

J. J. Duistermaat & J. A. C Kolk's Multidimensional Real Analysis 1:
Differentiation - A partial solutions manual

George Lydakis

December 3, 2024

Contents

1	Continuity	5
2	Differentiation	17

Chapter 1

Continuity

Exercise 6

Let A, B be any two subsets of \mathbb{R}^n .

- (i) Prove $\overset{\circ}{A} \subset \overset{\circ}{B}$ if $A \subset B$.
- (ii) Show $(A \cap B)^\circ = \overset{\circ}{A} \cap \overset{\circ}{B}$.
- (iii) From (ii) deduce $\overline{A \cup B} = \overline{A} \cup \overline{B}$.

Solution.

(i) Suppose $A \subset B$. Let then $x \in \overset{\circ}{A}$. It is then the case that there exists $r > 0$ such that $B_r(x) \subset A$. Then it also holds that $B_r(x) \subset B$. But then from the definition of the interior of a set, this means precisely that $x \in \overset{\circ}{B}$. Thus $\overset{\circ}{A} \subset \overset{\circ}{B}$.

(ii) Let first $x \in (A \cap B)^\circ$. Then there exists $r > 0$ such that $B_r(x) \subset A \cap B$. This means both that $B_r(x) \subset A$ and that $B_r(x) \subset B$. But then again from the definition of the interior of a set, this directly implies that $x \in \overset{\circ}{A}, x \in \overset{\circ}{B}$, and thus clearly $x \in \overset{\circ}{A} \cap \overset{\circ}{B}$. Therefore we have shown that $(A \cap B)^\circ \subset \overset{\circ}{A} \cap \overset{\circ}{B}$. In the other direction, let $x \in \overset{\circ}{A} \cap \overset{\circ}{B}$. Then $x \in \overset{\circ}{A}, x \in \overset{\circ}{B}$. This means that there exist $r_1, r_2 > 0$ such that $B_{r_1}(x) \subset A, B_{r_2}(x) \subset B$. Let $r = \min\{r_1, r_2\}$, in which case clearly $B_r(x) \subset B_{r_1}(x), B_r(x) \subset B_{r_2}(x)$, which implies that $B_r(x) \subset A, B_r(x) \subset B$. Then this also implies that $B_r(x) \subset A \cap B$. By the definition of the interior, $x \in (A \cap B)^\circ$.

(iii) For this we will need lemma 1.2.10 which states that for any $A \subset \mathbb{R}^n$, $(\overline{A})^c = (\overset{\circ}{A^c})$, $(\overset{\circ}{A})^c = \overline{A^c}$, i.e. that the complement of the closure of a set equals the interior of the complement, and that the complement of the interior of the set equals the closure of the set's complement. We will also need De Morgan's laws for sets. With that in mind we have that:

$$(\overline{A \cup B})^c = (A \cup B)^\circ = (A^c \cap B^c)^\circ = (\overset{\circ}{A^c}) \cap (\overset{\circ}{B^c}) = (\overline{A})^c \cap (\overline{B})^c = (\overline{A} \cup \overline{B})^c$$

, and then taking complements on both sides yields $\overline{A \cup B} = \overline{A} \cup \overline{B}$. Notice that on our third step we used the result from (ii).

Exercise 7

Let $n \in \mathbb{N} \setminus \{1\}$. Let F be a closed subset of \mathbb{R}^n with $\overset{\circ}{F} \neq \emptyset$. Prove that the boundary ∂F of F contains infinitely many points, unless $F = \mathbb{R}^n$ (in which case we have $\partial F = \emptyset$).

Solution.

We begin with the case where $F = \mathbb{R}^n$. Then clearly $F^c = \emptyset$, which implies that no boundary points exist (because then an open ball around them would have a non-empty intersection with the empty set). Thus $\partial F = \emptyset$.

Now let's examine the more interesting case where $F \neq \mathbb{R}^n$. Since the interior of F is non-empty, it contains at least one point, x . For this x it holds that there exists $r > 0$ such that $B_r(x) \subset F$. Furthermore, because $F \neq \mathbb{R}^n$, there exists at least one $y \notin F$. Let $S = \text{span}\{y - x\}$. Clearly, $\dim S = 1$. We know from Linear Algebra that:

$$\dim S + \dim S^\perp = \dim \mathbb{R}^n \implies \dim S^\perp = n - \dim S \geq n - 1 \geq 1$$

, since $n > 1$. Consequently, there exists at least one $z \in S^\perp, z \neq 0$ for which, by definition, $\langle z, y - x \rangle = 0$. Since z is orthogonal to $y - x$, we have that $z, y - x$ are linearly independent. Now we have that for any $\lambda \in [0, 1]$ the vector $w_\lambda = x + \frac{\lambda rz}{2\|z\|}$ is contained in $B_r(x)$.

Now consider, for any $\lambda \in [0, 1]$, the vector $y - w_\lambda$. We argue that there exists at least one boundary point of F that is of the form $b_\lambda = w_\lambda c(y - w_\lambda), c \in (0, 1)$. Indeed, consider the set $C = \{c \in (0, 1] | w_\lambda + c(y - w_\lambda) \in F\}$, which is bounded above. It is also non-empty, since $w_\lambda \in B_r(x)$, which implies that for a sufficiently small c , $w_\lambda + c(y - w_\lambda)$ remains in $B_r(x)$, and in fact we can find another sufficiently small radius r' such that $B_{r'}(w_\lambda + c(y - w_\lambda)) \subset B_r(x) \subset F$, thus at least one point of this form is in the interior of F .

Thus the set has a least upper bound, c_m . Firstly, c_m cannot be 1. If it was, then *all* points $w_\lambda + c(y - w_\lambda)$ would be in F , and then we could very easily construct a sequence of them that converges to y . But F is closed, so this is a contradiction. We now claim $b_\lambda = w_\lambda + c_m(y - w_\lambda)$ is a boundary point for F . If it was not, by negating the definition we obtain that $b_\lambda \in \overset{\circ}{F}$ or $b_\lambda \in \overset{\circ}{F}^c$. In the first case, there exists r' such that $B_{r'}(b_\lambda) \subset \overset{\circ}{F}$. But then for $c_{m'} = c_m + \frac{r'}{2}$, $w_\lambda + c_{m'}(y - w_\lambda)$ would be in F , and $c_{m'} > c_m$, contradicting the definition of c_m . In the second case, there exists r' such that $B_{r'}(b_\lambda) \subset \overset{\circ}{F}^c$. But then for $c_{m'} = c_m - \frac{r'}{2}$ we have that $w_\lambda + c_{m'}(y - w_\lambda) \notin F$, which means $c_{m'} < c_m$ is an upper bound for C , which contradicts the definition of c_m as the least upper bound of C .

Our conclusion is that for any $\lambda \in [0, 1]$, we can find a corresponding boundary point of F that has the form $w_\lambda + c(y - w_\lambda), c \in (0, 1)$. Now we claim that no two $b_{\lambda_1}, b_{\lambda_2}, \lambda_1 \neq \lambda_2$ can be equal. Indeed, suppose that this was the case. Then:

$$\begin{aligned} b_{\lambda_1} = b_{\lambda_2} &\implies w_{\lambda_1} + c_1(y - w_{\lambda_1}) = w_{\lambda_2} + c_2(y - w_{\lambda_2}) \implies \\ x + \frac{\lambda_1 rz}{\|z\|} + c_1(y - x - \frac{\lambda_1 rz}{\|z\|}) &= x + \frac{\lambda_2 rz}{\|z\|} + c_2(y - x - \frac{\lambda_2 rz}{\|z\|}) \implies \\ \frac{rz}{\|z\|}(\lambda_1 - \lambda_2 - c_1\lambda_1 + c_2\lambda_2) + (c_1 - c_2)(y - x) &= 0 \end{aligned}$$

But now recall that $z, y - x$ are linearly independent. As such, this equation is only satisfied for $c_1 = c_2, \lambda_1 - \lambda_2 - c_1\lambda_1 + c_1\lambda_2 = 0 \implies \lambda_1(1 - c_1) + \lambda_2(c_1 - 1) = 0$, and since $c_1, c_2 \neq 1$ this yields $\lambda_1 = \lambda_2$ which is a contradiction.

By the above we have proved that this mapping of $[0, 1]$ to boundary points of F is injective. As such, the cardinality of the set of boundary points of F is at least that of $[0, 1]$, which is equal to that of the real numbers, and thus there exist infinite boundary points.

Exercise 10

We define $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$f(x) = \begin{cases} \frac{x_1 x_2^2}{x_1^2 + x_2^6}, & x \neq 0 \\ 0, & x = 0 \end{cases}$$

Show that $\text{im}(f) = \mathbb{R}$ and prove that f is not continuous.

Solution.

Select any $y \in \mathbb{R}$. If $y = 0, f(0) = 0$, thus $0 \in \text{im}(f)$. If $y \neq 0$, we consider the equation:

$$\frac{x_1 x_2^2}{x_1^2 + x_2^6} = y \implies x_1 x_2^2 = y(x_1^2 + x_2^6)$$

Now, if we examine what happens when $x_1 = x_2^3$, we see that:

$$x_2^5 = y(x_2^6 + x_2^6) \implies x_2^5 = 2x_2^6 y \implies x_2^5(1 - 2x_2 y) = 0$$

Since $x_1 = x_2^3$, we examine the case where $x_2 \neq 0$ (otherwise we would not be able to consider this equation to begin with). For the above equation to hold, it must then be the case that:

$$1 - 2x_2 y = 0 \implies x_2 = \frac{1}{2y}$$

, which we can do since $y \neq 0$. In that case we then get that $x_1 = \frac{1}{8y^3}$. From this we conclude that for $y \neq 0$, $f(\frac{1}{8y^3}, \frac{1}{2y}) = y$, which means that any non-zero y belongs in the image of f as well. We've therefore shown that $\text{im}(f) = \mathbb{R}$.

For f to be continuous, it has to be continuous at 0. Consider then approaching 0 with a sequence such that $x_1 = x_2^3$ and $x_2 > 0$. Evaluated on this sequence, f equals:

$$f(x) = \frac{x_2^5}{x_2^6 + x_2^6} = \frac{1}{2x_2}$$

It is evident that as $x_2 \rightarrow 0$, this function tends to (positive) infinity, and thus the limit at $(0,0)$ cannot exist, which means that f is not continuous at 0, and consequently not continuous.

Exercise 11

A function $f : \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}$ is said to be positively homogeneous of degree $d \in \mathbb{R}$, if $f(tx) = t^d f(x)$ for all $x \neq 0$ and $t > 0$. Assume f to be continuous. Prove that f has an extension as a continuous function to \mathbb{R}^n precisely in the following cases: (i) if $d < 0$, then $f = 0$; (ii) if $d = 0$, then f is a constant function; (iii) if $d > 0$, then there is no further condition on f . In each case, indicate which value f has to assume at 0.

Solution. We examine the three cases one by one:

- If $d < 0$, we have the following. Suppose first that f is not the zero function, which means that for some $y \neq 0$, it holds that $f(y) \neq 0$. Consider then the sequence $i \rightarrow p_i = \frac{y}{i}$, which approaches 0 as $i \rightarrow \infty$. Then if a continuous extension of f exists, f would have to have a limit at 0. More specifically, the limit of $i \rightarrow f(p_i)$ would have to exist. We then have that:

$$\lim_{i \rightarrow \infty} f(p_i) = \lim_{i \rightarrow \infty} f\left(\frac{y}{i}\right) = \lim_{i \rightarrow \infty} \left(\frac{1}{i}\right)^d f(y) = \lim_{i \rightarrow \infty} i^{-d} f(y)$$

Clearly, this tends to infinity as $i \rightarrow \infty$, which means that the limit of f at 0 cannot exist, and thus f cannot be continuous.

Now let's assume that f is the zero function. In that case, clearly its limit at 0 will be 0, thus extending it by setting $f(0) = 0$ makes it continuous in \mathbb{R}^n .

- If $d = 0$, then $f(tx) = f(x)$ for $x \neq 0, t > 0$. First, let's assume that f is not constant, which means that for two $x_1 \neq x_2$, it holds that $f(x_1) \neq f(x_2)$. Then consider the two sequences $i \rightarrow p_i = \frac{x_1}{i}, i \rightarrow q_i = \frac{x_2}{i}$, both of which tend to 0 as $i \rightarrow \infty$. However:

$$\lim_{i \rightarrow \infty} f(p_i) = \lim_{i \rightarrow \infty} f\left(\frac{x_1}{i}\right) = \lim_{i \rightarrow \infty} f(x_1) = f(x_1)$$

, and by exactly the same procedure we can see that $\lim_{i \rightarrow \infty} f(q_i) = f(x_2)$, which means that the limit of f at 0 does not exist, and thus f cannot be continuous.

If we now assume that f is indeed constant, i.e., $f(x) = c, x \neq 0$, then it suffices to set $f(0) = c$ to make f continuous in \mathbb{R}^n .

- If $d > 0$, then $f(tx) = t^d f(x)$. Let S be the unit sphere in \mathbb{R}^n . We now have the following.

First, if f is 0 on S , then we observe that any non-zero $v \in \mathbb{R}^n$ can be written as $v = \frac{v'}{\|v\|}$, where $v' \in S$. Thus by the positive homogeneity (of degree d) of f we can obtain that $f(v) = 0$ for all non-zero v . Obviously, this means that the extension of f must then be valued 0 at 0 in order to be continuous.

Now suppose that f is not zero on at least one point $x \in S$. Form the sequence of points $i \rightarrow \frac{x}{i}$, which clearly converges to 0. Then because $f(\frac{x}{i}) = \frac{f(x)}{i^d}$, and $f(x) \neq 0$, in order for the extension of f to be continuous at 0 it has to be the case that $f(0) = \lim_{i \rightarrow \infty} f(\frac{x}{i}) = 0$. Now we need to prove that this is not only necessary but also sufficient for the extension to be 0 at 0.

Firstly, observe that S is a compact set. Therefore, f must obtain some maximum and minimum values on it. Call these $f(x_M), f(x_m)$, achieved on $x_M, x_m \in S$. Set $M = \max\{|f(x_M)|, |f(x_m)|\}$. Observe that since $f(x_m) \leq f(x) \leq f(x_M)$ for all $x \in S$, it holds that $|f(x)| < M$ for all $x \in S$. Now pick any $\epsilon > 0$ and set $\delta = (\frac{\epsilon}{M})^{\frac{1}{d}}$. Consider any $x \neq 0$ such that $\|x - 0\| < \delta$. We have that:

$$f(x) = f(\|x\|x') = \|x\|^d f(x')$$

, where $x' \in S$. Then $|f(x)| < ((\frac{\epsilon}{M})^{\frac{1}{d}})^d \cdot |f(x')| \leq \frac{\epsilon M}{M} = \epsilon$. But this precisely means that the limit of f as x tends to 0 is 0, and thus setting $f(0) = 0$ suffices to make the extension continuous.

Exercise 20

Suppose $(x_k)_{k \in \mathbb{N}}$ is a convergent sequence in \mathbb{R}^n with limit $a \in \mathbb{R}^n$. Show that $\{a\} \cup \{x_k | k \in \mathbb{N}\}$ is a compact subset of \mathbb{R}^n .

Solution.

Let the described set be called S . We need to show that S is both closed and bounded. First, set $M = \|a\| + 1$. Clearly, $\|a\| < M$. By the definition of the limit, there exists some N such that for $i > N$ it holds that $\|x_i - a\| < 1$, which by using the triangle inequality easily yields $\|x_i\| < 1 + \|a\| = M$. The set of $x_i, i \leq N$ contains a finite number of elements, and thus so does the set of the corresponding $\|x_i\|$. Set then $M' = \max\{\|x_i\|, i \leq N\} + 1$, and $L = \max\{M', M\}$. Then we have that for any $x \in S, \|x\| < L$, thus S is bounded.

Now we need to show that S is closed. Suppose that $i \rightarrow y_i$ is a convergent sequence in S . By the definition of S , it has to be the case that each y_i either equals some x_k or it equals a . Suppose y_i converges to a point $b \notin S$. Then for any $\epsilon > 0$ we can find $N > 0$ such that for $i > N, \|y_i - b\| < \epsilon$. If for some $N > 0$ it holds that $y_i = a$ for all $i > N$, then we can see that the sequence converges to $b \notin S$ but at the same time clearly converges to $a \in S$ (the distance of the terms from a becomes in fact zero). By the uniqueness of the limit, this is a contradiction. Thus for all $N > 0$ there have to exist an infinite number of $i > N$ for which $y_i \neq a$. These must then all be of the form x_k . Now assume that for some $N > 0$ all of these y_i equal x_k with $k \leq N$ (they don't have to equal the same x_k necessarily). But then these are countably many: they are at most N . Thus there exists a minimum distance of them from b . But this contradicts the hypothesis that y_i can get infinitely close to b .

What we have shown so far is that for all $N > 0$ there exists at least one y_i that equals some x_k with $k > N$. We can thus form a subsequence z_j of y_i by setting z_j be the first y_i which equals an x_k with $k > j$. This is now clearly a subsequence of x_k , and thus converges to a . But at the same time, it is a subsequence of y_i and thus converges to b , a contradiction. Therefore, y_i cannot converge to a $b \notin S$, which equivalently shows that S is closed.

Exercise 22

Show that $K \subset \mathbb{R}^n$ is compact if and only if every continuous function $K \rightarrow \mathbb{R}$ is bounded.

Solution.

\implies : Suppose K is compact. We then know that any continuous function defined on it has a maximum and a minimum value, which means that it is bounded.

\impliedby : Now suppose every continuous function from K to \mathbb{R} is bounded. Suppose K was not compact. In exercise 1.6.2 of Hubbard & Hubbard we proved that there exists then a continuous unbounded function from K to \mathbb{R} . Obviously, this contradicts our hypothesis of all continuous functions from K to \mathbb{R} being bounded, and thus K must be compact.

Exercise 23

Let $A \subset \mathbb{R}^n$ be arbitrary, let $K \subset \mathbb{R}^p$ be compact, and let $p : A \times K \rightarrow A$ be the projection with $p(x, y) = x$. Show that p is proper and deduce that it is a closed mapping.

Solution.

Firstly, a mapping $f : \mathbb{R}^n \rightarrow \mathbb{R}^p$ is said to be proper if the inverse image under f of every compact set in \mathbb{R}^p is compact in \mathbb{R}^n . Secondly, a mapping $f : \mathbb{R}^n \rightarrow \mathbb{R}^p$ is called closed if, for all $C \subset \mathbb{R}^n$ it holds that $f(C)$ is also closed.

In our case, p can be thought of as a mapping from a subset of \mathbb{R}^{n+p} to a subset of \mathbb{R}^n . As such, let S be a compact subset of A . We are interested in $p^{-1}(S)$. From p 's definition, we can see that it maps any element of the form (x, y) , $x \in S, y \in K$ to an element of S (namely, to x). Furthermore, if p maps an element (x, y) to an element $z \in S$, then from p 's definition it must hold that $x = z$ and that $y \in K$. From these observations we can conclude that $p^{-1}(S) = S \times K$. Both S, K are compact, which means that they are bounded, say by M_S, M_K respectively. For any element $x \in S \times K$, $\|x\|^2 = \|s\|^2 + \|k\|^2$, $s \in S, k \in K$. Therefore, $\|x\| = \sqrt{\|s\|^2 + \|k\|^2} < \sqrt{M_S^2 + M_K^2}$. We have thus found a bound for the norm of any element of $p^{-1}(S)$, which means that the set is bounded.

Now, to show that it is also closed, suppose $i \rightarrow s_i$ is any convergent sequence in $p^{-1}(S) = S \times K$ that converges to $(a_1, a_2, \dots, a_n, a_{n+1}, \dots, a_{n+p})$. We know that this implies that each of the coordinates of s_i converges to the respective coordinate of this limit. This in turn implies that the sequence $i \rightarrow (s_{i,1}, \dots, s_{i,n})$ of the first n coordinates converges to (a_1, \dots, a_n) and that the sequence $i \rightarrow (s_{i,n+1}, \dots, s_{i,n+p})$ converges to $(a_{n+1}, \dots, a_{n+p})$. All points of this first sequence belong in S , which is closed as a compact set, and thus $(a_1, \dots, a_n) \in S$. Similarly, all points of the second sequence belong in K , again a closed set due to being compact, and thus $(a_{n+1}, \dots, a_{n+p}) \in K$. But then $(a_1, a_2, \dots, a_n, a_{n+1}, \dots, a_{n+p}) \in S \times K = p^{-1}(S)$, meaning that $p^{-1}(S)$ is indeed closed, and thus compact.

Now, one can easily see that p is also continuous: whenever a sequence of points $i \rightarrow (x_i, y_i)$ converges to a point (x, y) , the sequence $p(x_i, y_i)$ converges trivially to x , which is also the value of p at (x, y) . But then from theorem 1.8.6 we have that f is a closed mapping due to being proper and continuous.

Exercise 25

A nonconstant polynomial function $f : \mathbb{C} \rightarrow \mathbb{C}$ is proper. Deduce that $f(\mathbb{C})$ is closed in \mathbb{C} .

Solution.

We have that $f(z) = \sum_{i=1}^m a_i z^i$, with $m \geq 1$ being the degree of the polynomial, and thus $a_m \neq 0$. Let $S \subset \mathbb{C}$ be a compact set. Let also $T = f^{-1}(S)$ be the inverse image of S under f . We need to show that T is also compact. Suppose, first, that T is not bounded. Therefore it contains a sequence $i \rightarrow z_i$ such that the norm of z_i goes to infinity as $i \rightarrow \infty$. Observe that for non-zero z we can write:

$$f(z) = a_m z^m \left(1 + \sum_{i=1}^{m-1} \frac{a_i}{a_m z^{m-i}} \right)$$

, where, crucially, $a_m \neq 0$. From this one can fairly easily draw the conclusion that if $i \rightarrow z_i$ is such that the norm of z_i is unbounded, the first term has a norm tending to infinity when $i \rightarrow \infty$ whereas the second tends to 1. Thus the sequence of $f(z_i)$ would also be unbounded, a contradiction because S is compact. Now suppose T is not closed. Then there exists a sequence $i \rightarrow z_i \in T$ that converges to $a \notin T$. Because f is continuous as a polynomial, we have that:

$$\lim_{i \rightarrow \infty} z_i = a \implies \lim_{i \rightarrow \infty} f(z_i) = f(a)$$

Because $a \notin T$, by definition it has to hold that $f(a) \notin S$. However, all of $f(z_i)$ by definition are in S , which means that because S is compact, their sequence contains a convergent subsequence that converges to a point in S . By the above however, this subsequence would have to converge to $f(a) \notin S$, a contradiction. Therefore T is closed.

Having shown that f is proper and because it is also continuous and \mathbb{C} is closed, we can draw the conclusion that f is a closed mapping and $f(\mathbb{C})$ is also closed.

Exercise 28

For every polynomial function $p : \mathbb{C} \rightarrow \mathbb{C}$ there exists $w \in \mathbb{C}$ with $|p(w)| = \inf\{|p(z)| \mid z \in \mathbb{C}\}$. Prove this using exercise 1.25. This is known as *Cauchy's Minimum Theorem*.

Solution.

First consider the case of p being a constant polynomial. Then clearly every $w \in \mathbb{C}$ fulfills the required condition, since p only achieves one value. Now, if p is not constant, exercise 25 applies. We thus know that $p(\mathbb{C})$ is a closed set in \mathbb{C} . We now form the set $S = \{|z| \mid z \in p(\mathbb{C})\}$, that is, the set of all possible magnitudes of $p(z)$. Notice that by the definition of the norm of complex numbers, this set has a lower bound, 0. Thus, by the Least Upper Bound axiom and its consequence for lower bounds, S has a *greatest* lower bound, $\inf S$. Now, by the characterization of supremum and infimum in Carothers, exercise 3, we have that for every $\epsilon > 0$, there exists $x \in S$ such that $x < \inf S + \epsilon \implies x - \inf S < \epsilon$. Form a sequence of x_i for $\epsilon_i = \frac{1}{i}$. By the definition of S , we can always find z_i such that $|p(z_i)| = x_i$.

Now clearly all of these magnitudes have to be less than $\inf S + 1$. Therefore, we can form a *compact* set of all z with $|p(z)| \leq \inf S + 1$. The sequence z_i belongs in this compact set, and thus has a convergent subsequence (Bolzano-Weierstrass theorem), z_{i_k} , that converges to $p(w)$ for some $w \in \mathbb{C}$. But because the norm is a continuous function, by applying it on the convergent sequence z_{i_k} we obtain that it must hold that:

$$\lim_{i_k \rightarrow \infty} |p(z_{i_k})| = |p(w)| = \inf S$$

This concludes the proof.

Exercise 29

A function $f : F \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ is said to be a contraction in F when f is Lipschitz continuous with a Lipschitz constant $\epsilon < 1$, in other words, when:

$$\|f(x) - f(x')\| \leq \epsilon \|x - x'\| < \|x - x'\|, (x, x' \in F)$$

The Contraction Lemma states the following:

Assume $F \subset \mathbb{R}^n$ is closed and $x_0 \in F$. Let $f : F \rightarrow F$ be a contraction with contraction factor $\leq \epsilon$. Then there exists a unique point $x \in F$ with

$$f(x) = x; \text{ furthermore } \|x - x_0\| \leq \frac{1}{1 - \epsilon} \|f(x_0) - x_0\|$$

The notation below is as in the Contraction Lemma. Define $g : F \rightarrow \mathbb{R}$ by $g(x) = \|x - f(x)\|$.

(i) Verify $|g(x) - g(x')| \leq (1 + \epsilon)\|x - x'\|$ for all $x, x' \in F$, and deduce that g is continuous on F .

If F is bounded the continuous function g assumes its minimum at a point p belonging to the compact set F .

(ii) Show $g(p) \leq g(f(p)) \leq \epsilon g(p)$, and conclude that $g(p) = 0$. That is, $p \in F$ is a fixed point of f .

If F is not bounded, set $F_0 = \{x \in F \mid g(x) \leq g(x_0)\}$. Then $F_0 \subset F$ is nonempty and closed.

(iii) Prove, for $x \in F_0$,

$$\|x - x_0\| \leq \|x - f(x)\| + \|f(x) - f(x_0)\| + \|f(x_0) - x_0\| \leq 2g(x_0) + \epsilon \|x - x_0\|$$

Hence $\|x - x_0\| \leq \frac{2g(x_0)}{1 - \epsilon}$ for $x \in F_0$, and therefore F_0 is bounded.

(iv) Show that f is a mapping of F_0 in F_0 by noting $g(f(x)) \leq \epsilon g(x) \leq g(x_0)$ for $x \in F_0$ and proceed as above to find a fixed point $p \in F_0$.

Solution.

(i) For any two $x, x' \in F$, we have that:

$$\begin{aligned} |g(x) - g(x')| &= \left| \|x - f(x)\| - \|x' - f(x')\| \right| \leq \| (x - f(x)) - (x' - f(x')) \| \\ &\leq \| (x - x') + (f(x') - f(x)) \| \leq \|x - x'\| + \|f(x') - f(x)\| \leq (1 + \epsilon)\|x - x'\| \end{aligned}$$

, where we used the triangle inequality and the Lipschitz continuity of f . This shows that g is Lipschitz continuous in F , and thus more specifically continuous.

(ii) f maps F to F , therefore $f(p) \in F$. Because g achieves its minimum value on p , it holds that $g(p) \leq g(x)$ for all $x \in F$. More particularly then, $g(p) \leq g(f(p))$. We also have that $g(f(p)) = \|f(p) - f(f(p))\| \leq \epsilon \|p - f(p)\| = \epsilon g(p)$, by using the Lipschitz continuity of f . But then the only way that $g(p) \leq g(f(p)) \leq \epsilon g(p)$ can be true with $\epsilon < 1$ is if $g(p) = 0$, which implies $p = f(p)$ by the definition of g , therefore that p is a fixed point of f .

(iii) Here we have the following:

$$\begin{aligned} \|x - x_0\| &= \|x - f(x) + f(x) - x_0 + f(x_0) - f(x_0)\| \\ &\leq \|x - f(x)\| + \|f(x) - f(x_0)\| + \|f(x_0) - x_0\| = g(x) + g(x_0) + \|f(x) - f(x_0)\| \\ &\leq 2g(x_0) + \epsilon \|x - x_0\| \end{aligned}$$

, where we used the triangle inequality, the defining property of F_0 and the Lipschitz continuity of f . From this we can deduce that $\|x - x_0\| \leq \frac{2g(x_0)}{1-\epsilon}$. But then by the triangle inequality, and since x_0 is fixed, $\|x\|$ is bounded, meaning that F_0 is bounded.

(iv) In order to show that f is a mapping of F_0 to F_0 , we need to show that for any $x \in F_0$, $f(x)$ is also in F_0 . This means that we need to show that $g(f(x)) \leq g(x_0)$. We have that:

$$g(f(x)) = \|f(x) - f(f(x))\| \leq \epsilon \|x - f(x)\| = \epsilon g(x) \leq \epsilon g(x_0) \leq g(x_0)$$

, where the last inequality holds because $\epsilon < 1$. We've thus shown that $f(x) \in F_0$, that is, that f is a mapping from F to F . Therefore F_0 is a compact set on which g is continuous, meaning that it achieves a minimum value on it. This means that we can directly apply part (ii) to show that there exists a fixed point p of f that belongs in F_0 .

We note that the uniqueness of the fixed point f is shown in the same way as in the contraction lemma: assume two fixed points $x, x', x \neq x'$, and then $0 < \|x - x'\| = \|f(x) - f(x')\| < \|x - x'\|$, which is a clear contradiction, and the inequality asserted by the lemma can be shown by:

$$\begin{aligned} \|x - x_0\| &= \|f(x_0) + x - f(x_0) - x_0\| \leq \|f(x_0) - x_0\| + \|x - f(x_0)\| = \|f(x) - f(x_0)\| + \|f(x_0) - x_0\| \\ &\leq \epsilon \|x - x_0\| + \|f(x_0) - x_0\| \implies \|x - x_0\|(1 - \epsilon) \leq \|f(x_0) - x_0\| \end{aligned}$$

Exercise 30

Let $A \subset \mathbb{R}^n$ be bounded and $f : A \rightarrow \mathbb{R}^p$ uniformly continuous. Show that f is bounded on A .

Solution.

Suppose that f is not bounded. Then there must exist a sequence of points $y_1, y_2, \dots \in \mathbb{R}^p$ such that $y_i = f(x_i), x_i \in A$ and $\|y_i\| > i$. Recall that A is bounded, which means that its closure, \bar{A} is bounded as well. Indeed, if this was not the case, there would exist boundary points with a norm greater than any $M > 0$, and by definition there would exist points inside A arbitrarily close to them, which means that they as well would have unbounded norms, a contradiction.

Since \bar{A} is bounded and is by definition closed, it is compact. Furthermore, all $x_i \in A \subset \bar{A}$, which means that by the Bolzano-Weierstrass theorem the sequence (x_n) has a convergent subsequence (x_{n_k}) . This sequence converges to a point $x \in \bar{A}$, and, by an easy extension of theorem 1.16 of Carothers, is equivalently Cauchy. In exercise 33 we prove that a uniformly continuous f maps Cauchy sequences to Cauchy sequences. This would mean that $(f(x_{n_k}))$ is equivalently convergent, and thus bounded. However, by the definition of $y_i = f(x_i)$, it cannot be the case that $(f(x_{n_k}))$ are bounded, since $\|f(x_{n_k})\| > n_k$. We arrive at a contradiction. Thus, f must be bounded on A .

Exercise 31

Let $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ be uniformly continuous.

(i) Show that $f + g$ is uniformly continuous.

(ii) Show that fg is uniformly continuous if f, g are bounded. Give a counterexample to show that the condition of boundedness cannot be dropped in general.

(iii) Assume $n = 1$. Is $g \circ f$ uniformly continuous?

Solution.

(i) f, g are both uniformly continuous, which means that for any $\epsilon > 0$, there exist $\delta_1, \delta_2 > 0$ such that for all $x, y \in \mathbb{R}^n$ with $\|x - y\| < \delta_1$ it holds that $|f(x) - f(y)| < \epsilon$ and for all x, y with $\|x - y\| < \delta_2$ it holds that $|g(x) - g(y)| < \epsilon$.

Pick then any $\epsilon > 0$. Then there exist $\delta_1, \delta_2 > 0$ such that $\|x - y\| < \delta_1 \implies |f(x) - f(y)| < \frac{\epsilon}{2}$ and $\|x - y\| < \delta_2 \implies |g(x) - g(y)| < \frac{\epsilon}{2}$. Now set $\delta = \min\{\delta_1, \delta_2\}$. For all x, y with $\|x - y\| < \delta$ it holds that:

$$|(f + g)(x) - (f + g)(y)| = |f(x) - f(y) + g(x) - g(y)| \leq |f(x) - f(y)| + |g(x) - g(y)| < \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon$$

By finding this choice of δ for any ϵ we have shown that $f + g$ is indeed uniformly continuous.

(ii) Pick any $\epsilon > 0$. We begin as follows, for any two $x, y \in \mathbb{R}^n$:

$$\begin{aligned} |(fg)(x) - (fg)(y)| &= |f(x)g(x) - f(x)g(y) + f(x)g(y) - f(y)g(y)| \\ &= |f(x)(g(x) - g(y)) + g(y)(f(x) - f(y))| \leq |f(x)| \cdot |g(x) - g(y)| + |g(y)| \cdot |f(x) - f(y)| \end{aligned}$$

f and g are both bounded, so there exist M_1, M_2 such that $|f(x)| < M_1, |g(y)| < M_2$ for all x, y . Set $M = \max\{M_1, M_2\}$. If $M = 0$, both functions are the zero function and fg is zero too, thus trivially uniformly continuous. If $M > 0$, pick any $\epsilon > 0$. Then, there exist $\delta_1, \delta_2 > 0$ such that:

$$\|x - y\| < \delta_1 \implies |f(x) - f(y)| < \frac{\epsilon}{2M}$$

$$\|x - y\| < \delta_2 \implies |g(x) - g(y)| < \frac{\epsilon}{2M}$$

Set $\delta = \min\{\delta_1, \delta_2\}$. Then for all x, y with $\|x - y\| < \delta$ we can see that:

$$|(fg)(x) - (fg)(y)| \leq |f(x)| \cdot |g(x) - g(y)| + |g(y)| \cdot |f(x) - f(y)| < M \cdot \frac{\epsilon}{2M} + M \cdot \frac{\epsilon}{2M} = \epsilon$$

, which proves that fg is uniformly continuous. Now for a counterexample, consider:

$$f(x_1, x_2) = x_1, g(x_1, x_2) = x_2$$

, both of which are clearly uniformly continuous as linear functions, but neither of which is bounded. Then $(fg)(x_1, x_2) = x_1 x_2$. Suppose now that fg is uniformly continuous. Pick $\epsilon = 1$. Then there must exist $\delta > 0$ such that for all $x, y \in \mathbb{R}^2, \|x - y\| < \delta \implies |(fg)(x) - (fg)(y)| < 1$.

For this δ , consider $x = (x_1, x_1), y = (x_1 + \frac{\delta}{2}, x_2 + \frac{\delta}{2})$. It is clear that $\|x - y\| < \delta$. Then it must hold that:

$$|(fg)(x) - (fg)(y)| < \epsilon \implies |x_1^2 - (x_1 + \frac{\delta}{2})^2| < \epsilon \implies |\delta x_1 + \frac{\delta^2}{4}| < \epsilon$$

If we pick, for example, $x_1 = \frac{\epsilon}{\delta}$, we can clearly see that this does not hold, and thus fg cannot be uniformly continuous.

(iii) Pick any $\epsilon > 0$. Then there exists $\delta > 0$ such that for all $x, y, \|x - y\| < \delta$ it holds that $|g(x) - g(y)| < \delta$. Now for this $\delta > 0$, because f is uniformly continuous, there exists $\delta' > 0$ such that for all x', y' with $\|x' - y'\| < \delta'$ it holds that $|f(x') - f(y')| < \delta$. But then this means that $|g(f(x')) - g(f(y'))| < \epsilon$, which means that if we choose δ' for the ϵ , we satisfy the definition of uniform continuity for $g \circ f$.

Exercise 32

Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be continuous and define $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ by $f(x) = g(x_1x_2)$. Show that f is uniformly continuous only if g is a constant function.

Solution.

In order to show this, we must show that f being uniformly continuous implies that g is continuous. We will do this by contradiction. However, we will first prove a useful lemma regarding uniformly continuous functions:

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is uniformly continuous if and only if for all sequences x_i, y_i such that:

$$\lim_{i \rightarrow \infty} \|x_i - y_i\| = 0$$

, it holds that:

$$\lim_{i \rightarrow \infty} \|f(x_i) - f(y_i)\| = 0$$

.

For this lemma we have the following:

\Rightarrow : Suppose f is uniformly continuous and the sequences x_i, y_i satisfy the condition stated above. Pick any $\epsilon > 0$. Then by the definition of uniform continuity, there exists $\delta > 0$ such that whenever $\|x - y\| < \delta$, it holds that $\|f(x) - f(y)\| < \epsilon$. For this $\delta > 0$, by the definition of the limit of a sequence, there exists $N > 0$ such that whenever $i > N$ it holds that $\|x_i - y_i\| < \delta$. Putting these together, for all $i > N$ it also holds that $\|f(x_i) - f(y_i)\| < \epsilon$. But this means precisely that the limit of the sequence $\|f(x_i) - f(y_i)\|$ is zero.

\Leftarrow : Suppose f is not uniformly continuous. Then there exists $\epsilon > 0$ for which for each $\delta > 0$ there exist at least two x, y such that $\|x - y\| < \delta$ and $\|f(x) - f(y)\| \geq \epsilon$. Construct then two sequences $i \rightarrow x_i, i \rightarrow y_i$ for which x_i, y_i are x, y such that the previous statement holds for $\delta = \frac{1}{i}$. We can easily see that $\lim_{i \rightarrow \infty} \|x_i - y_i\| = 0$. However, for this particular ϵ , it is *always* the case that $\|f(x_i) - f(y_i)\| \geq \epsilon$, which means that the limit of $\|f(x_i) - f(y_i)\|$ *cannot* be zero, which is a contradiction. Therefore f must be uniformly continuous.

Let us now proceed to the problem itself. As mentioned, suppose that g is not a constant function. This means that there exist x_1, x_2 such that $g(x_1) \neq g(x_2)$. This can be rewritten as $g(x_1) \neq g(x_1 + \eta)$ for some $\eta \neq 0$. Consider then the sequence $i \rightarrow (\frac{x_1}{i}, i)$ and call it a_i . Consider also the sequence $i \rightarrow (\frac{x_1 + \eta}{i}, i)$ and call it b_i . Notice that $b_i - a_i = (\frac{\eta}{i}, 0)$, and consequently that $\lim_{i \rightarrow \infty} \|b_i - a_i\| = 0$. By our previous lemma, this means that it must also hold that $\lim_{i \rightarrow \infty} \|f(b_i) - f(a_i)\| = 0$. We have that:

$$\|f(b_i) - f(a_i)\| = \left| g\left(\frac{x_1 + \eta}{i}i\right) - g\left(\frac{x_1}{i}i\right) \right| = |g(x_1 + \eta) - g(x_1)|$$

Clearly, the limit of this quantity as i tends to infinity is the quantity itself, which by our hypothesis is not zero. But then this directly contradicts the fact that f is uniformly continuous, and we thus arrive at a contradiction. Therefore, g must be constant.

Exercise 33

Let $A \subset \mathbb{R}^n$ and $f : A \rightarrow \mathbb{R}^p$ be uniformly continuous.

- (i) Show that $(f(x_k))_{k \in \mathbb{N}}$ is a Cauchy sequence in \mathbb{R}^p whenever $(x_k)_{k \in \mathbb{N}}$ is a Cauchy sequence in A .
- (ii) Using (i) prove that f can be extended in a unique fashion as a continuous function to \bar{A} .
- (iii) Give an example of a set $A \subset \mathbb{R}$ and a continuous function $A \rightarrow \mathbb{R}$ that does not take Cauchy sequences in A to Cauchy sequences in \mathbb{R} .

Solution.

(i) Suppose $(x_k)_{k \in \mathbb{N}}$ is a Cauchy sequence in A . This means that for any $\epsilon > 0$, there exists $N > 0$ such that for all $n_1, n_2 > N$ it holds that $\|x_{n_1} - x_{n_2}\| < \epsilon$. Pick now any $\epsilon > 0$. By the uniform continuity of f it must hold that there exists $\delta > 0$ such that for all $x, y \in A$ such that $\|x - y\| < \delta$, it holds that $\|f(x) - f(y)\| < \epsilon$. Because (x_k) is Cauchy, for this δ we can find $N > 0$ such that for $n_1, n_2 > N$ it holds that $\|x_{n_1} - x_{n_2}\| < \delta$. Thus we can obtain the corollary of the definition of uniform continuity, which is

that for all $n_1, n_2 > N$ it also holds that $\|f(x_{n_1}) - f(x_{n_2})\| < \epsilon$. But this means precisely that the sequence $(f(x_k))$ is also Cauchy.

(ii) First of all, if A is closed then it equals its closure, and thus the question is trivial: the extension of f to \bar{A} is itself, and since it's uniformly continuous there can be no other functions that satisfy this, otherwise the limit at at least one point would not match the value of f at that point.

Therefore, from now on assume that $A \neq \bar{A}$ and recall that $\bar{A} = A \cup \partial A$, which means that there exists at least one boundary point which does not belong in A . Let b be any such point. By the definition of a boundary point, there must exist a sequence of points of A that converges to b . Indeed, if one takes a sequence of open balls of radii $i \rightarrow \frac{1}{i}$ around b , these will always intersect A . Thus we can always select a point inside each such open ball to form the i -th point, x_i , of our constructed sequence. Notice that these points form a Cauchy sequence: after the N -th point, they can be at a distance of at most $\frac{1}{N}$, which means that for a given ϵ we can always find N_ϵ after which the points of the sequence are at least ϵ -close to each other.

Now by part (i) we have that f maps (x_i) to a sequence that is also Cauchy, and by an easy extension of theorem 1.16 of Carothers we know that Cauchy sequences in \mathbb{R}^n always converge. Therefore, $(f(x_i))$ converges. This means that there is precisely one *possible* value for the extension of f at b : the limit of $(f(x_i))$. This by itself, however, does not guarantee continuity, as we have not shown that *every* sequence in \bar{A} converging to b is mapped by an appropriate extension of f to a sequence which converges to this limit. For now let us commit a small abuse of notation by writing $f(b)$ to mean this unique potential value of the extension of f at b such that it *could* be continuous.

Let now $a, b, a \neq b$ be any two points in \bar{A} . By a similar argument as above, we can always find two sequences $a_n \rightarrow a, b_n \rightarrow b$ such that for all $n, a_n, b_n \in A$. Furthermore, assuming that if any of a, b are on the boundary and not in A , $f(a), f(b)$ refer to the limit described above, and otherwise to the value of f at the corresponding point, we observe that:

$$\begin{aligned} \|f(a) - f(b)\| &= \|f(a) - f(a_n) + f(b_n) - f(b) + f(a_n) - f(b_n)\| \\ &\leq \|f(a) - f(a_n)\| + \|f(b_n) - f(b)\| + \|f(a_n) - f(b_n)\| \end{aligned}$$

We've shown above that the first and third terms of the RHS correspond to quantities that can be made arbitrarily small by picking $n > N$ appropriately. Suppose now that we pick any $\epsilon > 0$. Then the uniform continuity of f guarantees that we can find $\delta > 0$ such that $\|x - y\| < \delta \implies \|f(x) - f(y)\| < \frac{\epsilon}{3}, x, y \in A$. Pick now any two $a, b \in \bar{A}$ such that $\|a - b\| < \frac{\delta}{2}$, and select sequences $a_n \rightarrow a, b_n \rightarrow b$ as described above. The continuity of the norm implies that:

$$\lim_{n \rightarrow \infty} \|a_n - b_n\| = \|a - b\|$$

, thus there exists $N > 0$ such that for $n > N$, $\|a_n - b_n\| - \|a - b\| < \frac{\delta}{2} \implies \|a_n - b_n\| < \delta$. But then by the uniform continuity of f , $\|f(a_n) - f(b_n)\| < \frac{\epsilon}{3}$. Thus, by picking N being the maximum of N_1, N_2, N_3 for which $n > N_1, n > N_2, n > N_3$ imply that each of the three terms of the RHS respectively are smaller than $\frac{\epsilon}{3}$, we obtain that $\|f(a) - f(b)\| < \epsilon$. But then we have shown precisely that the extension is uniformly continuous on all of \bar{A} .

Now, by exercise 13, because A is dense in \bar{A} , any two $f_1, f_2 : \bar{A} \rightarrow \mathbb{R}^p$ that are continuous have to be equal in A in order to be extensions of f , and by the corollary of said exercise, $f_1 = f_2$, which means that the extension is unique.

(iii) Consider the set $A = [0, 1)$ and the function $f(x) = \frac{1}{1-x}$, which is continuous in A as a rational function. Consider the sequence $x_i = 1 - \frac{1}{i}$, which converges to 1 and, by the theorem mentioned above (from Carothers), is thus Cauchy. Observe now that $f(x_i) = \frac{1}{1-(1-\frac{1}{i})} = i$. Obviously the sequence $f(x_i)$ does not converge (it is unbounded), and by the same theorem cannot be Cauchy.

Exercise 35

Assume that $K \subset \mathbb{R}^n$ and $L \subset \mathbb{R}^p$ are compact sets, and that $f : K \times L \rightarrow \mathbb{R}$ is continuous.

(i) Show that for every $\epsilon > 0$ there exists $\delta > 0$ such that

$$a, b \in K \text{ with } \|a - b\| < \delta, y \in L \implies |f(a, y) - f(b, y)| < \epsilon$$

Next, define $m : K \rightarrow \mathbb{R}$ by $m(x) = \max\{f(x, y) \mid y \in L\}$.

(ii) Prove that for every $\epsilon > 0$ there exists $\delta > 0$ such that $m(a) > m(b) - \epsilon$ for all $a, b \in K$ with $\|a - b\| < \delta$.

(iii) Deduce from (ii) that $m : K \rightarrow \mathbb{R}$ is uniformly continuous.

(iv) Show that $\max\{m(x) \mid x \in K\} = \max\{f(x, y) \mid x \in K, y \in L\}$.

Consider a continuous function f on \mathbb{R}^n that is defined on a product of n intervals in \mathbb{R} . It is a consequence of (iv) that the problem of finding maxima for f can be reduced to the problem of successively finding maxima for functions associated with f that depend on one variable only. Under suitable conditions the latter problem might be solved by using the differential calculus in one variable.

(v) Apply this method for proving that the function

$$f : [-\frac{1}{2}, 1] \times [0, 2] \rightarrow \mathbb{R} \text{ with } f(x) = \|x\|^2 e^{-x_1 - x_2^2}$$

attains its maximum value $e^{-\frac{1}{4}}$ at $\frac{1}{2}(-1, \sqrt{3})$.

Solution.

(i) Pick any $\epsilon > 0$. Because f is continuous, we have that there exists $\delta > 0$ such that:

$$\|(a, x) - (b, y)\| < \delta, a, b \in K, x, y \in L \implies |f(a, x) - f(b, y)| < \epsilon$$

Suppose now we apply the above implication for $x = y \in L$ and for $a, b \in K$ with $\|a - b\| < \delta$. It's important to note that because the domain of f is the entire Cartesian product of K, L the points $(a, y), (b, y)$ always belong in it. Now observe that $\|(a, x) - (b, y)\| = \|(a, x) - (b, x)\| = \|(a - b, 0)\| = \|a - b\|$

, by the definition of the norm in $K \times L$ given that we use the Euclidean norm in both K and L . Thus the above premise stands for any such two $a, b \in K$ and any $y \in L$, and we obtain the corollary:

$$|f(a, x) - f(b, y)| < \epsilon \implies |f(a, y) - f(b, y)| < \epsilon$$

, thus completing the proof.

(ii) First of all, we observe that because both K and L are compact sets, m is well defined (continuous functions achieve a maximum value on compact sets). By part (i), for a given $\epsilon > 0$ we can find $\delta > 0$ such that for all $y \in L$ and for all $a, b \in K$ with $\|a - b\| < \delta$ we have that $|f(a, y) - f(b, y)| < \epsilon$. Consequently:

$$-\epsilon < f(a, y) - f(b, y) < \epsilon \implies f(b, y) - \epsilon < f(a, y)$$

Now, by the definition of m , we have that $m(a) \geq f(a, y)$ for all $y \in L$. Thus $f(b, y) - \epsilon < m(a), y \in L$. Furthermore, again by the definition of m , we have that $m(b) = f(b, y_b)$ for some $y_b \in L$ (possibly not unique). In any case though, because the previously stated inequality holds for all $y \in L$, we obtain that $f(b, y_b) - \epsilon < m(a) \implies m(b) - \epsilon < m(a)$.

(iii) Pick any $\epsilon > 0$, and select the $\delta > 0$ such that for all $a, b \in K, \|a - b\| < \delta$ the implication of (ii) holds. This yields $m(a) > m(b) - \epsilon \implies m(a) - m(b) > -\epsilon$, and by exchanging the roles of a, b also that $m(b) > m(a) - \epsilon \implies \epsilon > m(a) - m(b)$. We have thus that $-\epsilon < m(a) - m(b) < \epsilon$, or, equivalently, $|m(a) - m(b)| < \epsilon$. But this is precisely the definition of uniform continuity for m , thus completing the proof.

(iv) Because both K, L are compact, $K \times L$ is compact. Indeed, both K, L are bounded, and by the definition of the norm in the Cartesian product of K, L , $K \times L$ will be bounded as well. Furthermore, any convergent sequence in $K \times L$ will result in the "coordinate-wise" corresponding sequences converging as well. Because K, L are closed, each of these will converge to a point in K, L respectively, and thus the ordered pair of these two will belong in $K \times L$, thus showing that the convergent sequence in $K \times L$ converges to a point inside it, i.e., $K \times L$ is closed, and thus indeed compact.

Therefore, there exist $x_M \in K, y_M \in L$ such that $f(x_M, y_M)$ is the maximum value of f . By the definition of m , $m(x_M) = f(x_M, y_M)$. Furthermore, we have that for any $x \in K, m(x) = f(x, y_x)$ for some $y_x \in L$, and as such $m(x) \leq f(x_M, y_M)$. We have shown that all $m(x)$ are at most equal to the maximum value of f , and also that for some $x, m(x)$ equals this maximum value. Therefore:

$$\max\{m(x) \mid x \in K\} = f(x_M, y_M) = \max\{f(x, y) \mid x \in K, y \in L\}$$

(v) Here $K = [-\frac{1}{2}, 1]$, $L = [0, 2]$, and these are both compact sets as closed intervals. We have that $m(x_1) = \max\{f(x_1, x_2) \mid x_2 \in L\}$ can be written as the maximum value of the function $g(x_2) = (x_1^2 + x_2^2)e^{-x_1 - x_2^2}$. Compute the derivative to see that:

$$g'(x_2) = 2x_2e^{-x_1 - x_2^2} + (-2x_2)(x_1^2 + x_2^2)e^{-x_1 - x_2^2} = 2x_2e^{-x_1 - x_2^2}(1 - x_1 - x_2^2)$$

Because $x_2 \in [0, 2]$, the sign of this expression depends entirely on $(1 - x_1 - x_2^2)$. The zeros of the derivative in $(0, 2)$ are then $x_2 = \sqrt{1 - x_1^2}$ (since $1 - x_1^2$ is always non-negative for $x_1 \in [-\frac{1}{2}, 1]$). Furthermore, one can easily see that the derivative is positive for $x_2 < \sqrt{1 - x_1^2}$ and negative for $x_2 > \sqrt{1 - x_1^2}$. Therefore $\sqrt{1 - x_1^2}$ is g 's maximum, since the derivative does not change sign again in $[0, 2]$. We have thus found that:

$$m(x_1) = f(x_1, \sqrt{1 - x_1^2}) = e^{-1 - x_1 + x_1^2}$$

Now to maximize this we again compute the derivative $m'(x_1) = (-1 + 2x_1)e^{-1 - x_1 + x_1^2}$, observe that its zeros are $\frac{1}{2}$ only, but at that point its sign changes from negative to positive, meaning that it corresponds to a minimum, and thus conclude that the maximum is attained at one of the endpoints $-\frac{1}{2}, 1$. A simple calculation leads to the conclusion that this happens at $x_1 = -\frac{1}{2}$, and thus from part (iv) that the maximum of f must be at $x_1 = -\frac{1}{2}, x_2 = \sqrt{1 - x_1^2} = \frac{\sqrt{3}}{2}$, where the value is $f(\frac{1}{2}(-1, \sqrt{3})) = e^{-\frac{1}{4}}$.

Chapter 2

Differentiation

Exercise 1

We study some properties of the operation of taking the trace; this will provide some background for the Euclidean norm. All the results obtained below apply, *mutatis mutandis*, to $\text{End}(\mathbb{R}^n)$ as well.

(i) Verify that $\text{tr} : \text{Mat}(n, \mathbb{R}) \rightarrow \mathbb{R}$ is a linear mapping, and that, for $A, B \in \text{Mat}(n, \mathbb{R})$,

$$\text{tr} A = \text{tr} A^T, \text{tr} AB = \text{tr} BA, \text{tr} BAB^{-1} = \text{tr} A$$

if $B \in \text{GL}(n, \mathbb{R})$. Deduce that the mapping tr vanishes on *commutators* in $\text{Mat}(n, \mathbb{R})$, that is, for $A, B \in \text{Mat}(n, \mathbb{R})$,

$$\text{tr}([A, B]) = 0, \text{ where } [A, B] = AB - BA.$$

Define a bilinear functional $\langle \cdot, \cdot \rangle : \text{Mat}(n, \mathbb{R}) \times \text{Mat}(n, \mathbb{R}) \rightarrow \mathbb{R}$ by $\langle A, B \rangle = \text{tr}(A^T B)$.

(ii) Show that $\langle A, B \rangle = \text{tr}(B^T A) = \text{tr}(AB^T) = \text{tr}(BA^T)$.

(iii) Let a_j be the j -th column vector of $A \in \text{Mat}(n, \mathbb{R})$, thus $a_j = Ae_j$. Verify, for $A, B \in \text{Mat}(n, \mathbb{R})$:

$$A^T B = (\langle a_i, b_j \rangle)_{1 \leq i, j \leq n}, \langle A, B \rangle = \sum_{1 \leq j \leq n} \langle a_j, b_j \rangle$$

(iv) Prove that $\langle \cdot, \cdot \rangle$ defines an inner product on $\text{Mat}(n, \mathbb{R})$, and that the corresponding norm equals the Euclidean norm on $\text{Mat}(n, \mathbb{R})$.

(v) Using part (ii), show that the decomposition from Lemma 2.1.4 is an orthogonal direct sum, in other words, $\langle A, B \rangle = 0$ if $A \in \text{End}^+(\mathbb{R}^n)$ and $B \in \text{End}^-(\mathbb{R}^n)$. Verify that $\dim \text{End}^\pm(\mathbb{R}^n) = \frac{1}{2}(n \pm 1)$. Define the bilinear functional $B : \text{Mat}(n, \mathbb{R}) \times \text{Mat}(n, \mathbb{R}) \rightarrow \mathbb{R}$ by $B(A, A') = \text{tr}(AA')$.

(vi) Prove that B satisfies $B(A_\pm, A_\pm) \geq 0$, for $0 \neq A_\pm \in \text{End}^\pm(\mathbb{R}^n)$.

Finally, we show that the trace is completely determined by some of the properties listed in part (i).

(vii) Let $\tau : \text{Mat}(n, \mathbb{R}) \rightarrow \mathbb{R}$ be a linear mapping satisfying $\tau(AB) = \tau(BA)$, for all $A, B \in \text{Mat}(n, \mathbb{R})$. Prove

$$\tau = \frac{1}{n} \tau(I) \text{tr}$$

Hint: Note that τ vanishes on commutators. Let $E_{ij} \in \text{Mat}(n, \mathbb{R})$ be given by having a single 1 at position (i, j) and zeros elsewhere. Then the statement is immediate from $[E_{ij}, E_{jj}] = E_{ij}$ for $i \neq j$, and $[E_{ij}, E_{ji}] = E_{ii} - E_{jj}$.

Solution.

(i) We know the following from Linear Algebra (the proofs are simple, and based on the definition of the trace of a matrix):

- For A, B square matrices, it holds that $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$.
- For A square matrix and $\lambda \in \mathbb{R}$, it holds that $\text{tr}(\lambda A) = \lambda \text{tr}(A)$.

Therefore, tr is indeed a linear mapping. For the properties requested, we have that:

- Taking the transpose of a matrix leaves the diagonal unchanged, therefore the sum of the diagonal elements is also unchanged. This equals the trace of the transpose, which means that $\text{tr}(A^T) = \text{tr}(A)$.
- This is something that we have also already seen in Linear Algebra, and the proof once more follows from the definition. We have that $\text{tr}(AB) = \sum_i (AB)_{i,i} = \sum_i \sum_j A_{ij} B_{ji}$, $\text{tr}(BA) = \sum_i (BA)_{i,i} = \sum_i \sum_j B_{ij} A_{ji}$, and these are of course equal, since they sum the exact same products.
- Using the above property, for an invertible matrix B we obtain that $\text{tr}(BAB^{-1}) = \text{tr}((BA)B^{-1}) = \text{tr}(B^{-1}(BA)) = \text{tr}(IA) = \text{tr}(A)$.

We additionally have that $\text{tr}([A, B]) = \text{tr}(AB - BA) = \text{tr}(AB) - \text{tr}(BA) = \text{tr}(AB) - \text{tr}(AB) = 0$.

(ii) Using the definition and the properties we proved in (i), we have that:

- $\text{tr}(B^T A) = \text{tr}((B^T A)^T) = \text{tr}(A^T B) = \langle A, B \rangle$
- $\text{tr}(AB^T) = \text{tr}(B^T A) = \langle A, B \rangle$
- $\text{tr}(BA^T) = \text{tr}(A^T B) = \langle A, B \rangle$

(iii) The (i, j) -th element of $A^T B$ equals the following (by the definition of matrix multiplication and the transpose):

$$(A^T B)_{i,j} = \sum_k A_{i,k}^T B_{k,j} = \sum_k A_{k,i} B_{k,j}$$

Each term of the last sum is nothing but the standard Euclidean inner product of a_i, b_j . Therefore, $(A^T B)_{i,j} = \langle a_i, b_j \rangle$.

We now also have that $\langle A, B \rangle = \text{tr}(A^T B) = \sum_{1 \leq j \leq n} (A^T B)_{j,j} = \sum_{1 \leq j \leq n} \langle a_j, b_j \rangle$.

(iv) To show that $\langle \cdot, \cdot \rangle$ defines an inner product for square matrices, we have that (by using results previously proved in the exercise):

- **Additivity in the first slot:** We have that $\langle A_1 + A_2, B \rangle = \text{tr}((A_1 + A_2)^T B) = \text{tr}(A_1^T + A_2^T)B = \text{tr}(A_1^T B + A_2^T B) = \text{tr}(A_1^T B) + \text{tr}(A_2^T B) = \langle A_1, B \rangle + \langle A_2, B \rangle$.
- **Homogeneity in the first slot:** We have that $\langle \lambda A, B \rangle = \text{tr}((\lambda A)^T B) = \text{tr}(\lambda A^T B) = \text{tr}(\lambda(A^T B)) = \lambda \text{tr}(A^T B)$.
- **Symmetry:** Note that in \mathbb{R}^n we care about symmetry instead of conjugate symmetry, thus: $\langle B, A \rangle = \text{tr}(B^T A) = \langle A, B \rangle$ by part (ii).
- **Positivity:** We have that $\langle A, A \rangle = \sum_{1 \leq j \leq n} \langle a_j, a_j \rangle \geq 0$, by the positivity of the Euclidean inner product for vectors.
- **Definiteness:** By the immediately preceding bullet, we have that $\langle A, A \rangle$ is a sum of non-negative quantities, each of which is zero iff the j -th column of A is zero (definiteness of the Euclidean inner product). This leads us to the conclusion that $\langle A, A \rangle$ is zero iff all of its columns are zero, which means that A is the zero matrix, which is the zero element for matrices, thus proving definiteness of $\langle \cdot, \cdot \rangle$.

Therefore, $\langle \cdot, \cdot \rangle$ does indeed define an inner product for square matrices. The associated norm is $\|A\| = \sqrt{\langle A, A \rangle} = \sqrt{\sum_{1 \leq j \leq n} \langle a_j, a_j \rangle} = \sqrt{\sum_{i,j} A_{i,j} A_{i,j}}$ which indeed equals the Euclidean norm for the matrix A .

(v) Lemma 2.1.4 states that the vector space $\text{End}(\mathbb{R}^n)$ of operators in \mathbb{R}^n can be written as a direct sum of the subspace $\text{End}^+(\mathbb{R}^n)$ of self-adjoint operators and the subspace $\text{End}^-(\mathbb{R}^n)$ of anti-adjoint operators (that is, operators for which $A^* = -A$). In particular:

$$\text{End}(\mathbb{R}^n) = \text{End}^+(\mathbb{R}^n) \oplus \text{End}^-(\mathbb{R}^n)$$

, since any operator T can be written as $T = \frac{1}{2}(T + T^*) + \frac{1}{2}(T - T^*)$.

Consider then $A \in \text{End}^+(\mathbb{R}^n), B \in \text{End}^-(\mathbb{R}^n)$. Using the defining properties of these operators and part (ii), we obtain that:

$$\langle A, B \rangle = \text{tr}(B^*A) = \text{tr}(-BA) = -\text{tr}(BA) = -\text{tr}(AB) = -\text{tr}(A^*B) = -\langle A, B \rangle$$

Clearly, this implies that $\langle A, B \rangle = 0$. Now let x, y be the dimensions of $\text{End}^+(\mathbb{R}^n), \text{End}^-(\mathbb{R}^n)$ respectively. Due to the direct sum and the fact that $\dim \text{End}(\mathbb{R}^n) = n^2$, we have that $x + y = n^2$. Recall also the isomorphism between operators and square matrices, and the fact that in \mathbb{R}^n self-adjoint operators are precisely those whose matrix wrt. the standard basis equals its transpose. Such a matrix is fully defined by choosing the elements on and above the diagonal. In total, these are $n + \frac{n^2-n}{2} = \frac{n(n+1)}{2}$.

In particular, such a matrix can be written uniquely as a sum of a linearly independent list of matrices consisting of either a one on the diagonal and zeros everywhere else, or a one in positions $(i, j), (j, i), i \neq j$, and zeros everywhere else. These are again in total $\frac{n(n+1)}{2}$. We conclude that the dimension x equals $\frac{n(n+1)}{2}$, and thus that $y = n^2 - \frac{n(n+1)}{2} = \frac{n(n-1)}{2}$.

(vi) Consider first a self-adjoint A . Then, we have that:

$$B(A, A) = \text{tr}(AA) = \text{tr}(A^*A) = \langle A, A \rangle$$

This is clearly greater than zero if A is not zero, by the properties of the inner product. For an anti-adjoint A we have that:

$$B(A, A) = \text{tr}(AA) = \text{tr}((-A^*)A) = -\text{tr}(A^*A) = -\langle A, A \rangle$$

By the exact same reasoning as before, for a non-zero A this is clearly less than zero.

(vii) Since τ is a linear mapping, it is uniquely determined by its values on a basis of $\text{Mat}(n, \mathbb{R})$. Observe that the E_{ij} given in the hint form such a basis. Consider the value of τ on the commutator $[A, B] = AB - BA$. By its linearity and its given defining property, it must hold that:

$$\tau([A, B]) = \tau(AB - BA) = \tau(AB) - \tau(BA) = \tau(AB) - \tau(AB) = 0$$

For $i \neq j$, it holds that $[E_{ij}, E_{jj}] = E_{ij}E_{jj} - E_{jj}E_{ij} = E_{ij} - 0 = E_{ij}$ (since E_{ij} maps everything to zero except for mapping e_i to e_j and E_{jj} maps everything to zero except for mapping e_j to e_j . By the above observation regarding commutators, we conclude that:

$$\tau(E_{ij}) = 0$$

For $i = j$, it holds that $[E_{ii}, E_{jj}] = E_{ii} - E_{jj}$. Again, we conclude that:

$$\tau(E_{ii} - E_{jj}) = 0 \implies \tau(E_{ii}) = \tau(E_{jj})$$

This means that τ maps every vector $E_{ij}, i \neq j$ to zero, and maps *all* vectors E_{ii} to the same value, let it be called a . By its linearity, observe that:

$$\tau(I) = \tau(E_{11} + E_{22} + \dots + E_{nn}) = na \implies a = \frac{1}{n}\tau(I)$$

Finally, observe that the trace is also a linear function, maps $E_{ij}, i \neq j$ to zero and maps each E_{ii} to 1 (since the trace equals the sum of the diagonal). But then it becomes clear that for all vectors E_{ij} it holds that

$$\tau(E_{ij}) = a \text{tr}(E_{ij}) = \frac{1}{n}\tau(I)\text{tr}(E_{ij})$$

, and since these are linear mappings, this also holds for any matrix, i.e. $\tau = \frac{1}{n}\tau(I)\text{tr}$.

Exercise 4

Suppose $A \in \text{End}(\mathbb{R}^n)$ is an orthogonal transformation, that is, $\|Ax\| = \|x\|$, for all $x \in \mathbb{R}^n$.

(i) Show that $A \in \text{Aut}(\mathbb{R}^n)$ and that A^{-1} is orthogonal.

(ii) Deduce from the polarization identity in Lemma 1.1.5.(iii) that

$$\langle A^*Ax, y \rangle = \langle Ax, Ay \rangle = \langle x, y \rangle, x, y \in \mathbb{R}^n$$

(iii) Prove $A^*A = I$ and deduce that $\det A = \pm 1$. Furthermore, using (i) show that A^* is orthogonal, thus $AA^* = I$; and also obtain $\langle Ae_i, Ae_j \rangle = \langle A^*e_i, A^*e_j \rangle = \langle e_i, e_j \rangle = \delta_{i,j}$, where e_1, \dots, e_n is the standard basis for \mathbb{R}^n . Conclude that the column and row vectors, respectively, in the corresponding matrix of A form an orthonormal basis for \mathbb{R}^n .

(iv) Deduce from (iii) that the corresponding matrix $A = (a_{ij}) \in \mathbf{GL}(n, \mathbb{R})$ is *orthogonal* with coefficients in \mathbb{R} and belongs to the *orthogonal group* $\mathbf{O}(n, \mathbb{R})$, which means that it satisfies $A^T A = I$; in addition, deduce

$$\sum_{1 \leq k \leq n} a_{ki} a_{kj} = \sum_{1 \leq k \leq n} a_{ik} a_{jk} = \delta_{ij}, 1 \leq i, j \leq n$$

Solution.

(i) We need to show that A is a bijection. Since A is an operator, it suffices to show that it is injective. Consider then an x such that $Ax = 0$. We then have that $\|Ax\| = 0$, which, since A is an isometry, implies $\|x\| = 0$, which means $x = 0$, therefore A is injective and thus bijective. Therefore, A^{-1} is well defined in \mathbb{R}^n .

Now consider any $y \in \mathbb{R}^n$. It holds that $A^{-1}y = x$ such that $Ax = y$. Since A is an isometry, we have that:

$$\|Ax\| = \|x\| \implies \|y\| = \|A^{-1}y\|$$

This means that A^{-1} is also an isometry (in the terminology of this book, an orthogonal transformation).

(ii) From the polarization identity we obtain that:

$$\begin{aligned} \langle Ax, Ay \rangle &= \frac{1}{4}(\|Ax + Ay\|^2 - \|Ax - Ay\|^2) = \frac{1}{4}(\|A(x + y)\|^2 - \|A(x - y)\|^2) \\ &\implies \langle Ax, Ay \rangle = \frac{1}{4}(\|x + y\|^2 - \|x - y\|^2) = \langle x, y \rangle \end{aligned}$$

By the definition of the adjoint, we also have that:

$$\langle A^*Ax, y \rangle = \langle Ax, (A^*)^*y \rangle = \langle Ax, Ay \rangle$$

(iii) Let x be any element of \mathbb{R}^n . Using the properties proved in (ii), we have that:

$$\begin{aligned} \langle A^*Ax - x, A^*Ax - x \rangle &= \langle A^*Ax, A^*Ax \rangle + \langle x, x \rangle - 2\langle A^*Ax, x \rangle = \\ &\langle x, A^*Ax \rangle + \langle x, x \rangle - 2\langle x, x \rangle = \langle x, x \rangle + \langle x, x \rangle - 2\langle x, x \rangle = 0 \end{aligned}$$

By the properties of the inner product, the only way this can be true is if $A^*Ax - x = 0 \implies A^*Ax = x$. Since x was selected arbitrarily, this means that $A^*A = I$. Because in \mathbb{R} A and A^* have the same eigenvalues with the same multiplicities, and thus equal determinants, this also implies that:

$$\det(A^*A) = \det(I) \implies \det(A^*)\det(A) = 1 \implies \det(A)^2 = 1 \implies \det(A) = \pm 1$$

Since $A^*A = I$, A^* is a left inverse of A , and from Linear Algebra we know that this means that A^* is the inverse of A . Therefore, $A^* = A^{-1}$, and from (i) we know that A^{-1} is an isometry/orthogonal.

For any two e_i, e_j of the standard basis of \mathbb{R}^n , we have that:

$$\langle Ae_i, Ae_j \rangle = \langle e_i, e_j \rangle = \delta_{ij}$$

, since the standard basis is orthonormal (and thus $\langle e_i, e_j \rangle = 0, i \neq j, \langle e_i, e_j \rangle = 1, i = j$). Also:

$$\langle A^*e_i, A^*e_j \rangle = \langle e_i, (A^*)^*A^*e_j \rangle = \langle e_i, AA^*e_j \rangle = \langle e_i, e_j \rangle = \delta_{ij}$$

From the first of those equations, we deduce that the columns of A are orthonormal vectors. Furthermore, since they are n in total, they must form a basis for \mathbb{R}^n . With respect to the standard basis, we know that the matrix of A^* equals the transpose of the matrix of A . From the second of those equations, we then have that the columns of this matrix are also orthonormal and n in total. Thus, the columns of the transpose of $\mathcal{M}(A)$, i.e. the rows of $\mathcal{M}(A)$, also form an orthonormal basis for \mathbb{R}^n .

(iv) This is immediately obvious from the definitions of the general linear group and orthogonal group of matrices. The given sum $\sum_{1 \leq k \leq n} a_{ki} a_{kj}$ is nothing but the Euclidean inner product of the i -th and j -th columns of $\mathcal{M}(A)$, which from (iii) equals δ_{ij} . Similarly, the given sum $\sum_{1 \leq k \leq n} a_{ik} a_{jk}$ is nothing but the Euclidean inner product of the i -th and j -th rows of $\mathcal{M}(A)$, which again from (iii) equals δ_{ij} .

Exercise 5

Let $\mathbf{SO}(3, \mathbb{R})$ be the special orthogonal group in \mathbb{R}^3 , consisting of the orthogonal matrices in $\text{Mat}(3, \mathbb{R})$ with determinant equal to 1. One then has, for every $R \in \mathbf{SO}(3, \mathbb{R})$, see Exercise 2.4,

$$R \in \mathbf{GL}(3, \mathbb{R}), R^T R = I, \det R = 1$$

Prove that for every $R \in \mathbf{SO}(3, \mathbb{R})$ there exist $a \in \mathbb{R}, 0 \leq a \leq \pi$ and $v \in \mathbb{R}^3$ with $\|v\| = 1$ with the following properties: R fixes v and maps the linear subspace N_v orthogonal to v into itself; in N_v the action of R is that of rotation by the angle a such that, for $0 < a < \pi$ and for all $y \in N_v$, one has $\det(v \ y \ Ry) > 0$.

We write $R = R_{a,v}$, which is referred to as the counterclockwise rotation in \mathbb{R}^3 by the angle a around the axis of rotation $\mathbb{R}v$.

Solution.

We begin by observing that since R is an orthogonal matrix, the corresponding linear transformation T must be an isometry. Indeed, for any $v \in \mathbb{R}^3$ we have that:

$$\|Tv\|^2 = \langle Tv, Tv \rangle = \langle v, T^* Tv \rangle = \langle v, R^T R v \rangle = \langle v, I v \rangle = \|v\|^2$$

Since \mathbb{R}^3 has odd dimension, T must have a real eigenvalue λ . Furthermore, because $\|Tv\| = \|v\|$, λ can only be 1 or -1. We claim that it cannot be the case that the only real eigenvalue is -1, because then the determinant of T could not be 1. We prove this by contradiction. Recall that the determinant of T equals the product of the eigenvalues of T each raised to the corresponding multiplicity. Besides λ , T can only have at most two more eigenvalues, λ_1, λ_2 . If this is the case, because T is the complexification of T , it must hold that $\lambda_2 = \overline{\lambda_1}$. But then if $\lambda = -1, \lambda \cdot \lambda_1 \cdot \lambda_2 = -1 \cdot |\lambda_1|^2$, and this cannot be equal to 1. The only other possibility is for λ_1, λ_2 to both be real. But then they would have to be either 1 or -1. If $\lambda = \lambda_1 = \lambda_2 = -1$ their product is -1, contradiction.

Therefore, 1 is an eigenvalue of T and there must exist a non-zero corresponding eigenvector v with norm 1. Since T is an isometry, it is normal. Again from Linear Algebra, we know then that the orthogonal complement of v , N_v is invariant under T . Note also that N_v must have dimension 2, since $\mathbb{R}^3 = \text{span}\{v\} \oplus N_v$. Since T is an isometry, it must be the case that $T|_{N_v}$ is also an isometry. In a two-dimensional real vector space, we know that we can always find an orthonormal basis with respect to which the matrix of the isometry is either of the form:

$$\begin{pmatrix} \cos a & -\sin a \\ \sin a & \cos a \end{pmatrix}$$

, with $a \in [0, \pi]$, or of the form $\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$ or $\begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$. In our case, these last two can be ruled out because they would lead to T having a determinant of -1. Therefore, since N_v has dimension 2 we can choose an orthonormal basis v_1, v_2 for N_v with respect to which $\mathcal{M}(T)$ has the first of the three forms described above. Putting this together with the facts that $\mathbb{R}^3 = \text{span}\{v\} \oplus N_v$, $Tv = v$ and N_v is the orthogonal complement of $\text{span}\{v\}$ we obtain that the matrix of T with respect to the orthonormal basis v, v_1, v_2 must be of the form:

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos a & -\sin a \\ 0 & \sin a & \cos a \end{pmatrix}$$

It becomes clear that for any vector in N_v , T rotates the vector by the angle a , and that this is counterclockwise, by the geometric definition of measuring angles between vectors. Also, from the definition of invariant subspaces, T maps N_v into itself. Therefore, we only need to show that $\det(v \ y \ Ry)$ for all $y \in N_v$. We have that (where everything is with respect to our basis v, v_1, v_2):

$$(v \ y \ Ry) = \begin{pmatrix} 1 & y_1 & y_1 \\ 0 & y_2 & \cos a y_2 - \sin a y_3 \\ 0 & y_3 & \sin a y_2 + \cos a y_3 \end{pmatrix}$$

Observe that for this matrix, the terms of the determinant sum that take the first column element from the second or third row are zero. Therefore, the determinant is:

$$\begin{aligned} \det(v \ y \ Ry) &= \text{sign}(1, 2, 3)1 \cdot y_2 \cdot (\sin a y_2 + \cos a y_3) + \text{sign}(1, 3, 2)1 \cdot y_3 \cdot (\cos a y_2 - \sin a y_3) \\ &= y_2^2 \sin a + y_2 y_3 \cos a - y_2 y_3 \cos a + y_3^2 \sin a = \sin a (y_2^2 + y_3^2) \end{aligned}$$

Because $a \in [0, \pi]$, this determinant is always non-negative, and is zero whenever $y_2, y_3 = 0$, which means that y is a multiple of v . Since $y \in N_v$ this can only happen if $y = 0$ (which also makes intuitive sense, because then $Ry = 0$).

Exercise 6

Assume a differentiable function $f : \mathbb{R} \rightarrow \mathbb{R}$ has a zero, so $f(x) = 0$ for some $x \in \mathbb{R}$. Suppose $x_0 \in \mathbb{R}$ is a first approximation to this zero and consider the tangent line to the graph $\{(x, f(x)) | x \in \text{dom}(F)\}$ of f at $(x_0, f(x_0))$; this is the set

$$\{(x, y) \in \mathbb{R}^2 | y - f(x_0) = f'(x_0)(x - x_0)\}$$

Then determine the intercept x_1 of that line with the x -axis, in other words $x_1 \in \mathbb{R}$ for which $-f(x_0) = f'(x_0)(x_1 - x_0)$, or, under the assumption that $A^{-1} = f'(x_0) \neq 0$,

$$x_1 = F(x_0) \text{ where } F(x) = x - Af(x)$$

It seems plausible that x_1 is nearer the required zero x than x_0 , and that iteration of this procedure will get us nearer still. In other words, we hope that the sequence $(x_k)_{k \in \mathbb{N}_0}$ with $x_{k+1} = F(x_k)$ converges to the required zero x , for which then $F(x) = x$. Now we formalize this heuristic argument.

Let $x_0 \in \mathbb{R}^n$, $\delta > 0$, and let $V = V(x_0; \delta)$ be the closed ball in \mathbb{R}^n of center x_0 and radius δ ; furthermore, let $f : V \rightarrow \mathbb{R}^n$. Assume $A \in \text{Aut}(\mathbb{R}^n)$ and a number ϵ with $0 \leq \epsilon < 1$ to exist such that:

- (i) the mapping $F : V \rightarrow \mathbb{R}^n$ with $F(x) = x - A(f(x))$ is a contraction with contraction factor $\leq \epsilon$;
- (ii) $\|Af(x_0)\| \leq (1 - \epsilon)\delta$

Prove that there exists a unique $x \in V$ with

$$f(x) = 0; \text{ and also } \|x - x_0\| \leq \frac{1}{1 - \epsilon} \|Af(x_0)\|$$

Solution.

Since V is closed and F is a contraction, we would like to use the contraction lemma. However, this requires F to map V to itself. For this we should therefore show that for any $x \in V$ it also holds that $F(x) \in V$, which is equivalent to showing that $\|F(x) - x_0\| \leq \delta$. We know that F is a contraction, which means that:

$$\begin{aligned} \|F(x) - F(x_0)\| &\leq \epsilon \|x - x_0\| \implies \|F(x) - x_0 + A(f(x_0))\| \leq \epsilon \|x - x_0\| \\ \implies \left| \|F(x) - x_0\| - \|A(f(x_0))\| \right| &\leq \epsilon \|x - x_0\| \implies \|F(x) - x_0\| - \|A(f(x_0))\| \leq \epsilon \|x - x_0\| \\ \implies \|F(x) - x_0\| &\leq \epsilon \delta + (1 - \epsilon)\delta \implies \|F(x) - x_0\| \leq \delta, \end{aligned}$$

which is precisely the condition we require for the contraction lemma. By applying it we have that there exists a unique $x \in V$ such that:

$$F(x) = x \implies x - A(f(x)) = x \implies A(f(x)) = 0$$

Because A is invertible, the only way this can hold is if $f(x) = 0$. Furthermore, again from the contraction lemma we know that:

$$\|x - x_0\| \leq \frac{1}{1 - \epsilon} \|F(x_0) - x_0\| \implies \|x - x_0\| \leq \frac{1}{1 - \epsilon} \|Af(x_0)\|$$

Exercise 8

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable.

(i) Prove that f is constant if $Df = 0$.

Hint: Recall that the result is known if $n = 1$ and use directional derivatives.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^p$ be differentiable.

(ii) Prove that f is constant if $Df = 0$.

(iii) Let $L \in L(\mathbb{R}^n, \mathbb{R}^p)$, and suppose $Df(x) = L$, for every $x \in \mathbb{R}^n$. Prove the existence of $c \in \mathbb{R}^p$ with $f(x) = Lx + c$, for every $x \in \mathbb{R}^n$.

Solution.

(i) Suppose f is not constant. Then there exist $x, y, x \neq y$ such that $f(x) \neq f(y)$. This means also that $y - x \neq 0$. Then let $g : \mathbb{R} \rightarrow \mathbb{R}$ be $g(t) = f(t(y - x) + x)$, in which case $g(0) = f(x), g(1) = f(y)$. Because f is differentiable, so is g (composition of differentiable functions). Then, apply the Mean Value Theorem in the interval $[0, 1]$ to obtain that there exists $c \in (0, 1)$ such that:

$$g'(c) = \frac{g(1) - g(0)}{1 - 0} = f(y) - f(x) \neq 0$$

But then observe that:

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

, which equals precisely the directional derivative of f along the vector $y - x$ at $c(y - x) + x$. Because this equals $[Df(c(y - x) + x)](y - x)$, and $Df = 0$, it must be zero, which is a contradiction.

(ii) Again, suppose f is not constant. Then there exist $x = (x_1, \dots, x_n), y = (y_1, \dots, y_n), x \neq y$ such that $f(x) = (z_1, \dots, z_p) \neq (w_1, \dots, w_p) = f(y)$. As a consequence, $z_j \neq w_j$ for at least one j . Then consider the function $g : \mathbb{R}^n \rightarrow \mathbb{R}, g(x_1, \dots, x_n) = f_j(x_1, \dots, x_n)$. It is the case that:

$$[Dg](x) = [Df](x)_{j, \cdot}$$

, that is, the derivative of g at any x equals the j -th row of the derivative of f at x . Therefore, D_g is also zero (a row of zeros). But then by part (i), g is constant, and thus $g(z) = f_j(z)$ must be equal to $g(w) = f_j(w)$, a contradiction. Therefore f is constant.

(iii) Consider the function $g : \mathbb{R}^n \rightarrow \mathbb{R}^p, g(x) = f(x) - Lx$. At any given x , we know that the derivative of Lx is L (clearly, the “best linear approximation” of a linear approximation is itself). Furthermore, we are given that $Df = L$. By the addition rule for derivatives, we have firstly that g is differentiable and secondly that:

$$[Dg](x) = [Df](x) - L = L - L = 0$$

But then by part (ii), g is constant, therefore there exists $c \in \mathbb{R}^p$ such that $g(x) = c$ for all $x \in \mathbb{R}^n$. This means of course that $f(x) - Lx = c \implies f(x) = Lx + c$.

Exercise 9

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be homogeneous of degree 1, in the sense that $f(tx) = tf(x)$ for all $x \in \mathbb{R}^n$ and $t \in \mathbb{R}$. Show that f has directional derivatives at 0 in all directions. Prove that f is differentiable at 0 if and only if f is linear, which is the case if and only if f is additive.

Solution.

Observe first that $f(\mathbf{0}) = f(0 \cdot \mathbf{0}) = 0f(\mathbf{0}) = 0$. Now, pick any $v \in \mathbb{R}^n$ and examine the definition of the directional derivative of f along v at 0:

$$\lim_{h \rightarrow 0} \frac{f(0 + hv) - f(0)}{h} = \lim_{h \rightarrow 0} \frac{f(hv) - 0}{h} = \lim_{h \rightarrow 0} \frac{hf(v)}{h} = f(v)$$

, where we used the fact that f is homogeneous. We observe then that this limit always evaluates to some number, namely, to the corresponding $f(v)$. This means that all directional derivatives of f at 0 exist. Now we examine the second claim of the exercise. We observe first that if f is linear, it is clearly additive, whereas if it is additive, since we are given that it is homogeneous, it is then also linear. This means that it suffices to prove that “ f is differentiable at 0 if and only if it is additive” We have that:

\Rightarrow : Suppose first that f is differentiable at 0. Then $[Df(0)]$ is a unique linear transformation, with the property that the directional derivative of f along any v at 0 equals $[Df(0)](v)$. Pick any two $a, b \in \mathbb{R}^n$. We then have that:

$$[Df(0)](a + b) = f(a + b)$$

, since the directional derivative of f at 0 along any v equals $f(v)$. Additionally, $[Df(0)]$ is linear, and therefore:

$$[Df(0)](a + b) = [Df(0)](a) + [Df(0)](b) = f(a) + f(b)$$

, again by using the same property for the directional derivatives of f at 0. But then we have shown that $f(a + b) = f(a) + f(b)$, which means precisely that f is additive.

\Leftarrow : Now suppose f is additive. We know from the first part of the exercise that all directional derivatives of f at 0 exist. If e_1, \dots, e_n are the vectors of the standard basis of \mathbb{R}^n , then we know that the corresponding directional derivatives of f at 0 are equal to $f(e_1), \dots, f(e_n)$. Let then L be the linear transformation such that $L(e_i) = f(e_i)$ (linear transformations are completely defined by their values on a basis). Consider any vector $h = \sum_{i=1}^n a_i e_i$ which tends to zero. Then:

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{1}{\|h\|} (f(0 + h) - f(0) - L(h)) &= \lim_{h \rightarrow 0} \frac{1}{\|h\|} (f(\sum_{i=1}^n a_i e_i) - f(0) - L(\sum_{i=1}^n a_i e_i)) \\ &= \lim_{h \rightarrow 0} \frac{1}{\|h\|} (\sum_{i=1}^n a_i f(e_i) - 0 - \sum_{i=1}^n a_i L(e_i)) = 0 \end{aligned}$$

, where we used the linearity of both f and L . This limit being equal to 0 means precisely that f is differentiable at 0.

Exercise 10

Let g be as in Example 1.3.11, that is:

$$g(x) = \begin{cases} \frac{x_1 x_2^2}{x_1^2 + x_2^4}, & x \neq 0 \\ 0, & x = 0 \end{cases}$$

Then we define $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$f(x) = x_2 g(x) = \begin{cases} \frac{x_1 x_2^3}{x_1^2 + x_2^4}, & x \neq 0 \\ 0, & x = 0 \end{cases}$$

Show that $D_v f(0) = 0$ for all $v \in \mathbb{R}^2$; and deduce that $Df(0) = 0$, if f were differentiable at 0. Prove that $D_v f(x)$ is well-defined for all $v, x \in \mathbb{R}^2$, and that $v \mapsto D_v f(x)$ belongs to $\mathcal{L}(\mathbb{R}^2, \mathbb{R})$, for all $x \in \mathbb{R}^2$. Nevertheless, verify that f is not differentiable at 0 by showing

$$\lim_{x_2 \rightarrow 0} \frac{|f(x_2^2, x_2)|}{\|(x_2^2, x_2)\|} = \frac{1}{2}$$

Solution.

Pick any $v \in \mathbb{R}^2, v \neq 0$ and examine the definition of the directional derivative of f along v at 0:

$$\lim_{h \rightarrow 0} \frac{f(0 + hv) - f(0)}{h} = \lim_{h \rightarrow 0} \frac{\frac{(hv_1)(hv_2)^3}{(hv_1)^2 + (hv_2)^4} - 0}{h} = \lim_{h \rightarrow 0} \frac{h^4 v_1 v_2^3}{h^3 v_1^2 + h^4 v_2^4} = \lim_{h \rightarrow 0} \frac{h v_1 v_2^3}{v_1^2 + h v_2^4}$$

If $v_1 = 0$, it has to be the case that $v_2 \neq 0$, in which case the limit clearly evaluates to 0. If $v_2 = 0$, it has to be the case that $v_1 \neq 0$, in which case again the limit evaluates to 0. If both $v_1, v_2 \neq 0$, then the denominator tends to v_1^2 and the numerator tends to 0, which again means that the limit evaluates to 0. Therefore, the limit can indeed be shown to equal 0 for all $v \in \mathbb{R}^2 \setminus 0$. Thus, if f were differentiable at 0, it would have to be the case that $Df(0) = 0$, since it would have to equal a linear transformation which maps all vectors to zero.

We have already shown that $D_v f(x)$ is zero whenever $x = 0$. For $x \neq 0$, observe that f is defined as a rational function with a non-zero numerator. This is then clearly a differentiable function, which means $Df(x)$ exists, and then we know that the directional derivative along v at x is nothing but $[Df(x)](v)$, so obviously the function $v \mapsto D_v f(x)$ is linear.

However, we will now show that the zero matrix/function does not fulfill the definition of being the derivative of f at 0, that is, we will show that:

$$\lim_{h \rightarrow 0} \frac{1}{\|h\|} ((f(0 + h) - f(0)) - Lh) \neq 0$$

Consider approaching 0 with a sequence of points $i \rightarrow (\frac{1}{i^2}, \frac{1}{i}), i > 0$. This corresponds to the limit:

$$\lim_{x_2 \rightarrow 0} \frac{1}{\|(x_2^2, x_2)\|} (f(x_2^2, x_2) - 0 - 0h) = \lim_{x_2 \rightarrow 0} \frac{1}{\sqrt{x_2^4 + x_2^2}} \left(\frac{x_2^2 x_2^3}{x_2^4 + x_2^4} \right) = \lim_{x_2 \rightarrow 0} \frac{x_2}{2x_2 \sqrt{1 + x_2^2}} = \frac{1}{2}$$

Since this is not zero, the limit as $h \rightarrow 0$ cannot be zero either, thus showing that f is not differentiable at 0 (since we showed that if it were, it would have to hold that $Df(0) = 0$).

Exercise 11

Define $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$f(x) = \begin{cases} \frac{x_1^3}{\|x\|^2}, & x \neq 0 \\ 0, & x = 0 \end{cases}$$

Prove that f is continuous on \mathbb{R}^2 . Show that directional derivatives of f at 0 in all directions do exist. Verify, however, that f is not differentiable at 0. Compute

$$D_1 f(x) = 1 + x_2^2 \frac{x_1^2 - x_2^2}{\|x\|^4}, D_2 f(x) = -\frac{2x_1^3 x_2}{\|x\|^4} \quad (x \in \mathbb{R}^2 \setminus \{0\})$$

In particular, both partial derivatives of f are well-defined on all of \mathbb{R}^2 . Deduce that at least one of these partial derivatives has to be discontinuous at 0. To see this explicitly, note that

$$D_1 f(tx) = D_1 f(x), D_2 f(tx) = D_2 f(x) \quad (t \in \mathbb{R} \setminus \{0\}, x \in \mathbb{R}^2)$$

This means that the partial derivatives assume constant values along every line through the origin under omission of the origin. More precisely, in every neighborhood of 0 the function $D_1 f$ assumes every value in $[0, \frac{9}{8}]$ while $D_2 f$ assumes every value in $[-\frac{3}{8}\sqrt{3}, \frac{3}{8}\sqrt{3}]$. Accordingly, in this case both partial derivatives are discontinuous at 0.

Solution.

For $x \neq 0$, f is continuous as a fraction of two continuous functions. In order for f to be continuous at 0 it has to be the case that $\lim_{x \rightarrow 0} f(x) = f(0) = 0$. We have that:

$$\left| \frac{x_1^3}{\|x\|^2} \right| = \left| \frac{x_1^3}{x_1^2 + x_2^2} \right| = \left| \frac{x_1}{1 + (\frac{x_2}{x_1})^2} \right| \leq |x_1|$$

, since we are dividing $|x_1|$ with a quantity that is always at least 1. Now, as $x \rightarrow 0$, $|x_1| \rightarrow 0$, which means that the LHS is upper-bounded by a quantity that tends to zero, which means that it must itself tend to zero. This means precisely that $\lim_{x \rightarrow 0} f(x) = f(0) = 0$. Therefore f is continuous on \mathbb{R}^2 .

Now, for computing the directional derivative of f at 0 along v we set

$$g(h) = f(0 + hv) = f(hv) = \frac{h^3 v_1^3}{h^2 v_1^2 + h^2 v_2^2} = \frac{h v_1^3}{v_1^2 + v_2^2}$$

and we observe that $D_v f(0) = g'(0) = \frac{v_1^3}{v_1^2 + v_2^2}$. If f were differentiable at 0 it would have to be the case that $Df(0) = \begin{pmatrix} 1 & 0 \end{pmatrix}$, which is obtained by substituting $v = (1, 0), v = (0, 1)$ in the formula above. It would also have to hold that:

$$\begin{aligned} \lim_{h \rightarrow 0} \frac{1}{\|h\|} (f(0 + h) - f(0) - Df(0)h) &= 0 \implies \lim_{h \rightarrow 0} \frac{1}{\|h\|} \left(\frac{h_1^3}{h_1^2 + h_2^2} - 0 - \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} \right) = 0 \\ \implies \lim_{h \rightarrow 0} \frac{1}{\|h\|} \left(\frac{h_1^3}{h_1^2 + h_2^2} - h_1 \right) &= 0 \implies \lim_{h \rightarrow 0} \frac{1}{\|h\|} \left(\frac{h_1^3 - h_1 h_2^2}{h_1^2 + h_2^2} \right) = 0 \end{aligned}$$

Consider what happens when h approaches 0 with $h_2 = 0$. Then the quantity inside the limit evaluates to:

$$\frac{h_1^3 - 0}{|h_1| \cdot h_1^2} = \frac{h_1}{|h_1|},$$

which tends to 1 for $h \rightarrow 0^+$ and to -1 for $h \rightarrow 0^-$. This shows that the limit above does not exist, and therefore f cannot be differentiable at 0.

Now let us compute $D_1 f(x), D_2 f(x)$ for $x \neq 0$. By standard differentiation rules:

$$D_1 f(x) = \frac{\partial}{\partial x_1} \left(\frac{x_1^3}{x_1^2 + x_2^2} \right) = \frac{3x_1^2(x_1^2 + x_2^2) - 2x_1^4}{(x_1^2 + x_2^2)^2} = \frac{x_1^4 + 3x_1^2 x_2^2}{\|x\|^4} = \frac{x_1^4 + \|x\|^4 - \|x\|^4 + 3x_1^2 x_2^2}{\|x\|^4}$$

$$= 1 + \frac{x_1^4 - (x_1^2 + x_2^2)^2 + 3x_1^2x_2^2}{\|x\|^4} = 1 + \frac{-2x_1^2x_2^2 - x_2^4 + 3x_1^2x_2^2}{\|x\|^4} = 1 + \frac{x_2^2(x_1^2 - x_2^2)}{\|x\|^4}$$

One can similarly obtain that $D_2f(x) = -\frac{2x_1^3x_2}{\|x\|^4}$ (calculations omitted for brevity).

For both of these formulas one can notice that if we substitute (x_1, x_2) by (tx_1, tx_2) , the powers are arranged in such a way that all t 's cancel out, the result being that $D_1f(tx) = D_1f(x)$, $D_2f(tx) = D_2f(x)$.

Now consider a neighborhood around 0 containing an open ball $B_\epsilon(0)$. Observe that if we select a point v inside this ball such that $v = (\sqrt{3}x, x)$ we have that:

$$D_1f(v) = 1 + v^2 \frac{3v^2 - v^2}{(3v^2 + v^2)^2} = 1 + \frac{2v^4}{16v^4} = \frac{9}{8}$$

But recall that $D_1f(0) = 1$, which means that if D_1f were continuous at 0, its limit as $x \rightarrow 0$ would have to be 0 also. However, the existence of a point v with the above property for *any* ϵ such that $D_1f(v) = \frac{9}{8}$ implies that the limit cannot be zero. Therefore, D_1f is discontinuous at 0.

Similarly, if D_2f were continuous at 0, its limit would have to equal $D_2f(0) = 0$. By the same argument as above, observe that for any $\epsilon > 0$ we can find $v \in B_\epsilon(0)$, $v = (\sqrt{3}x, x)$, and then:

$$D_2f(v) = -2 \frac{(\sqrt{3}x)^3x}{(3x^2 + x^2)^2} = -\frac{2\sqrt{3} \cdot 3x^4}{16x^4} = -\frac{3\sqrt{3}}{8}$$

The existence of this point shows that the limit of D_2f as $x \rightarrow 0$ cannot be zero, and thus D_2f is not continuous at 0.

Exercise 12

Theorem 2.3.4 states that:

Let $U \subset \mathbb{R}^n$ be an open set, let $a \in U$ and $f : U \rightarrow \mathbb{R}^p$. Then f is differentiable at a if f is partially differentiable in a neighborhood of a and all its partial derivatives are continuous at a .

Let the notation be as in the above theorem, but with $n \geq 2$. Show that the conclusion of the theorem remains valid under the weaker assumption that $n - 1$ partial derivatives of f exist in a neighborhood of a and are continuous at a while the remaining partial derivative merely exists at a .

Hint: Write

$$f(a + h) - f(a) = \sum_{n \geq j \geq 2} (f(a + h^{(j)}) - f(a + h^{(j-1)})) + f(a + h^{(1)}) - f(a)$$

Apply the method of the theorem to the sum over j and the definition of derivative to the remaining difference.

Background: As a consequence of this result we see that at least two different partial derivatives of f must be discontinuous at a if f fails to be differentiable at a , compare with Example 2.3.5.

Solution.

WLOG we will assume that the partial derivative of f with respect to the first variable merely exists at a , and all others exist in a neighborhood of a and are continuous at a (this is because in all other cases everything below is the same up to a “permutation” of variables).

We begin by restating the definition of $h^{(j)}$ that was used in the proof of theorem 2.3.4:

$$h^{(j)} = \sum_{1 \leq k \leq j} h_k e_k \in \mathbb{R}^n, n \geq j \geq 0$$

We also recall that $h^{(j)} = h^{(j-1)} + h_j e_j$. Essentially, $h^{(j)}$ is a vector that lies in the span of the first j vectors of the standard basis of \mathbb{R}^n . The letter h serves to remind us of the fact that it will be used as an infinitesimal “perturbation” vector. Now, as given in the hint, observe that:

$$f(a + h) - f(a) = \sum_{n \geq j \geq 2} (f(a + h^{(j)}) - f(a + h^{(j-1)})) + f(a + h^{(1)}) - f(a),$$

where we see that almost all terms of the sum lead to cancellations of these “partially perturbed” $f(a+h^{(j)})$. Now, as in theorem 2.3.4, we define $g_j : \mathbb{R} \rightarrow \mathbb{R}$ as $g_j(t) = f(a + h^{(j-1)} + te^j)$. We can thus rewrite:

$$f(a+h) - f(a) = \sum_{n \geq j \geq 2} (g_j(h_j) - g_j(0)) + g_1(h_1) - g_1(0)$$

For $j = 2, \dots, n$, the j -th partial derivative of f exists in a neighborhood of a , which means that there exists an open interval in \mathbb{R} containing 0 such that g_j is differentiable on it. Consider a sufficiently small h such that each h_j is contained in the corresponding interval for g_j . As in the proof of 2.3.4, we can apply then the mean value theorem to each g_j to obtain that there exist $\tau_j \in [0, h_j]$ such that:

$$D_j f(a + h^{(j-1)} + \tau_j e_j) = g'_j(\tau_j) = \frac{f(a + h^{(j-1)} + h_j e^j) - f(a + h^{(j-1)})}{h_j}$$

One can thus write that:

$$f(a+h) - f(a) = \sum_{n \geq j \geq 2} [D_j f(a + h^{(j-1)} + \tau_j e_j)] h_j + (g_1(h_1) - g_1(0))$$

We recall now that the remaining partial derivative at a ($D_1 f(a)$) exists. Equivalently, g_1 is differentiable at 0, which, from Hadamard’s lemma means that there exists $\phi_1 \in \mathbb{R} \rightarrow \mathcal{L}(\mathbb{R}, \mathbb{R})$ such that $g_1(t) - g_1(0) = \phi_1(t)t$. By a slight abuse of notation, $\phi_1(t)$ can be identified with a real number also denoted as $\phi_1(t)$. Then let $\phi : \mathbb{R}^n \rightarrow \mathcal{L}(\mathbb{R}^n, \mathbb{R})$ be:

$$\phi(a+h) = (\phi_1(h_1) \quad D_2 f(a + h^{(1)} + \tau_2 e_2) \quad \dots \quad D_n f(a + h^{(n-1)} + \tau_n e_n))$$

By the equations above one can see that $f(a+h) - f(a) = \phi(a+h)(h)$, and because $D_2 f, \dots, D_n f$ are all continuous at the indicated points and ϕ_1 is continuous at 0, ϕ is continuous at a , which is precisely condition (ii) of the Hadamard lemma, which we know is equivalent to f being differentiable at a .

Exercise 15

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a C^1 function and $k > 0$, and assume that $|D_j f(x)| \leq k$ for $1 \leq j \leq n$ and $x \in \mathbb{R}^n$.

(i) Prove that f is Lipschitz continuous with $\sqrt{n}k$ as a Lipschitz constant.

Next, let $g : \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}$ be a C^1 function and $k > 0$, and assume that $|D_j g(x)| \leq k$, for $1 \leq j \leq n$ and $x \in \mathbb{R}^n \setminus \{0\}$.

(ii) Show that if $n \geq 2$, then g can be extended to a continuous function defined on all of \mathbb{R}^n . Show that this is false if $n = 1$ by giving a counterexample.

Solution.

(i) Suppose $x = (x_1, x_2, \dots, x_n), y = (y_1, y_2, \dots, y_n), x, y \in \mathbb{R}^n$. We need to show that $|f(y) - f(x)| \leq \sqrt{n}k\|y - x\|$. Let $h_i = y_i - x_i$, in which case we have that:

$$\begin{aligned} |f(y_1, \dots, y_n) - f(x_1, \dots, x_n)| &= |f(x_1 + h_1, \dots, x_n + h_n) - f(x_1, \dots, x_n)| \\ &= |f(x_1 + h_1, \dots, x_n + h_n) - f(x_1, x_2 + h_2, \dots, x_n + h_n) + f(x_1, x_2 + h_2, \dots, x_n + h_n) - \dots \\ &\quad + f(x_1, \dots, x_n + h_n) - f(x_1, \dots, x_n)| \end{aligned}$$

By writing out like this, we observe that we now have a sum of differences, each of which only differs in one variable of f . Because f is C^1 , by the mean value theorem we can write:

$$\begin{aligned} f(x_1, \dots, x_{i-1}, x_i + h_i, \dots, x_n + h_n) - f(x_1, \dots, x_{i-1}, x_i, \dots, x_n + h_n) &= \\ h_i D_i(f)(x_1, \dots, x_{i-1}, b_i, x_{i+1} + h_{i+1}, \dots, x_n + h_n), \end{aligned}$$

for some $b_i \in (x_i, x_i + h_i)$. If we call c_i the point at which $D_i(f)$ is evaluated, we can rewrite the difference we started with as:

$$|f(y) - f(x)| = \left| \sum_{i=1}^n h_i D_i(f)(c_i) \right| \leq k \left| \sum_{i=1}^n h_i \right| \leq k\sqrt{n} \|h\| = k\sqrt{n} \|y - x\|,$$

where we used the inequality given for $D_i(f)$ and the bound for the l -1 norm that we have seen in e.g. exercise 18, section 3.2 of Carothers.

(ii) If we can show that for any two sequences $(x_n), (y_n) \rightarrow 0$ the corresponding $(g(x_n)), (g(y_n))$ converge to same real number L , we will have shown that a continuous extension of g to all of \mathbb{R}^n exists (by defining its value at 0 as L). We can restrict our attention to the unit ball $B_1(0)$. We would like to be able to use the mean value theorem, but we observe that because $0 \in B_1(0)$, it may not always be the case that for $a, b \in B_1(0) \setminus \{0\}$ the segment $[a, b]$ lies entirely within the domain of g .

Motivated by this, consider partitioning $B_1(0) \setminus \{0\}$ into 2^n “open quadrants”, for all possible combinations of the signs of the coordinates. For example, $Q_1 = \{(x_1, \dots, x_n) \in \mathbb{R}^n | x_1 > 0, x_2 > 0 \dots x_n > 0\}$, $Q_2 = \{(x_1, \dots, x_n) \in \mathbb{R}^n | x_1 < 0, x_2 > 0 \dots x_n > 0\}$. Notice that these are disjoint, and that if we take the union of all Q_i and of all the planes where at least coordinate equals zero we get $B_1(0) \setminus \{0\}$. Notice also that each Q_i is open, and that since all partial derivatives of g are bounded, the proof of (i) applies on each Q_i (instead of on all of \mathbb{R}^n).

Therefore, for any two sequences $x_k, y_k \rightarrow 0$ inside the same quadrant, we have that $|g(y_k) - g(x_k)| \leq k\sqrt{n} \|y_k - x_k\|$. Firstly, this shows that g 's values are bounded inside each quadrant (since the RHS is also bounded). Secondly, the convergence of the two sequences to zero shows that as $k \rightarrow \infty$, $|f(x_k) - f(y_k)| \rightarrow 0$. Putting together these two facts leads to the conclusion that $(f(x_k)), (f(y_k))$ both converge to the same real number.

Now assume that $(x_k) \subset Q_i$, $x_k \rightarrow 0$ and that $y_k \rightarrow 0$ is a sequence that lies on the subset plane which is such that at least one coordinate is zero, and that all non-zero coordinates have the same signs as those of Q_i (intuitively, one of the “borders” of Q_i). Crucially for the remainder of the proof, as long as $n \geq 2$ this subset never degenerates to a single point, which is exactly what happens when $n = 1$. Notice now that the segment formed by the points x_k, y_k lies entirely inside the domain of g (in fact, excluding y_k it lies entirely inside Q_i). Thus, the mean value theorem again applies, and we obtain a similar bound on $|g(y_k) - g(x_k)|$, showing that in fact $g(y_k) \rightarrow L_i$ as well, where L_i is the unique limit obtained previously for Q_i . The final part of the proof is the observation that given any two quadrants, we can find a “path” between them that consists of subsets of planes as described above, and potentially other quadrants. This allows us to transitively reason that all L_i must in fact be equal, thus completing the proof.

Lastly, a concrete counterexample for $n = 1$ is the function $g(x) = -1, x < 0, g(x) = 1, x > 0$, which has zero partial derivatives everywhere, but cannot be extended to a continuous function in all of \mathbb{R} .

Exercise 18

Let $U = \mathbb{R}^2 \setminus \{(0, x_2) | x_2 \geq 0\}$ and define $f : U \rightarrow \mathbb{R}$ by

$$f(x) = \begin{cases} x_2^2, & x_1 > 0 \text{ and } x_2 \geq 0 \\ 0, & x_1 < 0 \text{ and } x_2 < 0 \end{cases}$$

(i) Show that $D_1 f = 0$ on all of U but that f is not independent of x_1 .

Now suppose that $U \subset \mathbb{R}^2$ is an open set having the property that for each $x_2 \in \mathbb{R}$ the set $\{x_1 \in \mathbb{R} | (x_1, x_2) \in U\}$ is an interval.

(ii) Prove that $D_1 f = 0$ on all of U implies that f is independent of x_1 .

Solution.

(i) Let $x \in \mathbb{R}^2$. If $x_1 < 0$, then $f(x) = 0$ and thus $D_1 f(x) = 0$. If $x_1 > 0$, then $f(x) = x_2^2$, and thus again $D_1 f(x) = 0$. However, observe that $f(-1, 1) = 0$ and $f(1, 1) = 1$, which means that when varying only x_1 , the value of f varies, so it cannot be independent of x_1 . A one-dimensional analogue of this would be a function that is stepwise constant (e.g. $f : \mathbb{R} \setminus \{0\}, f(x) = 1, x < 1, f(x) = 2, x > 0$), and as such has derivative 0 in its entire domain but cannot be said to be independent (in this case, constant) of x .

(ii) We need to show that for any three $x_1, y_1, c \in \mathbb{R}$ such that $(x_1, c), (y_1, c) \in U$ it holds that $f(x_1, c) = f(y_1, c)$. In other words, that when x_2 is kept constant, f 's value cannot change. We know that for a

given $c, \{x_1 \in \mathbb{R} | (x_1, c) \in U\}$ is an interval. Let it be called I , with endpoints a, b which it may or may not include. Define $g : I \rightarrow \mathbb{R}, g(x) = f(x, c)$. Then it is the case that g is differentiable, and that $g'(x) = D_1 f(x, c) = 0$. In other words, g has a derivative of 0 and its domain is an interval. From calculus 1, we know that this means that g is constant (one can prove this by contradiction and the mean value theorem), in other words that $g(x) = d \implies f(x, c) = d$ for all $x \in I$, which is exactly what we wanted to show.

Exercise 22

Let $U \subset \mathbb{R}^n$ and $V \subset \mathbb{R}^p$ be open and let $f : U \rightarrow V$ be a C^1 mapping. Define $Tf : U \times \mathbb{R}^n \rightarrow \mathbb{R}^p \times \mathbb{R}^p$, the *tangent mapping* of f to be the mapping given by

$$Tf(x, h) = (f(x), Df(x)h)$$

Let $V \subset \mathbb{R}^p$ be open and let $V \rightarrow \mathbb{R}^q$ be a C^1 mapping. Show that the chain rule takes the natural form

$$T(g \circ f) = Tg \circ Tf : U \times \mathbb{R}^n \rightarrow \mathbb{R}^q \times \mathbb{R}^q$$

Solution.

The chain rule states that for $x \in U$:

$$[D(g \circ f)](x)(h) = [D(g)](f(x))Df(x)(h)$$

By using the tangent mapping definition given in the exercise, we have that for $x \in U, h \in \mathbb{R}^n$:

$$T(g \circ f)(x, h) = ((g \circ f)(x), [D(g \circ f)](x)(h)) = (g(f(x)), [D(g)](f(x))Df(x)(h))$$

Consider now what the RHS of the given equality describes: we take the result of Tf on (x, h) , which belongs in $V \times \mathbb{R}^p$ (a point in V and a direction vector), and pass it as an argument to Tg , which indeed accepts arguments belonging in $V \times \mathbb{R}^p$ (a point in V and a direction vector). This is therefore well-defined. Furthermore:

$$(Tg \circ Tf)(x, h) = Tg(Tf(x, h)) = Tg(f(x), Df(x)h) = (g(f(x)), [D(g)](f(x))Df(x)(h)),$$

which yields that $T(g \circ f) = Tg \circ Tf$.

Exercise 24

Let U be a convex open subset of \mathbb{R}^n and let $f : U \rightarrow \mathbb{R}^p$ be a differentiable mapping. Prove that the following assertions are equivalent.

- (i) The mapping f is Lipschitz continuous on U with Lipschitz constant k .
- (ii) $\|Df(x)\| \leq k$ for all $x \in U$, where $\|\cdot\|$ denotes the operator norm.

Solution.

(i) \implies (ii): Suppose first that f is Lipschitz continuous. Recall the definition of differentiability at a particular x :

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

with $Df(x)$ being the (unique) linear operator that satisfies this. Take any $\epsilon > 0$. Then there exists $\delta > 0$ such that whenever $\|h\| < \delta$, it holds that:

$$\left\| \frac{f(x+h) - f(x) - Df(x)(h)}{\|h\|} - 0 \right\| \leq \epsilon \implies \|f(x+h) - f(x) - Df(x)(h)\| \leq \epsilon \|h\|$$

By the triangle inequality we have that:

$$\left| -\|f(x+h) - f(x)\| + \|Df(x)(h)\| \right| \leq \epsilon \|h\| \implies \|Df(x)(h)\| < \epsilon \|h\| + \|f(x+h) - f(x)\|$$

Now, Lipschitz continuity guarantees that $\|f(x+h) - f(x)\| \leq k\|x+h-x\| = k\|h\|$. We therefore have that $\|Df(x)(h)\| < \epsilon\|h\| + k\|h\|$. Note now that for computing the operator norm of $Df(x)$ we are interested in h that have unit norm. Pick then any such h , and observe that for any $\epsilon > 0$ there exists $\delta > 0$ such that $h' = \delta h$ satisfies the above inequality, which means:

$$\|Df(x)(\delta h)\| < \epsilon\|\delta h\| + k\|\delta h\| \implies \|Df(x)(h)\| < (\epsilon + k)$$

Because this holds for *any* $\epsilon > 0$, we must necessarily have that $\|Df(x)(h)\| \leq k$, and since this is true for any unit-length vector, we can conclude that $\|Df(x)\| \leq k$ as well.

(ii) \implies (i): First, pick any vector $h \in \mathbb{R}^n$ and observe that:

$$\|Df(x)\| \leq k \implies \|Df(x)(h/\|h\|)\| \leq k \implies \|Df(x)(h)\| \leq k\|h\|$$

Given that U is convex and open, we can now simply apply theorem 2.5.3 to conclude that f is Lipschitz continuous on U with Lipschitz constant k .

Exercise 25

Let $U \subset \mathbb{R}^n$ be open and $a \in U$, and let $f : U \setminus \{a\} \rightarrow \mathbb{R}^p$ be differentiable. Suppose there exists $L \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^p)$ with $\lim_{x \rightarrow a} Df(x) = L$. Prove that f is differentiable at a with $Df(a) = L$.

Hint: Apply the Mean Value Theorem to $x \mapsto f(x) - Lx$.

Solution.

We first note that for the exercise to make sense it has to be the case that $n > 1$, otherwise easy counterexamples can be found, and we cannot safely draw the needed intermediate conclusion (by exercise 15 part 2, and by observing that $\lim_{x \rightarrow a} Df(x) = L$ implies that $D_j f(x)$ are bounded for all j) that f is continuous at a and thus $f(a)$ is well-defined.

Consider any open half-ball $H_\delta(a) \subset U \setminus \{a\}$, and any two x, y inside it. This set is convex, and as such we can apply the mean value theorem on f to obtain the existence of c on the segment connecting x, y such that:

$$\begin{aligned} f(y) - f(x) - L(y-x) &= Df(c)(y-x) - L(y-x) \implies \|f(y) - f(x) - L(y-x)\| = \|Df(c)(y-x) - L(y-x)\| \\ &\implies \|f(y) - f(x) - L(y-x)\| \leq \|Df(c) - L\| \cdot \|y-x\|, \end{aligned}$$

where we use the Frobenius norm for the linear function $Df(c) - L$.

Now we note that because $\lim_{x \rightarrow a} Df(x) = L$, for any given $\epsilon > 0$, we can find δ such that for $\|x - a\| < \delta$ it holds that $\|Df(c) - L\| < \epsilon/2$. For an arbitrary ϵ , pick then this corresponding δ , and consider any two x, y with $\|y - x\| < \delta$. If these lie inside *some* half-ball $H_\delta(a)$, then by applying the inequality above we draw the conclusion that $\|f(y) - f(x) - L(y-x)\| < \frac{\epsilon}{2} \cdot \|y-x\|$. If they don't, this means that they are colinear. Crucially, because $n \geq 1$, the “perpendicular bisector” (used loosely here, since it may well be more than one-dimensional) of the segment connecting them is not empty. More specifically, there exists z such that $\|z - x\| = \|z - y\| < \min\{\delta_1, \delta_2\}$, where δ_1, δ_2 are obtained by using the limit definition above for $\epsilon/4$ on half-balls containing y, z (for δ_1) and z, x (for δ_2). This means that:

$$\begin{aligned} \|f(y) - f(x) - L(y-x)\| &\leq \|f(y) - f(z) + f(z) - f(x) - L(y-z+z-x)\| \\ &\leq \|f(y) - f(z) - L(y-z)\| + \|f(z) - f(x) - L(z-x)\| \\ &\leq \|Df(c_1) - L\| \cdot \|z-y\| + \|Df(c_2) - L\| \cdot \|z-x\| < \frac{\epsilon}{4}\|y-x\| + \frac{\epsilon}{4}\|y-x\| = \frac{\epsilon}{2}\|y-x\| \end{aligned}$$

Notice then that if we let $y \rightarrow a$, and by recalling that limits preserve non-strict inequalities, we have that conclusion that for any $\epsilon > 0$, there exists $\delta > 0$ such that for $\|x - a\| < \delta$, $\|f(a+h) - f(a) - L(h)\| \leq \frac{\epsilon}{2}\|h\| < \epsilon\|h\|$, which is precisely the definition of $Df(a) = L$.

Exercise 32

Let $f : \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}$ be a differentiable function.

(i) Let $x \in \mathbb{R}^n \setminus \{0\}$ be a fixed vector and define $g : \mathbb{R}_+ \rightarrow \mathbb{R}$ by $g(t) = f(tx)$. Show that $g'(t) = Df(tx)(x)$.

Assume f to be positively homogeneous of degree $d \in \mathbb{R}$, that is, $f(tx) = t^d f(x)$, for all $x \neq 0, t \in \mathbb{R}_+$.

(ii) Prove the following, known as Euler's identity:

$$Df(x)(x) = \langle x, \nabla f(x) \rangle = df(x), x \in \mathbb{R}^n \setminus \{0\}$$

(iii) Conversely, prove that f is positively homogeneous of degree d if we have $Df(x)(x) = df(x)$, for every $x \in \mathbb{R}^n \setminus \{0\}$. **Hint:** Calculate the derivative of $t \mapsto t^{-d}g(t), t \in \mathbb{R}_+$.

Solution.

(i) Set $h : \mathbb{R}^+ \rightarrow \mathbb{R}^n, h(t) = tx$. Then h is differentiable, and it is the case that $Dh(t)(s) = x$. Also, $g(t) = f(h(t))$, so as a composition of differentiable functions g is also differentiable, and:

$$Dg(t)(s) = ([D(f)]h(t))([Dh(t)](s)) = D(f)(tx)(x)$$

(ii) We know that for any differentiable $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$\langle v, \nabla f(x) \rangle = [Df(x)](v)$$

Therefore $Df(x)(x) = \langle x, \nabla f(x) \rangle$. By setting $g(t) = f(tx) = t^d f(x)$, we have that $g'(t) = dt^{d-1}f(x)$. Therefore, $g'(1) = df(x)$, while from part (i) we have that $g'(1) = D(f)(x)(x)$. This implies that $Df(x)(x) = df(x) = \langle x, \nabla f(x) \rangle$.

(iii) First, we have that by the linearity of the derivative, $df(tx) = Df(tx)(tx) = tDf(tx)(x)$. Then, by using the hint, we define $h(t) = t^{-d}g(t)$ with g as in (i) and we observe that this is a differentiable function with:

$$\begin{aligned} h'(t) &= g'(t)t^{-d} + (-d)t^{-d-1}g(t) = t^{-d}Df(tx)(x) + (-d)t^{-d-1}f(tx) \\ &= t^{-d}\frac{df(tx)}{t} - dt^{-d-1}f(tx) = 0 \end{aligned}$$

But then this means that h is a constant function, therefore that $c = t^{-d}g(t) = t^{-d}f(tx)$ for some $c \in \mathbb{R}$. Also, by setting $t = 1$ we have that $c = f(x)$, which means that $t^{-d}f(tx) = f(x) \implies f(tx) = t^d f(x)$, i.e. that f is positively homogeneous of degree d .

Exercise 44

Let $A \in \mathbb{R}^{n \times n}$. We write $A = (a_1 \cdots a_n)$ if $a_j \in \mathbb{R}^n$ is the j -th column vector of A , $1 \leq j \leq n$. This means that we identify $\mathbb{R}^{n \times n}$ with $\mathbb{R}^n \times \cdots \mathbb{R}^n$. Furthermore, we denote by $A^\#$ the complementary matrix as in Cramer's rule (so $A^\# = (a_{ij}^\# = ((-1)^{i-j} \det A_{ji}))$, where A_{ij} is the matrix obtained by deleting the i -th row and j -th column).

(i) Consider the function $\det : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ as an element of $\mathcal{L}^n(\mathbb{R}^n, \mathbb{R})$ via the identification above. Now prove that its derivative $D(\det)(A) \in \mathcal{L}(\mathbb{R}^n, \mathbb{R})$ is given by

$$D(\det)(A) = \text{tr} \circ A^\#,$$

where tr denotes the trace of a matrix. More explicitly:

$$D(\det)(A)H = \text{tr}(A^\#H) = \langle (A^\#)^T, H \rangle, H \in \mathbb{R}^{n \times n}$$

Here we use the result from exercise 2.1 (ii). In particular, verify that $\nabla \det A = (A^\#)^T$, the cofactor matrix, and

$$(D \det)(I) = \text{tr} \quad (*), (D \det)(A) = \det A (\text{tr} \circ A^{-1}) \quad (**),$$

for A invertible, by means of Cramer's rule.

(ii) Let $X \in \mathbb{R}^{n \times n}$ and define $e^X \in \mathbb{R}^{n \times n}$ as in Example 2.4.10 (i.e. $e^X = I + X + \frac{1}{2}X^2 + \frac{1}{3!}X^3 + \cdots$). Show that for $t \in \mathbb{R}$:

$$\frac{d}{dt} \det(e^{tX}) = \frac{d}{ds} \bigg|_{s=0} \det(e^{(s+t)X}) = \frac{d}{ds} \bigg|_{s=0} \det(e^{sX}) \det(e^{tX}) = \text{tr} X \det(e^{tX})$$

Deduce by solving this differential equation that (compare also with Formula 5.33)

$$\det(e^{tX}) = e^{t \cdot \text{tr} X}$$

(iii) Assume that A is an invertible $n \times n$ matrix. Using $\det(A + H) = \det(A) \det(I + A^{-1}H)$, deduce $(**)$ from $(*)$ in (i). Now also derive $(\det A)A^{-1} = A^\#$ from (i).

(iv) Assume that $A : \mathbb{R} \rightarrow \mathbf{GL}(n, \mathbb{R})$ (i.e. A always yields invertible matrices) is a differentiable mapping. Derive from part (i) the following formula for the derivative of $\det \circ A : \mathbb{R} \rightarrow \mathbb{R}$:

$$\frac{1}{\det \circ A} (\det \circ A)' = \text{tr}(A^{-1} \circ DA),$$

in other words (compare with part (ii))

$$\frac{d}{dt} (\log \det A)(t) = \text{tr}(A^{-1} \frac{dA}{dt})(t), t \in \mathbb{R}$$

(v) Suppose that the mapping $A : \mathbb{R} \rightarrow \mathbb{R}^{n \times n}$ is differentiable and that $F : \mathbb{R} \rightarrow \mathbb{R}^{n \times n}$ is continuous, while we have the differential equation

$$\frac{dA}{dt}(t) = F(t)A(t), t \in \mathbb{R}$$

In this context, the *Wronskian* $w \in C^1(\mathbb{R})$ of A is defined as $w(t) = \det A(t)$. Use part (i), exercise 2.1 (i) and Cramer's rule to show that for $t \in \mathbb{R}$

$$\frac{dw}{dt}(t) = \text{tr} F(t)w(t), \text{ and deduce } w(t) = w(0)e^{\int_0^t \text{tr} F(\tau) d\tau}$$

Solution.

(i) We have shown in exercise 2.15 of Spivak that:

$$D(\det)(a_1, \dots, a_n)(h_1, \dots, h_n) = \sum_{i=1}^n \det \begin{pmatrix} a_1 \\ \vdots \\ h_i \\ \vdots \\ a_n \end{pmatrix}$$

We therefore have but to show that the RHS of the equation given in (i) equals the RHS of this equation. We have that $(\text{tr} \circ A^\#)(H) = \text{tr}(A^\# H)$. Due to the trace we are only interested in the diagonal elements of the argument. In particular, we have that:

$$\text{tr}(A^\# H) = \sum_{i=1}^n \sum_{j=1}^n (-1)^{i-j} \det A_{ji} h_{ji}$$

We now manipulate the expression we have from Spivak using the property of row linearity:

$$D(\det)(a_1, \dots, a_n)(h_1, \dots, h_n) = \sum_{i=1}^n \sum_{j=1}^n h_{ij} \det \begin{pmatrix} a_1 \\ \vdots \\ a_{i-1} \\ e_j \\ \vdots \\ a_n \end{pmatrix}$$

Now we observe that we can apply row-wise Laplace expansion on the i -th row of each of the summands for which the outer index equals i to rewrite it as:

$$\det \begin{pmatrix} a_1 \\ \vdots \\ a_{i-1} \\ e_j \\ \vdots \\ a_n \end{pmatrix} = (-1)^{i+j} \det A_{ij},$$

where we used the fact that this row contains exactly one non-zero element (at column j), and that deleting this row and column yields the same result as deleting the same row and column from A . Note that $(-1)^{i+j} = (-1)^{i-j}$, and thus from this derivation we obtain that:

$$D(\det)(a_1, \dots, a_n)(h_1, \dots, h_n) = \text{tr}(H A^\#) = \text{tr}(A^\# H),$$

by using the commutative property of the trace. This concludes the proof, and now we can of course write $D(\det)(A)H = \langle (A^\#)^T, H \rangle$ by the definition of inner product for matrices we saw in exercise 2.1. From this we additionally conclude that $\nabla \det A = (A^\#)^T$: the gradient is precisely the vector whose inner product with each “direction vector” (here, the matrix H) yields the directional derivative along H , i.e. $D(\det)(A)H$.

Now if we were to evaluate at $A = I$, we first have that $A^\# = I$. To see why this is the case, consider erasing column i and row i : we are simply left with a smaller identity matrix, which of course has determinant 1. On the other hand, erasing row i , column $j, i \neq j$ will result in an either upper or lower triangular matrix that has at least one zero on the diagonal, and such a matrix will have determinant zero.

Thus:

$$D(\det)(I)(H) = \text{tr}(IH) = \text{tr}(H)$$

For A invertible, we know by Cramer’s rule that $\det A \cdot I = AA^\# = A^\# A$, and since A^{-1} is well defined, $(\det A)A^{-1} = A^\#$. Substituting into the trace expression from before:

$$D(\det)(A)(H) = \text{tr}(A^\# H) = \text{tr}((\det A)A^{-1}H) = \det A \text{tr}(A^{-1}H)$$

(ii) Call $f(t) = e^{tX}$, $g(A) = \det A$, in which case $f : \mathbb{R} \rightarrow \mathbb{R}^{n \times n}$ and $g : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$, and in which case we are interested in $(g \circ f)'$. By the chain rule, we have that:

$$D(g \circ f)(t)(h) = [D(g)(f(t))][D(f)(t)](h)$$

We know from Example 2.4.10 that $[D(f)(t)](h) = he^{tX}X$, and by also using the expression we obtained in part (i) for A invertible (as is e^{tX}) we get:

$$D(g \circ f)(t)(h) = [D(g)(f(t))](he^{tX}X) = \det(e^{tX}) \operatorname{tr}((e^{tX})^{-1}he^{tX}X) = \det(e^{tX}) \operatorname{tr}(hX) = \det e^{tX} \operatorname{tr}(X)h,$$

which is equivalent to saying that $\frac{d}{dt} \det(e^{tX}) = \operatorname{tr}(X) \det(e^{tX})$.

Now notice that the function $g \circ f$ maps from \mathbb{R} to \mathbb{R} , and therefore this differential equation is of the general form:

$$\frac{d}{dt}(g \circ f)(t) = a(g \circ f)(t),$$

whose solution is well known to be $(g \circ f)(t) = Ce^{at}$, which implies that $\det(e^{tX}) = Ce^{\operatorname{tr}(X)t}$. To determine C , we set $t = 0$, in which case e^{tX} is the identity matrix, and thus the LHS is 1 and the RHS is C , meaning that $C = 1$.

(iii) We begin by writing out the limit that corresponds to (*), i.e. to $(D \det)(I) = \operatorname{tr}$:

$$\lim_{H \rightarrow 0} \frac{1}{\|H\|} (\det(I + H) - \det(I) - \operatorname{tr}(H)) = 0 \implies \lim_{H \rightarrow 0} \frac{1}{\|H\|} (\det(I + H) - 1 - \operatorname{tr}(H)) = 0$$

Now, for the derivative of \det at an invertible A , we are interested in finding the linear transformation \mathcal{L} such that:

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

Using the provided hint:

$$\lim_{H \rightarrow 0} \frac{1}{\|H\|} (\det(A + H) - \det(A) - \mathcal{L}(H)) = \lim_{H \rightarrow 0} \frac{1}{\|H\|} (\det(A) \det(I + A^{-1}H) - \det(A) - \mathcal{L}(H))$$

This then motivates us to set $U = A^{-1}H$, which means $AU = H$, and also that since A^{-1} is kept constant as $H \rightarrow 0$, it is also the case that $U \rightarrow 0$. Therefore, the above limit can be rewritten as:

$$\lim_{U \rightarrow 0} \frac{1}{\|AU\|} (\det(A) \det(I + U) - \det(A) - \mathcal{L}(AU)) = \lim_{U \rightarrow 0} \left(\frac{1}{\|AU\|} \det(A) (\det(I + U) - 1) - \frac{1}{\|AU\|} \mathcal{L}(AU) \right)$$

Now we have that $\|U\| = \|A^{-1}AU\| \leq \|A^{-1}\| \cdot \|AU\|$, we observe that for the first term of the sum:

$$\frac{1}{\|AU\|} \|\det(A) (\det(I + U) - 1)\| \leq \frac{1}{\|A^{-1}\| \cdot \|U\|} \|\det(A) (\det(I + U) - 1)\|$$

But now the right side is known to tend to zero if one were to subtract $\det(A) \operatorname{tr}(U)$ from the quantity inside the magnitude of the numerator. This therefore motivates us to have $\mathcal{L}(AU) = \det(A) \operatorname{tr}(U)$, which, substituting U , yields $\mathcal{L}(H) = \det(A) \operatorname{tr}(A^{-1}H)$, and since a linear transformation that makes the limit zero is *the* derivative of \det at A , we have indeed arrived at (**).

(iv) We have shown in (i) that for A invertible:

$$D(\det)(A)(H) = \det(A) \operatorname{tr}(A^{-1}H)$$

Let then $f(t) = \det(A(t))$, $f : \mathbb{R} \rightarrow \mathbb{R}$. By using the chain rule:

$$D(f)(t)(h) = D(\det)(A(t))(D(A)(t)(h)) \implies D(f)(t)(h) = \det(A(t)) \operatorname{tr}((A(t))^{-1}D(A)(t)(h))$$

Because $A(t)$ is always invertible, we can divide both sides by it, and by simplifying the notation to match the exercise we obtain:

$$\frac{1}{\det A} (\det \circ A)' = \operatorname{tr}(A^{-1} \circ DA)$$

Now, by recalling the derivative of the (natural) logarithm, we can also rewrite this as:

$$\frac{d}{dt}(\log \det A)(t) = \operatorname{tr}(A^{-1} \frac{dA}{dt})(t)$$

(v) We begin by differentiating $w(t)$ and using part (i):

$$D(w)(t)(h) = D(\det)(A(t))(D(A)(t)(h)) \implies D(w)(t)(h) = \operatorname{tr}(A(t)^{\#} D(A)(t)(h))$$

Note now that by simplifying the notation and substituting from the given differential equation we obtain that:

$$\begin{aligned} \frac{dw}{dt}(t) &= \operatorname{tr}(A(t)^{\#} F(t) A(t)) \implies \frac{dw}{dt}(t) = \operatorname{tr}(A t^{\#} A(t) F(t)) \implies \frac{dw}{dt}(t) = \operatorname{tr}(\det A(t) I F(t)) \\ &\implies \frac{dw}{dt}(t) = \det A(t) \operatorname{tr}(F(t)) = w(t) \operatorname{tr}(F(t)), \end{aligned}$$

where we have used properties of the trace (shown in exercise 2.1) and Cramer's rule. This is a differential equation somewhat similar to (ii), except that $w(t)$ is multiplied with a function of t instead of a constant. Because F is continuous, the solution to this differential equation is of the form:

$$w(t) = c e^{\int_0^t \operatorname{tr} F(\tau) d\tau},$$

and by substituting $t = 0$ we have $w(0) = c$, where $w(0)$ depends on $A(0)$ (an "initial condition" for the differential equation).