2 \end{pmatrix}$$

Our intuition about reflection tells us that it must be the case that $M^2 = I$. In order not to compute the matrix explicitly, we observe that:

$$\begin{aligned} T_v(T_v(u)) &= T_v(u) - 2(v \cdot T_v(u))v = u - 2(v \cdot u)v - 2(v \cdot (u - 2(v \cdot u)v))v \\ &= u - 2(v \cdot u)v - 2(v \cdot u - 2(v \cdot u)(v \cdot v))v = u - 2(v \cdot u)v - 2(v \cdot u - 2v \cdot u)v \\ &= u \end{aligned}$$

, thereby confirming that T_v^2 is the identity function, and thus that it must also be the case that $M^2 = I$.

Exercise 28

Let A be an $n \times m$ matrix. Show that $|A|^2 = \text{tr}(A^T A)$.

Solution.

Let:

$$A = \begin{pmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{pmatrix}$$

Then:

$$A^T = \begin{pmatrix} a_{11} & \dots & a_{n1} \\ \vdots & \ddots & \vdots \\ a_{1m} & \dots & a_{nm} \end{pmatrix}$$

By definition:

$$|A|^2 = \sum_{i=1}^n \sum_{j=1}^m a_{ij}^2$$

Carrying out the matrix multiplication $A^T A$ (an $m \times m$ matrix) we obtain that:

$$A^T A = \begin{pmatrix} \sum_{k=1}^n a_{k1}a_{k1} & \cdots & \sum_{k=1}^n a_{k1}a_{km} \\ \vdots & \ddots & \vdots \\ \sum_{k=1}^n a_{km}a_{k1} & \cdots & \sum_{k=1}^n a_{km}a_{km} \end{pmatrix}$$

Then, the trace is equal to:

$$\text{tr}(A^T A) = \sum_{i=1}^m \sum_{k=1}^n a_{ki}^2$$

Thus it is indeed true that $\text{tr}(A^T A) = |A|^2$.

Another proof, given that $A^T A$ is an $m \times m$ matrix:

$$\text{tr}(A^T A) = \sum_{i=1}^m (A^T A)_{ii} = \sum_{i=1}^m \sum_{k=1}^n (A^T)_{ik} A_{ki} = \sum_{i=1}^m \sum_{k=1}^n A_{ki} A_{ki} = \sum_{i=1}^m \sum_{k=1}^n A_{ki}^2 = |A|^2$$

, where we used the definition of matrix multiplication and the definition of the transpose.

1.5 Limits and continuity

Exercise 1

For each of the following subsets, state whether it is open or closed (or both or neither) and say why.

- (a) $\{x \in \mathbb{R} | 0 < x \leq 1\}$ as a subset of \mathbb{R}
- (b) $\{(x, y) \in \mathbb{R}^2 | \sqrt{x^2 + y^2} < 1\}$ as a subset of \mathbb{R}^2
- (c) the interval $(0, 1]$ as a subset of \mathbb{R}
- (d) $\{(x, y) \in \mathbb{R}^2 | \sqrt{x^2 + y^2} \leq 1\}$ as a subset of \mathbb{R}^2
- (e) $\{x \in \mathbb{R} | 0 \leq x \leq 1\}$ as a subset of \mathbb{R}
- (f) $\{(x, y, z) \in \mathbb{R}^3 | \sqrt{x^2 + y^2 + z^2} \leq 1, x, y, z \neq 0\}$ as a subset of \mathbb{R}^3
- (g) the empty set as a subset of \mathbb{R}

Solution.

(a) This set is not open. To see why this is the case, consider its point $x = 1$. Any open ball (in this case, open interval) around it must include points of the form $1 + \epsilon, \epsilon > 0$, which do not belong in the set. Furthermore, the set is not closed. Its complement is $\{x \in \mathbb{R} | x \leq 0 \text{ or } x > 1\}$. Again, any open ball around $x = 0$ must include points of the form $0 + \epsilon, \epsilon > 0$, which do not belong in the set. Therefore the complement is not open, which means the set is not closed.

(b) This set is open (it is the unit disk in \mathbb{R}^2 , boundary excluded). To prove this formally, name the set S and consider any point $P = (x, y) \in S$. Then $\|P\| = \sqrt{x^2 + y^2} < 1$, therefore set $r = 1 - \|P\| > 0$. We claim that the open ball $B_r(P)$ is completely contained in S . Indeed, if $P' \in B_r(P)$, then $\|P'\| = \|P + (P' - P)\| \leq \|P\| + \|P' - P\| < \|P\| + r = \|P\| + 1 - \|P\| = 1$, where we've used the fact that P' is contained in the open ball, and thus by definition $\|P' - P\| < r$. The proven inequality yields that $P' \in S$, thus that $B_r(P) \subset S$, therefore S is open.

The set is not closed since its complement contains the unit circle, and taking any open ball around a point on the unit circle will include points whose magnitude is less than 1, and thus belong in S .

(c) This is identical to (a).

(d) This set, let it again be called S , is closed and not open. Let us first examine why it is closed. Its complement is $S' = \{(x, y) \in \mathbb{R}^2 | \sqrt{x^2 + y^2} > 1\}$. Consider any point $P \in S'$. Then $\|P\| > 1$, therefore set $r = \|P\| - 1$ and take $B_r(P)$ the open ball of radius r around it. For any point $P' \in B_r(P)$ it holds that $\|P'\| = \|P' - P + P\| \geq \| \|P\| - \|P' - P\| \|$. It holds that $\|P' - P\| < \|P\| - 1$, therefore $\|P\| - \|P' - P\| > 1$, thus $\|P'\| > 1$, which means $P' \in S'$, which means that $B_r(P) \subset S'$. Thus S' is open, which means S is closed.

S is not open because any open ball around, for example, $(1, 0) \in S$ must include points that have a norm of $\|(1, 0)\| + \epsilon > 1$, which thus do not belong in S .

(e) This set is closed and not open. It is closed because its complement is $\{x \in \mathbb{R} | x < 0 \text{ or } x > 1\}$, and for any point x in it the open ball $B_{|x|}(x)$ if $x < 0$ or $B_{|x|-1}(x)$ if $x > 1$ is completely contained in the complement, which means that the complement is open and thus the set is closed.

The set is not open because any ball $B_r(0)$ must include points of the form $-\epsilon, \epsilon > 0$, which thus do not belong in the set.

(f) This set, S , is neither closed nor open. It is not open because, for example, any open ball around $P = (\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}) \in S$ must include points that have norm $\|P\| + \epsilon > 1$, and thus do not belong in the set. The set is also not closed. Its complement is

$$S' = \{(x, y, z) \in \mathbb{R}^3 | \sqrt{x^2 + y^2 + z^2} > 1 \text{ or at least one of } x, y, z \text{ is zero} \}$$

This includes, for example, the point $(0, 0, 0)$. Any open ball $B_r((0, 0, 0))$ of radius $r < 1$ includes some points of the form $(\frac{\epsilon}{\sqrt{3}}, \frac{\epsilon}{\sqrt{3}}, \frac{\epsilon}{\sqrt{3}})$, which clearly have norm less than 1, and all of their coordinates are non-zero, thus they belong in S and not in S' . Consequently, S' is not open, and S is not closed.

(g) The empty set is both closed and open. It is open because the statement “for any point in the subset, there exists an open ball around it such that it is fully contained in the set” is vacuously true. It is closed because its complement is \mathbb{R} , and this is clearly an open set, since any open ball around any of its points belongs in \mathbb{R} .

Exercise 2

For each of the following subsets, state whether it is open or closed (or both or neither), and say why.

- (a) (x, y) -plane in \mathbb{R}^3
- (b) $\mathbb{R} \subset \mathbb{C}$
- (c) the line $x = 5$ in the (x, y) -plane
- (d) $(0, 1) \subset \mathbb{C}$
- (e) $\mathbb{R}^n \subset \mathbb{R}^n$
- (f) the unit sphere in \mathbb{R}^3

Solution.

Let us first outline a methodology for proving “open-ness”/“closed-ness”.

- For one, showing that a set C is open is equivalent to showing that $\mathbb{R}^n - C$ is closed. Indeed, if $\mathbb{R}^n - C$ is closed, C is by definition open, since it is its complement. On the other hand, if C is open, then its complement $\mathbb{R}^n - C$ has an open complement (C) and thus is closed.
- Secondly, showing that a set is closed is equivalent to showing that every convergent sequence in it converges to a point that is also in it.
- Thirdly, for any continuous function f , and any convergent sequence x_i that converges to x_0 , we know that the sequence $f(x_i)$ converges to $f(x_0)$.
- Suppose now we want to show that C is closed. Take any convergent sequence x_i in it converging to x_0 . If we can find a continuous f such that $f(x_i)$ all have some property that implies x_i belong in C , and if f maintains this property at the limit, then x_0 also belongs in the set, and thus C is closed.

The following lemma will also be useful: the only subsets of \mathbb{R}^n that are both open and closed are \emptyset and \mathbb{R}^n . We will prove this by contradiction. Suppose there is a set $X, X \neq \emptyset, X \neq \mathbb{R}^n$ that is both open and closed. Then, there exists at least one $x \in X$ and at least one $y \notin X$. Consider all vectors of the form:

$$x + \lambda(y - x), \lambda \in [0, 1]$$

Now let S be the set of all $\lambda \in [0, 1]$ such that $x + \lambda(y - x) \in X$. Since $y \notin X$, this set has an upper bound, 1. By the completeness of the real numbers, this means that it has a supremum, λ_0 . It is either the case that $\lambda_0 \in S$ or $\lambda_0 \notin S$.

- Suppose $\lambda_0 \in S$, which means $x + \lambda_0(y - x) \in X$, and by necessity $\lambda_0 < 1$. Consider the sequence of points:

$$i \rightarrow y + f(i)(x - y), i = 1, 2, \dots,$$

$$f(1) = 0, f(i) = \frac{1 - \lambda_0}{1 + \frac{1}{i}}, i > 1$$

We can fairly easily see that as i approaches infinity, this sequence converges to $y + (1 - \lambda_0)(x - y) = x + \lambda_0(y - x) \in X$. Additionally, by rewriting we can see that:

$$i \rightarrow x - x + y + \frac{1 - \lambda_0}{1 + \frac{1}{i}}(x - y) = x + (y - x)\left(1 - \frac{1 - \lambda_0}{1 + \frac{1}{i}}\right) = x + \frac{\frac{1}{i} + \lambda_0}{1 + \frac{1}{i}}(y - x)$$

It is true that $\frac{\frac{1}{i} + \lambda_0}{1 + \frac{1}{i}} > \lambda_0 \iff \lambda_0 + \frac{1}{i} > \lambda_0 + \frac{\lambda_0}{i} \iff 1 > \lambda_0$, which by the definition of λ_0 means that every term of the sequence (including obviously $0 \rightarrow y$) belongs in $\mathbb{R}^n - X$. We have thus found a sequence in $\mathbb{R}^n - X$ that converges to a point outside of it, which means it is not closed, and thus X cannot be open, contradiction.

- Suppose $\lambda_0 \notin S$, which means $x + \lambda_0(y - x) \notin X$. If $x + \lambda_0(y - x)$ is not in the closure of X , then there exists an open ball $B_r(x + \lambda_0(y - x))$ such that it has no intersection with X . But then this ball includes points of the form $x + \lambda'(y - x)$, $\lambda' < \lambda_0$ that do not belong in X , and this contradicts the definition of λ_0 as the supremum of S . Thus $x + \lambda_0(y - x) \in \overline{X}$, $x + \lambda_0(y - x) \notin X$. But then for any open ball around $x + \lambda_0(y - x)$, we can find a point in it that belongs in X , and by decreasing the radii of these balls we can form a convergent sequence of points in X that converges to $x + \lambda_0(y - x) \notin X$, which means that X is not closed, contradiction.

Both cases result in contradictions, thus no set other than the empty set and \mathbb{R}^n can be both open and closed. Thus from now on showing that a —non-empty, not equal to \mathbb{R}^n — set is open directly implies it is not closed, and vice versa.

With that in mind, we have that:

(a) This set is defined as $C = \{(x, y, z) \in \mathbb{R}^3 | z = 0\}$. Consider any convergent subsequence $v_i = (x_i, y_i, 0) \in C$ that converges to $v_0 = (x_0, y_0, z_0)$. Consider the function $f(x, y, z) = z$, which is continuous as a polynomial. We know then that:

$$\lim_{i \rightarrow \infty} f(v_i) = f(v_0) \implies f(v_0) = \lim_{i \rightarrow \infty} 0 = 0 \implies z_0 = 0$$

Clearly, this means that $v_0 \in C$, and we've therefore shown that C is a closed set.

(b) As a subset of \mathbb{C} , we can think of \mathbb{R} as $\mathbb{R} = \{z \in \mathbb{C} | \text{Im}\{z\} = 0\}$. Consider the function $f(z) = \text{Im}\{z\} : \mathbb{C} \rightarrow \mathbb{C}$. For any $z_0 \in \mathbb{C}$, we have the following. Given any $\epsilon > 0$, set $\delta = \epsilon$, and then for every z such that $\|z - z_0\| < \delta$:

$$\begin{aligned} \|f(z) - f(z_0)\|^2 &= |\text{Im}\{z - z_0\}|^2 \leq |\text{Re}\{z - z_0\}|^2 + |\text{Im}\{z - z_0\}|^2 = \|z - z_0\|^2 < \delta^2 \\ \implies \|f(z) - f(z_0)\| &< \epsilon \end{aligned}$$

, which means that f is continuous. Now pick any convergent sequence $z_j = x_j + 0i \in \mathbb{R}$ that converges to z_0 . Then:

$$\lim_{j \rightarrow \infty} f(z_j) = f(z_0) \implies \lim_{j \rightarrow \infty} \text{Im}\{z_j\} = f(z_0) \implies \text{Im}\{z_0\} = f(z_0) = \lim_{j \rightarrow \infty} 0 = 0$$

This means that $z_0 \in \mathbb{R}$, and thus \mathbb{R} is closed.

(c) We have $C = \{(x, y) \in \mathbb{R}^2 | x = 5\}$. Consider the function $f(x, y) = x$, which is continuous as a polynomial. Then, consider any convergent sequence $v_i = (5, y_i) \in C$ converging to $v_0 = (x_0, y_0)$. We have that:

$$\lim_{i \rightarrow \infty} f(v_i) = f(v_0) \implies \lim_{i \rightarrow \infty} 5 = f(v_0) \implies x_0 = 5$$

, which means that $v_0 \in C$, which means that C is closed.

(d) We have that $C = \{z \in \mathbb{C} | 0 < \operatorname{Re}\{z\} < 1, \operatorname{Im}\{z\} = 0\}$ and its complement is $\mathbb{C} - (0, 1) = \{z \in \mathbb{C} | \operatorname{Re}\{z\} \leq 0 \text{ or } \operatorname{Re}\{z\} \geq 1 \text{ or } \operatorname{Im}\{z\} \neq 0\}$. Consider any convergent sequence $z_i \in \mathbb{C} - (0, 1)$ converging to z_0 , and let $f(z) = \operatorname{Im}\{z\}(\operatorname{Re}\{z\} - 1)\operatorname{Re}\{z\}$. By the same reasoning we applied in (b) for the imaginary part as a function we can also show that the real part is continuous. Therefore f is continuous as a product of continuous functions. Observe furthermore that f is zero for any $z \in \mathbb{C} - (0, 1)$. We therefore have that:

$$\lim_{z_i \rightarrow z_0} f(z) = f(z_0) \implies f(z_0) = \lim_{z_i \rightarrow z_0} 0 \implies f(z_0) = 0$$

But this means that at least one of the factors of $f(z_0)$ is zero, which means that by definition it belongs in $\mathbb{C} - (0, 1)$, thus this set is closed and $(0, 1)$ is open.

(e) $\mathbb{R}^n \subset \mathbb{R}^n$ is both open and closed. Indeed, any convergent sequence in it converges trivially to a point also in it, thus it is closed. Additionally, for any point in it we can always —trivially— find an open ball around it that also belongs in \mathbb{R}^n , and thus \mathbb{R}^n is also open.

(f) The unit sphere in \mathbb{R}^3 is the set $C = \{(x, y, z) \in \mathbb{R}^3 | \sqrt{x^2 + y^2 + z^2} = 1\}$. Consider any convergent sequence v_i in it, converging to $v_0 = (x_0, y_0, z_0)$. Consider also the function $f(x, y, z) = \sqrt{x^2 + y^2 + z^2}$, which is continuous as a composition of continuous functions. Observe that $f(v_i) = 1$ for all v_i . Then, we have that:

$$\lim_{v_i \rightarrow v_0} f(v_i) = f(v_0) \implies \lim_{v_i \rightarrow v_0} 1 = f(v_0) \implies 1 = f(v_0)$$

But this means that $1 = \sqrt{x_0^2 + y_0^2 + z_0^2}$, thus that $v_0 \in C$, which means that C is closed.

Exercise 3

Prove the following statements for open subsets of \mathbb{R}^n :

- (a) Any union of open sets is open.
- (b) A finite intersection of open sets is open.
- (c) An infinite intersection of open sets is not necessarily open.

Solution.

(a) Consider the open sets $S_i, i \in I$, where I is an arbitrary index set. Let $S = \bigcup_i S_i$. Let x be any point in S . By definition, x belongs in at least one S_i . Since S_i is open, there exists an open ball $B_r(x) \subset S_i$. Because $S_i \subset S$, it is also true that $B_r(x) \subset S$, thus for any point in S we can find an open ball around it that is entirely contained in S , therefore S is indeed open.

(b) Consider the open sets S_1, S_2, \dots, S_n . Let $S = \bigcap_i S_i$. Suppose $x \in S$. This means that $x \in S_i$ for all i . Since all of these sets are open, there exist open balls $B_{r_i}(x) \subset S_i$ for all i . Let r_m be the smallest of these radii. This is well defined because the set $\{r_i\}$ is finite, and thus has a minimum element. Clearly, $B_{r_m}(x)$ is a subset of all S_i , which means that it is also a subset of their intersection S . We have thus found an open ball around x that is completely contained in S , thus S is open.

(c) For any real number $r > 0$, consider the set $S_r = B_r(0)$. Let $S = \bigcap S_r$, i.e. the intersection of all of these sets. Observe that this set contains 0: for any positive r , 0 belongs in $B_r(0)$, and thus it belongs in S as well. Observe also that the set contains no other points. Indeed, suppose $v \neq 0$ is in S . Note that $\|v\| > 0$, thus there exists $\epsilon > 0$ such that $0 < \epsilon < \|v\|$. Consequently, we have that v is contained in any $B_r(0)$, more specifically that it is contained in $B_\epsilon(0)$, which is a contradiction. Therefore $S = \{0\}$, and this is a closed set (any convergent sequence in it, which is just the trivial sequence 0, converges to $0 \in S$).

Exercise 4

- (a) Show that the interior of A is the biggest open set contained in A .
- (b) Show that the closure of A is the smallest closed set that contains A .
- (c) Show that the closure of a set A is A plus its boundary: $\overline{A} = A \cup \partial A$.
- (d) Show that the boundary is the closure minus the interior: $\partial A = \overline{A} - \mathring{A}$.

Solution.

(a) The interior of A is the set of $x \in A$ such that there exists an open ball $B_r(x) \subset A$. First of all, note that by definition the interior of A is indeed contained in A . Secondly, we need to prove that it is open. Select any $x \in \mathring{A}$. Then there exists $B_r(x) \subset A$. Now select any $y \in B_r(x)$. Let $r_y = r - \|y - x\|$. This is clearly a strictly positive quantity: $y \in B_r(x)$, thus $\|y - x\| < r$. Now select any $z \in B_{r_y}(y)$. It holds that:

$$\|z - y\| < r_y \implies \|z - y\| < r - \|y - x\| \implies \|y - x\| + \|z - y\| < r$$

By the triangle inequality, $\|z - x\| \leq \|z - y\| + \|y - x\| < r$. This means that $z \in B_r(x)$. Consequently, $z \in A$, which means that $B_{r_y}(y) \subset A$, which means that y belongs in the interior of A . We have thus proved that for $x \in \mathring{A}$, $B_r(x) \subset \mathring{A}$, thus proving that \mathring{A} is an open set.

Finally, we need to prove that \mathring{A} is the *biggest* open set contained in A . Select then any open $S \subset A$. For any $x \in S$, there exists $r > 0$ such that $B_r(x) \subset S \subset A$. But since $B_r(x) \subset A$, x by definition belongs in the interior of A . Thus $S \subset \mathring{A}$ for any open S , which means that \mathring{A} is indeed the biggest open set contained in A .

(b) The closure of A is the set of all x such that for all $r > 0$, $B_r(x) \cap A \neq \emptyset$. First of all, note for any $x \in A$, and for any $r > 0$ it holds that $B_r(x) \cap A$ contains at least x , thus A is indeed contained in \overline{A} . We now need to prove that the closure is indeed closed.

Select any convergent sequence $x_i \in \overline{A}$ that converges to x_0 . Since $x_i \rightarrow x_0$, there exists $M > 0$ such that $\|x_i - x_0\| < r$ for every $i > M$. But then this means that for every $r > 0$, the open ball $B_r(x_0)$ contains some elements of A , namely all $x_i, i > M$. Therefore, for every $r > 0$, $B_r(x_0) \cap A \neq \emptyset$, which means $x_0 \in \overline{A}$, thus the closure of A is closed.

Now we need to show that \overline{A} is the *smallest* closed set that contains A . Select then any closed $S \supset A$, and let $x \in \overline{A}$. By definition, for any $r > 0$, $B_r(x) \cap A \neq \emptyset$. More specifically, for any positive integer i , we can find $x_i \in A$ such that $x_i \in B_i(x)$. This sequence of x_i clearly converges to x . Because $x_i \in A$, $x_i \in S$, and because S is closed, the limit of the sequence, x , must also be contained in S . Thus for any $x \in \overline{A}$, $x \in S$, which means that the closure of A is indeed contained in all closed sets that contain A .

(c) The boundary of A is the set of all x such that every neighborhood of x intersects both A and the complement of A . More precisely, in \mathbb{R}^n this means that any open ball around x intersects both A and $\mathbb{R}^n - A$.

Let $x \in \overline{A}$. Then for all $r > 0$, $B_r(x)$ such that $B_r(x) \cap A \neq \emptyset$. If $x \in A$, clearly $x \in A \cup \partial A$. If $x \notin A$, then the open ball $B_r(x)$ contains at least one element not in A , namely, x itself. In addition, $B_r(x) \cap A \neq \emptyset$, thus $B_r(x)$ indeed intersects both A and its complement. Therefore $x \in \partial A$, thus $x \in A \cup \partial A$, thus $\overline{A} \subset A \cup \partial A$. Now let $x \in A \cup \partial A$. If $x \in A$, $x \in \overline{A}$. If $x \in \partial A$, then for all $r > 0$, the open ball $B_r(x)$ intersects both A and $\mathbb{R}^n - A$, thus it intersects \overline{A} , thus $x \in \overline{A}$. Therefore $A \cup \partial A \subset \overline{A}$. We conclude that $\overline{A} = A \cup \partial A$.

(d) We need to show that $\partial A = \overline{A} \cap (\mathbb{R}^n - \mathring{A})$.

Suppose $x \in \partial A$. Then clearly $x \in \overline{A}$ because any open ball around x intersects both A and A 's complement. Additionally, x cannot belong in \mathring{A} , because that would imply that at least one $B_r(x)$ is completely contained in A , i.e. $B_r(x) \cap (\mathbb{R}^n - A)$ would be empty. Thus $x \in \mathbb{R}^n - \mathring{A}$, thus $x \in \overline{A} \cap (\mathbb{R}^n - \mathring{A})$.

Now suppose $x \in \overline{A} \cap (\mathbb{R}^n - \mathring{A})$. Since $x \in \overline{A}$, any open ball around x intersects A . Since $x \in (\mathbb{R}^n - \mathring{A})$, all open balls around x must intersect $\mathbb{R}^n - A$, otherwise at least one of them would be completely contained in A and thus x would belong in \mathring{A} . Because each neighborhood of x contains an open ball around it, and this open ball was just shown to have non-empty intersections with both A and $\mathbb{R}^n - A$, the neighborhood also intersects both these sets, thus $x \in \partial A$. Therefore $\overline{A} \cap (\mathbb{R}^n - \mathring{A}) \subset \partial A$.

We have thus shown that $\partial A = \overline{A} \cap (\mathbb{R}^n - \mathring{A})$.

Exercise 5

For each of the following subsets of \mathbb{R} and \mathbb{R}^2 , state whether it is open or closed (or both or neither), and prove it.

- (a) $\{(x, y) \in \mathbb{R}^2 | 1 < x^2 + y^2 < 2\}$
- (b) $\{(x, y) \in \mathbb{R}^2 | xy \neq 0\}$
- (c) $\{(x, y) \in \mathbb{R}^2 | y = 0\}$
- (d) $\{\mathbb{Q} \subset \mathbb{R}\}$

Solution.

We first prove the following useful lemma: If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is continuous and $K \subset \mathbb{R}^m$ is closed, then the set $f^{-1}(K) = \{x \in \mathbb{R}^n | f(x) \in K\}$ is also closed. To do this, suppose x_i is a convergent sequence in $f^{-1}(K)$ that converges to x_0 . Since f is continuous, we have that:

$$\lim_{x_i \rightarrow x_0} f(x_i) = f(x_0)$$

By the definition of the inverse image, all of $f(x_i)$ belong in K . Additionally, we just showed that the sequence $f(x_i)$ converges to $f(x_0)$. Because K is closed, $f(x_0)$ also belongs in K . But then again by the

definition of the inverse image, x_0 must belong in $f^{-1}(K)$. Since x_0 is the limit of the sequence x_i , which was selected arbitrarily, we have equivalently shown that $f^{-1}(K)$ is indeed a closed set. With that in mind we have that:

(a) The complement of this set is $\{(x, y) \in \mathbb{R}^2 | x^2 + y^2 \leq 1 \text{ or } x^2 + y^2 \geq 2\}$. Let $f(x, y) = x^2 + y^2$, and consider the set $K = (-\infty, 1] \cup [2, \infty)$. This is easily shown to be a closed set (its complement is the open interval $(1, 2)$). By applying the lemma above on our —continuous— f and closed K we have that $f^{-1}(K)$ is closed. But $f^{-1}(K)$ are precisely the (x, y) such that $x^2 + y^2 \leq 1$ or $x^2 + y^2 \geq 2$, which is the complement of the given set, thus the set itself is open (and by a previous result cannot be closed).

(b) The complement of this set is $\{(x, y) \in \mathbb{R}^2 | xy = 0\}$. Consider the continuous function $f(x, y) = xy$. Consider also the set $K = \{0\}$. This is trivially a closed set (its complement is the union of two open intervals). Therefore, once again $f^{-1}(K)$ is also closed, and this equals precisely all (x, y) such that $xy = 0$, i.e. the complement of the given set, thus the set itself is open (and cannot be closed).

(c) Let $f(x, y) = y$ and $K = \{0\}$. K is closed and f is continuous, thus $f^{-1}(K) = \{(x, y) \in \mathbb{R}^2 | y = 0\}$ is also closed (and cannot be open).

(d) We will take as a given the fact that $\sqrt{2}$ is not a rational number. Consider first the sequence of numbers $x_i = \frac{\sqrt{2}}{i}, i = 1, 2, \dots$. These are all clearly irrational numbers, i.e. they all belong in the complement of \mathbb{Q} . However, the limit of this sequence can easily be shown to be 0, which is a rational number. Therefore the complement of \mathbb{Q} is not a closed set, thus \mathbb{Q} is not open.

If we now consider the sequence x_i such that the i -th number in the sequence equals the number formed by the first i decimals of $\sqrt{2}$, these are all clearly rational numbers (they have a finite number of digits). However, the limit of the sequence is clearly $\sqrt{2}$, which is irrational. Therefore \mathbb{Q} is not closed either.

Exercise 7

For each of the following formulas, find its natural domain, and show whether the natural domain is open, closed or neither.

- (a) $\sin \frac{1}{xy}$
- (b) $\log \sqrt{x^2 - y}$
- (c) $\log \log x$
- (d) $\arcsin \frac{3}{x^2 + y^2}$
- (e) $\sqrt{e^{\cos xy}}$
- (f) $\frac{1}{xyz}$

Solution.

(a) This formula “requires” $xy \neq 0$. Therefore the natural domain is $\{(x, y) \in \mathbb{R}^2 | xy \neq 0\}$. This is precisely part (b) of exercise 5, thus the natural domain is open and not closed.

(b) This formula “requires” $\sqrt{x^2 - y} > 0$ and $x^2 - y \geq 0$. The two constraints restrict the natural domain to $\{(x, y) \in \mathbb{R}^2 | x^2 - y > 0\}$. The complement of this set is $\{(x, y) \in \mathbb{R}^2 | x^2 - y \leq 0\}$, which can easily be shown to be a closed set by considering the continuous function $f(x, y) = x^2 - y$ and the closed set $K = (-\infty, 0]$. Thus the natural domain itself is open.

(c) Here we require $x > 0$ and $\log x > 0$. By taking exponents in the second inequality we obtain $x > 1$, thus the natural domain is $(1, \infty)$, which is clearly open and not closed.

(d) Here we require firstly $x^2 + y^2 \neq 0$ and secondly $-1 \leq \frac{3}{x^2 + y^2} \leq 1$. The second inequality can be rewritten as $x^2 + y^2 \geq -3$ and simultaneously $3 \leq x^2 + y^2$. Clearly putting all of the constraints together yields $3 \leq x^2 + y^2$, and the corresponding set $\{(x, y) \in \mathbb{R}^2 | 3 \leq x^2 + y^2\}$ is easily shown to be closed and not open.

(e) The quantity in the square root is clearly non-negative, thus the natural domain is the entire \mathbb{R}^2 , which is both open and closed.

(f) Here we require $xyz \neq 0$, which is similar to (a) and yields a natural domain that is open and not closed.

Exercise 9

Suppose $\sum_{i=1}^{\infty} x_i$ is a convergent series in \mathbb{R}^n . Show that the triangle inequality applies:

$$\left\| \sum_{i=1}^{\infty} x_i \right\| \leq \sum_{i=1}^{\infty} \|x_i\|$$

Solution.

To begin, note that the LHS equals the norm of the vector to which the series converges, and thus is always well defined. The RHS, though, may diverge, in which case it trivially holds that the LHS is “less than infinity”. From now on assume that the LHS also converges. We can rewrite this inequality as:

$$\lim_{n \rightarrow \infty} \left\| \sum_{i=1}^n x_i \right\| \leq \lim_{n \rightarrow \infty} \sum_{i=1}^n \|x_i\|$$

Because both limits exist, we can write this also as:

$$\lim_{n \rightarrow \infty} \left(\left\| \sum_{i=1}^n x_i \right\| - \sum_{i=1}^n \|x_i\| \right) \leq 0$$

Suppose that this is not true, i.e. that the limit equals a positive number $c > 0$. Then, by the definition of the limit of a sequence, for every $\epsilon > 0$ there must exist $M > 0$ such that for every $n > M$ it holds that:

$$\begin{aligned} \left| \left\| \sum_{i=1}^n x_i \right\| - \sum_{i=1}^n \|x_i\| - c \right| < \epsilon &\implies -\epsilon < \left\| \sum_{i=1}^n x_i \right\| - \sum_{i=1}^n \|x_i\| - c < \epsilon \\ \implies c - \epsilon < \left\| \sum_{i=1}^n x_i \right\| - \sum_{i=1}^n \|x_i\| < c + \epsilon \end{aligned}$$

Because $c > 0$, there exists $0 < \epsilon < c$, otherwise c would be the supremum of the set $\{0\}$. Thus we can always choose ϵ and find the corresponding M such that the left-most quantity above is strictly positive. Note, however, that for this *finite* M , the triangle inequality holds:

$$\left\| \sum_{i=1}^n x_i \right\| \leq \sum_{i=1}^n \|x_i\| \implies \left\| \sum_{i=1}^n x_i \right\| - \sum_{i=1}^n \|x_i\| \leq 0$$

However, the immediately preceding argument showed that this quantity is strictly positive. We get a contradiction, thus $c \leq 0$, which shows that the triangle inequality holds at the limit as well.

Alternative proof

We first prove that the norm as a function in \mathbb{R}^n is continuous. Let then $f(x) = \|x\|$. Consider any sequence x_i that converges to x_0 . We need to show that $\lim_{i \rightarrow \infty} f(x_i) = f(x_0) \implies \lim_{i \rightarrow \infty} \|x_i\| = \|x_0\|$. Pick any $\epsilon > 0$. Then there exists $M > 0$ such that for every $i > M$ it holds that $\|x_i - x_0\| < \epsilon$. We also have that:

$$|\|x_i\| - \|x_0\|| \leq \|x_i - x_0\| < \epsilon$$

But this means that the sequence of $f(x_i)$ indeed converges to $f(x_0)$, thus the norm is continuous. Now because the given series $\sum_{i=1}^{\infty} x_i$ converges, $f(\sum_{i=1}^n x_i)$ converges too, and in fact it converges to $f(\sum_{i=1}^{\infty} x_i)$. Namely:

$$\lim_{n \rightarrow \infty} f\left(\sum_{i=1}^n x_i\right) = \lim_{n \rightarrow \infty} \left\| \sum_{i=1}^n x_i \right\| = \left\| \lim_{n \rightarrow \infty} \sum_{i=1}^n x_i \right\|$$

Limits preserve non-strict inequalities. Therefore:

$$\left\| \sum_{i=1}^n x_i \right\| \leq \sum_{i=1}^n \|x_i\| \implies \left\| \lim_{n \rightarrow \infty} \sum_{i=1}^n x_i \right\| \leq \lim_{n \rightarrow \infty} \sum_{i=1}^n \|x_i\|$$

Exercise 10

Let A be an $n \times n$ matrix and define:

$$e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k = I + A + \frac{1}{2}A^2 + \frac{1}{3!}A^3 + \dots$$

- (a) Show that the series converges for all A , and find a bound for $\|e^A\|$ in terms of $\|A\|$ and n .
- (b) Compute explicitly e^A for $A = \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}$, $A = \begin{pmatrix} 0 & a \\ 0 & 0 \end{pmatrix}$, $A = \begin{pmatrix} 0 & a \\ -a & 0 \end{pmatrix}$.
- (c) Prove the following or find counterexamples:
1. $e^{A+B} = e^A e^B$ for all A, B
 2. $e^{A+B} = e^A e^B$ for A, B that satisfy $AB = BA$
 3. $e^{2A} = (e^A)^2$ for all A

Solution.

(a) Firstly, let's consider the case in which A is nilpotent. From linear algebra we know that this means that $A^n = 0$, that any power $k > n$ of A is also zero and that potentially some previous powers are zero as well. In any case, it becomes clear that the sum is a sum of a finite number of non-zero terms, which obviously converges. From now on we thus assume that A is not nilpotent, which means $\|A^k\|$ is never zero. Consider the series:

$$b_n = \sum_{k=0}^n \frac{1}{k!} \|A^k\|$$

This is a series in \mathbb{R} . We can therefore use the well known ratio test to see if it converges. We thus examine the limit:

$$\lim_{n \rightarrow \infty} \left| \frac{b_{n+1}}{b_n} \right| = \lim_{n \rightarrow \infty} \left| \frac{\frac{1}{(n+1)!} \|A^{n+1}\|}{\frac{1}{n!} \|A^n\|} \right| = \lim_{n \rightarrow \infty} \frac{\|A^{n+1}\|}{(n+1) \|A^n\|}$$

Recall that it holds that $\|AB\| \leq \|A\| \cdot \|B\|$. Thus:

$$\|A^{n+1}\| \leq \|A^n\| \cdot \|A\| \implies \frac{\|A^{n+1}\|}{\|A^n\|} \leq \|A\| \implies \frac{\|A^{n+1}\|}{(n+1) \|A^n\|} \leq \frac{\|A\|}{n+1}$$

Clearly, the RHS here goes to zero as n goes to infinity, since the norm of A is constant. This means that the same is true for the LHS as well, thus the ratio we examine in the test is less than 1, and the series converges absolutely. But this means that the “vectorized” version of the given series converges as well (absolute convergence implies convergence), thus the matrix formulation of it converges too.

Additionally, from exercise 9 we know that the triangle inequality applies “at infinity” as well, i.e.

$$\begin{aligned} \|e^A\| &\leq \|I\| + \|A\| + \frac{1}{2}\|A^2\| + \frac{1}{3!}\|A^3\| + \dots \leq \sqrt{n} + \|A\| + \frac{1}{2}\|A\|^2 + \frac{1}{3!}\|A\|^3 + \dots = \\ &\sqrt{n} - 1 + 1 + \|A\| + \frac{1}{2}\|A\|^2 + \frac{1}{3!}\|A\|^3 + \dots = \sqrt{n} - 1 + e^{\|A\|} \end{aligned}$$

(b) For $A = \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}$, each power of it is of the form $A^k = \begin{pmatrix} a^k & 0 \\ 0 & b^k \end{pmatrix}$. Thus:

$$e^A = \begin{pmatrix} 1 + a + \frac{1}{2!}a^2 + \frac{1}{3!}a^3 + \dots & 0 \\ 0 & 1 + b + \frac{1}{2!}b^2 + \frac{1}{3!}b^3 + \dots \end{pmatrix} = \begin{pmatrix} e^a & 0 \\ 0 & e^b \end{pmatrix}$$

For $A = \begin{pmatrix} 0 & a \\ 0 & 0 \end{pmatrix}$, all powers of A starting from 2 are zero. Thus:

$$e^A = I + A = \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix}$$

For $A = \begin{pmatrix} 0 & a \\ -a & 0 \end{pmatrix}$ we can observe that A maps the first vector of the standard basis to $-a$ times the second, and that it maps the second vector of the standard basis to a times the first. Therefore powers of A intuitively alternate between, one, mapping e_1 to a multiple of e_2 and vice versa, or two, mapping e_1 to a multiple of e_1 and e_2 to a multiple of e_2 . This leads to the conclusion—which can also be easily proved by induction—that:

$$A^k = \begin{pmatrix} (-1)^{\frac{k}{2}} a^k & 0 \\ 0 & (-1)^{\frac{k}{2}} a^k \end{pmatrix}, \text{ for } k \text{ even}$$

$$A^k = \begin{pmatrix} 0 & (-1)^{\frac{k+1}{2}} a^k \\ (-1)^{\frac{k+1}{2}} a^k & 0 \end{pmatrix}, \text{ for } k \text{ odd}$$

The terms already begin to remind us of the series for sine and cosine. Writing out the series for e^A we have that:

$$e^A = \begin{pmatrix} 1 - \frac{1}{2!}a^2 + \frac{1}{4!}a^4 + \dots & a - \frac{1}{3!}a^3 + \frac{1}{5!}a^5 + \dots \\ -a + \frac{1}{3!}a^3 - \frac{1}{5!}a^5 + \dots & 1 - \frac{1}{2!}a^2 + \frac{1}{4!}a^4 + \dots \end{pmatrix} = \begin{pmatrix} \cos(a) & \sin(a) \\ -\sin(a) & \cos(a) \end{pmatrix}$$

(c) 1. This is not always true. Consider $A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, B = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$. From (b) we have that:

$$e^A = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, e^B = \begin{pmatrix} \cos(-1) & \sin(-1) \\ -\sin(-1) & \cos(-1) \end{pmatrix} \implies e^A e^B = \begin{pmatrix} \cos(-1) - \sin(-1) & \sin(-1) + \cos(-1) \\ -\sin(-1) & \cos(-1) \end{pmatrix}$$

Also $A + B = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}$. Observe here that each power of $A + B$ that is at least 2 equals the zero matrix. Therefore, we see that:

$$e^{A+B} = I + A + B = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$$

, which clearly does not equal $e^A e^B$.

2. If it is true that $AB = BA$, the following identity holds (this can be proved by induction):

$$(A + B)^n = \sum_{k=0}^n \binom{n}{k} A^k B^{n-k}$$

Thus, we have that:

$$\begin{aligned} e^{A+B} &= \sum_{n=0}^{\infty} \frac{(A+B)^n}{n!} = \sum_{n=0}^{\infty} \frac{1}{n!} \sum_{k=0}^n \binom{n}{k} A^k B^{n-k} = \sum_{n=0}^{\infty} \frac{1}{n!} \sum_{k=0}^n \frac{n!}{(n-k)!k!} A^k B^{n-k} \\ &= \sum_{n=0}^{\infty} \sum_{k=0}^n \frac{1}{(n-k)!k!} A^k B^{n-k} \end{aligned}$$

This is the *Cauchy product* of the series e^A and e^B . We will now try to prove Merten's theorem for series of matrices. Merten's theorem states that if one of two series (of numbers) a_n, b_n converges absolutely and the other converges, then their Cauchy product c_n converges to the product of the limits l_1, l_2 of the two series. Let us try to set up the proof in the same way Apostol does for numbers. Below, A_n, B_n indicate convergent series of (square) matrices that converge to A, B respectively, C_n indicates their Cauchy product, and without loss of generality, A_n converges absolutely.

First, define the partial sums:

$$\mathbf{A}_n = \sum_{k=0}^n A_k, \mathbf{B}_n = \sum_{k=0}^n B_k, \mathbf{C}_n = \sum_{k=0}^n C_k$$

As a reminder, $C_k = \sum_{i=0}^k A_k B_{k-i}$. Now we define $D_n = B - \mathbf{B}_n$. Intuitively, D_n measures the difference between the sum of the first n terms of B_n and the limit of this series. We also define $E_n = \sum_{k=0}^n A_k D_{n-k}$. This is harder to interpret intuitively, but we will see its usefulness soon. By definition of the Cauchy product, we now have that:

$$\mathbf{C}_p = \sum_{n=0}^p \sum_{k=0}^n A_k B_{n-k}$$

We would now like to have the same top indices in both sums, as this would make some further processing easier. To do this, define:

$$f_n(k) = \begin{cases} A_k B_{n-k} & n \geq k \\ 0 & n < k \end{cases}$$

Observe that this now allows us to write:

$$\mathbf{C}_p = \sum_{n=0}^p \sum_{k=0}^p f_n(k)$$

, since we are just summing some zero terms whenever $n < k$, which in the first version of \mathbf{C}_p would result in $n - k$ being negative and B_{n-k} being undefined. Now we do the following manipulations:

$$\mathbf{C}_p = \sum_{n=0}^p \sum_{k=0}^p f_n(k) = \sum_{k=0}^p \sum_{n=0}^p f_n(k) = \sum_{k=0}^p \sum_{n=k}^p A_k B_{n-k} = \sum_{k=0}^p A_k \sum_{m=0}^{p-k} B_m$$

, where the second equality comes from rearranging the sum, the third comes from eliminating zero terms and the fourth comes from renaming indices. Proceeding:

$$\mathbf{C}_p = \sum_{k=0}^p A_k \mathbf{B}_{p-k} = \sum_{k=0}^p A_k (B - \mathbf{D}_{p-k}) = \mathbf{A}_p B - E_n$$

, where we used the definition of the partial sums of B_n , the definition of D_n as the difference of the limit from each partial sum and the definition of E_n , whose motivation now is more clear.

Now observe that B is a constant (the limit of the series B_n) and multiplies the partial sum \mathbf{A}_n . This partial sum converges to A . If we can show that E_n converges to 0, we'll thus have shown that \mathbf{C}_p converges to AB , and of course the limit of the partial sum \mathbf{C}_p is the limit of the Cauchy product series.

The sequence D_n converges to the zero matrix (easy to see by its definition). Therefore, the norms of its terms are bounded, so we can find $M > 0$ such that $\|D_n\| \leq M$ for all n . Because A_n converges absolutely, $K = \sum_{n=0}^{\infty} \|A_n\|$ is well defined.

Now given an $\epsilon > 0$, choose $N > 0$ such that for $n > N$, $\|D_n\| < \frac{\epsilon}{2K}$ and at the same time $\sum_{n=N+1}^{\infty} \|A_n\| < \frac{\epsilon}{2M}$. Note we can satisfy the second inequality precisely because for any ϵ we can pick N such that the partial sum up to A_N has a difference less than ϵ from the series' limit, and thus the remaining terms of the sum, $\sum_{n=N+1}^{\infty} A_n$, must sum to less than ϵ . Now for $p > 2N$ we have that:

$$\begin{aligned} \|E_p\| &\leq \sum_{k=0}^N \|A_k D_{p-k}\| + \sum_{k=N+1}^p \|A_k D_{p-k}\| \\ &\leq \sum_{k=0}^N \|A_k\| \cdot \|D_{p-k}\| + \sum_{k=N+1}^p \|A_k\| \cdot \|D_{p-k}\| \leq \frac{\epsilon}{2K} \sum_{k=0}^N \|A_k\| + M \sum_{k=N+1}^p \|A_k\| \end{aligned}$$

For the first inequality we used the triangle inequality, and for the second the well-known inequality regarding matrix norms. For the third inequality, note that for the first term of the sum, $0 \leq k \leq N$, $p > 2N \implies p - k > N$. Thus our upper bound $\frac{\epsilon}{2K}$ applies to the terms of that sum. For the second term, the "global" upper bound for D_n applies. Now we have that $\sum_{k=0}^N \|A_k\| < K$. Lastly, by our previous choices $\sum_{k=N+1}^p \|A_k\| < \sum_{k=N+1}^{\infty} \|A_k\| < \frac{\epsilon}{2M}$. Putting these all together we obtain:

$$\|E_p\| < \frac{\epsilon}{2K} K + M \frac{\epsilon}{2M} = \epsilon$$

Thus we've shown that E_n tends to the zero matrix as n tends to infinity, thus as previously mentioned that the limit of the Cauchy product is indeed AB .

Finally, as previously mentioned, this can be directly applied to e^{A+B} , since both e^A and e^B converge and e^{A+B} is their Cauchy product. Therefore it indeed holds that $e^{A+B} = e^A e^B$ whenever $AB = BA$.

3. Clearly, setting $B = A$ in (2) yields that $AB = BA$ holds for all A and thus we obtain:

$$e^{2A} = e^{A+B} = e^A e^B = e^A e^A = (e^A)^2$$

Exercise 11

Let $\phi : (0, \infty) \rightarrow (0, \infty)$ be a function such that $\lim_{u \rightarrow 0} \phi(u) = 0$.

- (a) Show that the sequence $i \rightarrow a_i$ in \mathbb{R}^n converges to a iff for any $u > 0$, there exists $N > 0$ such that for $n > N$ we have that $\|a_n - a\| < \phi(u)$.
(b) Find an analogous statement for limits of functions.

Solution.

(a) \implies : Suppose first that a_i converges to a . Then, for any $u > 0$, we know that $\phi(u) > 0$. This means that there exists $N > 0$ such that for every $n > N$ it holds that:

$$\|a_i - a\| < \phi(u)$$

\Leftarrow : Suppose now that for every $u > 0$ there exists $N > 0$ such that for $n > N$ it holds that $\|a_n - a\| < \phi(u)$. Pick $\epsilon > 0$. Because $\lim_{u \rightarrow 0} \phi(u) = 0$, there exists $\delta > 0$ such that for every $0 < u < \delta$ it holds that $|\phi(u)| < \epsilon$. For these $\phi(u)$, we know that there exists $N > 0$ such that for every $n > N$ we have that $\|a_n - a\| < \phi(u)$. Since $\phi(u) < \epsilon$, for this arbitrarily selected $\epsilon > 0$ we were able to find $N > 0$ such that $\|a_n - a\| < \epsilon$ for $n > N$.

This means that a_i indeed converges to a .

(b) When considering limits of functions instead of sequences, an analogous statement is:

Let $\phi : (0, \infty) \rightarrow (0, \infty)$ be a function such that $\lim_{u \rightarrow 0} \phi(u) = 0$ and $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Then the limit of f as x approaches x_0 equals L iff for any $\epsilon > 0$ there exists $\delta > 0$ such that $\|x - x_0\| < \delta$ implies $\|f(x) - L\| < \phi(\epsilon)$.

Exercise 12

Let u be a strictly positive function of $\epsilon > 0$, such that $u(\epsilon) \rightarrow 0$ as $\epsilon \rightarrow 0$, and let U be a subset of \mathbb{R}^n . Prove that the following statements are equivalent:

1. A function $f : U \rightarrow \mathbb{R}$ has the limit a at x_0 if $x_0 \in \overline{U}$ and if for all $\epsilon > 0$, there exists $\delta > 0$ such that when $\|x - x_0\| < \delta$ and $x \in U$, then $\|f(x) - a\| < \epsilon$.
2. A function $f : U \rightarrow \mathbb{R}$ has the limit a at x_0 if $x_0 \in \overline{U}$ and if for all $\epsilon > 0$, there exists $\delta > 0$ such that when $\|x - x_0\| < \delta$ and $x \in U$, then $\|f(x) - a\| < u(\epsilon)$.

Solution.

1 \implies 2: Select any $\epsilon > 0$. Then we know that $\epsilon' = u(\epsilon) > 0$. Because f has the limit a at $x_0 \in \overline{U}$, there exists δ such that when $x \in U$, $\|x - x_0\| < \delta$, it holds that $\|f(x) - a\| < \epsilon' = u(\epsilon)$, thereby completing the proof in this direction.

2 \implies 1: Select any $\epsilon > 0$. Because $u(x)$ tends to 0 as x tends to 0, there exists $\delta > 0$ such that for $|x| < \delta$, it holds that $|u(x)| < \epsilon$. Select then any of those x , and call $u(x) = \epsilon'$. Then by our premise (2), there exists $\delta' > 0$ such that when $x \in U$, $\|x - x_0\| < \delta'$, it holds that $\|f(x) - a\| < \epsilon' = u(x) < \epsilon$. But this is precisely statement (1), and thus the proof is completed in this direction as well.

Exercise 13

Prove the converse of Proposition 1.5.17 (i.e. prove that if every convergent sequence in a set $C \subset \mathbb{R}^n$ converges to a point in C , then C is closed).

Solution.

We equivalently need to prove that if every convergent sequence in $C \subset \mathbb{R}^n$ converges to a point in C , $\mathbb{R}^n - C$ is open. Suppose that this is not the case. Then there exists at least one point $x \in \mathbb{R}^n - C$ such that for every $r > 0$ the open ball $B_r(x)$ is not a subset of $\mathbb{R}^n - C$. Consequently, each $B_r(x)$ contains at least one point y_r that is *not* in $\mathbb{R}^n - C$, and therefore is in C . Consider the sequence of y_{r_i} formed by selecting $r_i = \frac{1}{i}, i = 1, 2, \dots$. Each of those points is in C . However, it is clear that the limit of this sequence is x : for any $r > 0$, we have but to set $M = \lceil \frac{1}{r} \rceil$, and then for $i > M$, $y_{r_i} \in B_{r_i}(x), r_i < r \implies \|y_{r_i} - x\| < r_i < r$. Therefore, a sequence of points in C converges to a point outside of it, which is a contradiction. Therefore, $\mathbb{R}^n - C$ is open and C is closed.

Exercise 14

State whether the following limits exist, and prove it.

- (a) $\lim_{(x,y) \rightarrow (1,2)} \frac{x^2}{x+y}$
- (b) $\lim_{(x,y) \rightarrow (0,0)} \frac{\sqrt{|x|}y}{x^2+y^2}$
- (c) $\lim_{(x,y) \rightarrow (0,0)} \frac{\sqrt{|xy|}}{\sqrt{x^2+y^2}}$
- (d) $\lim_{(x,y) \rightarrow (1,2)} x^2 + y^3 - 3$

Solution.

(a) This is a quotient of two continuous functions (polynomials), none of which is zero at $(1, 2)$. Therefore, by the rules of limit calculation, we obtain that the limit exists, and in fact:

$$\lim_{(x,y) \rightarrow (1,2)} \frac{x^2}{x+y} = \frac{1}{3}$$

(b) As we know, for the limit to exist at a point, it has to be the same no matter “how we approach” (recall that the definition features universal quantifiers) the point. Consider, therefore, approaching $(0, 0)$ by moving on the curve $(x, \sqrt{|x|})$. Then the limit would have to be:

$$\lim_{(x,\sqrt{|x|}) \rightarrow (0,0)} \frac{\sqrt{|x|}\sqrt{|x|}}{x^2 + |x|} = \lim_{(x,\sqrt{|x|}) \rightarrow (0,0)} \frac{|x|}{|x|(|x| + 1)} = \lim_{(x,\sqrt{|x|}) \rightarrow (0,0)} \frac{1}{|x| + 1}$$

, which evaluates to 1. However, if we now repeat the exact same procedure for $y = -\sqrt{|x|}$, we see that the limit becomes -1 (only the final sign of the fraction changes).

(c) We apply the same logic as above: consider approaching $(0, 0)$ when $y = x$. Then the limit would evaluate to:

$$\lim_{(x,x) \rightarrow (0,0)} \frac{\sqrt{|x^2|}}{\sqrt{x^2 + x^2}} = \lim_{(x,x) \rightarrow (0,0)} \frac{|x|}{\sqrt{2}|x|} = \frac{1}{\sqrt{2}}$$

However, if we were to approach $(0, 0)$ when $y = 2x$, then this would instead evaluate to $\sqrt{\frac{2}{5}}$. Therefore the limit does not exist.

(d) This is a continuous function (a polynomial), and thus the limit at $(1, 2)$ is simply:

$$\lim_{(x,y) \rightarrow (1,2)} x^2 + y^3 - 3 = 1 + 8 - 3 = 6$$

Exercise 17

Prove the following theorem (1.5.16): Let $i \rightarrow a_i, i \rightarrow b_i$ be two sequences in \mathbb{R}^n and let $i \rightarrow c_i$ be a sequence in \mathbb{R} . Then:

- (1) If $i \rightarrow a_i, i \rightarrow b_i$ converge to a, b respectively, then $i \rightarrow a_i + b_i$ converges to $a + b$.
- (2) if $i \rightarrow a_i$ and $i \rightarrow c_i$ converge to a, c , then $i \rightarrow c_i a_i$ converges to ca .
- (3) If $i \rightarrow a_i, i \rightarrow b_i$ converge to a, b then the sequence $i \rightarrow \langle a_i, b_i \rangle$ converges to $\langle a, b \rangle$.
- (4) If $i \rightarrow a_i r$ is bounded and $i \rightarrow c_i$ converges to 0, then $\lim_{i \rightarrow \infty} c_i a_i = 0$.

Solution.

(1) Select any $\epsilon > 0$. We then know that there exist positive integers N_1, N_2 such that for $n > N_1, \|a_n - a\| < \epsilon$ and for $n > N_2, \|b_n - b\| < \epsilon$. Set $N = \max\{N_1, N_2\}$. Then for $n > N$ both of the inequalities hold and:

$$\|(a_n - a) + (b_n - b)\| \leq \|a_n - a\| + \|b_n - b\| < 2\epsilon$$

The function $\phi(\epsilon) = 2\epsilon$ tends to 0 as $\epsilon > 0$ tends to 0, and therefore by the theorem we proved in exercise 11, $a_n + b_n$ indeed converges to $a + b$.

(2) Select any $\epsilon > 0$. Then there exist N_1, N_2 such that for $n > N_1, \|a_n - a\| < \epsilon$ and for $n > N_2, |c_n - c| < \epsilon$. Set $N = \max\{N_1, N_2\}$. Then, for $n > N$ we have that:

$$\begin{aligned} \|c_n a_n - ca\| &= \|c_n a_n - c_n a + c_n a - ca\| \leq \|c_n a_n - c_n a\| + \|c_n a - ca\| = |c_n| \cdot \|a_n - a\| + \|a\| \cdot |c_n - c| \\ &\leq |c_n| \epsilon + \|a\| \epsilon \end{aligned}$$

By the triangle inequality we also have that:

$$||c_n| - |c|| \leq |c_n - c| < \epsilon \implies |c_n| < |c| + \epsilon$$

Therefore:

$$\|c_n a_n - ca\| \leq (|c| + \epsilon) \epsilon + \|a\| \epsilon$$

Because $|c|, \|a\|$ are non-negative constants, the RHS is a function of ϵ that tends to 0 as ϵ tends to 0. Therefore, by using the same theorem we have that $c_n a_n$ indeed converges to ca .

(3) As previously, select any $\epsilon > 0$. We then know that there exist positive integers N_1, N_2 such that for $n > N_1, \|a_n - a\| < \epsilon$ and for $n > N_2, \|b_n - b\| < \epsilon$. Set $N = \max\{N_1, N_2\}$. Then for $n > N$:

$$\begin{aligned} |\langle a_n, b_n \rangle - \langle a, b \rangle| &= |\langle a_n, b_n \rangle + \langle a_n, b \rangle - \langle a_n, b \rangle - \langle a, b \rangle| = |\langle a_n, b_n - b \rangle + \langle a_n - a, b \rangle| \\ &\leq \|a_n\| \cdot \|b_n - b\| + \|b\| \cdot \|a_n - a\| < \epsilon \|a_n\| + \|b\| \epsilon \end{aligned}$$

, where we used the Cauchy-Schwarz inequality. As above, observe that $\|a_n\| < \|a\| + \epsilon$, which makes the RHS a function of ϵ that tends to 0 as ϵ tends to 0, which guarantees that $\langle a_n, b_n \rangle$ converges to $\langle a, b \rangle$.

(4) Select any $\epsilon > 0$. Then:

$$\|c_n a_n\| \leq |c_n| \cdot \|a_n\|$$

a_i is bounded, which means

$$\|a_i\| < M$$

for some $M \geq 0$. In addition, there exists $N > 0$ such that for $n > N, |c_n - 0| < \epsilon$. But then for $n > N$ our first inequality becomes:

Since M is a non-negative constant, the RHS is again a function of ϵ tending to 0 as ϵ tends to 0, and thus $c_n a_n$ indeed converges to 0 (the zero vector).

Exercise 18

Prove that if a sequence $i \rightarrow a_i$ converges to a , then any subsequence converges to the same limit.

Solution.

A subsequence b_i of a_i is a sequence such that $b_i = a_{f(i)}$, where $f : \mathbb{Z}^+ \rightarrow \mathbb{Z}^+$ and $f(i) < f(j)$ whenever $i < j$ (also often denoted as a_{k_i}). Pick then any $\epsilon > 0$. Then there exists $N > 0$ such that when $i > N, \|a_i - a\| < \epsilon$. Because b_i is an infinite sequence, it must be the case that for some $k, f(k) > N$. If this were not true, then the set of all $f(i)$ would have at most N different values, since $f(i) < f(j)$ whenever $i < j$. This means that for this k , and for all $i > k$, it holds that $f(i) > f(k) > N$. But then:

$$\|b_i - a\| = \|a_{f(i)} - a\| < \epsilon$$

, which completes the proof that b_i indeed converges to a .

Exercise 19

Let $U \subset \text{Mat}(2, 2)$ be the set of matrices A such that $I - A$ is invertible.

(a) Show that U is open, and find a sequence in U that converges to I .

(b) Consider the mapping $f : U \rightarrow \text{Mat}(2, 2)$ given by $f(A) = (A^2 - I)(A - I)^{-1}$. Does $\lim_{A \rightarrow I} f(A)$ exist? If so, what is the limit?

(c) Let $B = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$, and let $V \subset \text{Mat}(2, 2)$ be the set of matrices A such that $A - B$ is invertible.

Again, show that V is open, and that B can be approximated by elements of V .

(d) Consider the mapping $g : V \rightarrow \text{Mat}(2, 2)$ given by

$$g(A) = (A^2 - B^2)(A - B)^{-1}$$

Does $\lim_{A \rightarrow B} g(A)$ exist? If so, what is the limit?

Solution.

(a) Consider any element A of U . We then have that $I - A$ is invertible, and hence $\|(I - A)^{-1}\|$ is not zero. Let then H be any 2×2 matrix such that $\|H\| < \frac{1}{\|(I - A)^{-1}\|}$. Then it holds that:

$$\|H\| \cdot \|(I - A)^{-1}\| < 1$$

, and since $\|AB\| \leq \|A\| \cdot \|B\|$ for any two square matrices, we have that $\|(I - A)^{-1}H\| < 1$. By corollary 1.5.39, we know that $I - H(I - A)^{-1}$ is invertible. Then we can write:

$$(I - (I - A)^{-1}H)^{-1}(I - A)^{-1} = ((I - A)(I - (I - A)^{-1}H))^{-1} = (I - A - H)^{-1} = (I - (A + H))^{-1}$$

, and thus we conclude that $I - (A + H)$ is invertible. But then this means that $A + H$ belongs in U . Thus, for any given A we can find a neighborhood of A (by picking H simply by constraining its norm) such that $A + H$ is in U , thus U is open.

Now consider the sequence $i \rightarrow A_i = \begin{pmatrix} 1 - \frac{1}{i+1} & 0 \\ 0 & 1 - \frac{1}{i+1} \end{pmatrix}$. For any i , we have that $I - A_i = \begin{pmatrix} \frac{1}{i+1} & 0 \\ 0 & \frac{1}{i+1} \end{pmatrix}$. Clearly, for any $i > 0$, this diagonal matrix has non-zero diagonal entries, and thus its determinant is not zero, meaning that it is invertible. Therefore $A_i \in U$ for all i . We now have to show that the sequence of A_i converges to I . We have that:

$$\|A_i - I\| = \left\| \begin{pmatrix} -\frac{1}{i+1} & 0 \\ 0 & -\frac{1}{i+1} \end{pmatrix} \right\| = \sqrt{\frac{2}{(i+1)^2}} = \frac{\sqrt{2}}{i+1}$$

From this we easily conclude that for any $\epsilon > 0$ we can pick a sufficiently large i such that $\|A_i - I\| < \epsilon$. Thus A_i indeed converges to I .

(b) For any $A \in U$, we know that $I - A$, and thus $A - I$, are invertible. We furthermore observe that we can write $(A^2 - I) = (A + I)(A - I)$. Thus:

$$f(A) = (A^2 - I)(A - I)^{-1} = (A + I)(A - I)(A - I)^{-1} = A + I$$

From this it's fairly straightforward to show that the limit of f as A approaches I is $2I$. All we have to do is observe that $\|A + I - 2I\| = \|A - I\|$, and then use the fact that A approaches I to bound this quantity by any given ϵ .

(c) Let A be any element of V . Since $A - B$ is invertible, $\|(A - B)^{-1}\| \neq 0$. Now pick any 2×2 matrix H such that $\|H\| < \frac{1}{\|(A - B)^{-1}\|}$. By the same procedure as in (a), we can obtain that $\|-(A - B)^{-1}H\| < 1$. Now by corollary 1.5.39 we have that $I + (A - B)^{-1}H$ is invertible, and then:

$$(I + (A - B)^{-1}H)^{-1}(A - B)^{-1} = ((A - B)(I + (A - B)^{-1}H))^{-1} = (A - B + H)^{-1}$$

, which means that $(A + H) - B$ is invertible, thus $A + H$ is in V . Thus, we have found a neighborhood of matrices $A + H$ around A such that $A + H \in V$, which means V is open.

Now for approximating B with elements of V , it suffices to consider the sequence:

$$i \rightarrow A_i = \begin{pmatrix} 1 - \frac{1}{i+1} & 0 \\ 0 & -1 - \frac{1}{i+1} \end{pmatrix}$$

For all A_i it is again clear that $A_i - B$ is invertible, as a diagonal matrix with non-zero diagonal entries. Furthermore, as $i \rightarrow \infty$, we will again have that $A - B$ approaches B : the proof is quite similar to (a), based on computing the norm of $A_i - B$ and showing that for any $\epsilon > 0$, we can find i sufficiently large to bound it by ϵ , precisely because the norm is a decreasing function of i .

(d) For the limit to exist, it has to be the case that for any sequence of A_i approaching B , $(A_i^2 - B^2)(A_i - B)^{-1}$ converges to the same matrix. Consider first the sequence of A_i defined in (c), i.e.

$$A_i = \begin{pmatrix} 1 - \frac{1}{i+1} & 0 \\ 0 & -1 - \frac{1}{i+1} \end{pmatrix}$$

We have that:

$$A_i^2 = \begin{pmatrix} 1 - \frac{2}{i+1} + \frac{1}{(i+1)^2} & 0 \\ 0 & 1 + \frac{2}{i+1} + \frac{1}{(i+1)^2} \end{pmatrix} \Rightarrow A_i^2 - B^2 = \begin{pmatrix} -\frac{2}{i+1} + \frac{1}{(i+1)^2} & 0 \\ 0 & \frac{2}{i+1} + \frac{1}{(i+1)^2} \end{pmatrix}$$

Also:

$$A_i - B = \begin{pmatrix} -\frac{1}{i+1} & 0 \\ 0 & -\frac{1}{i+1} \end{pmatrix} \Rightarrow (A_i - B)^{-1} = \begin{pmatrix} -(i+1) & 0 \\ 0 & -(i+1) \end{pmatrix}$$

Thus:

$$(A_i^2 - B^2)(A_i - B)^{-1} = \begin{pmatrix} 2 - \frac{1}{i+1} & 0 \\ 0 & -2 - \frac{1}{i+1} \end{pmatrix}$$

Fairly straightforwardly, this sequence of $g(A_i)$ converges to $\begin{pmatrix} 2 & 0 \\ 0 & -2 \end{pmatrix}$ as $i \rightarrow \infty$.

Now consider the sequence $i \rightarrow C_i = \begin{pmatrix} 1 - \frac{1}{i+1} & \frac{1}{i+1} \\ 0 & -1 - \frac{1}{i+1} \end{pmatrix}$. This sequence can easily be shown to converge to B in the same way as we approached (a). We have that:

$$C_i^2 = \begin{pmatrix} 1 - \frac{2}{i+1} + \frac{1}{(i+1)^2} & \frac{1}{i+1} - \frac{1}{(i+1)^2} - \frac{1}{i+1} - \frac{1}{(i+1)^2} \\ 0 & 1 + \frac{2}{i+1} + \frac{1}{(i+1)^2} \end{pmatrix} \Rightarrow C_i^2 - B^2 = \begin{pmatrix} -\frac{2}{i+1} + \frac{1}{(i+1)^2} & -\frac{2}{i+1} + \frac{1}{(i+1)^2} \\ 0 & \frac{2}{i+1} + \frac{1}{(i+1)^2} \end{pmatrix}$$

By inspection, we have that $(C_i - B)^{-1} = \begin{pmatrix} -(i+1) & i+1 \\ 0 & -(i+1) \end{pmatrix}$. Thus:

$$(C_i^2 - B^2)(C_i - B)^{-1} = \begin{pmatrix} 2 - \frac{1}{i+1} & -2 + \frac{1}{i+1} + \frac{2}{i+1} \\ 0 & -2 - \frac{1}{i+1} \end{pmatrix} = \begin{pmatrix} 2 - \frac{1}{i+1} & -2 + \frac{3}{i+1} \\ 0 & -2 - \frac{1}{i+1} \end{pmatrix}$$

Now, however, we observe that if we subtract the matrix $\begin{pmatrix} 2 & -2 \\ 0 & -2 \end{pmatrix}$ from $(C_i^2 - B^2)(C_i - B)^{-1}$, we end up with a matrix whose norm can easily be shown to tend to zero as $i \rightarrow \infty$, from which we can also obtain that this sequence of $g(C_i)$ tends to $\begin{pmatrix} 2 & -2 \\ 0 & -2 \end{pmatrix}$ as C_i tends to B , which shows that the limit $\lim_{A \rightarrow B} g(A)$ does not exist.

Exercise 21

For the following functions, can you choose a value for f at $(0,0)$ to make the function continuous at the origin?

(a) $f(x, y) = \frac{1}{x^2 + y^2 + 1}$

(b) $f(x, y) = \frac{\sqrt{x^2 + y^2}}{|x| + |y|^{\frac{1}{3}}}$

(c) $f(x, y) = (x^2 + y^2) \log(x^2 + 2y^2)$

(d) $f(x, y) = (x^2 + y^2) \log(|x + y|)$

Solution.

(a) This function is a rational function with the denominator being always positive. As such, its limit as (x, y) approaches the origin is simply $\frac{1}{0+0+1} = 1$, and therefore setting $f(0, 0) = 1$ makes it trivially continuous at $(0, 0)$.

(b) Consider approaching $(0, 0)$ with a sequence of (x, y) such that $y = x$. Evaluated on this sequence, $f(x, y)$ equals:

$$f(x, y) = \frac{\sqrt{x^2 + x^2}}{|x| + |x|^{\frac{1}{3}}} = \frac{\sqrt{2}|x|}{|x| + |x|^{\frac{1}{3}}} = \frac{\sqrt{2}|x|}{|x|(1 + |x|^{-\frac{2}{3}})} = \frac{\sqrt{2}}{1 + \frac{1}{|x|^{\frac{2}{3}}}}$$

Now, from calculus we know that this fraction approaches zero as x approaches 0. Now consider instead approaching $(0, 0)$ with a sequence of (x, y) such that $y = x^3$. Then for this sequence, $f(x, y)$ equals:

$$f(x, y) = \frac{\sqrt{x^2 + x^6}}{|x| + |x^3|^{\frac{1}{3}}} = \frac{\sqrt{x^2(1 + x^4)}}{|x| + |x|} = \frac{|x|\sqrt{1 + x^4}}{2|x|} = \frac{\sqrt{1 + x^4}}{2}$$

, where we used the fact that x is not zero. Now this final expression is continuous at 0, and thus its limit as x approaches 0 is $\frac{1}{2}$. But then this means that no matter what value we assign to $f(0, 0)$, the limit at $(0, 0)$ can never exist, since it would have to be the same on all possible sequences converging to $(0, 0)$.

(c) For (x, y) such that $0 < x^2 + 2y^2 < 1$, we have that $\log(x^2 + 2y^2) < 0$. Additionally, we have that $x^2 + y^2 > 0$, thus $f(x, y) < 0$. It is furthermore true that $0 < x^2 + y^2 \leq x^2 + 2y^2 \implies \log(x^2 + y^2) \leq \log(x^2 + 2y^2)$. This also implies that $(x^2 + y^2) \log(x^2 + y^2) \leq (x^2 + y^2) \log(x^2 + 2y^2) = f(x, y)$.

Now observe that from calculus (L'Hopital's rule) we know that $\lim_{z \rightarrow 0} z \log(z) = 0$, and here we have that as $(x, y) \rightarrow (0, 0)$, $x^2 + y^2 \rightarrow 0$. Thus the limit of the LHS above at $(0, 0)$ is 0, and then by the sandwich theorem, the limit of $f(x, y)$ at $(0, 0)$ is also 0. Thus it suffices to define $f(0, 0) = 0$ to make the function continuous at the origin.

(d) The function is not defined for $y = -x$. However, we can examine its behavior when $y = -x + e^{-\frac{1}{x^2}}$, $x \neq 0$. As $x \rightarrow 0$, the exponent here tends to negative infinity, which makes $e^{-\frac{1}{x^2}}$ tend to 0, and thus y tends to 0. Thus this is indeed a sequence tending to $(0, 0)$. We then have that:

$$\begin{aligned} f(x, y) &= (x^2 + (-x + e^{-\frac{1}{x^2}})^2) \log(|x - x + e^{-\frac{1}{x^2}}|) = (x^2 + x^2 + e^{-\frac{2}{x^2}} - 2xe^{-\frac{1}{x^2}}) \left(-\frac{1}{x^2}\right) \\ &= -2 - \frac{e^{-\frac{2}{x^2}}}{x^2} + \frac{2e^{-\frac{1}{x^2}}}{x} \end{aligned}$$

The last two terms can again via L'Hopital's rule and calculus procedures be shown to converge to 0 as x approaches 0. Thus, $f(x, y)$ approaches -2 on this sequence of (x, y) . If, however, we set instead $y = x$, $x > 0$ and approach $(0, 0)$, we end up with a case very similar to (c) where we used L'Hopital's rule for $z \log(z)$, and we conclude that $f(x, y)$ approaches 0. Thus the limit does not exist and no value at $(0, 0)$ can make f continuous at $(0, 0)$.

Exercise - Unlisted; Arose from a discussion of the Bolzano - Weierstrass theorem

Prove the following generalization of the nested interval theorem: If B_0, B_1, \dots are bounded, non-empty boxes in \mathbb{R}^n such that $B_i = [a_{i1}, b_{i1}] \times [a_{i2}, b_{i2}] \times \dots \times [a_{in}, b_{in}]$ with $B_0 \supset B_1 \supset B_2 \supset \dots$, then their intersection $\cap_{i=1}^{\infty} (B_i) \neq \emptyset$. Furthermore, if each “box side” $s_{ij} = b_{ij} - a_{ij}$ for $j = 1, 2, \dots, n$ is such that $\lim_{i \rightarrow \infty} s_{ij} = 0$, then the intersection contains precisely one point.

Solution.

We examine each coordinate separately. Consider the j -th coordinate. For every box B_i , it is the case that $a_{ij} \leq b_{ij}$. It is furthermore the case that $B_{i+1} \subset B_i$ implies $a_{ij} \leq a_{i+1,j} \leq b_{i+1,j} \leq b_{ij}$. We then have that the intervals $[a_{ij}, b_{ij}]$ are non-empty and bounded (since the boxes themselves are non-empty and bounded), and they also form a nested sequence $[a_{i+1,j}, b_{i+1,j}] \subset [a_{ij}, b_{ij}]$. Therefore, by the nested interval theorem, the intersection $\cap_{i=1}^{\infty} [a_{ij}, b_{ij}]$ is not empty. This means that there exists at least one $x_j \in [a_{ij}, b_{ij}]$ for all i .

Because everything stated above holds for all coordinates, we have at least one such x_j for all j . Thus, if we take $x = (x_1, \dots, x_n)$, we observe that $x \in B_i$ for all i , which means that the intersection of the boxes is not empty.

Now if each box side $s_{ij}, j = 1, 2, \dots, n$ is such that its length $b_{ij} - a_{ij}$ goes to zero as i goes to infinity, the nested interval theorem in \mathbb{R} dictates that the intersection $\cap_{i=1}^{\infty} [a_{ij}, b_{ij}]$ contains *precisely one* point x_j . Once again, this holds for all j . Then, if we take $x = (x_1, \dots, x_n)$ we have that this x belongs in the intersection of all boxes. If this intersection contained at least one other point $y = (y_1, \dots, y_n)$, then for some k it has to be the case that $y_k \neq x_k$. But then $y_k \in [a_{ik}, b_{ik}]$ for all i , meaning that the corresponding interval intersection $\cap_{i=1}^{\infty} [a_{ik}, b_{ik}]$ contains at least two points, x_k, y_k , which contradicts the nested interval theorem. Therefore, the box intersection contains precisely one point.

1.6 Five big theorems

Exercise 1

Show that a set is bounded if it is contained in a ball centered anywhere; it does not have to be centered at the origin.

Solution.

Suppose that a set $S \subset \mathbb{R}^n$ is contained in a ball $B_r(x_0)$ centered at x_0 . Consider any point x inside this ball. It holds that $\|x - x_0\| < r$. We observe then that $\|x\| = \|x - x_0 + x_0\| \leq \|x - x_0\| + \|x_0\| \leq \|x_0\| + r$. More specifically, for all elements of S the above inequality holds, meaning that they are contained in the ball $B_{r+\|x_0\|}(0)$, which means precisely that S is a bounded set.

Exercise 2

Let $A \subset \mathbb{R}^n$ be a subset that is not compact. Show that there exists a continuous unbounded function on A .

Solution.

Because A is not compact, it must either be unbounded or not closed. Suppose first that A is unbounded. Let $f : A \rightarrow \mathbb{R}$ be the function $f(x) = \|x\|$. We know that this is a continuous function everywhere, thus it is also continuous in A . However, because A is unbounded, for every $R > 0$ there exists $x \in A$ such that $\|x\| > R$. But this means precisely that $|f(x)| > R$, which means that f is unbounded.

Now suppose that A is bounded but is not closed (the only case we have not covered). Because A is not closed, there must exist at least one sequence $x_i \in A$ that converges to a point $y \notin A$. Let then $f : A \rightarrow \mathbb{R}$ be the function $f(x) = \frac{1}{\|x - y\|}$. Because both the nominator and the denominator are continuous functions, f is also continuous on all points for which $\|x - y\| \neq 0$, that is, on all $x \neq y$. More specifically, because $y \notin A$, f is continuous on A . We now need to show that f is unbounded. Consider thus any $M > 0$ and set $\epsilon = \frac{1}{M} > 0$. Because $x_i \rightarrow y$, for this ϵ there must exist $N > 0$ such that for all $i > N$, $\|x_i - y\| < \epsilon$. Consequently, $\frac{1}{\epsilon} < \frac{1}{\|x_i - y\|} \implies M < f(x_i)$. Since these x_i belong in A , we've shown that for any $M > 0$ we can find $x \in A$ such that $|f(x)| > M$, which means that f is unbounded on A .

Exercise 7

Show that if f is differentiable on a neighborhood of $[a, b]$ and we have $f'(a) < m < f'(b)$ then there exists $c \in (a, b)$ such that $f'(c) = m$. Now the second part is also

Solution.

We consider the function $g : [a, b] \rightarrow \mathbb{R}, g(x) = f(x) - (m(x - a) + f(a))$, which, geometrically, is a kind of signed distance between the graph of f from the line passing through $(a, f(a))$ with slope m . Firstly, observe that this is a continuous and differentiable function (since f is differentiable), and that $[a, b]$ is compact.

We now recall that any continuous function on a compact set has a global maximum and a global minimum. Let then c be the global minimum of g . We know that if $c \in (a, b)$, then it must be the case that $g'(c) = 0$. It is the case that $g'(x) = f'(x) - m$, which implies $f'(c) = m$, which is what we require. Now suppose that $c = a$. This means that it is more specifically true that $g(a + \epsilon) \geq g(a)$ for every $\epsilon > 0, \epsilon \leq b - a$. Thus:

$$\begin{aligned} g(a + \epsilon) \geq g(a) &\implies f(a + \epsilon) - (m(a + \epsilon - a) + f(a)) \geq f(a) - (m(a - a) + f(a)) \\ &\implies f(a + \epsilon) - f(a) - m\epsilon \geq 0 \implies \frac{f(a + \epsilon) - f(a)}{\epsilon} \geq m \end{aligned}$$

More specifically, because f is differentiable at a , the limit of the left side as $\epsilon \rightarrow 0$ is well defined and equals $f'(a)$, and at the limit the inequality still holds, which would mean $f'(a) \geq m$ which contradicts the hypothesis. Suppose now that $c = b$. Then by the same reasoning, $g(b - \epsilon) \geq g(b)$ for $\epsilon > 0, \epsilon \leq b - a$. Thus:

$$\begin{aligned} g(b - \epsilon) \geq g(b) &\implies f(b - \epsilon) - (m(b - \epsilon - a) + f(a)) \geq f(b) - (m(b - a) + f(a)) \\ &\implies f(b - \epsilon) - f(b) + m\epsilon \geq 0 \implies \frac{f(b - \epsilon) - f(b)}{-\epsilon} \leq m \end{aligned}$$

Again, because f is differentiable at b , and by substituting $h = -\epsilon$ we can observe that the LHS equals the quantity whose limit as $h \rightarrow 0$ is $f'(b)$, meaning that it has to hold that $f'(b) \leq m$ which again contradicts the hypothesis.

Therefore, $c \in (a, b)$ which directly implies that $f'(c) = m$ as requested.

Exercise 9

If $a, b \in \mathbb{C}, a, b \neq 0$ and $j \geq 1$, find $p_0 > 0$ such that if $0 < p < p_0$ and $u = p(\cos(\theta) + i\sin(\theta))$, then there exist j values of θ such that $a + bu^j$ is between 0 and a .

Solution.

This exercise will be easier to solve if we think about complex numbers using exponential notation, that is, using the fact that for non-zero, $z = \rho e^{i\theta}, \rho > 0, \theta \in [0, 2\pi)$. Using this notation we have:

$$a = \rho_a e^{i\theta_a}, b = \rho_b e^{i\theta_b}, u = p e^{i\theta}$$

We want $a + bu^j$ to be on the line segment between 0 and a . This means that $a + bu^j = \lambda a$ for some $\lambda \in (0, 1)$. Consequently:

$$bu^j = a(\lambda - 1) \implies \rho_b e^{i\theta_b} (p e^{i\theta})^j = (\lambda - 1) \rho_a e^{i\theta_a} \implies (p^j \rho_b) e^{i(j\theta + \theta_b)} = \rho_a (1 - \lambda) (-1) e^{i\theta_a}$$

Notice that we wrote $\lambda - 1$ as $(-1)(1 - \lambda)$ since we want the “radius” part of the exponential notation to be positive. Recall also that $e^{i\pi} = -1$, and that for two complex numbers to be equal, it must hold that when written in exponential notation their radii are equal and their angles differ by a multiple of 2π . Therefore we have first that:

$$p^j \rho_b = \rho_a (1 - \lambda) \implies p^j = \frac{\rho_a (1 - \lambda)}{\rho_b} \implies p = \left(\frac{\rho_a}{\rho_b} (1 - \lambda) \right)^{\frac{1}{j}}$$

Notice that this is a decreasing function of $\lambda \in (0, 1)$, and that for $\lambda = 0$ we have that $p = \left(\frac{\rho_a}{\rho_b} \right)^{\frac{1}{j}}$, which we can see constitutes an upper bound for p in order for the above radii to be able to be equal. Thus we

propose $p_0 = \left(\frac{\rho_a}{\rho_b}\right)^{\frac{1}{j}}$ and we must now check whether there exist j values of θ making equality possible. By equating the angles of the above complex numbers we obtain:

$$j\theta + \theta_b = \pi + \theta_a + 2k\pi \implies \theta = \frac{\theta_a - \theta_b + 2k\pi + \pi}{j} = \frac{\theta_a - \theta_b + \pi}{j} + \frac{2k\pi}{j}$$

, with k integer. This represents a unique angle as long as it is in $[0, 2\pi)$, which means that:

$$0 \leq \frac{\theta_a - \theta_b + \pi}{j} + \frac{2k\pi}{j} < 2\pi \implies 2k\pi < 2\pi j - \theta_a + \theta_b - \pi \implies k < j + \frac{\theta_b - \theta_a - \pi}{2\pi}$$

, and also that

$$\theta_b - \theta_a - \pi \leq 2k\pi \implies \frac{\theta_b - \theta_a - \pi}{2\pi} \leq k$$

Now observe that $-\frac{3}{2} < \frac{\theta_b - \theta_a - \pi}{2\pi} < \frac{1}{2}$, and thus that if this quantity x is positive, k must be at least one and at most $j + x < j + \frac{1}{2}$, thus the possible values are $0, 1, \dots, j$, which are j in total. If this quantity is in $(-1, 0]$ then k must be at least 0 and has to be at most $j + x > j - 1$, thus the possible values are $0, 1, \dots, j-1$, again j in total. Lastly, if it is in $(-\frac{3}{2}, -1]$ then k has to be at least -1 and at most $j + x > j - 2$ which means that the possible values are $-1, 0, \dots, j-2$, again j in total.

In all cases, there do indeed exist j unique choices of θ as long as $p < \left(\frac{\rho_a}{\rho_b}\right)^{\frac{1}{j}}$.

1.7 Derivatives in several variables as linear transformations

Exercise 4

Example 1.7.2 may lead you to expect that if $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable at a , then $f(a+h) - f(a) - f'(a)h$ will be comparable to h^2 . It is not true that once you get rid of the linear term you always have a term that includes h^2 . Using the definitions, check whether the following functions are differentiable at 0.

a. $f(x) = |x|^{3/2}$

b. $f(x) = \begin{cases} x \log|x|, & x \neq 0 \\ 0, & x = 0 \end{cases}$

c. $f(x) = \begin{cases} \frac{x}{\log|x|}, & x \neq 0 \\ 0, & x = 0 \end{cases}$

In each case, if f is differentiable at 0, is $f(0+h) - f(0) - f'(0)h$ comparable to h^2 ?

Solution.

a. Suppose first that $h \rightarrow 0^+$. Then $f(0+h) = h^{3/2}$, and $f(0) = 0$. We have that:

$$\begin{aligned} \lim_{h \rightarrow 0^+} \frac{1}{h} (f(0+h) - f(0) - mh) = 0 &\implies \lim_{h \rightarrow 0^+} \frac{1}{h} (h^{3/2} - mh) = 0 \implies \lim_{h \rightarrow 0^+} (h^{1/2} - m) = 0 \\ &\implies m = 0 \end{aligned}$$

Therefore, if the derivative at 0 exists, it has to be 0. Let's see what happens when $h \rightarrow 0^-$:

$$\lim_{h \rightarrow 0^-} \frac{1}{h} (f(0-h) - f(0) - mh) = 0 \implies \lim_{h \rightarrow 0^-} \frac{1}{h} ((-h)^{3/2} - mh) = 0$$

Perform a change of variables setting $u = -h$, which implies $u \rightarrow 0^+$. Then the above becomes:

$$\lim_{u \rightarrow 0^+} \frac{1}{-u} (u^{3/2} + mu) = 0 \implies \lim_{u \rightarrow 0^+} (-u^{1/2} - m) = 0 \implies m = 0$$

We conclude that the derivative of f at 0 is indeed 0, and we observe that here we didn't have a square term, but a square root one instead.

(b) First let's examine what happens as $h \rightarrow 0^+$. Then $f(0+h) = h \log h$, and:

$$\lim_{h \rightarrow 0^+} \frac{1}{h}(f(0+h) - f(0) - mh) = 0 \implies \lim_{h \rightarrow 0^+} \frac{1}{h}(h \log h - 0 - mh) = 0 \implies \lim_{h \rightarrow 0^+} (\log h - m) = 0$$

We observe that as h approaches 0, the logarithm function tends to negative infinity. Therefore, there can be no m that satisfies this equation, and the derivative of f at 0 does not exist. Here the problem is that the remaining term does not converge at all.

(c) First let $h \rightarrow 0^+$. Then $f(0+h) = \frac{h}{\log h}$. Then:

$$\lim_{h \rightarrow 0^+} \frac{1}{h}(f(0+h) - f(0) - mh) = 0 \implies \lim_{h \rightarrow 0^+} \frac{1}{h}\left(\frac{h}{\log h} - 0 - mh\right) = 0 \implies \lim_{h \rightarrow 0^+} \left(\frac{1}{\log h} - m\right) = 0$$

Now, as h approaches 0, the logarithm tends to negative infinity, and thus the fraction tends to 0 (note: we assume that f is not defined at $x = 1$). Thus it must be the case that $m = 0$. If we now let $h \rightarrow 0^-$, then by following the same procedure as above we will end up with a term that exhibits the same behavior, which again leads us to conclude that $m = 0$, meaning that f is indeed differentiable at 0 with $f'(0) = 0$.

Exercise 6

Calculate the partial derivatives $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}$ for the \mathbb{R}^m -valued functions

a. $f(x, y) = \begin{pmatrix} \cos x \\ x^2 y + y^2 \\ \sin(x^2 - y) \end{pmatrix}$

b. $f(x, y) = \begin{pmatrix} \sqrt{x^2 + y^2} \\ xy \\ \sin^2(xy) \end{pmatrix}$

Solution.

a. We have that $\frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \frac{\partial f_3}{\partial x} \end{pmatrix}$, with:

$$\frac{\partial f_1}{\partial x} = \frac{\partial(\cos x)}{\partial x} = -\sin x, \quad \frac{\partial f_2}{\partial x} = \frac{\partial(x^2 y + y^2)}{\partial x} = 2xy, \quad \frac{\partial f_3}{\partial x} = \frac{\partial \sin(x^2 - y)}{\partial x} = 2x \cos(x^2 - y)$$

Similarly, $\frac{\partial f}{\partial y} = \begin{pmatrix} \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial y} \\ \frac{\partial f_3}{\partial y} \end{pmatrix}$, with:

$$\frac{\partial f_1}{\partial y} = \frac{\partial(\cos x)}{\partial y} = 0, \quad \frac{\partial f_2}{\partial y} = \frac{\partial(x^2 y + y^2)}{\partial y} = x^2 + 2y, \quad \frac{\partial f_3}{\partial y} = \frac{\partial \sin(x^2 - y)}{\partial y} = -\cos(x^2 - y)$$

b. Again, $\frac{\partial f}{\partial x} = \begin{pmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \frac{\partial f_3}{\partial x} \end{pmatrix}$, with:

$$\frac{\partial f_1}{\partial x} = \frac{\partial(\sqrt{x^2 + y^2})}{\partial x} = x(x^2 + y^2)^{-1/2}, \quad \frac{\partial f_2}{\partial x} = \frac{\partial(xy)}{\partial x} = y, \quad \frac{\partial f_3}{\partial x} = \frac{\partial \sin^2(xy)}{\partial x} = 2y \sin(xy) \cos(xy)$$

For the partial of f_1 , note that for $(x, y) = (0, 0)$ this cannot be defined due to the non-differentiability of the square root function at zero.

Similarly, $\frac{\partial f}{\partial y} = \begin{pmatrix} \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial y} \\ \frac{\partial f_3}{\partial y} \end{pmatrix}$, with:

$$\frac{\partial f_1}{\partial y} = \frac{\partial(\sqrt{x^2 + y^2})}{\partial y} = y(x^2 + y^2)^{-1/2}, \quad \frac{\partial f_2}{\partial y} = \frac{\partial(xy)}{\partial y} = x, \quad \frac{\partial f_3}{\partial y} = \frac{\partial \sin^2(xy)}{\partial y} = 2x \sin(xy) \cos(xy)$$

, where, as above, we observe that the partial of f_1 at $(0, 0)$ cannot be defined due to the non-differentiability of the square root function at zero.

Exercise 7

Write the answers to exercise 1.7.6 in the form of the Jacobian matrix.

Solution.

a. The corresponding Jacobian is:

$$J_f(x, y) = \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \\ \frac{\partial f_3}{\partial x} & \frac{\partial f_3}{\partial y} \end{pmatrix} = \begin{pmatrix} -\sin x & 0 \\ 2xy & x^2 + 2y \\ 2x \cos(x^2 - y) & -\cos(x^2 - y) \end{pmatrix}$$

b. Whenever $(x, y) \neq (0, 0)$, the Jacobian is:

$$J_f(x, y) = \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} \\ \frac{\partial f_3}{\partial x} & \frac{\partial f_3}{\partial y} \end{pmatrix} = \begin{pmatrix} x(x^2 + y^2)^{-1/2} & y(x^2 + y^2)^{-1/2} \\ y & x \\ 2y \sin(xy) \cos(xy) & 2x \sin(xy) \cos(xy) \end{pmatrix}$$

Exercise 10

Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ be a function.

a. Prove that if f is affine, then for any $a, v \in \mathbb{R}^2$,

$$f(a_1 + v_1, a_2 + v_2) = f(a_1, a_2) + D_f(a)v$$

b. Prove that if f is not affine, this is not true.

Solution.

a. By the definition of affine function, it must hold that $g(x, y) = f(x, y) - f(0, 0)$ is a linear function. Consequently, g corresponds to matrix multiplication with some matrix M . Pick any $a \in \mathbb{R}^2$. Then let us examine whether there exists a linear transformation L such that:

$$\lim_{h \rightarrow 0} \frac{1}{\|h\|} (f(a_1 + h_1, a_2 + h_2) - f(a_1, a_2) - Lh) = 0$$

Because $f(a_1, a_2) = g(a_1, a_2) + f(0, 0)$, $f(a_1 + h_1, a_2 + h_2) = g(a_1 + h_1, a_2 + h_2) + f(0, 0)$, and because g is linear, we have that their difference equals $g(h_1, h_2) = M \begin{pmatrix} h_1 \\ h_2 \end{pmatrix}$. Observe then that if we simply set L to be the linear transformation whose matrix wrt. the standard basis equals M , the limit becomes zero. We know then that f is differentiable at a , and that $D_f(a) = L = M$ for *any* a . More specifically, this means that $D_f(a) = D_f(a')$ for any two $a, a' \in \mathbb{R}^2$. Therefore, using the linearity of g we conclude that for any two $a, v \in \mathbb{R}^2$:

$$\begin{aligned} f(a_1 + v_1, a_2 + v_2) &= g(a_1 + v_1, a_2 + v_2) + f(0, 0) = g(a_1 + v_1, a_2 + v_2) + (f(a_1, a_2) - g(a_1, a_2)) \\ &= g(v_1, v_2) + f(a_1, a_2) = D_f(a)v + f(a_1, a_2) \end{aligned}$$

b. Suppose $f(x, y) = (x^2, 0)$, which is clearly not an affine function. Select $a = (0, 0)$, in which case $D_f(a) = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$ (as can easily be obtained by computing partials at 0). Then for any $v \in \mathbb{R}^2$ it would have to hold that:

$$f(v_1, v_2) = f(0, 0) + D_f(a)v = (0, 0)$$

, which is obviously not true for e.g. $v = (1, 0)$.

Exercise 13

Show that if $f(x) = |x|$, then for any number m ,

$$\lim_{h \rightarrow 0} (f(0+h) - f(0) - mh) = 0$$

, but

$$\lim_{h \rightarrow 0} \frac{1}{h} (f(0+h) - f(0) - mh) = 0$$

is never true: there is no number m such that mh is a “good approximation” to $f(h) - f(0)$ in the sense of Definition 1.7.5.

Solution.

We know that f is continuous, and thus is more specifically continuous at $x = 0$. The same is true for the function mh . Therefore, for any number m :

$$\lim_{h \rightarrow 0} (f(0+h) - f(0) - mh) = f(0+0) - f(0) - m \cdot 0 = 0$$

Let us examine the second limit for $h \rightarrow 0^+$, in which case $f(0+h) = h$:

$$\lim_{h \rightarrow 0^+} \frac{1}{h} (f(0+h) - f(0) - mh) = \lim_{h \rightarrow 0^+} \frac{1}{h} (h - 0 - mh) = \lim_{h \rightarrow 0^+} (1 - m) = 1 - m$$

Now, for $h \rightarrow 0^-$ we have that $f(0+h) = -h$ and:

$$\lim_{h \rightarrow 0^-} \frac{1}{h} (f(0+h) - f(0) - mh) = \lim_{h \rightarrow 0^+} \frac{1}{h} (-h - 0 - mh) = \lim_{h \rightarrow 0^+} (-1 - m) = -1 - m$$

For the limit to exist for some m , it would have to be the case that for that m , $1 - m = -1 - m \implies 2 = 0$, which is of course impossible. Therefore, there is no m for which mh is a “good approximation” to $|x|$ near 0. Intuitively, this makes sense: the “corner” formed by the graph of this function means that a linear function that approximates it very well “on one side” will approximate it very badly on the other.

Exercise 15

- Define what it means for a mapping $F : \text{Mat}(n, m) \rightarrow \text{Mat}(k, l)$ to be differentiable at a point $A \in \text{Mat}(n, m)$.
- Consider the function $F : \text{Mat}(n, m) \rightarrow \text{Mat}(n, n)$ given by $F(A) = AA^T$. Show that F is differentiable, and compute the derivative $[DF(A)]$.

Solution.

- We will apply the “best linear approximation” definition of the derivative. Namely, a mapping such as the one described here is differentiable at a point $A \in \text{Mat}(n, m)$ (that is, A is an $n \times m$ matrix) if and only if there exists a linear transformation $L : \text{Mat}(n, m) \rightarrow \text{Mat}(k, l)$ such that:

$$\lim_{H \rightarrow 0} \frac{1}{\|H\|} ((F(A+H) - F(A)) - L(H)) = \mathbf{0}$$

Here we should observe that $\dim \text{Mat}(n, m) = n \cdot m$, $\dim \text{Mat}(k, l) = k \cdot l$, which means that, with respect to the standard bases of these vector space, L is isomorphic to a matrix $M(L)$ of dimensions $(k \cdot l) \times (n \cdot m)$. Furthermore the “zero” which the limit equals is the zero $k \times l$ matrix.

- This reminds us of the squaring function for real numbers. As such, we guess that $[DF(A)](H) = AH^T + HA^T$. We know have that:

$$\begin{aligned} F(A+H) - F(A) - [DF(A)](H) &= (A+H)(A+H)^T - AA^T - AH^T + HA^T \\ &= (A+H)(A^T + H^T) - AA^T - AH^T - HA^T = AA^T + AH^T + HA^T + HH^T - AA^T - AH^T - HA^T = HH^T \end{aligned}$$

Now:

$$\frac{\|HH^T\|}{\|H\|} \leq \frac{\|H\|^2}{\|H\|} = \|H\|$$

, which means that as $\|H\| \rightarrow 0$, the LHS here tends to zero as well. But then this means that $\frac{1}{\|H\|} \cdot \|(F(A+H) - F(A)) - [DF(A)](H)\|$ tends to zero, and thus $\frac{1}{\|H\|} \cdot ((F(A+H) - F(A)) - L(H))$ tends to zero too. This means that $[DF(A)]$ satisfies the definition of the derivative of F , and, since it is a linear function, it is indeed the unique linear function which satisfies this definition (AKA, it is *the* derivative of F).

Exercise 18

Let U be an open subset of \mathbb{R}^n and let $f : U \rightarrow \mathbb{R}$ be differentiable at $a \in U$. Show that if v is a *unit* vector making an angle θ with the gradient $\nabla f(a)$, then

$$[DF(a)]v = \|\nabla f(a)\| \cos \theta$$

Why does this justify saying that $\nabla f(a)$ points in the direction in which f increases the fastest, and that $\|\nabla f(a)\|$ is this fastest rate of increase?

Solution.

We know that given any vector v , we can express the directional derivative of f along v as:

$$[Df(a)](v)$$

Because $f : U \rightarrow \mathbb{R}$, it holds that $[Df(a)] \in \mathbb{R}^{1 \times n}$, and $\nabla f(a) = [Df(a)]^T$. Now (e.g. by the Riesz representation theorem) we can see that:

$$[Df(a)](v) = \langle \nabla f(a), v \rangle = \|\nabla f(a)\| \cdot \|v\| \cdot \cos \theta$$

, and if we assume v to be a unit vector, this simplifies to the desired expression. Now observe that $\cos \theta$ is in the range $[-1, 1]$, and becomes 1 or -1 precisely when v is colinear with $\nabla f(a)$. This is also precisely when the magnitude of this vector becomes highest, and as such the gradient can be thought of as the direction of fastest change in the function. Clearly, this rate of the increase is the magnitude of $[Df(a)](v)$, which equals the magnitude of the gradient when $\cos \theta = 1$.

1.8 Rules for computing derivatives

Exercise 5

The following “proof” of part 5 of Theorem 1.8.1 (derivative of the product of $f : U \rightarrow \mathbb{R}, g : U \rightarrow \mathbb{R}^m$) is correct as far as it goes, but it is not a complete proof. Why not?

“**Proof**”: By part 3 of Theorem 1.8.1, we may assume that $m = 1$ (i.e. that g is scalar valued). Then

$$\begin{aligned} [Dfg(a)](h) &= [(D_1fg)(a), \dots, (D_nfg)(a)]h \\ &= [f(a)(D_1g)(a) + (D_1f)(a)g(a), \dots, f(a)(D_ng)(a) + (D_nf)(a)g(a)]h \\ &= f(a)[(D_1g)(a), \dots, (D_ng)(a)]h + [(D_1f)(a), \dots, (D_nf)(a)]g(a)h \\ &= f(a)[Dg(a)]h + ([Df(a)](h))g(a) \end{aligned}$$

Solution.

The problem here is that as we know, the Jacobian matrix does indeed equal the matrix of the derivative of fg at a , but only if the limit-based definition of $D(fg)(a)$ is satisfied. It is possible for all partial derivatives at a to exist, and thus for the Jacobian to be well-defined, yet at the same time for it to not satisfy the limit-based definition. Therefore the “proof” shown here does not establish differentiability of fg at a .

Exercise 6

- a. Prove the rule for differentiating dot products (part 7 of Theorem 1.8.1) directly from the definition of the derivative.
- b. Let $U \subset \mathbb{R}^3$ be open. Show by a similar argument that if $f, g : U \rightarrow \mathbb{R}^3$ are both differentiable at a , then so is the cross product $f \times g : U \rightarrow \mathbb{R}^3$. Find the formula for this derivative.

Solution.

(a) We begin by examining the definition of the derivative for the proposed $D(\langle f, g \rangle)(a)$:

$$\lim_{h \rightarrow 0} \frac{1}{\|h\|} (\langle f(a+h), g(a+h) \rangle - \langle f(a), g(a) \rangle - \langle Df(a)(h), g(a) \rangle - \langle f(a), Dg(a)(h) \rangle) = 0$$

We work with the numerator to obtain the following:

$$\begin{aligned} & \langle f(a+h), g(a+h) \rangle - \langle f(a), g(a) \rangle - \langle Df(a)(h), g(a) \rangle - \langle f(a), Dg(a)(h) \rangle \\ = & \langle f(a+h), g(a+h) \rangle - \langle f(a), g(a) \rangle + \langle f(a), g(a+h) \rangle - \langle f(a), g(a+h) \rangle - \langle Df(a)(h), g(a) \rangle - \langle f(a), Dg(a)(h) \rangle \\ = & \langle f(a+h) - f(a), g(a+h) \rangle + \langle f(a), g(a+h) - g(a) - Dg(a)(h) \rangle - \langle Df(a)(h), g(a) \rangle \\ = & \langle f(a+h) - f(a), g(a+h) \rangle + \langle Df(a)(h), g(a+h) \rangle - \langle Df(a)(h), g(a+h) \rangle \\ & + \langle f(a), g(a+h) - g(a) - Dg(a)(h) \rangle - \langle Df(a)(h), g(a) \rangle \\ = & \langle f(a+h) - f(a) - Df(a)(h), g(a+h) \rangle + \langle Df(a)(h), g(a+h) - g(a) \rangle + \langle f(a), g(a+h) - g(a) - Dg(a)(h) \rangle \end{aligned}$$

Now observe that if we compute the norm of this expression over $\|h\|$, let it be known as N , we obtain that:

$$\begin{aligned} N = & \left\| \left\langle \frac{f(a+h) - f(a) - Df(a)(h)}{\|h\|}, g(a+h) \right\rangle + \left\langle \frac{Df(a)(h)}{\|h\|}, g(a+h) - g(a) \right\rangle \right. \\ & \left. + \left\langle f(a), \frac{g(a+h) - g(a) - Dg(a)(h)}{\|h\|} \right\rangle \right\| \leq \left\| \frac{f(a+h) - f(a) - Df(a)(h)}{\|h\|} \right\| \cdot \|g(a+h)\| \\ & + \left\| \frac{Df(a)(h)}{\|h\|} \right\| \cdot \|g(a+h) - g(a)\| + \|f(a)\| \cdot \left\| \frac{g(a+h) - g(a) - Dg(a)(h)}{\|h\|} \right\| \end{aligned}$$

, where we used the triangle inequality and the Cauchy-Schwarz inequality. Now, since both f, g are differentiable at a , they are also continuous at a , and thus bounded in a neighborhood around a . By all of this we easily conclude that the first and third terms of the RHS sum clearly tend to 0 as $h \rightarrow 0$. As for the second term of the sum, we have already seen that the quotient $\frac{\|Df(a)(h)\|}{\|h\|}$ is bounded, whereas $g(a+h) - g(a) \rightarrow 0$ as $h \rightarrow 0$. Therefore, this term also tends to 0 as $h \rightarrow 0$. We conclude then that $N \rightarrow 0$ as $h \rightarrow 0$, and thus the original limit is zero as well, meaning that the proposed formula for $D(\langle f, g \rangle)(a)$ is indeed correct.

(b) We first write out the formula for the cross product, assuming that $f(x) = (f_1(x), f_2(x), f_3(x))$, $g(x) = (g_1(x), g_2(x), g_3(x))$, where each $f_i, g_i : U \rightarrow \mathbb{R}$ is differentiable at a :

$$(f \times g)(a) = \begin{pmatrix} f_2(a)g_3(a) - f_3(a)g_2(a) \\ -f_1(a)g_3(a) + f_3(a)g_1(a) \\ f_1(a)g_2(a) - f_2(a)g_1(a) \end{pmatrix}$$

Observe that each of the components here is differentiable at a (by the sum and product rules for derivatives). It is therefore the case that $(f \times g)$ is also differentiable at a , and:

$$D(f \times g)(a) = \begin{pmatrix} D(f \times g)_1(a) \\ D(f \times g)_2(a) \\ D(f \times g)_3(a) \end{pmatrix},$$

where:

$$D(f \times g)_1(a) = \begin{pmatrix} D_1 f_2(a)g_3(a) + D_1 g_3(a)f_2(a) - D_1 f_3(a)g_2(a) - D_1 g_2(a)f_3(a) \\ D_2 f_2(a)g_3(a) + D_2 g_3(a)f_2(a) - D_2 f_3(a)g_2(a) - D_2 g_2(a)f_3(a) \\ D_3 f_2(a)g_3(a) + D_3 g_3(a)f_2(a) - D_3 f_3(a)g_2(a) - D_3 g_2(a)f_3(a) \end{pmatrix}^T,$$

and the rest is computed similarly to result in a 3×3 matrix for $D(f \times g)(a)$.

Exercise 13

Let $U \subset \mathbb{R}^{n \times n}$ be the set of matrices A such that $A + A^2$ is invertible. Compute the derivative of the map $F : U \rightarrow \mathbb{R}^{n \times n}$ given by $F(A) = (A + A^2)^{-1}$.

Solution.

Consider the functions $F_1, F_2 : U \rightarrow \mathbb{R}^{n \times n}$, $F_1(A) = A + A^2$, $F_2(A) = A^{-1}$. We then have that $F(A) = F_2(F_1(A))$. From previous results and the addition rule for derivatives, we have that $DF_1(A)(H) = IH + AH + HA = H + AH + HA$ and that $DF_2(A)(H) = -A^{-1}HA^{-1}$. By an application of the chain rule we then obtain that:

$$DF(A)(H) = [DF_2(F_1(A))](DF_1(A)(H)) = -(A + A^2)^{-1}(H + AH + HA)(A + A^2)^{-1}$$