

Chapter 8

PERMUTATIONS, DETERMINANTS, AND THE GEOMETRY OF LINEAR TRANSFORMATIONS

8.1 Permutations

8.1.1 The permutation group

The concept of a *transformation group* is fundamentally important to modern mathematics, and to geometry in particular. In this section we introduce a basic example of a transformation group: *the Permutation Group*. As we shall explain, the permutation group plays an essential role in the computations of area, volumes and their higher dimensional generalizations.

Definition 78 (Permutation). *A permutation of $\{1, 2, \dots, n\}$ is a function σ from this set onto itself.*

Recall that “onto” means that for every j in $\{1, 2, \dots, n\}$, there is an i with $\sigma(i) = j$. We can specify a permutation σ of $\{1, 2, \dots, n\}$ by listing the assignments it makes:

$$\begin{array}{ccccccccc} 1 & 2 & 3 & \cdots & n \\ \downarrow & \downarrow & \downarrow & \cdots & \downarrow \\ \sigma(1) & \sigma(2) & \sigma(3) & \cdots & \sigma(n) \end{array}$$

For example, if $n = 3$, and $\sigma(1) = 2$, $\sigma(2) = 3$ and $\sigma(3) = 1$,

$$\sigma = \begin{array}{ccc} 1 & 2 & 3 \\ \downarrow & \downarrow & \downarrow \\ 2 & 3 & 1 \end{array}.$$

The arrows do not really tell us much; we can remember that the top row is inputs, and the bottom row is outputs. Let's shorten the notation to

$$\sigma = \begin{array}{ccc} 1 & 2 & 3 \\ 2 & 3 & 1 \end{array}$$

The generalization of this way of writing permutations to higher values of n is plain, and we use it freely.

There are exactly $n!$ permutations of $\{1, 2, \dots, n\}$: Consider any permutation σ of $\{1, 2, \dots, n\}$. There are n choices for the value of $\sigma(1)$. Make this choice, and then, $\sigma(1)$ being taken, there are $n - 1$ choices remaining for value of $\sigma(2)$. Next, there are $n - 2$ choices for $\sigma(3)$, the value to be assigned to 3. Continuing in this way, there are $n(n - 1)(n - 2) \cdots 1 = n!$ choices to make, and each one leads to a distinct permutation.

Example 118 (Permutations of $\{1, 2, 3\}$). *There are six permutations of $\{1, 2, 3\}$:*

$$\begin{array}{ccc} \sigma_a = & \begin{array}{ccc} 1 & 2 & 3 \\ 1 & 2 & 3 \end{array} & \sigma_b = \begin{array}{ccc} 1 & 2 & 3 \\ 2 & 1 & 3 \end{array} \\ \\ \sigma_c = & \begin{array}{ccc} 1 & 2 & 3 \\ 1 & 3 & 2 \end{array} & \sigma_d = \begin{array}{ccc} 1 & 2 & 3 \\ 2 & 3 & 1 \end{array} \\ \\ \sigma_e = & \begin{array}{ccc} 1 & 2 & 3 \\ 3 & 1 & 2 \end{array} & \sigma_f = \begin{array}{ccc} 1 & 2 & 3 \\ 3 & 2 & 1 \end{array} \end{array} \quad (8.1)$$

Since permutations of $\{1, 2, \dots, n\}$ are functions from this set into itself, we can compose them: If σ_1 and σ_2 are two permutations of $\{1, 2, \dots, n\}$, then $\sigma_2 \circ \sigma_1$ is defined by

$$\sigma_2 \circ \sigma_1(i) = \sigma_2(\sigma_1(i)) \quad , \quad \text{for each } i = 1, \dots, n . \quad (8.2)$$

Example 119 (Composing permutations). *Let us compute $\sigma_d \circ \sigma_b$ where σ_d and σ_b are the permutations given in (8.1). From (8.1) we see that*

$$\begin{aligned} \sigma_d \circ \sigma_b(1) &= \sigma_d(\sigma_b(1)) = \sigma_d(2) = 3 \\ \sigma_d \circ \sigma_b(2) &= \sigma_d(\sigma_b(2)) = \sigma_d(1) = 2 \\ \sigma_d \circ \sigma_b(3) &= \sigma_d(\sigma_b(3)) = \sigma_d(3) = 1 \end{aligned}$$

Thus,

$$\sigma_d \circ \sigma_b = \begin{array}{ccc} 1 & 2 & 3 \\ 3 & 2 & 1 \end{array} = \sigma_f .$$

In this example, the composition product of two permutations was another permutation. In fact, the composition of two permutations is *always* a permutation: Consider any two permutations σ_2 and σ_1 of $\{1, 2, \dots, n\}$. To see that $\sigma_2 \circ \sigma_1$ is also a permutation, we just need to check that for each j in $\{1, 2, \dots, n\}$, there is an i such that $\sigma_2 \circ \sigma_1(i) = j$. But since σ_2 is a permutation, there is a k so that $\sigma_2(k) = j$. And since σ_1 is a permutation, there is an i so that $\sigma_1(i) = k$. Then $\sigma_2 \circ \sigma_1(i) = \sigma_2(\sigma_1(i)) = \sigma_2(k) = j$. We have found the i for which $\sigma_2 \circ \sigma_1(i) = j$, so $\sigma_2 \circ \sigma_1$ is a permutation.

The permutation σ_a at the upper left of the list in (8.1) is called the *identity* permutation since it just sends each element of $\{1, 2, 3\}$ to itself. This has an obvious generalization to other values of n . Moreover, every permutation σ has an inverse, σ^{-1} , which simply sends any j in back to the integer i in $\{1, 2, \dots, n\}$ from whence it came. (Since $\{1, 2, \dots, n\}$ is a finite set, and since σ is onto, it is also one-to-one. Indeed, if $\sigma(i) = \sigma(j)$ for $i \neq j$, it would have spent two of n “shots” at hitting a single target, which would preclude hitting all n . So σ is necessarily one-to-one from $\{1, 2, \dots, n\}$ onto itself, and hence invertible.)

The inverse too is a map of $\{1, 2, \dots, n\}$ onto itself, and hence a permutation. (It is just the original map “in reverse”).

Definition 79 (Permutation group). Let \mathcal{S}_n denote the set of all $n!$ permutations of $\{1, \dots, n\}$, equipped with the composition product $\sigma_1 \circ \sigma_2$. This is the permutation group on $\{1, \dots, n\}$.

The term “group” has a precise technical meaning in mathematics; It is a generalization of the more concrete notion of a “transformation group” which is what the permutations group is: A *transformation group on a set X* is a set \mathcal{G} of invertible functions from X to X such that whenever $g \in \mathcal{G}$, then $g^{-1} \in \mathcal{G}$, and such that whenever $g_1, g_2 \in \mathcal{G}$, then $g_1 \circ g_2 \in \mathcal{G}$. Note that, as a consequence of the definition, \mathcal{G} contains the identity transformation $i(x) = x$ for all $x \in X$. Since \mathcal{S}_n contains *all* invertible transformations from $\{1, \dots, n\}$ into itself, it is the largest transformation group on $\{1, \dots, n\}$.

8.1.2 The character of a permutation

In this subsection, we define a function χ on \mathcal{S}_n with values in $\{-1, 1\}$, called the *character*, that is essential to the theory of determinants. The definition of χ depends on another function which measures the “degree of mixing” of a permutation σ , or in other words, “how far σ is from the identity permutation”.

Consider once more the list (8.1) of permutations. Except for the identity permutation, σ_a , all of these permutations “mix things up” to some extent. In fact, we have arranged these permutations in an order that reflects a measure of “how much mixing” is involved in each one, starting from no mixing in the identity transformation $\begin{array}{ccc} 1 & 2 & 3 \\ 1 & 2 & 3 \end{array}$ at the upper left, to the most mixing in the “order reversing” permutation $\begin{array}{ccc} 1 & 2 & 3 \\ 3 & 2 & 1 \end{array}$ at the lower right. Now, you may well ask: In what sense is the order reversing permutation “farthest from the identity”? After all, it does send 2 to 2, and

other permutations have no such fixed points. To answer this question, we must explain how we quantitatively measure “the degree of mixing” .

We shall quantify the the degree of mixing of a permutation σ by counting “the number of pairs it puts out of order”: Consider the set $P := \{(i, j) : 1 \leq i, j \leq n\}$ of all distinct ordered pairs chosen from $\{1, \dots, n\}$, which is a set of $n(n-1)$ elements. Define the disjoint sets

$$P_{\text{up}} = \{(i, j) : 1 \leq i < j \leq n\} \quad \text{and} \quad P_{\text{down}} = \{(i, j) : 1 \leq j < i \leq n\} .$$

P_{up} is the set of all “increasing” pairs and P_{down} is the set of all “decreasing” pairs. Note that both of these sets consist of $n(n-1)/2$ ordered pairs, and $P = P_{\text{up}} \cup P_{\text{down}}$.

For any $\sigma \in \mathcal{S}_n$, the function f_σ from P into itself defined by

$$f_\sigma(i, j) = (\sigma(i), \sigma(j))$$

is invertible. In fact, it is a permutation of the elements of P .

Since f_σ is one-to-one and onto, each pair that f_σ moves out of P_{up} into P_{down} must be replaced by a pair that f_σ moves out of P_{down} into P_{up} so that *the number of pairs that f_σ moves out of P_{up} into P_{down} coincides with the number of pairs it moves out of P_{down} into P_{up}* : It is simply the number of pairs that f_σ “swaps” between P_{down} and P_{up} .

Definition 80 (Definition (Degree of mixing)). *The degree of mixing of a permutation σ of $\{1, 2, \dots, n\}$ is the number of pairs of integers (i, j) in $\{1, 2, \dots, n\}$ with*

$$i < j \quad \text{and} \quad \sigma(i) > \sigma(j) . \quad (8.3)$$

This number is denoted $D(\sigma)$. In terms of the notation introduced in the preceding paragraph, $D(\sigma)$ is the number of pairs that f_σ swaps between P_{down} and P_{up} . The more “reversed” pairs, the more mixing there is.

Example 120 (Computing the degree of mixing). *Let us compute $D(\sigma)$ for each of the six permutations of $\{1, 2, 3\}$. There are exactly three pairs (i, j) with $i < j$, namely*

$$(1, 2) \quad (1, 3) \quad (2, 3) .$$

To compute the degree of mixing of σ , we look at

$$(\sigma(1), \sigma(2)) \quad (\sigma(1), \sigma(3)) \quad (\sigma(2), \sigma(3)) ,$$

and count the number of these pairs that are “out of order”. You can easily check that

$$D(\sigma_a) = 0 \quad D(\sigma_b) = 1 \quad D(\sigma_c) = 1 \quad D(\sigma_d) = 2 \quad D(\sigma_e) = 2 \quad D(\sigma_f) = 3 .$$

Thus, with this definition of the degree of mixing, the order reversing permutation $\begin{smallmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{smallmatrix}$ *has the highest degree of mixing among all permutations of $\{1, 2, 3\}$.*

Lemma 18 (Degree of mixing and inverses). *For each $\sigma \in \mathcal{S}_n$,*

$$D(\sigma^{-1}) = D(\sigma) .$$

Proof: Let f_σ be the invertible function on pairs induced by σ , as explained above. Then evidently $(f_\sigma)^{-1} = f_{\sigma^{-1}}$. However many pairs f_σ swaps between P_{up} and P_{down} , $f_{\sigma^{-1}}$ swaps the same same number back again in undoing the effects of f_σ . \square

Example 121 (A graphical view of the degree of mixing). *There is a very good graphical way to compute $D(\sigma)$ that I have learned from Roe Goodman: For example, start from*

$$\begin{array}{ccc} 1 & 2 & 3 \\ 1 & 2 & 3 \end{array}$$

which denotes the identity permutation. Now, to represent any other permutation σ on $\{1, 2, 3\}$, simply draw in arrow running from i in the top row to $\sigma(i)$ in the bottom row for $i = 1, 2, 3$. Do this for each of the permutations in the previous example, and note that $D(\sigma)$ is exactly the number of intersection points of the three arrows. With a bit of thought, you will be able to see that this generalizes to any value of n , and that the number of intersection points always counts the number of out-of-order pairs.

The definition of $D(\sigma)$ is useful because of the way it interacts with the composition product: Consider the following question:

- Given two permutations σ_1 and σ_2 of $\{1, 2, \dots, n\}$, what can we say about $D(\sigma_2 \circ \sigma_1)$?

Lemma 19 (Degree of mixing and composition). *For any $\sigma_1, \sigma_2 \in \mathcal{S}_n$,*

$$D(\sigma_2 \circ \sigma_1) = D(\sigma_2) + D(\sigma_1) - 2c \quad (8.4)$$

where c is a non-negative integer.

Proof: First, in applying σ_1 , we reverse the order of $D(\sigma_1)$ pairs. Then, applying σ_2 after that, we reverse the order of $D(\sigma_2)$ pairs. So the number of pairs that are reversed by $\sigma_2 \circ \sigma_1$ is no more than $D(\sigma_1) + D(\sigma_2)$.

However, some of the pairs that σ_2 reverses may have already been put out of order by σ_1 . In this case, σ_2 puts them back in order. An extreme case is when $\sigma_2 = (\sigma_1)^{-1}$. Then σ_2 undoes all of the mixing done by σ_1 , and $D(\sigma_2 \circ \sigma_1) = 0$.

So we conclude that $0 \leq D(\sigma_2 \circ \sigma_1) \leq D(\sigma_1) + D(\sigma_2)$. We can say more: Suppose that when σ_2 is applied, c pairs that had been put out of order by σ_1 are “reordered” when we apply σ_2 . Then,

- Of the $D(\sigma_1)$ pairs reversed by σ_1 , exactly $D(\sigma_1) - c$ are still reversed after applying σ_2 .
- Of the $D(\sigma_2)$ pair reversals created by σ_2 , c are “used up” undoing reversals created by σ_1 , and so exactly $D(\sigma_2) - c$ new reversals are created.

Adding things up, $D(\sigma_2 \circ \sigma_1) = (D(\sigma_1) - c) + (D(\sigma_2) - c) = D(\sigma_1) + D(\sigma_2) - 2c$, which proves (8.4). \square

We now come to our first application of Lemma 19. Note that that whatever c is in (8.4), $2c$ is always an even integer, and so $(-1)^{2c} = 1$, and

$$(-1)^{D(\sigma_2 \circ \sigma_1)} = (-1)^{D(\sigma_1)} (-1)^{D(\sigma_2)} . \quad (8.5)$$

Definition 81 (Character of a Permutation). *The character $\chi(\sigma)$ of a permutation σ is defined by*

$$\chi(\sigma) = (-1)^{D(\sigma)} . \quad (8.6)$$

where $D(\sigma)$ is given by (1.7). A permutation σ is called an even permutation if $\chi(\sigma) = 1$, and an odd permutation if $\chi(\sigma) = -1$.

The key property of the character function is that $\chi(\sigma_2 \circ \sigma_1) = \chi(\sigma_2)\chi(\sigma_1)$. That is, *the character of a product equals the product of the characters*. This follows directly from (8.5). In Example 121 σ_a , σ_d and σ_e are even permutations whereas σ_b , σ_c and σ_f are odd permutations.

If you want to determine $\chi(\sigma)$ for a given permutation σ you need not compute $D(\sigma)$ first, and then apply the definition (8.6). There are some general rules for particular kinds of permutations.

Definition 82 (Pair Permutations). *For each $i < j$ in $\{1, 2, \dots, n\}$ the pair permutation $\sigma_{i,j}$ is defined by*

$$\sigma_{i,j}(i) = j \quad , \quad \sigma_{i,j}(j) = i \quad \text{and} \quad \sigma_{i,j}(k) = k \quad \text{for} \quad k \neq i, j . \quad (8.7)$$

It is called an adjacent pair permutation in case $j = i + 1$ for $i < n$, or if $(i, j) = (n, 1)$; i.e., if j follows i in the cyclic order on $\{1, \dots, n\}$.

Example 122. For $n = 4$, $\sigma_{2,4} = \begin{array}{cccc} 1 & 2 & 3 & 4 \\ 1 & 4 & 3 & 2 \end{array}$.

Notice that each pair permutation is its own inverse – applying it twice swaps the reversed pair back into place.

Next notice that for each adjacent pair permutation $\sigma_{i,i+1}$, $D(\sigma_{i,i+1}) = 1$, and hence $\chi(\sigma_{i,i+1}) = -1$. What about general pair permutations?

- For any $i < j$, $\sigma_{i,j}$ can be written as the product of $2k - 1$ adjacent pair permutations where $k = j - i$.

Therefore, since the character of a product is the product of the characters,

$$\chi(\sigma_{i,j}) = (-1)^{2k-1} = -1$$

for every pair permutation, adjacent or not.

To justify the claim about σ_{ij} with $i < j$, write $j = i + k$. Then one can “move” i to the right of j using k adjacent pair permutations. One can then move j back to the i th spot with $k - 1$ pair permutations. Only $k - 1$ are required, because the last pair permutation used to move i into the j th place already moved j one place to the left.

We summarize the discussion in the following theorem:

Theorem 73 (Properties of the character). *For any two permutations σ_1 and σ_2 of $\{1, 2, \dots, n\}$,*

$$\chi(\sigma_2 \circ \sigma_1) = \chi(\sigma_2)\chi(\sigma_1) . \quad (8.8)$$

Moreover, for any pair permutation $\sigma_{i,j}$,

$$\chi(\sigma_{i,j}) = -1 . \quad (8.9)$$

The theorem gives us a convenient way to compute $\chi(\sigma)$: Bring the sequence $(1, 2, \dots, n)$ into the order $(\sigma(1), \sigma(2), \dots, \sigma(n))$ by swapping pairs; that is, by pair permutations. Then $\chi(\sigma)$ is the product of the characters of these pairs permutations, so it is $(-1)^\ell$, where ℓ is the number of pair permutations you used.

Example 123 (Computing $\chi(\sigma)$ counting pair permutations). Consider $\sigma = \begin{smallmatrix} 1 & 2 & 3 & 4 \\ 4 & 1 & 3 & 2 \end{smallmatrix}$. We can

transform $(1, 2, 3, 4)$ to $(4, 1, 3, 2)$ using pair permutations as follows:

$$(1, 2, 3, 4) \rightarrow (4, 2, 3, 1) \rightarrow (4, 1, 3, 2)$$

or as well by

$$(1, 2, 3, 4) \rightarrow (1, 2, 4, 3) \rightarrow (1, 4, 2, 3) \rightarrow (4, 1, 2, 3) \rightarrow (4, 1, 3, 2)$$

In the first case we used 2 pair permutations, and in the second case we used 4. Either way, we see $\chi(\sigma) = (-1)^2 = (-1)^4 = 1$, so σ is even.

You might wonder why we did not *define* $\chi(\sigma)$ to be $(-1)^\ell$ where ℓ is the number of “pair swaps” required to produce σ . The point is this: Suppose you could write some σ as a product of 7 pair permutations, and also 242 pair permutations. Then you would have $\chi(\sigma) = (-1)^7 = -1$ and $\chi(\sigma) = (-1)^{242} = 1$, and *both* cannot be right. If this happened, $\chi(\sigma) = (-1)^\ell$ would not be a well defined function.

Evidently, our analysis above implies that for any given permutation σ , if there is a way to write σ as a product of an odd number of pair permutations, then *every* way of writing σ as a product of pair permutations uses an odd number of them. This fact is not obvious! We know it is true because we have proved Lemma 19.

At this point we have covered as much of the theory of the permutation group as we shall use in explaining the theory of determinants. However, the permutation group is such a fundamental example of a transformation group, and the notion of a transformation group is so essential to modern analysis and geometry, that it is worthwhile to go somewhat further with the theory of permutations, and to study \mathcal{S}_n as a metric space. We do this in the next subsection.

8.1.3 The permutation group as a metric space

Definition 83 (Distance in \mathcal{S}_n). Let ϱ be the function on $\mathcal{S}_n \times \mathcal{S}_n$ given by

$$\varrho(\sigma_1, \sigma_2) = D(\sigma_1^{-1} \circ \sigma_2) .$$

This function is called the length function or distance function on \mathcal{S}_n .

It is not hard to see that the length function we have just defined is a metric on \mathcal{S}_n . That is, it satisfies the three requirements of a metric:

$$(1) \text{ For all } \sigma_1, \sigma_2 \in \mathcal{S}_n, \varrho(\sigma_1, \sigma_2) \geq 0, \text{ and } \varrho(\sigma_1, \sigma_2) = 0 \iff \sigma_1 = \sigma_2.$$

(2) For all $\sigma_1, \sigma_2 \in \mathcal{S}_n$, $\varrho(\sigma_1, \sigma_2) = \varrho(\sigma_2, \sigma_1)$.

(3) For all $\sigma_1, \sigma_2, \sigma_3 \in \mathcal{S}_n$, $\varrho(\sigma_1, \sigma_3) \leq \varrho(\sigma_1, \sigma_2) + \varrho(\sigma_2, \sigma_3)$.

To see that this is the case, note for (1) that ϱ is defined to be a non-negative integer, and $D(\sigma_1^{-1} \circ \sigma_2) = 0$ if and only if there are “no crossings” in $\sigma_1^{-1} \circ \sigma_2$, which is the case if and only if $\sigma_1^{-1} \circ \sigma_2$ is the identity, which is the case if and only if $\sigma_1 = \sigma_2$. For (2), Note that

$$(\sigma_1^{-1} \circ \sigma_2)^{-1} = \sigma_2^{-1} \circ \sigma_1$$

and since, by Lemma 18, D is unaffected by taking inverses,

$$\varrho(\sigma_1, \sigma_2) = D(\sigma_1^{-1} \circ \sigma_2) = D((\sigma_1^{-1} \circ \sigma_2)^{-1}) = D(\sigma_2^{-1} \circ \sigma_1) = \varrho(\sigma_2, \sigma_1) .$$

Finally, for (3) we use (8.4) and the fact that, due to the associative nature of composition,

$$\sigma_1^{-1} \circ \sigma_3 = (\sigma_1^{-1} \circ \sigma_2) \circ (\sigma_2^{-1} \circ \sigma_3) .$$

Thus, by (8.4), since $2c \geq 0$ for all non-negative integers c ,

$$\varrho(\sigma_1, \sigma_3) = D(\sigma_1^{-1} \circ \sigma_3) = D((\sigma_1^{-1} \circ \sigma_2) \circ (\sigma_2^{-1} \circ \sigma_3)) \leq D(\sigma_1^{-1} \circ \sigma_2) + D(\sigma_2^{-1} \circ \sigma_3) = \varrho(\sigma_1, \sigma_2) + \varrho(\sigma_2, \sigma_3) .$$

We now explain how one can think of $\varrho(\sigma_1, \sigma_2)$ as the *length of the shortest path in \mathcal{S}_n from σ_1 to σ_2* . Given any $\sigma \in \mathcal{S}_n$, consider the set of permutations

$$\{\sigma \circ \tau : \tau \text{ is an adjacent pair permutation}\}$$

We call this set the set of the *nearest neighbors* of σ in \mathcal{S}_n . In terms of the graphical representation discussed in Example 121, the diagram representing σ differs from the diagram representing any of its nearest neighbors only in having the tails of two adjacent arrows swapped. (We use the cyclic order on $\{1, \dots, n\}$ in which 1 and n are adjacent; 1 follows n .)

Now think of “moving” from σ to $\sigma \circ \tau$, where τ is an adjacent pair transposition, as a “step” from σ to one of its nearest neighbors. By a *path in \mathcal{S}_n from σ_1 to σ_2* , we mean a sequence of such steps starting at σ_1 and ending at σ_2 .

Definition 84 (Paths in \mathcal{S}_n). *For any σ_1 and σ_2 in \mathcal{S}_n , a path from σ_1 to σ_2 is a sequence $\{\tau_1, \dots, \tau_\ell\}$ of adjacent pair permutations such that*

$$\sigma_2 = \sigma_1 \circ \tau_1 \cdots \circ \tau_\ell .$$

For example, if $\{\tau_1, \tau_2, \tau_3\}$ is a path from σ_1 to σ_2 , then the sequences of steps

$$\sigma_1 \longrightarrow \sigma_1 \circ \tau_1 \longrightarrow \sigma_1 \circ \tau_1 \circ \tau_2 \longrightarrow \sigma_1 \circ \tau_1 \circ \tau_2 \circ \tau_3 = \sigma_2$$

is a sequence of “one step moves between nearest neighbors” that starts at σ_1 and ends at σ_2 .

Theorem 74 (The metric in \mathcal{S}_n as a minimal path length). *For each $\sigma_1, \sigma_2 \in \mathcal{S}_n$, there is a path from σ_1 to σ_2 , and*

$$\varrho(\sigma_1, \sigma_2) = \min\{ \ell : \text{there exists a path } \{\tau_1, \dots, \tau_\ell\} \text{ from } \sigma_1 \text{ to } \sigma_2 \} .$$

Theorem 74 says that for each $\sigma_1, \sigma_2 \in \mathcal{S}_n$, there is a way to get from σ_1 to σ_2 by making a finite number of steps from one nearest neighbor to another, and that $\varrho(\sigma_1, \sigma_2)$ is the *least* number of such steps in which this can be done. The following lemma is the key to the proof.

Lemma 20 (Reduction lemma). *For all $\sigma \in \mathcal{S}_n$ except the identity, there is some k with $1 \leq k \leq n-1$ such that $\sigma(k) > \sigma(k+1)$. For any such k , let τ be the adjacent pair permutation $\tau = \sigma_{k,k+1}$. Then*

$$D(\sigma \circ \tau) = D(\sigma) - 1 .$$

Proof: Suppose for each $i = 1, \dots, n-1$, $\sigma(i+1) > \sigma(i)$. Then

$$\sigma(1) < \sigma(2) < \dots < \sigma(n) .$$

The only order preserving permutation is the identity, and since σ is not the identity, there is some $k \in \{1, \dots, n-1\}$ such that $\sigma(k) > \sigma(k+1)$. Let τ denote any adjacent pair permutation $\sigma_{k,k+1}$ for some such value of k .

Define the following sets of ordered pair (i, j) :

$$\begin{aligned} A &:= \{ (i, j) : i < k, j > k+1 \} \\ B &:= \{ (i, j) : j = k \text{ or } k+1, j > k+1 \} \\ C &:= \{ (i, j) : i < k, j = k \text{ or } k+1 \} . \end{aligned}$$

The sets A, B, C are disjoint from each other and from $\{(k, k+1)\}$, and $A \cup B \cup C \cup \{(k, k+1)\}$ is the set of all ordered pairs (i, j) with $i < j$.

Note that for $(i, j) \in A$, $(\sigma(i), \sigma(j)) = (\sigma \circ \tau(i), \sigma \circ \tau(j))$. Hence the image of A under f_σ is the same as the image of A under $f_{\sigma \circ \tau}$, and so σ and $\sigma \circ \tau$ reverse the same number of pairs in A .

Note that for $(i, j) \in B$,

$$(\sigma(i), \sigma(k)) = (\sigma \circ \tau(i), \sigma \circ \tau(k+1)) \quad \text{and} \quad (\sigma(i), \sigma(k+1)) = (\sigma \circ \tau(i), \sigma \circ \tau(k)) .$$

Hence the image of B under f_σ is the same as the image of B under $f_{\sigma \circ \tau}$, and so σ and $\sigma \circ \tau$ reverse the same number of pairs in B .

Note that for $(i, j) \in C$,

$$(\sigma(k), \sigma(j)) = (\sigma \circ \tau(k+1), \sigma \circ \tau(j)) \quad \text{and} \quad (\sigma(k+1), \sigma(j)) = (\sigma \circ \tau(k), \sigma \circ \tau(j)) .$$

Hence the image of C under f_σ is the same as the image of C under $f_{\sigma \circ \tau}$, and so σ and $\sigma \circ \tau$ reverse the same number of pairs in C .

Finally, by the choice of k , σ reverses $(k, k+1)$, but then by the definition of τ , $\sigma \circ \tau$ does not. Hence $\sigma \circ \tau$ reverses exactly one fewer pair than does σ . \square

Proof of Theorem 74: First, suppose that $\{\tau_i, \dots, \tau_\ell\}$ is a path of length ℓ from σ_1 to σ_2 . Then $\sigma_2 = \sigma_1 \circ \tau_1 \circ \dots \circ \tau_\ell$. Therefore, $\sigma_1^{-1} \sigma_2 = \tau_1 \circ \dots \circ \tau_\ell$ and so

$$D(\sigma_1^{-1} \sigma_2) = D(\tau_1 \circ \dots \circ \tau_\ell) .$$

Then by Lemma 19

$$\begin{aligned} D(\tau_1 \circ \cdots \circ \tau_\ell) &\leq D(\tau_1) + D(\tau_2 \circ \cdots \circ \tau_\ell) \\ &= 1 + D(\tau_2 \circ \cdots \circ \tau_\ell) \end{aligned}$$

since $D(\tau) = 1$ for any adjacent pair permutation. Proceeding inductively, we find

$$D(\tau_1 \circ \cdots \circ \tau_\ell) \leq \ell .$$

Hence, any path from σ_1 to σ_2 takes at least $D(\sigma_1^{-1}\sigma_2)$ steps.

On the other hand, by Lemma 20, as long as $\sigma_1 \neq \sigma_2$, or what is the same, $D(\sigma_2^{-1} \circ \sigma_1) \neq 0$, there exists an adjacent pair permutation τ_1 such that

$$D(\sigma_2^{-1} \circ \sigma_1 \circ \tau_1) = D(\sigma_2^{-1} \circ \sigma_1) - 1 .$$

Next as long as $D(\sigma_2^{-1} \circ \sigma_1 \circ \tau_1) \neq 0$, there exists an adjacent pair permutation τ_2 such that

$$\begin{aligned} D(\sigma_2^{-1} \circ \sigma_1 \circ \tau_1 \circ \tau_2) &= D(\sigma_2^{-1} \circ \sigma_1 \circ \tau_1) - 1 \\ &= D(\sigma_2^{-1} \circ \sigma_1) - 2 . \end{aligned}$$

Continuing this way, we find a sequence $\{\tau_1, \dots, \tau_{D(\sigma_2^{-1} \circ \sigma_1)}\}$ adjacent pair permutations such that

$$D(\sigma_2^{-1} \circ \sigma_1 \circ \tau_1 \circ \cdots \circ \tau_{D(\sigma_2^{-1} \circ \sigma_1)}) = 0 .$$

But this means that $\sigma_2^{-1} \circ \sigma_1 \circ \tau_1 \circ \cdots \circ \tau_{D(\sigma_2^{-1} \circ \sigma_1)}$ is the identity, and therefore,

$$\sigma_2 = \sigma_1 \circ \tau_1 \circ \cdots \circ \tau_{D(\sigma_2^{-1} \circ \sigma_1)} .$$

Hence, there exists a path from σ_1 to σ_2 of length $D(\sigma_2^{-1} \circ \sigma_1)$. Note that $(\sigma_2^{-1} \circ \sigma_1)^{-1} = \sigma_1^{-1} \circ \sigma_2$, and then by Lemma 18,

$$D(\sigma_2^{-1} \circ \sigma_1) = D(\sigma_1^{-1} \circ \sigma_2) .$$

Therefore, there exists a path from σ_1 to σ_2 consisting of $D(\sigma_1^{-1} \circ \sigma_2)$ steps.

By what we have proved above, this is the least number of steps taken in any path from σ_1 to σ_2 . \square

8.2 Algebraic properties of the determinant

8.2.1 The determinant formula

We are going to break down the formula for the determinant into “building blocks”. The building blocks will be two simple functions that we will combine to form the determinant function. The first one is the character function on the permutations. Here is the second one:

Definition 85 (The function $A \mapsto \sigma(A)$). *For any $n \times n$ matrix A , and any permutation σ on $\{1, \dots, n\}$, define the number $\sigma(A)$ by*

$$\sigma(A) := A_{\sigma(1),1} A_{\sigma(2),2} \cdots A_{\sigma(n),n} = \prod_{j=1}^n A_{\sigma(j),j} . \quad (8.10)$$

Definition 86 (The determinant function). *The determinant function $\det(A)$ on the set of $n \times n$ matrices is defined by*

$$\det(A) = \sum_{\sigma \in \mathcal{S}_n} \chi(\sigma) \sigma(A) . \quad (8.11)$$

Let us first check that this definition gives us what we expect for $n = 2$ and $n = 3$.

Example 124 (2×2 determinants). *Consider the general 2×2 matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$. There are only two permutations of $\{1, 2\}$ to consider, namely*

$$\sigma_1 = \begin{array}{cc} 1 & 2 \\ 1 & 2 \end{array} \quad \text{and} \quad \sigma_2 = \begin{array}{cc} 1 & 2 \\ 2 & 1 \end{array} .$$

Clearly $\chi(\sigma_1) = 1$ and $\chi(\sigma_2) = -1$. Hence $\det(A) = A_{1,1}A_{2,2} - A_{2,1}A_{1,2} = ad - bc$, which is the usual formula.

Example 125 (3×3 determinants). *Consider a general 3×3 matrix A . We have already worked out a list of the six permutations of $\{1, 2, 3\}$ in (8.1) of the previous section, and computed the characters of each of them. In the 3×3 case then, the definition (8.11) leads to*

$$\begin{aligned} \det(A) &= A_{1,1}A_{2,2}A_{3,3} + A_{2,1}A_{3,2}A_{1,3} + A_{3,1}A_{1,2}A_{2,3} \\ &\quad - A_{2,1}A_{1,2}A_{3,3} - A_{1,1}A_{3,2}A_{2,3} - A_{3,1}A_{2,2}A_{1,3} . \end{aligned}$$

This too is reassuring – the formula (8.11) leads us to the usual formula for 3×3 determinants.

Theorem 75 (Characteristic properties of the determinant). *Let \det be the numerically valued function on the $n \times n$ matrices defined by (8.11). Then:*

- (1) $\det(A)$ changes sign when any two rows of A are interchanged.
- (2) $\det(A)$ is linear in each row of A .
- (3) $\det(I_{n \times n}) = 1$, where $I_{n \times n}$ denotes the $n \times n$ identity matrix.

Moreover, these three properties characterize the determinant: the function $A \mapsto \det(A)$ is the only function on the $n \times n$ matrices with the three properties (1), (2) and (3).

Proof: To prove (1), suppose that B is obtained from A by interchanging the k th and ℓ th rows of A . Then we have to show that $\det(B) = -\det(A)$.

To see this, note that $B_{i,j} = A_{\sigma_{k,\ell}(i),j}$, and hence, for any permutation σ ,

$$\sigma(B) = (\sigma \circ \sigma_{k,\ell})(A)$$

Since $\sigma_{k,\ell}$ is a pair permutation,

$$\chi(\sigma \circ \sigma_{k,\ell}) = -\chi(\sigma) .$$

Therefore,

$$\det(B) = \sum_{\sigma} \chi(\sigma) \sigma(B) = - \sum_{\sigma} \chi(\sigma \circ \sigma_{k,\ell}) (\sigma \circ \sigma_{k,\ell})(A) . \quad (8.12)$$

Let τ denote the permutation $\tau := \sigma \circ \sigma_{k,\ell}$. Since $\sigma_{k,\ell}$ is its own inverse, $\sigma = \tau \circ \sigma_{k,\ell}$. That is, the map $\sigma \mapsto \tau := \sigma \circ \sigma_{k,\ell}$ is a one-to-one map of the set of permutations on $\{1, \dots, n\}$. Hence

$$\sum_{\sigma} \chi(\sigma \circ \sigma_{k,\ell})(\sigma \circ \sigma_{k,\ell})(A) = \sum_{\tau} \chi(\tau)\tau(A) = \det(A). \quad (8.13)$$

(In the sum on the middle, we are summing over τ instead of σ , but τ is just a “dummy” variable; we are summing over *all* permutations of on $\{1, \dots, n\}$. Hence $\sum_{\tau} \chi(\tau)\tau(A) = \det(A)$). Combining (8.12) and (8.13) we have $\det(B) = -\det(A)$, and this proves (1).

To prove (2), we have to show that if

$$\mathbf{r}_i = \alpha \mathbf{v} + \beta \mathbf{w}$$

then

$$\det \begin{pmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \alpha \mathbf{v} + \beta \mathbf{w} \\ \vdots \\ \mathbf{r}_n \end{pmatrix} = \alpha \det \begin{pmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \mathbf{v} \\ \vdots \\ \mathbf{r}_n \end{pmatrix} + \beta \det \begin{pmatrix} \mathbf{r}_1 \\ \mathbf{r}_2 \\ \vdots \\ \mathbf{w} \\ \vdots \\ \mathbf{r}_n \end{pmatrix}. \quad (8.14)$$

This is true since each product $\sigma(A) = A_{\sigma(1),1} A_{\sigma(2),2} \cdots A_{\sigma(n),n}$ contains exactly one factor coming from the i th row, and hence is a linear function of the entries of the i th row. By definition, $\det(A)$ is a linear combination of the $\sigma(A)$. A linear combination of linear functions is linear, and so the determinant is a linear function of the entries of the i th row.

To prove (3), note that by the definition of $\sigma(I)$, $\sigma(I) = 1$ if σ is the identity permutation, and $\sigma(I) = 0$ otherwise. Hence (3) follows from the formula the determinant.

To prove the uniqueness, let f be any function on the $n \times n$ matrices that has properties (1), (2) and (3). Let A be any $n \times n$ matrix. Recall that by subtracting multiples of one row from another, and perhaps swapping rows, in a finite number of steps we can transform A into a matrix B that is in row echelon form. Since f has properties (1) and (2), by what we have proved so far, $\det(B) = (-1)^\ell \det(A)$ where ℓ is the number of row swaps used in transforming A into B .

Since B is a square matrix in row echelon form, either all of its diagonal entries are non-zero, or else the bottom row (at least) of B is $\mathbf{0}$. Then, since f has property (2), multiplying the bottom row of B by 2 doubles $f(B)$. But since the bottom row of B is $\mathbf{0}$, multiplying the bottom row of B by 2 does not change B , and therefore does not change $f(B)$. The only number that is its own double is 0, and so $f(B) = 0$ whenever any of its diagonal entries is zero.

If none of the diagonal entries of B is zero, we can do further row operations to “clean out” the part of B above the diagonal, leaving us with a diagonal matrix D . Since f has properties (1) and (2), This does not change the value of f , and so

$$f(A) = (-1)^\ell f(B) = (-1)^\ell f(D)$$

where D is the diagonal matrix whose j th diagonal entry is $B_{j,j}$. But then by the linearity of f in

each row, we can pull out these factors from each row, leaving us with

$$f(A) = (-1)^\ell \prod_{j=1}^n B_{j,j} f(I_{n \times n}) . \quad (8.15)$$

Therefore, given that $f(I_{n \times n}) = 1$, we have $f(A) = (-1)^\ell \prod_{j=1}^n B_{j,j}$ under the assumption that none of the diagonal entries of B is zero. However, when any of the diagonal entries is zero, the product of the diagonal entries is zero, and also as we have explained above, $f(A) = 0$. Hence (8.15) is valid without any restrictions.

This gives us a formula for $f(A)$, showing how to compute it in terms of a reduction of A to row-echelon form. Since for *any* function f with properties (1), (2) and (3) this formula gives the value of f , there is at most one such function. Since $\det(A)$ is such a function, $f(A) = \det(A)$. \square

Example 126 (Computing determinants using row operations). *The formula (8.15) that we have encountered in the proof of Theorem 75 is very useful for computing determinants, especially for larger values of n . All one has to do is to reduce A to row echelon form with a finite sequence of row operations, keeping track of the number of row swaps that is used. For example, consider the matrix*

$$A = \begin{bmatrix} 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix} .$$

Then subtracting multiples of one row from another, we transform

$$A \rightarrow \begin{bmatrix} 1 & 2 & 4 \\ 0 & 1 & 5 \\ 0 & 2 & 12 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 2 & 4 \\ 0 & 1 & 5 \\ 0 & 0 & 2 \end{bmatrix}$$

By (3), the determinant of the upper triangular matrix on the right is 2. But since our row operations did not change the value of the determinant, this is also the value of $\det(A)$. Hence $\det(A) = 2$. You can readily check that this is what the usual formula gives as well.

8.2.2 Algebraic properties of the determinant function

Theorem 76 (Determinants and invertibility). *Let A be any $n \times n$ matrix. Then A is invertible if and only if $\det(A) \neq 0$.*

Proof: This follows immediately from the formula (8.15) that we have derived in the proof of Theorem 75. As we have seen, the rank of A is n if and only if all of the diagonal entries of B are non zero. Thus, by (8.15), the rank of A is n if and only if $\det(A) \neq 0$. But the rank of A is n if and only if A is invertible. \square

The uniqueness part of Theorem 75 has an important consequence:

Theorem 77 (Product property of the determinant). *Let A and B be any $n \times n$ matrices. Then*

$$\det(AB) = \det(A) \det(B) .$$

Proof: If $\det(B) = 0$, then B is not invertible. When linear transformations are not invertible, they are *neither* one to one nor onto, so the transformation generated by B is not one to one, and hence neither is AB . So AB is not invertible, and thus $\det(AB) = 0$.

It remains to consider the case $\det(B) \neq 0$. Fix such a matrix B , and define a function f_B on the $n \times n$ matrices by

$$f_B(A) = \frac{\det(AB)}{\det(B)} .$$

It is easy to see that swapping two rows of A swaps the same two rows of AB , and subtracting a multiple of one row of A from another results in and subtracting the same multiple of one row of AB from another – the same rows also. So f_B has properties (1) and (2).

Next, by the previous theorem, if A is not invertible, so that neither is AB , $\det(AB) = 0$, and hence by the definition of f_B , $f_B(A) = 0$. Also by the definition of f_B , $f_B(I_{n \times n}) = 1$. Thus, f_B has the property (3'). By uniqueness part of Theorem 76, this means $f_B(A) = \det(A)$. But by the definition of f_B , this means $\det(AB) = \det(A) \det(B)$, □

It may seem that we have focused on the rows as opposed to the columns in our definition, but this is not the case:

Theorem 78 (Invariance of the determinant under the transpose). *Let A be an $n \times n$ matrix and A^T its transpose. Then*

$$\det(A^T) = \det(A) .$$

In particular, since the transpose operation swaps rows and columns, $\det(A)$ is linear in the columns of A , and changes sign when two columns of A are swapped.

Proof: Let τ be any permutation on $\{1, \dots, n\}$. For any n numbers a_1, \dots, a_n ,

$$\prod_{j=1}^n a_j = \prod_{j=1}^n a_{\tau(j)} :$$

The only difference between the products on the left and the right is that we are doing the multiplication in a different order, but since multiplication is commutative, the order does not matter.

Therefore, for any two permutations σ, τ be any permutation on $\{1, \dots, n\}$, and any $n \times n$ matrix A ,

$$\sigma(A) = \prod_{j=1}^n A_{\sigma(j), j} = \prod_{j=1}^n A_{\sigma(\tau(j)), \tau(j)} .$$

Now taking $\tau = \sigma^{-1}$, we have

$$\sigma(A) = \prod_{j=1}^n A_{j, \tau(j)} = \prod_{j=1}^n A_{\tau(j), j}^T = \tau(A^T) .$$

since the character of the identity permutation is 1, for $\tau = \sigma^{-1}$, $\chi(\sigma \circ \tau) = 1$ and so $\chi(\tau) = \chi(\sigma)$. Therefore,

$$\det(A) = \sum_{\sigma} \chi(\sigma) \sigma(A) = \sum_{\tau} \chi(\tau) \tau(A^T) = \det(A^T) .$$

This proves the theorem. □

8.3 Geometric properties of the determinant

In our study of integration in \mathbb{R}^2 and \mathbb{R}^3 , we have seen that the determinant of a 2×2 matrix A gives the *area magnification factor* of the linear transformation associated to A , and that the determinant of a 3×3 matrix A gives the *volume magnification factor* of the linear transformation associated to A .

In this section, we present another proof of these results that makes use of the algebraic properties of the determinant that we have proved in the previous section. This approach has the advantage that it yields analogous results $n \times n$ matrices. The key to all of this is a *factorization result* for $m \times n$ matrices called the *singular value decomposition*.

The Singular Value Decomposition Theorem allows us to express *any* $m \times n$ matrix as the product of three simple matrices: One will be an $m \times n$ diagonal matrix; i.e., a matrix whose i, j th entry is zero when $i \neq j$. Moreover, the diagonal entries will be non-negative. Such matrices describe very simple transformations!

The other two factors in the decomposition will be *orthogonal matrices*. These too are very simple once one is familiar with their basic properties. However, these basic properties do require an introduction. That is provided in the next subsection.

8.3.1 Orthogonal matrices

Definition 87 (Orthogonal matrix). *An $n \times n$ matrix $U = [\mathbf{u}_1, \dots, \mathbf{u}_n]$ is an orthogonal matrix if and only if $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ is an orthonormal basis of \mathbb{R}^n .*

Theorem 79 (Characteristic properties of orthogonal matrices). *Let U be an $n \times n$ matrix. Then the following are equivalent:*

- (1) U is orthogonal.
- (2) $\|U\mathbf{x}\| = \|\mathbf{x}\|$ for all $\mathbf{x} \in \mathbb{R}^n$.
- (3) $\mathbf{x} \cdot \mathbf{y} = (U\mathbf{x}) \cdot (U\mathbf{y})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,
- (4) U is invertible, and the inverse of U is the transpose of U ; i.e., $U^{-1} = U^T$.
- (5) The rows of U , $\{\mathbf{r}_1, \dots, \mathbf{r}_n\}$, are an orthonormal basis of \mathbb{R}^n .

Proof: Suppose that (1) is true so that $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ is an orthonormal basis of \mathbb{R}^n . Then for any $\mathbf{x} \in \mathbb{R}^n$, $U\mathbf{x} = \sum_{j=1}^n x_j \mathbf{u}_j$ and

$$\|U\mathbf{x}\|^2 = \left\| \sum_{j=1}^n x_j \mathbf{u}_j \right\|^2 = \sum_{j=1}^n x_j^2 = \|\mathbf{x}\|^2.$$

Therefore, (1) \Rightarrow (2).

Now assume that (2) is true.

$$\|U(\mathbf{x} + \mathbf{y})\|^2 = (U\mathbf{x} + U\mathbf{y}) \cdot (U\mathbf{x} + U\mathbf{y}) = \|U\mathbf{x}\|^2 + 2(U\mathbf{x}) \cdot (U\mathbf{y}) + \|U\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2(U\mathbf{x}) \cdot (U\mathbf{y}) + \|\mathbf{y}\|^2.$$

Also by assumption, $\|U(\mathbf{x} + \mathbf{y})\|^2 = \|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2\mathbf{x} \cdot \mathbf{y} + \|\mathbf{y}\|^2$. Comparing calculations, we see $\|\mathbf{x}\|^2 + 2(U\mathbf{x}) \cdot (U\mathbf{y}) + \|\mathbf{y}\|^2 = \|\mathbf{x}\|^2 + 2\mathbf{x} \cdot \mathbf{y} + \|\mathbf{y}\|^2$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. Cancelling $\|\mathbf{x}\|^2$ and $\|\mathbf{y}\|^2$ off both sides, we see that (2) \Rightarrow (3).

On the other hand, when (3) is true,

$$\mathbf{u}_i \cdot \mathbf{u}_j = (U\mathbf{e}_i) \cdot (U\mathbf{e}_j) = \mathbf{e}_i \cdot \mathbf{e}_j = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}.$$

Thus (3) \Rightarrow (1). At this point we have proved that (1), (2) and (3) are equivalent.

Now assume that U is orthogonal, and hence (2), is true. To show that U is invertible, suppose that $U\mathbf{x} = U\mathbf{y}$, and so $0 = \|U(\mathbf{x} - \mathbf{y})\| = \|\mathbf{x} - \mathbf{y}\|$, using (2). This implies that $\mathbf{x} = \mathbf{y}$, and so the linear transformation described by U is one-to-one. By the Fundamental Theorem of Linear algebra, it is also onto, and hence U is invertible.

Moreover, combining (3) with the fundamental property of the transpose, for all \mathbf{x}, \mathbf{y}

$$\mathbf{x} \cdot \mathbf{y} = U\mathbf{x} \cdot U\mathbf{y} = \mathbf{x} \cdot U^T U\mathbf{y}.$$

Therefore,

$$\mathbf{x} \cdot (I - U^T U)\mathbf{y} = 0$$

for all \mathbf{x}, \mathbf{y} . Taking $\mathbf{x} = (I - U^T U)\mathbf{y}$, we see $\|(I - U^T U)\mathbf{y}\| = 0$ for all \mathbf{y} , and hence $(I - U^T U)\mathbf{y} = \mathbf{0}$ for all \mathbf{y} . This means that $U^T U\mathbf{y} = \mathbf{y}$, and since U is invertible, this means that $U^{-1} = U^T$. Thus when U is orthogonal, (4) is true.

Next, suppose that (4) is true. Then, by the fundamental property of the transpose, for all \mathbf{x}, \mathbf{y}

$$U\mathbf{x} \cdot U\mathbf{y} = \mathbf{x} \cdot U^T U\mathbf{y} = \mathbf{x} \cdot U^{-1} U\mathbf{y} = \mathbf{x} \cdot \mathbf{y}$$

so that (4) \Rightarrow (3), and hence that U is orthogonal. Hence, (4) is equivalent to U being orthogonal.

Finally, suppose that U is orthogonal. Then since (2) is true, for all \mathbf{x} ,

$$\|\mathbf{x}\| = \|UU^{-1}\mathbf{x}\| = \|U(U^{-1}\mathbf{x})\| = \|U^{-1}\mathbf{x}\|.$$

This shows that U^{-1} has the property in (2), and by what we have proved above this means that U is orthogonal. But since U is orthogonal, $U^{-1} = U^T$, and so U^T is orthogonal. Hence the columns of U^T are orthonormal. But the columns of U^T are the rows of U , and so (5) is true for all orthogonal matrices U .

Conversely, suppose that (5) is true. Then U^T has orthonormal columns and is therefore orthogonal. But by what we just proved, every orthogonal matrix also has orthonormal rows. Hence the rows of U^T are orthonormal, and this means that the columns of U are orthonormal. Hence U is orthogonal whenever (5) is true. \square

One consequence of Theorem 79 is that the set of $n \times n$ orthogonal matrices, regarded as a set of transformations of \mathbb{R}^n is a *transformation group*. Indeed, whenever U is an orthogonal matrix, its inverse U^T is also an orthogonal matrix by Theorem 79: Since both the columns and rows of U are orthonormal, and since the inverse of U is U^T , both the columns and rows of U^{-1} are orthonormal,

and hence U^{-1} is orthogonal. Furthermore, if U and V are any two $n \times n$ orthogonal matrices, then for all $\mathbf{x} \in \mathbb{R}^n$,

$$\|(VU)\mathbf{x}\| = \|V(U\mathbf{x})\| = \|U\mathbf{x}\| = \|\mathbf{x}\| .$$

Hence the product VU has property (2) and by Theorem 79 is orthogonal. Thus, the product of any two orthogonal matrices is again an orthogonal matrix. Since the matrix product corresponds to the composition of the corresponding linear transformation of \mathbb{R}^n , this shows that the set of all orthogonal matrices, regarded as a set of transformations of \mathbb{R}^n is a *transformation group*.

Definition 88 (The orthogonal group on \mathbb{R}^n). *The set of all $n \times n$ orthogonal matrices is called the orthogonal group on \mathbb{R}^n , and is denoted by $O(n)$.*

Theorem 80 (Determinants of orthogonal matrices). *Let $U \in O(n)$. Then*

$$\det(U) = \pm 1 .$$

The set of matrices $U \in O(n)$ such that $\det(U) = 1$, regarded as a set of transformations of \mathbb{R}^n , forms a transformation group on \mathbb{R}^n .

Proof: For all $U \in O(n)$, $I = U^T U$. This

$$1 = \det(I) = \det(U^T U) = \det(U^T) \det(U) = (\det(U))^2 ,$$

where we have used Theorems 77 and 78. The only solutions of the equation $x^2 = 1$ are ± 1 , and so $\det(U) = \pm 1$.

Now suppose $\det(U) = 1$. Then by Theorem 79 and Theorem 78,

$$\det(U^{-1}) = \det(U^T) = \det(U) = 1 .$$

Hence the inverse of U has the same property. Next, let $V, U \in O(n)$ be such that $\det(V) = \det(U) = 1$. Then by Theorem 77,

$$\det(VU) = \det(V) \det(U) = 1 .$$

Hence the subset of $O(n)$ consisting of matrices with unit determinant is closed under taking inverses and products. It is therefore a transformation group, and, as such, a *subgroup* of $O(n)$.

Definition 89 (The special orthogonal group on \mathbb{R}^n). *The subset of $O(n)$ consisting of orthogonal matrices U with $\det(U) = 1$ is called the special orthogonal group on \mathbb{R}^n , and is denoted by $SO(n)$.*

Example 127 (Two dimensional orthogonal matrices). *Let $U = [\mathbf{u}_1, \mathbf{u}_2] \in O(2)$. Then \mathbf{u}_1 is a unit vector. Hence*

$$\mathbf{u}_1 = (\cos \theta, \sin \theta)$$

for some θ . Since \mathbf{u}_2 must be a unit vector orthogonal to \mathbf{u}_1 , there are only two choices for \mathbf{u}_2 :

$$\mathbf{u}_2 = (-\sin \theta, \cos \theta) \quad \text{or else} \quad \mathbf{u}_2 = (\sin \theta, -\cos \theta) .$$

Thus either we have

$$U = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \quad \text{or else} \quad U = \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix} . \quad (8.16)$$

Note that

$$\det \left(\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \right) = 1 \quad \text{and} \quad \det \left(\begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix} \right) = -1 .$$

Thus, the matrices in $SO(2)$ are precisely the matrices on the left in (8.16), and we recognize these as the two dimensional rotation matrices.

The matrices on the right in (8.16) reflection matrices. Indeed, let $\mathbf{u} = (\cos(\theta/2), \sin(\theta/2))$. Then the Householder reflection in \mathbb{R}^2 given by \mathbf{u} has the matrix

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 2 \begin{bmatrix} \cos^2(\theta/2) & \cos(\theta/2)\sin(\theta/2) \\ \cos(\theta/2)\sin(\theta/2) & \sin^2(\theta/2) \end{bmatrix} = \begin{bmatrix} \cos \theta & \sin \theta \\ \sin \theta & -\cos \theta \end{bmatrix} ,$$

by the double-angle formulas. Thus the matrices in $O(n)$ that are not in $SO(n)$ are precisely the reflection matrices.

Example 128 (Three dimensional orthogonal matrices). Let $U = [\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3] \in O(3)$, so that, by definition, $\{\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3\}$ is an orthonormal basis of \mathbb{R}^3 .

We have seen in Chapter One that there is a linear transformation \mathbf{f} from \mathbb{R}^3 to \mathbb{R}^3 that is the composition of at most 3 Householder reflections such that $(\mathbf{e}_j) = \mathbf{u}_j$ for $j = 1, 2, 3$. This means that the matrix representing \mathbf{f} is the matrix U , and hence U is the product of at most three matrices representing Householder reflections. We have seen that whenever $\mathbf{h}_{\mathbf{u}}$ is a Householder reflection, and $H_{\mathbf{u}} := [\mathbf{h}_{\mathbf{u}}(\mathbf{e}_1), \mathbf{h}_{\mathbf{u}}(\mathbf{e}_2), \mathbf{h}_{\mathbf{u}}(\mathbf{e}_3)]$ is the 3×3 matrix representing $\mathbf{h}_{\mathbf{u}}$,

$$\det(H_{\mathbf{u}}) = \det([\mathbf{h}_{\mathbf{u}}(\mathbf{e}_1), \mathbf{h}_{\mathbf{u}}(\mathbf{e}_2), \mathbf{h}_{\mathbf{u}}(\mathbf{e}_3)]) = \mathbf{h}_{\mathbf{u}}(\mathbf{e}_1) \cdot \mathbf{h}_{\mathbf{u}}(\mathbf{e}_2) \times \mathbf{h}_{\mathbf{u}}(\mathbf{e}_3) = -1 .$$

Now suppose U is not the identity matrix. Then U is the product of either 1, 2 or 3 Householder reflection matrices. Since the determinant of each of these is -1 , by Theorem 77, $\det(U) = 1$ if and only if U is the product of exactly 2 Householder reflection matrices.

As we have seen in Chapter Two, the product of any two Householder reflections is a rotation: Each Householder reflection leaves a plane through the origin - the plane of reflection - unchanged. The two planes of reflection meet in a line through the origin which is left unchanged by the composition of the two reflections. This line is the axis of rotation. Thus, $SO(3)$ consists of the 3×3 rotation matrices. Every matrix $U \in O(3)$ that is not in $SO(3)$ is the product of some matrix in $SO(3)$ and a Householder reflection matrix.

8.3.2 Orthogonal matrices, area, volume and shape

If $U \in O(2)$, the action of U on \mathbb{R}^2 preserves the area of subsets of \mathbb{R}^2 : We have proved that the area magnification factor of the linear transformation given by any 2×2 matrix A is $|\det(A)|$. Since $|\det(U)| = 1$ for all $U \in O(2)$, the area magnification factor associated to U is 1.

However, the action of a matrix in $O(2)$ preserves much more than the area of subsets of \mathbb{R}^2 ; it preserve the distances between each pair of points in the set - there is no distortion of the shape of the set; all that changes is the way it is situated in the plane. For example, Consider the cat-shaped set in the unit square in \mathbb{R}^2 shown below:

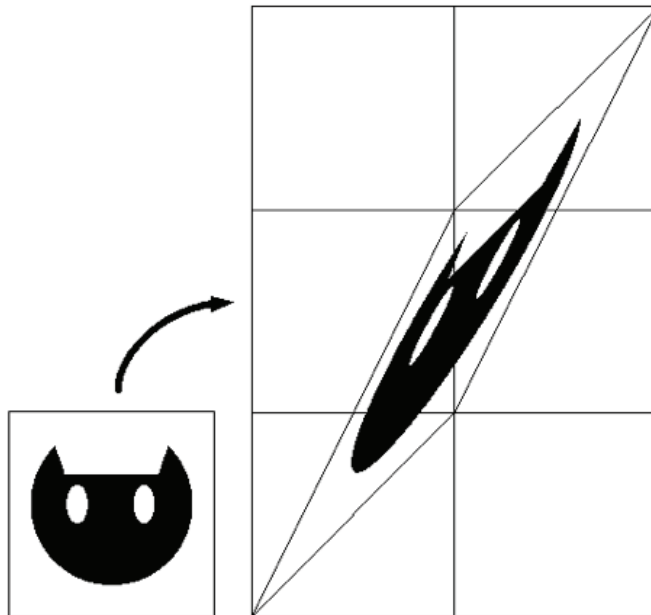


Consider the matrix $U \in O(2)$ given by $U = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$. This matrix U acts on \mathbb{R}^2 by counter-clockwise rotation through the angle $\pi/2$. The image of the cat shaped set under the action of U is therefore



The orthogonal transformation has changed the orientation of the set in the plane but has not distorted it in any way.

Next consider the 2×2 matrix A given by $A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$. Note that $\det(A) = 1$, and so the action of A on \mathbb{R}^2 preserves the area of subsets of \mathbb{R}^2 , but it strongly distorts the shapes of subsets of \mathbb{R}^2 . Here is picture showing the original cat shaped set, and also its image under the action of A :



While this transformation preserves area, it does not preserve distances: Notice that after the transformation, the distance between the tips of the cat's ears is greater than it was before, while the distance from the tip of the left ear to the center of the cat's face is less than it was. Also, the transformation changes angles: After the transformation, the angle in the lower left corner of the

bounding box is much less than it was, while the angle in the lower right corner of the bounding box is much greater than it was.

We now show that transformations of \mathbb{R}^n that preserve distances automatically preserve angles, areas and volumes as well. This brings an important definition, and an important characterization of orthogonal matrices:

Definition 90 (Euclidean transformation). *Let \mathbf{f} be a function from \mathbb{R}^n to \mathbb{R}^n with the property that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,*

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| = \|\mathbf{x} - \mathbf{y}\| . \quad (8.17)$$

Then \mathbf{f} is a Euclidean transformation.

That is, Euclidean transformations preserve the distances between points; there is no “stretching” or “compression” associated to a Euclidean transformation.

Theorem 81 (Euclidean transformation and orthogonal matrices). *A function \mathbf{f} from \mathbb{R}^n to \mathbb{R}^n is a Euclidean transformation if and only if there is an $\mathbf{x}_0 \in \mathbb{R}^n$ and a $U \in O(n)$ so that for all $\mathbf{x} \in \mathbb{R}^n$,*

$$\mathbf{f}(\mathbf{x}) = \mathbf{x}_0 + U\mathbf{x} . \quad (8.18)$$

Proof: Let \mathbf{f} be a Euclidean transformation. Define the transformation \mathbf{g} by

$$\mathbf{g}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{0})$$

and note that since for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y}) = \mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})$, \mathbf{g} is also a Euclidean transformation, and $\mathbf{g}(\mathbf{0}) = (\mathbf{0})$. Thus, for all $\mathbf{x} \in \mathbb{R}^n$,

$$\|\mathbf{x}\| = \|\mathbf{x} - \mathbf{0}\| = \|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{0})\| = \|\mathbf{g}(\mathbf{x})\|$$

so that the transformation \mathbf{g} preserves the lengths of vectors. But then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$.

$$\begin{aligned} 2\mathbf{x} \cdot \mathbf{y} &= \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2 - \|\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{g}(\mathbf{x})\|^2 + \|\mathbf{g}(\mathbf{y})\|^2 - \|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{y})\|^2 \\ &= 2\mathbf{g}(\mathbf{x}) \cdot \mathbf{g}(\mathbf{y}) . \end{aligned} \quad (8.19)$$

Thus,

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{g}(\mathbf{x}) \cdot \mathbf{g}(\mathbf{y}) ,$$

and the transformation \mathbf{g} preserves dot products. Therefore, if we define $\mathbf{u}_j = \mathbf{g}(\mathbf{e}_j)$ for $j = 1, \dots, n$, $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ is an orthonormal basis of \mathbb{R}^n .

Now, for any $\mathbf{x} \in \mathbb{R}^n$,

$$\begin{aligned} \mathbf{g}(\mathbf{x}) &= \sum_{j=1}^n (\mathbf{g}(\mathbf{x}) \cdot \mathbf{u}_j) \mathbf{u}_j = \sum_{j=1}^n (\mathbf{g}(\mathbf{x}) \cdot \mathbf{g}(\mathbf{e}_j)) \mathbf{u}_j \\ &= \sum_{j=1}^n (\mathbf{x} \cdot \mathbf{e}_j) \mathbf{u}_j = \sum_{j=1}^n x_j \mathbf{u}_j \end{aligned}$$

However, if we define the orthogonal matrix U by $U = [\mathbf{u}_1, \dots, \mathbf{u}_n]$, then by the rules of matrix multiplication,

$$U\mathbf{x} = \sum_{j=1}^n x_j \mathbf{u}_j .$$

Therefore, $\mathbf{g}(x) = U\mathbf{x}$ for all \mathbf{x} . But then by the definition of \mathbf{g} , if we define $\mathbf{x}_0 := \mathbf{f}(\mathbf{0})$, it follows that for all \mathbf{x} , $\mathbf{f}(\mathbf{x}) = \mathbf{x}_0 + U\mathbf{x}$.

Conversely, suppose that the transformation \mathbf{f} is given by $\mathbf{f}(\mathbf{x}) = \mathbf{x}_0 + U\mathbf{x}$ where $U \in O(n)$, Then

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| = \|U\mathbf{x} - U\mathbf{y}\| = \|U(\mathbf{x} - \mathbf{y})\| = \|\mathbf{x} - \mathbf{y}\|$$

since the final equality is true since $U \in O(n)$. Thus, \mathbf{f} is a Euclidean transformation. \square

We now claim that Euclidean transformations of \mathbb{R}^n preserve n -dimensional volume. At least in 2 dimensions, this is intuitively clear: By the previous theorem, a Euclidean transformation is either a reflection of a rotation, followed by a translation. None of these operations affects area, which is what we mean by 2-dimensional volume.

The main difficulty in justifying our claim is to make precise sense of what we mean by volume in \mathbb{R}^n . To discuss the notion of the volume of sets in \mathbb{R}^n in full generality would take us into a study of Lebesgue measure, and that is beyond the scope of our discussion. However, there is a well defined notion of the n -dimensional volume of sets in \mathbb{R}^n that is valid for all closed bounded sets in \mathbb{R}^n . The n -dimensional volume function is *additive*, meaning that if $A = B \cup C$ is a union of two disjoint sets for which the volume is defined, then

$$\text{volume}(A) = \text{volume}(B) + \text{volume}(C) .$$

Suppose that A is any cube of side length h in \mathbb{R}^n ; i.e., any set of the form

$$\mathbf{x}_0 + \sum_{j=1}^n s_j \mathbf{u}_j \quad \text{where} \quad 0 \leq s_j \leq h \quad \text{for all} \quad j = 1, \dots, n ,$$

where \mathbf{x}_0 is one corner of the cube. Then the n -dimensional volume function is defined so that

$$\text{volume}(A) = h^n .$$

Then by the additivity that we have just discussed, any set A that can be decomposed as a disjoint union of N cubes of side length h satisfies

$$\text{volume}(A) = Nh^n .$$

Now let $U \in O(n)$, Then the image of any cube of side length h in \mathbb{R}^n under the transformation given by U is again a cube of side length h : The transformation preserves dot products, and so it preserves the lengths and orthogonality of the edges. Hence if A is any disjoint union of N cubes of side length h , so is its image under the transformation U ; i.e., the set $U(A) := \{ U\mathbf{x} : \mathbf{x} \in A \}$. That is,

$$\text{volume}(U(A)) = Nh^3 = \text{volume}(A) .$$

In the theory of the Lebesgue measure, one proves that every subset A of \mathbb{R}^n for which the volume can be defined and is finite, A can be approximated up to a small error in the volume by

a finite disjoint union of cubes of some small side length h . Because of this, roughly speaking, the volume of A is the limit as h tends to 0 of h^n times the number of disjoint cubes of side length h that can be packed into A . Therefore, one has that:

- Whenever $U \in O(n)$, and whenever $A \subset \mathbb{R}^n$ has a well defined volume, the image of A under the transformation U ; i.e., $U(A)$, has the same volume as A .

8.3.3 The Spectral Theorem as a factorization theorem

Before introducing the singular value decomposition, we explain how a theorem with which we are already familiar, namely the Spectral Theorem, provides a factorization of an arbitrary $n \times n$ symmetric matrix A into the product of three factors,

$$A = U\Lambda U^T,$$

where $U = [\mathbf{u}_1, \dots, \mathbf{u}_n]$ is a matrix whose columns are an orthonormal basis for \mathbb{R}^n consisting of eigenvectors of A – the Spectral Theorem assures us that such a basis exists – and Λ is an $n \times n$ diagonal matrix whose j th diagonal entry is the eigenvalue λ_j corresponding to the eigenvector \mathbf{u}_j . Finally, U^T is the transpose of U .

As we shall see, each of these factors, U , Λ and U^T , has a simple geometric interpretation that helps us understand the geometric nature of the linear transformation associated to A .

Theorem 82 (Diagonalization of symmetric matrices). *Let A be a symmetric $n \times n$ matrix. Let $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ be an orthonormal basis of \mathbb{R}^n consisting of eigenvectors of A , recalling that such a basis always exists. Let Λ be the $n \times n$ diagonal matrix whose j th diagonal entry of Λ is the eigenvalue of A corresponding to \mathbf{u}_j . Then*

$$A = U\Lambda U^T.$$

Conversely, given any $n \times n$ diagonal matrix Λ , and any $U = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in O(n)$, $U\Lambda U^T$ is an $n \times n$ symmetric matrix such that for each j , \mathbf{u}_j is an eigenvector of $U\Lambda U^T$ with eigenvalue λ_j , where λ_j is the j th diagonal entry in Λ .

Proof: We claim that $U^T A U = \Lambda$. Once this is shown, we shall have our factorization of A since multiplying both sides on the left by U and on the right by U^T and using $U^T U = U U^T = I$, we get $A = U\Lambda U^T$.

To justify the claim, recall that for any matrix B , the i, j th entry is given by $B_{i,j} = \mathbf{e}_i \cdot B \mathbf{e}_j$. Hence, using the fundamental property of the transpose and the fact that $U \mathbf{e}_k = \mathbf{u}_k$ for all k ,

$$(U^T A U)_{i,j} = \mathbf{e}_i \cdot U^T A U \mathbf{e}_j = (U \mathbf{e}_i) \cdot A (U \mathbf{e}_j) = \mathbf{u}_i \cdot A \mathbf{u}_j = \lambda_j \mathbf{u}_i \cdot \mathbf{u}_j = \begin{cases} \lambda_j & i = j \\ 0 & i \neq j \end{cases}.$$

This shows that $U^T A U$ equals the diagonal matrix Λ .

For the converse, since $U^T \mathbf{u}_j = \mathbf{e}_j$,

$$(U\Lambda U^T) \mathbf{u}_j = U(\Lambda(U^T \mathbf{u}_j)) = U(\Lambda \mathbf{e}_j) = \lambda_j U \mathbf{e}_j = \lambda_j \mathbf{u}_j.$$

This proves that each \mathbf{u}_j is an eigenvector of $U\Lambda U^T$ with eigenvalue λ_j . Note that

$$(U\Lambda U^T)^T = (U^T)^T \Lambda^T U^T = U\Lambda U^T$$

so $U\Lambda U^T$ is symmetric. □

8.3.4 The Singular Value Decomposition

The following theorem generalizes the factorization of symmetric matrices that is provided by Theorem 82 to general matrices.

Theorem 83 (The Singular Value Decomposition Theorem). *Let A be an $m \times n$ matrix. Then there exist matrices U , V and Σ such that*

$$A = U\Sigma V^T$$

and

$$(1) U \in O(m)$$

$$(3) V \in O(n)$$

(3) Σ is an $m \times n$ diagonal matrix whose j th diagonal entry is σ_j where $\sigma_j \geq \sigma_{j+1} \geq 0$ for all $j = 1, \dots, \min\{m, n\} - 1$.

In any such factorization of A , the matrix Σ is always the same. In particular, the numbers $\{\sigma_1, \dots, \sigma_{\min\{m, n\}}\}$ are uniquely determined by A . We call these numbers the singular values of A .

For example, if $m = 3$ and $n = 4$, the matrix Σ has the form

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \end{bmatrix}.$$

If $m = 4$ and $n = 3$, the matrix Σ has the form

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \\ 0 & 0 & 0 \end{bmatrix}.$$

In both cases,

$$\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq 0.$$

Before we prove this theorem, let us see what it tells us about the geometry of the general invertible linear transformation from \mathbb{R}^n to \mathbb{R}^n . Let A be any invertible $n \times n$ matrix. Let $A = U\Sigma V^T$ be the factorization of it provided by the singular value decomposition. Then Σ is a diagonal matrix whose j th diagonal entry is σ_j . Since U and V are invertible, the fact that A is invertible implies that Σ is invertible. This in turn implies that $\sigma_j > 0$ for each $j = 1, \dots, n$.

In particular, if $n = 2$, and A is invertible, then

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix},$$

where $\sigma_1 \geq \sigma_2 > 0$.

The transformation given by Σ is very easy to understand geometrically: Let $D \subset \mathbb{R}^2$ be the closed unit disc; i.e., the set of points (x, y) such that $x^2 + y^2 \leq 1$. Let \widehat{D} denote the image of D under the transformation given by Σ , and regard \widehat{D} as a set in the u, v plane. Since

$$(u, v) = \Sigma(x, y) = (\sigma_1 x, \sigma_2 y) ,$$

$$\frac{u^2}{\sigma_1^2} + \frac{v^2}{\sigma_2^2} \leq 1 \quad \Longleftrightarrow \quad x^2 + y^2 \leq 1 .$$

This tells us that \widehat{D} is an ellipse centered on the origin whose major axis has length σ_1 , and whose minor axis has length $2\sigma_2$.

Now that we know what the image of the unit disc is under Σ , we ask:

- *What is the image of the unit disc under a general invertible linear transformation from \mathbb{R}^2 to \mathbb{R}^2 ?*

To answer this question, let $A = U\Sigma V^T$ be a singular value decomposition of A . The matrices U and V are both 2×2 orthogonal matrices, and hence by what we have seen in Example 127, they are either 2×2 rotations or 2×2 reflections. Since by Theorem 79, V^T is also an orthogonal matrix, it too is either a rotation or a reflection.

Now, the image of the unit disc under either a rotation or a reflection is again the unit disc, though points in it will generally get moved around. Here is a picture showing the original unit disc with the vectors \mathbf{e}_1 and \mathbf{r}_2 drawn in.

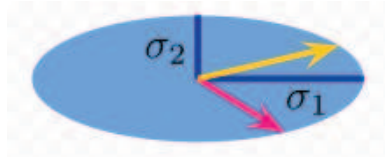


The next picture shows the image under V^T of the unit disc: The image is still the unit disc, but the vectors \mathbf{e}_1 and \mathbf{e}_2 have been rotated or reflected into new positions – in the picture it is a rotation.

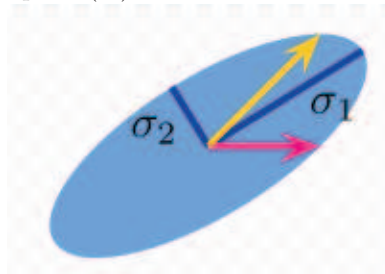


Now apply Σ : This distorts the unit disc into an ellipse whose major axis lies along the first coordinate axis and has length $2\sigma_1$, and whose minor axis lies along the second coordinate axis and

has length $2\sigma_2$. Here is a picture showing the action of Σ on the unit disk:



Now apply U : Since U is either a rotation or reflection, it either rotates or reflects the ellipse, changing only the orientation of the major and minor axes, and not their lengths. Here is a picture showing the action of U on the ellipse $\Sigma(D)$ where D is the unit disc.



Again, in this picture, U is a rotation, and the lengths of the major and minor axes of the ellipse are $2\sigma_1$ and $2\sigma_2$. This reasoning applies to *any* invertible 2×2 matrix A , and hence we conclude, that:

- The image of the unit disc in \mathbb{R}^2 under the action of any 2×2 invertible matrix A is an ellipse, and from the major and minor axes of this ellipse, one can read off the singular values of A .

These conclusions are readily extended to general $n \times n$ matrices. In fact, we have: .

Theorem 84 (Geometric consequences of the singular value decomposition). *Let A be any $m \times n$ matrix. Then the image of the unit ball in \mathbb{R}^n under the action of A is an m -dimensional ellipsoid in \mathbb{R}^m . The directions of the principle axes of the ellipsoid are the columns of U , and the length of the axis with direction \mathbf{u}_j is $2\sigma_j$.*

Finally, when $A = U\Sigma V^T$ is a singular value decomposition of an $n \times n$ matrix, by Theorems 77 and 78,

$$|\det(A)| = |\det(U\Sigma V^T)| = |\det(U)| |\det(\Sigma)| |\det(V)| = |\det(\Sigma)|$$

since $U, V \in O(n)$ so that their determinants are ± 1 , and since Σ is diagonal with all of its entries non-negative

$$\det \Sigma = \prod_{j=1}^n \sigma_j .$$

Now, consider an n -dimensional cube in \mathbb{R}^n whose edges are parallel to the coordinate axes, and whose sides have length h . Since the action of Σ on \mathbb{R}^n is simply stretch or compress distances along lines parallel to the coordinate axes, the image of such a cube under Σ is a rectangular box with edges parallel to the axes and side lengths $\sigma_1 h, \dots, \sigma_n h$. The transformation associated to Σ is often called a *scale transformation* since its action effectively changes the scale along the coordinate axes.

It is easy to understand the effect of the scale transformation Σ on volume. By what we have explained about volume in \mathbb{R}^n , the volume of such a box is the product of the edge lengths; i.e., it is

$$h^n \prod_{j=1}^n \sigma_j = h^n |\det(A)|.$$

Thus, the transformation Σ has the volume magnification factor $|\det(A)|$. But since $A = U\Sigma V^T$ and neither U nor V^T affect volume, this is also the volume magnification factor of A . We summarize:

- For any invertible $n \times n$ matrix A , whenever $E \subset \mathbb{R}^n$ has a well defined volume, the image of E under the transformation A ; i.e., $A(E)$, has the volume

$$\text{volume}(A(E)) = |\det(A)| \text{volume}(E).$$

Proof of Theorem 83: To prove the existence of the singular value decomposition $A = U\Sigma V^T$, it suffices to find an orthonormal bases $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ of \mathbb{R}^m and an orthonormal bases $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ of \mathbb{R}^n and numbers $\sigma_1 \geq \dots \geq \sigma_r > 0$, $r \leq \min\{m, n\}$, such that

$$\mathbf{u}_i \cdot A\mathbf{v}_j = \begin{cases} \sigma_i & i = j \leq r \\ 0 & \text{otherwise} \end{cases}. \quad (8.20)$$

To see this, suppose we have (8.20). Define $U = [\mathbf{u}_1, \dots, \mathbf{u}_m]$ and $V = [\mathbf{v}_1, \dots, \mathbf{v}_n]$. Then U is in $O(m)$ and $V \in O(n)$. Let us compute the i, j th entry of $U^T A V$. This is given by

$$(U^T A V)_{i,j} = \mathbf{e}_i \cdot U^T A V \mathbf{e}_j = (U \mathbf{e}_i) \cdot A(V \mathbf{e}_j) = \mathbf{u}_i \cdot \mathbf{v}_j. \quad (8.21)$$

Define Σ to be the $m \times n$ matrix with

$$\Sigma_{i,j} = \begin{cases} \sigma_i & i = j \leq r \\ 0 & \text{otherwise} \end{cases}.$$

Comparing this with (8.21) and using the assumption (8.20) we conclude that for all i and j , $(U^T A V)_{i,j} = \Sigma_{i,j}$. This of course means that

$$U^T A V = \Sigma.$$

Now, multiplying both sides on the left by U and on the right by V^T and using the fact that $U U^T = I_{m \times m}$ and $V V^T = I_{n \times n}$, we have that $A = U \Sigma V^T$, as desired.

To complete the proof, we now construct the orthonormal bases and the numbers σ_j that figure in (8.20)

Let A be any $m \times n$ matrix. Form the $(m+n) \times (m+n)$ matrices

$$B := \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix} \quad \text{and} \quad M := \begin{bmatrix} -I_{m \times m} & 0 \\ 0 & I_{n \times n} \end{bmatrix}.$$

More explicitly, the 0 entry in the upper left of B denotes the $m \times m$ zero matrix, and the 0 entry in the lower right of B denotes the $n \times n$ zero matrix. Likewise, the 0 entry in the upper right of M

denotes the $m \times n$ zero matrix, and the 0 entry in the lower left of M denotes the $n \times m$ zero matrix. A little thought about the rules of matrix multiplication shows that

$$MB = \begin{bmatrix} 0 & -A \\ A^T & 0 \end{bmatrix} \quad \text{and} \quad BM = \begin{bmatrix} 0 & A \\ -A^T & 0 \end{bmatrix}.$$

That is,

$$BM = -MB.$$

This has the following consequence: If $\mathbf{w} \in \mathbb{R}^{m+n}$ is an eigenvector of B with eigenvalue λ , then

$$B(M\mathbf{w}) = -M(B\mathbf{w}) = -M(\lambda\mathbf{w}) = -\lambda M\mathbf{w}.$$

Since M is orthogonal, $M\mathbf{w} \neq \mathbf{0}$, and so $M\mathbf{w}$ is an eigenvector of B with eigenvalue $-\lambda$.

Thus, again using the fact that $M \in O(m+n)$, if $\{\mathbf{w}_1, \dots, \mathbf{w}_{m+n}\}$ is any orthonormal basis of \mathbb{R}^{m+n} consisting of eigenvectors of B , so is $\{M\mathbf{w}_1, \dots, M\mathbf{w}_{m+n}\}$.

Moreover, the eigenvalue associated to the j th vector in the first basis is minus the eigenvalue associated to the j th vector in the second basis. Since the set of eigenvectors of B is the set of roots of the characteristic polynomial $pB(t) := \det(B - tI)$, it must be that the non-zero eigenvalues come in pairs, one positive and one negative, so that for some $r \leq \min\{m, n\}$, there are exactly r strictly positive eigenvalues, and r strictly negative eigenvalues in the spectrum of B , together with $n + m - 2r$ zero eigenvalues.

Let

$$\sigma_1 \geq \dots \geq \sigma_r > 0$$

be the r strictly positive eigenvalues arranged in non-decreasing order. Let $\{\mathbf{w}_1, \dots, \mathbf{w}_r\}$ be an orthonormal set of eigenvectors of B with $B\mathbf{w}_j = \sigma_j \mathbf{w}_j$ for each $j = 1, \dots, r$. We could obtain such a set by selecting the appropriate vectors from the first orthonormal basis introduced above, and then adjusting the indexing as needed.

Next define the vector $\mathbf{w}_{r+1}, \dots, \mathbf{w}_{2r}$ by

$$\{\mathbf{w}_{r+j}, \dots, \mathbf{w}_{2r}\} = \{M\mathbf{w}_1, \dots, M\mathbf{w}_r\}.$$

Since M is orthogonal, this set is orthonormal. Since B is symmetric, eigenvectors with distinct eigenvalues are orthogonal. Since every vector in the first set has a positive eigenvalue, and every vector in the second set has a negative eigenvalue, the combined set

$$\{\mathbf{w}_1, \dots, \mathbf{w}_r, \mathbf{w}_{r+1}, \dots, \mathbf{w}_{2r}\}$$

is orthonormal.

Next, for each $j = 1, \dots, r$, define the vectors $\mathbf{u}_j \in \mathbb{R}^m$ and $\mathbf{v}_j \in \mathbb{R}^n$ by

$$\mathbf{w}_j = \frac{1}{\sqrt{2}}(\mathbf{u}_j, \mathbf{v}_j).$$

That is, \mathbf{u}_j consists of the first m entries of $\sqrt{2}\mathbf{w}_j$, while \mathbf{v}_j consists of the last n entries of $\sqrt{2}\mathbf{w}_j$.

We now claim that $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$ is an orthonormal subset of \mathbb{R}^m , and that $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ is an orthonormal subset of \mathbb{R}^n . To see this, we note that, by construction,

$$\mathbf{w}_r + j\mathbf{j} = \frac{1}{\sqrt{2}}(-\mathbf{u}_j, \mathbf{v}_j) .$$

Since \mathbf{w}_j and \mathbf{w}_{r+j} are orthogonal.

$$0 = (\mathbf{u}_j, \mathbf{v}_j) \cdot (-\mathbf{u}_j, \mathbf{v}_j) = \|\mathbf{v}_j\|^2 - \|\mathbf{u}_j\|^2 .$$

Since

$$1 = \|\mathbf{w}_j\|^2 = \frac{1}{2}(\mathbf{u}_j, \mathbf{v}_j) \cdot (\mathbf{u}_j, \mathbf{v}_j) = \frac{1}{2}(\|\mathbf{u}_j\|^2 + \|\mathbf{v}_j\|^2) ,$$

we see that $\|\mathbf{u}_j\| = \|\mathbf{v}_j\| = 1$ for $j = 1, \dots, r$.

Next, for $i \neq j$, \mathbf{w}_i is orthogonal to both \mathbf{w}_j and \mathbf{w}_{r+j} . Then means

$$0 = (\mathbf{u}_i, \mathbf{v}_i) \cdot (\mathbf{u}_j, \mathbf{v}_j) = \mathbf{u}_i \cdot \mathbf{u}_j + \mathbf{v}_i \cdot \mathbf{v}_j \quad \text{and} \quad 0 = (\mathbf{u}_i, \mathbf{v}_i) \cdot (-\mathbf{u}_j, \mathbf{v}_j) = -\mathbf{u}_i \cdot \mathbf{u}_j + \mathbf{v}_i \cdot \mathbf{v}_j .$$

Adding and subtracting, we obtain that $\mathbf{u}_i \cdot \mathbf{u}_j = 0$ and $\mathbf{v}_i \cdot \mathbf{v}_j = 0$. This proves that these sets are orthonormal.

Next, in case $r < m$, we take any extension of $\{\mathbf{u}_1, \dots, \mathbf{u}_r\}$ to an orthonormal basis $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ of \mathbb{R}^m , and we define

$$\{\mathbf{w}_{2r+1}, \dots, \mathbf{w}_{r+m}\} = \{(\mathbf{u}_{r+1}, \mathbf{0}), \dots, (\mathbf{u}_m, \mathbf{0})\} ,$$

where $\mathbf{0}$ denotes the zero vector in \mathbb{R}^n so that each of the vectors in $\{\mathbf{w}_{2r+1}, \dots, \mathbf{w}_{r+m}\}$ is in \mathbb{R}^{m+n} . Each of these vectors is orthogonal to every eigenvector of B with a non-zero eigenvalue, and hence it must be in the zero eigenspace of B ; i.e., the null space of B .

Likewise, in $r < n$, we take any extension of $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$ to an orthonormal basis $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ of \mathbb{R}^n , and we define

$$\{\mathbf{w}_{m+r+1}, \dots, \mathbf{w}_{m+n}\} = \{(\mathbf{0}, \mathbf{v}_{m+r+1}), \dots, (\mathbf{0}, \mathbf{v}_{m+n})\} ,$$

where $\mathbf{0}$ denotes the zero vector in \mathbb{R}^m so that each of the vectors in $\{\mathbf{w}_{m+r+1}, \dots, \mathbf{w}_{m+n}\}$ is in \mathbb{R}^{m+n} . Each of these vectors is orthogonal to every eigenvector of B with a non-zero eigenvalue, and hence it must be in the zero eigenspace of B ; i.e., the null space of B .

We now claim that the orthonormal bases $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ that we have just constructed, together with the numbers $\{\sigma_1, \dots, \sigma_r\}$, satisfy (8.20). As explained at the beginning of the proof, this completes the proof of the existence of the singular value decomposition.

To do this, we compute

$$B\mathbf{w}_j = \frac{1}{\sqrt{2}}B(\mathbf{u}_j, \mathbf{v}_j) = \frac{1}{\sqrt{2}}(A\mathbf{v}_j, A^t\mathbf{u}_j) \quad \text{and} \quad B\mathbf{w}_{r+j} = \frac{1}{\sqrt{2}}B(-\mathbf{u}_j, \mathbf{v}_j) = \frac{1}{\sqrt{2}}(A\mathbf{v}_j, -A^t\mathbf{u}_j) .$$

We also have that for $j \leq r$.

$$B\mathbf{w}_j = \sigma_j \mathbf{w}_j = \sigma_j \frac{1}{\sqrt{2}}(\mathbf{u}_j, \mathbf{v}_j) \quad \text{and} \quad B\mathbf{w}_{r+j} = -\sigma_j \mathbf{w}_{r+j} = \sigma_j \frac{1}{\sqrt{2}}(\mathbf{u}_j, -\mathbf{v}_j) .$$

Therefore,

$$\sigma_j = \mathbf{w}_j \cdot B\mathbf{w}_j = \frac{1}{2}(\mathbf{u}_j, \mathbf{v}_j) \cdot (A\mathbf{v}_j, A^t\mathbf{u}_j) = \frac{1}{2}\mathbf{u}_j \cdot A\mathbf{v}_j + \frac{1}{2}\mathbf{v}_j \cdot A^t\mathbf{u}_j = \mathbf{u}_j \cdot A\mathbf{v}_j ,$$

where we have used the fundamental property of the transpose. An even easier computation shows $\mathbf{u}_i \cdot \mathbf{v}_j = 0$ if either $i > r$ or $j > r$.

Next, by similar computations

$$0 = \mathbf{w}_i \cdot B\mathbf{w}_j = \frac{1}{2}\mathbf{u}_i \cdot A\mathbf{v}_j + \frac{1}{2}\mathbf{u}_j \cdot A\mathbf{v}_i ,$$

and

$$0 = \mathbf{w}_i \cdot B\mathbf{w}_{r+j} = \frac{1}{2}\mathbf{u}_i \cdot A\mathbf{v}_j - \frac{1}{2}\mathbf{u}_j \cdot A\mathbf{v}_i ,$$

Combining these equations, we conclude $\mathbf{u}_i \cdot A\mathbf{v}_j = 0$ whenever $i \neq j$. This proves (8.20). \square

8.4 Exercises

1 Consider the following permutations

$$\sigma_1 = \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ 3 & 1 & 4 & 5 & 6 & 2 \end{array} \quad \sigma_2 = \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ 4 & 3 & 6 & 5 & 2 & 1 \end{array} \quad \sigma_3 = \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ 4 & 5 & 6 & 1 & 2 & 3 \end{array}$$

(a) Compute $D(\sigma_j)$ and $\chi(\sigma_j)$ for $j = 1, 2, 3$.

(b) For each $j = 1, 2, 3$, find a way to write σ_j as a product of pair permutations.

(c) Compute the value of $\chi(\sigma_1 \circ (\sigma_2 \circ \sigma_3)^{-1})$.

2 Consider the following permutations

$$\sigma_1 = \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ 2 & 4 & 6 & 1 & 3 & 5 \end{array} \quad \sigma_2 = \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ 5 & 1 & 6 & 4 & 2 & 3 \end{array} \quad \sigma_3 = \begin{array}{cccccc} 1 & 2 & 3 & 4 & 5 & 6 \\ 4 & 1 & 5 & 2 & 6 & 3 \end{array}$$

(a) Compute $D(\sigma_j)$ and $\chi(\sigma_j)$ for $j = 1, 2, 3$.

(b) For each $j = 1, 2, 3$, find a way to write σ_j as a product of pair permutations.

(c) Compute the value of $\chi(\sigma_1 \circ (\sigma_2 \circ \sigma_3)^{-1})$.

3 Let σ_1 , σ_2 and σ_3 be the permutations defined in Exercise 1. Compute the distances $\varrho(\sigma_1, \sigma_2)$, $\varrho(\sigma_2, \sigma_3)$ and $\varrho(\sigma_3, \sigma_1)$. Also, find geodesics from σ_1 to σ_2 , from σ_2 to σ_3 , and from σ_3 to σ_1 .

4 Let σ_1 , σ_2 and σ_3 be the permutations defined in Exercise 2. Compute the distances $\varrho(\sigma_1, \sigma_2)$, $\varrho(\sigma_2, \sigma_3)$ and $\varrho(\sigma_3, \sigma_1)$. Also, find geodesics from σ_1 to σ_2 , from σ_2 to σ_3 , and from σ_3 to σ_1 .

5 The *order reversing permutation* σ_* in \mathcal{S}_n is the permutation defined by

$$\sigma_* := \begin{array}{cccccc} 1 & 2 & \dots & n-1 & n \\ n & n-1 & \dots & 2 & 1 \end{array} .$$

In other words,

$$\sigma_*(k) = n - k + 1 \quad \text{for all} \quad k = 1, \dots, n .$$

(a) Show that $D(\sigma_*) = n(n-1)/2$, and that for all $\sigma \in \mathcal{S}_n$, $D(\sigma) < n(n-1)/2$ unless $\sigma = \sigma_*$.

(b) Prove that

$$\max\{ \varrho(\sigma_1, \sigma_2) : \sigma_1, \sigma_2 \in \mathcal{S}_n \} = n(n-1)/2 .$$

In other words, any two permutations in \mathcal{S}_n are connected by a path of at most $n(n-1)/2$ steps, and there exist pairs of permutations such that the shortest path connecting them has this many steps. This is often expressed by saying that *the diameter of \mathcal{S}_n is $n(n-1)/2$* .

6 Show that the set \mathcal{A}_n consisting of all even permutations in \mathcal{S}_n is a transformation group on $\{1, \dots, n\}$. \mathcal{A}_n is called the *alternating group of order n* . Show that there are exactly $n!/2$ permutations in \mathcal{A}_n , and show that the set of all odd permutations is not a transformation group.

7 For each $\sigma \in \mathcal{S}_n$, define the $n \times n$ matrix P_σ by

$$P_\sigma := [\mathbf{e}_{\sigma(1)}, \mathbf{e}_{\sigma(2)}, \dots, \mathbf{e}_{\sigma(n)}] ; \quad (8.22)$$

that is, the j th column of P_σ is $\mathbf{e}_{\sigma(j)}$. The $n!$ matrices P_σ with $\sigma \in \mathcal{S}_n$ are called the *permutation matrices*. Prove that for all $\sigma \in \mathcal{S}_n$, $\det(P_\sigma) = \chi(\sigma)$.

Chapter 9

FLUX AND CIRCULATION, DIVERGENCE AND CURL

9.1 Flows and flux

9.1.1 Vector fields and motion

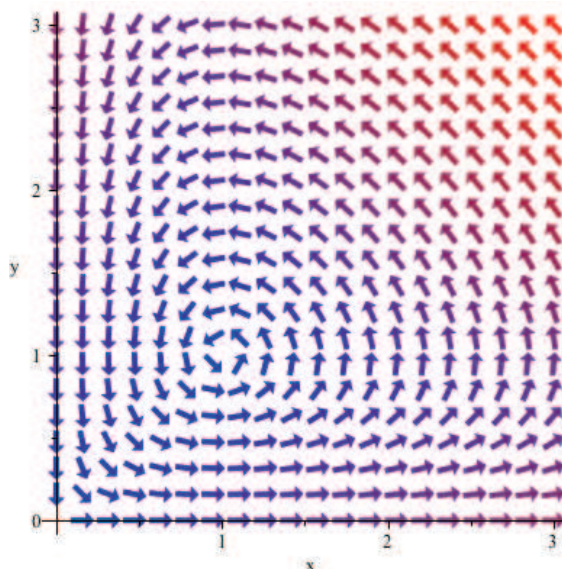
Let $\mathbf{F}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$ be a function from \mathbb{R}^n to \mathbb{R}^n . We have studied such functions already, but now our point of view will be slightly different. To go along with this new point of view, we introduce some new terminology:

Definition 91 (Vector field). *A vector field on an open set $U \subset \mathbb{R}^n$, possibly \mathbb{R}^n itself, is a function \mathbf{F} defined on U with values in \mathbb{R}^n . The vector field is said to be continuous, differentiable or continuously differentiable if the function \mathbf{F} is continuous, differentiable or continuously differentiable.*

We can represent vector fields on \mathbb{R}^2 in an informative graphical manner, and so our first examples concern the case $n = 2$. Let us look at the specific vector field

$$\mathbf{F}(x, y) = (x(1 - y), y(x - 1)) .$$

Here is a plot of this vector field for $0 \leq x, y \leq 3$:



The arrows at each point show the direction of the vector field at each point on a grid, and color is used to indicate the length of these vectors: As the colors range from blue to red, the length of the vectors ranges from shortest to longest.

Think of the arrows as describing motion. A point particle at \mathbf{x} will move to $\mathbf{x} + \mathbf{F}(\mathbf{x})dt$ in an infinitesimal time step dt , and then move on from there following the arrow at that point, and so on.

To make this more precise, take a starting point $\mathbf{x}_0 = (x_0, y_0)$ and fix a small time step h . Define a sequence of points $\{\mathbf{x}_n^{(h)}\}$ by $\mathbf{x}_n^{(h)} = \mathbf{x}_0$ and for $n \geq 1$,

$$\mathbf{x}_n^{(h)} = \mathbf{x}_{n-1}^{(h)} + h\mathbf{F}(\mathbf{x}_{n-1}^{(h)}) . \quad (9.1)$$

Run this sequence until the point leaves the region where the vector field \mathbf{F} is defined, or forever if \mathbf{F} is defined everywhere, or the sequence never leaves the the region where \mathbf{F} is defined.

Given the sequence $\mathbf{x}_n^{(h)}$, define the continuous function $\mathbf{x}^{(h)}(t)$ by “connecting the dots”; that is,

$$\mathbf{x}^{(h)}(t) = \mathbf{x}_{n-1}^{(h)} + (t/h - (n-1))(\mathbf{x}_n^{(h)} - \mathbf{x}_{n-1}^{(h)}) \quad \text{for} \quad (n-1)h \leq t \leq nh .$$

The resulting curve $\mathbf{x}^{(h)}(t)$ moves along by following the arrows provided by the vector field, updating the information from the vector field every time-step h . As h tends to zero, it follows the arrows provided by the vector field more and more accurately. Indeed, since (9.1) can be written

$$\frac{\mathbf{x}_n^{(h)} - \mathbf{x}_{n-1}^{(h)}}{h} = \mathbf{F}(\mathbf{x}_{n-1}^{(h)}) ,$$

one would expect a limiting curve $\mathbf{x}(t) := \lim_{h \rightarrow 0} \mathbf{x}^{(h)}(t)$ to exist and satisfy

$$\mathbf{x}'(t) = \mathbf{F}(\mathbf{x}(t)) . \quad (9.2)$$

When \mathbf{F} is continuously differentiable, this is indeed the case. Let us take this for granted for the moment – after all, it is quite intuitive – and focus instead on what we can learn from the curves satisfying (9.2).

Definition 92 (Flow lines). *Curves $\mathbf{x}(t)$ satisfying (9.2) on any open interval $t_0 < t < t_1$ are called flow curves of the vector field \mathbf{F} ; they describe the motion of a point particle that is “going with the flow” described by the vector field \mathbf{F} .*

In some cases, it is possible to explicitly solve (9.2) to find the flow lines. We begin with the simplest (but important) example:

Example 129 (Solving for flow lines). *Consider the constant vector field defined on all of \mathbb{R}^2 by*

$$\mathbf{F}(x, y) = (a, b)$$

where $a, b \in \mathbb{R}$ are constant. Then $\mathbf{x}(t)$ is a flow curve for \mathbf{F} if and only if $\mathbf{x}'(t) = (a, b)$. By the Fundamental Theorem of Calculus, this means

$$\mathbf{x}(t) = \mathbf{x}_0 + t(a, b) .$$

The flow curves are simply straight lines with direction vector given by (a, b) .

The next example is less simple, but also important:

Example 130 (Solving for flow curves). *Consider the vector field \mathbf{F} on \mathbb{R}^2 given by $\mathbf{F}(x, y) = (-y, x)$.*

Note that $\mathbf{F}(\mathbf{x})$ is a linear function of \mathbf{x} ; if we define the matrix $A := \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, we have

$$\mathbf{F}(x) = A\mathbf{x} .$$

If $\mathbf{x}(t)$ is any flow curve for \mathbf{F} , then $\mathbf{x}'(t) = A\mathbf{x}(t)$. Differentiating again, we get

$$\mathbf{x}''(t) = A^2\mathbf{x}(t) = -\mathbf{x}(t)$$

since $A^2 = -I$. Thus, whenever $\mathbf{x}(t) = (x(t), y(t))$ is a flow curve of \mathbf{F} , we have

$$x''(t) = -x(t) \quad \text{and} \quad y''(t) = -y(t) .$$

One way to solve these equations is to take $x(t) = \alpha \cos(t) + \beta \sin(t)$. Then $x(0) = \alpha$ and $x'(0) = \beta$. But $\mathbf{x}'(t) = A\mathbf{x}(t)$ says that $x'(t) = -y(t)$. Hence $\beta = x'(0) = -y(0)$. That is, $x(t) = x_0 \cos(t) - y_0 \sin(t)$. Then since $y(t) = -x'(t)$, $y(t) = x_0 \sin(t) + y_0 \cos(t)$. Thus,

$$\mathbf{x}(t) = (x_0 \cos(t) - y_0 \sin(t), x_0 \sin(t) + y_0 \cos(t)) \tag{9.3}$$

is a flow curve for \mathbf{F} through \mathbf{x}_0 , and it is defined for all t .

In fact, it is the unique flow curve for \mathbf{F} through \mathbf{x}_0 . To see this, suppose that $\mathbf{y}(t)$ is any other, and define $\mathbf{z}(t) = \mathbf{x}(t) - \mathbf{y}(t)$. Note that $\mathbf{z}(0) = \mathbf{x}(0) - \mathbf{y}(0) = \mathbf{x}_0 - \mathbf{x}_0 = \mathbf{0}$, since by hypothesis both curves are initially at \mathbf{x}_0 . Also,

$$\mathbf{z}'(t) = \mathbf{x}'(t) - \mathbf{y}'(t) = A(\mathbf{x}(t) - \mathbf{y}(t)) = A\mathbf{z}(t) .$$

Next, a simple calculation shows that $(a, b) \cdot A(a, b) = 0$ for all (a, b) . Hence

$$\frac{d}{dt} \|\mathbf{z}(t)\|^2 = 2\mathbf{z}(t) \cdot A\mathbf{z}(t) = 0 .$$

Thus, since $\|\mathbf{z}(0)\| = 0$, $\|\mathbf{z}(t)\| = 0$ for all t , and this means that $\mathbf{y}(t) = \mathbf{x}(t)$ for all t . The two flow curves are in fact the same, which proves the uniqueness of the flow curves for \mathbf{F} .

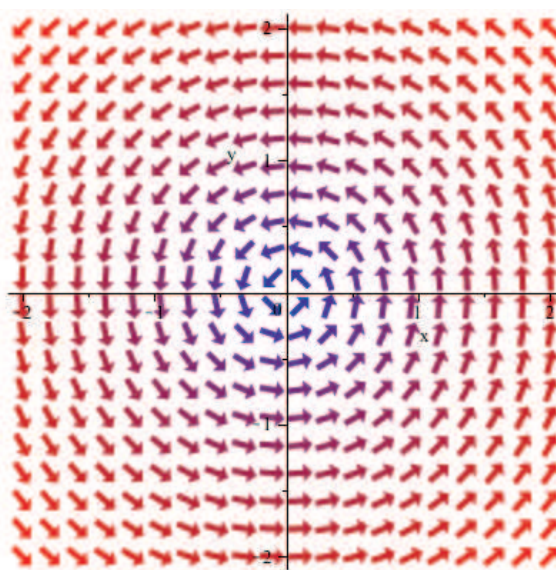
There is still more to be learned from this example if we write (9.3) another way. Define the t -dependent matrix

$$[R(t)] := \begin{bmatrix} \cos t & -\sin t \\ \sin t & \cos t \end{bmatrix} . \quad (9.4)$$

Then (9.3) is equivalent to

$$\mathbf{x}(t) = [R(t)]\mathbf{x}_0 . \quad (9.5)$$

You recognize $[R(t)]$ as a rotation matrix. The vector field describes circular motion about the origin at a constant speed that is equal to the distance from the origin. The flow exists for all times t . Here is a plot of the vector field:



Example 131 (Solving for flow lines). Consider the vector field $\mathbf{F}(x, y) = (x^2, 1)$. Then (9.2) gives us the system

$$\begin{aligned} x'(t) &= x^2(t) \\ y'(t) &= 1 . \end{aligned}$$

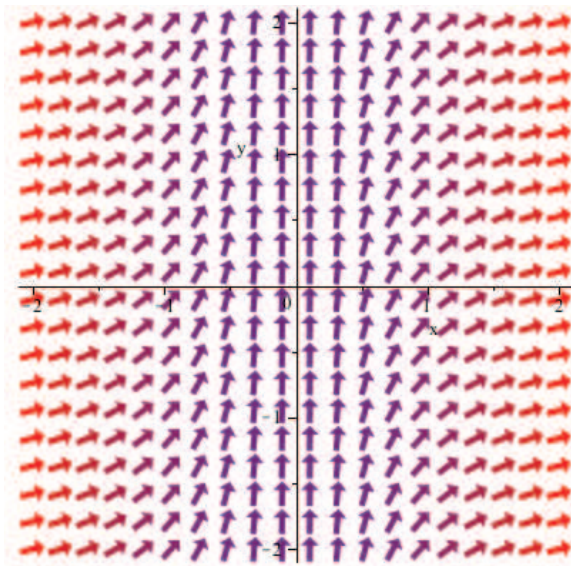
This system may also be solved by single variable calculus methods since the rate of change of $x(t)$ depends only on $x(t)$ and the rate of change of $y(t)$ is even constant. Indeed, by the Fundamental Theorem of Calculus, $y(t) = y(0) + \int_0^t y'(s)ds = y_0 + t$. Likewise, we have $1 = \frac{x'(t)}{x^2(t)} = -\frac{d}{dt} \frac{1}{x(t)}$. Therefore, by the Fundamental Theorem of Calculus,

$$\frac{1}{x(0)} - \frac{1}{x(t)} = \int_0^t 1ds = t ,$$

and so

$$\mathbf{x}(t) = \frac{x_0}{1 - tx_0} .$$

For $t = 1/x_0$, we would be dividing by zero. Indeed, notice $x(t)$ is well defined for $|t| < 1/|x_0|$, but that $\lim_{t \rightarrow 1/x_0} x(t) = \infty$. This is an example where, although the vector field is very simple, the flow only exists a finite time before it “explodes” to infinity.



For positive x , the vector field pushes point particles further off to the right, faster and faster, so that although the vector field itself is finite everywhere, it accelerates a point particle to infinite speed in a finite time.

Example 132 (Solving for flow lines). Consider the vector field $\mathbf{F}(x, y) = (x^{1/3}, 1)$. Then (9.2) gives us the system

$$\begin{aligned} x'(t) &= x^{1/3}(t) \\ y'(t) &= 1. \end{aligned}$$

Let $\mathbf{x}_0 = (0, 0)$. Then since $\mathbf{F}(0, t) = (0, 1)$ for all t , the curve $\mathbf{x}(t) = (0, t)$ is a flow curve of \mathbf{F} with $\mathbf{x}(0) = \mathbf{x}_0$. But there is another flow curve of \mathbf{F} through \mathbf{x}_0 . For $x(t) > 0$, the equation $x'(t) = x^{1/3}(t)$ is equivalent to

$$1 = \frac{x'(t)}{x^{1/3}(t)} = \frac{3}{2} \frac{d}{dt} x^{2/3}(t).$$

Hence, by the Fundamental Theorem of Calculus, $t = \frac{3}{2} x^{2/3}(t)$, so that

$$x(t) = \left(\frac{2}{3} t \right)^{3/2}.$$

Hence the curve $\tilde{\mathbf{x}}(t)$ given by $\tilde{\mathbf{x}}(t) = (0, t)$ for $t \leq 0$, and by

$$\tilde{\mathbf{x}}(t) = \left(\left(\frac{2}{3} t \right)^{3/2}, t \right)$$

for $t > 0$ is another flow curve of \mathbf{F} passing through \mathbf{x}_0 .

9.1.2 Flow transformations

To understand the flow described by a vector field, it is best to look at the set of flow curves as a whole instead of individually. The basic idea is that a “nice” vector field \mathbf{F} describes a one parameter set $\{\Phi_t : t \in \mathbb{R}\}$ of transformations of \mathbb{R}^n through the following simple rule: *for each t , and each \mathbf{x}_0 , $\Phi_t(\mathbf{x}_0)$ is the point one gets to by following the flow lines described by \mathbf{F} for a time t .* It is useful to picture this for $n = 2$ in terms an actual flow of liquid across the plane in which at each time t , a particle being carried along with the flowing liquid has the velocity $\mathbf{F}(\mathbf{x})$ as it passes through the point \mathbf{x} . As the particle moves, it traces out a flow curve. Other particles, placed elsewhere, would trace out other flow curves. The *flow transformation* tells where the flowing liquid takes any particle, started anywhere, in time t . It is a fundamental object of study in fluid mechanics, aerodynamics, and many other fields of science, including electrodynamics, for less obvious reasons.

In our applications of the flow curve concept, we need to know something not only about the existence, but also the *uniqueness* of flow curves. We have given a “recipe” for constructing flow curves through a limiting process. But maybe a different “recipe” would yield different curves. Indeed, we have seen in Example 132 that this can happen: Our recipe would yield the flow curve $\mathbf{x}(t) = (0, t)$, but as we have seen, there is another flow curve for this vector field \mathbf{F} that passes through $(0, 0)$ at time $t = 0$. If flow curves are not unique, our simple recipe for the flow transformation Φ_t does not define a function: Which flow curve do you follow?

Also, we have seen in Example 131 that flow curves might “blow up” in a finite time. So for some vector fields \mathbf{F} , it may not be possible to follow a flow curve for time t if t is too large.

On the other hand, in our first two example, the vector field was “nice” and had unique flow curves through each point that existed for all times t . How do we recognize “nice” vector fields?

As we shall see, a vector field \mathbf{F} behaves nicely whenever it is continuously differentiable and has a Jacobian matrix $[D_{\mathbf{F}}(\mathbf{x})]$ that has a bounded Frobenius norm, meaning that there is some finite number L such that

$$\|[D_{\mathbf{F}}(\mathbf{x})]\|_{\text{F}} \leq L \quad \text{for all} \quad \mathbf{x} \in \mathbb{R}^2. \quad (9.6)$$

The condition (9.6) has the following consequence:

Lemma 21. *Let \mathbf{F} be a vector field that satisfies (9.6). Then for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$,*

$$\|\mathbf{F}(\mathbf{x}) - \mathbf{F}(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Proof: Consider the parameterized line segment $\mathbf{z}(t)$ defined by $\mathbf{z}(t) = (1 - t)\mathbf{x} + t\mathbf{y}$ for which $\mathbf{z}(0) = \mathbf{x}$, $\mathbf{z}(1) = \mathbf{y}$ and $\mathbf{z}'(t) = \mathbf{y} - \mathbf{x}$ for all t . Then by the Fundamental Theorem of Calculus and the Chain Rule,

$$\mathbf{F}(\mathbf{y}) - \mathbf{F}(\mathbf{x}) = \int_0^1 \frac{d}{dt} \mathbf{F}(\mathbf{z}(t)) = \int_0^1 [D_{\mathbf{F}}(\mathbf{z}(t))](\mathbf{y} - \mathbf{x}) dt.$$

Then by the triangle inequality for integrals,

$$\|\mathbf{F}(\mathbf{y}) - \mathbf{F}(\mathbf{x})\| \leq \int_0^1 \|[D_{\mathbf{F}}(\mathbf{z}(t))](\mathbf{y} - \mathbf{x})\| dt \leq \int_0^1 L\|\mathbf{y} - \mathbf{x}\| dt = L\|\mathbf{y} - \mathbf{x}\|.$$

□

Theorem 85 (Uniqueness of flow curves). *Let \mathbf{F} be a vector field that satisfies (9.6). Let $\mathbf{x}(t)$ and $\mathbf{y}(t)$ be flow curves for \mathbf{F} defined for $-a < t < a$ for some $a > 0$. Let \mathbf{x}_0 denote $\mathbf{x}(0)$, and let \mathbf{y}_0 denote $\mathbf{y}(0)$. Then for all $-a < t < a$,*

$$e^{-|t|L}\|\mathbf{x}_0 - \mathbf{y}_0\| \leq \|\mathbf{x}(t) - \mathbf{y}(t)\| \leq e^{|t|L}\|\mathbf{x}_0 - \mathbf{y}_0\| . \quad (9.7)$$

Both inequalities in (9.7) tells us something important about flow lines: From the inequality on the right, we see that if $\mathbf{x}_0 = \mathbf{y}_0$, then $\mathbf{x}(t) = \mathbf{y}(t)$ for all $-a < t < a$, so that the two flow curves are the same. In other words, when \mathbf{F} satisfies (9.6), there is at most one flow curve through each point $\mathbf{x}_0 \in \mathbb{R}^n$.

But that is not all: fix any $\epsilon > 0$, and define $\delta(\epsilon) = e^{|t|L}\epsilon$. Then

$$\|\mathbf{y}_0 - \mathbf{x}_0\| < \delta(\epsilon) := e^{|t|L}\epsilon \quad \Rightarrow \quad \|\mathbf{x}(t) - \mathbf{y}(t)\| < \epsilon . \quad (9.8)$$

In other words, if the initial points \mathbf{x}_0 and \mathbf{y}_0 are sufficiently close, then the flow curves through them will be close at time t .

The inequality on the left in (9.7) tells the *flow curves never cross*: If $\mathbf{x}(0) \neq \mathbf{y}(0)$, then it is impossible to have $\mathbf{x}(t) = \mathbf{y}(t)$ for any t .

Proof of Theorem 85: We compute:

$$\frac{d}{dt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 = 2(\mathbf{x}(t) - \mathbf{y}(t)) \cdot (\mathbf{x}'(t) - \mathbf{y}'(t)) = 2(\mathbf{x}(t) - \mathbf{y}(t)) \cdot (\mathbf{F}(\mathbf{x}(t)) - \mathbf{F}(\mathbf{y}(t))) .$$

Then by the Cauchy-Schwarz inequality and Lemma 21,

$$\begin{aligned} \pm \frac{d}{dt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 &\leq 2\|\mathbf{x}(t) - \mathbf{y}(t)\|\|\mathbf{F}(\mathbf{x}(t)) - \mathbf{F}(\mathbf{y}(t))\| \\ &\leq 2L\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 , \end{aligned} \quad (9.9)$$

where we have taken advantage of the fact that we can multiply through by -1 before applying the Cauchy-Schwarz inequality since, after all, we are going to take absolute values. In other words,

$$-2L\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 \leq \frac{d}{dt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 \leq 2L\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 . \quad (9.10)$$

The inequality on the right in (9.10) says that

$$\frac{d}{dt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 - 2L\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 \leq 0 .$$

Multiplying through by e^{-2Lt} , we obtain $\frac{d}{dt}(e^{-2Lt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2) \leq 0$. It follows that

$$\varphi(t) := (e^{-2Lt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2)$$

is a non-increasing function of t .

Likewise, the inequality on the left in (9.10) says that

$$\frac{d}{dt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 + 2L\|\mathbf{x}(t) - \mathbf{y}(t)\|^2 \geq 0 .$$

Multiplying through by e^{2Lt} , we obtain $\frac{d}{dt}(e^{2Lt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2) \geq 0$. It follows that

$$\psi(t) := (e^{2Lt}\|\mathbf{x}(t) - \mathbf{y}(t)\|^2)$$

is a non-decreasing function of t . Since $\varphi(0) = \psi(0) = \|\mathbf{x}(0) - \mathbf{y}(0)\|^2$, we have that for $t > 0$ $\varphi(t) \leq \|\mathbf{x}(0) - \mathbf{y}(0)\|^2 \leq \psi(t)$, while for $t < 0$, $\psi(t) \leq \|\mathbf{x}(0) - \mathbf{y}(0)\|^2 \leq \varphi(t)$. From the definitions of φ and ψ , either way we have

$$e^{-2L|t|} \|\mathbf{x}(t) - \mathbf{y}(t)\|^2 \leq \|\mathbf{x}(0) - \mathbf{y}(0)\|^2 \leq e^{2L|t|} \|\mathbf{x}(t) - \mathbf{y}(t)\|^2 .$$

Taking squares roots and rearranging terms gives us (9.7). \square

We now state a theorem whose proof we postpone until later. This theorem guarantees the existence of flow curves under the conditions we have showed guarantee uniqueness.

Theorem 86 (Existence of flow curves). *Let \mathbf{F} be a vector field on \mathbb{R}^n that satisfies (9.6). Then for all $\mathbf{x}_0 \in \mathbb{R}^n$, there exists a flow curve $\mathbf{x}(t)$ for \mathbf{F} that is defined for all $-\infty < t < \infty$ and such that $\mathbf{x}(0) = \mathbf{x}_0$.*

We may use the flow curves of a vector field \mathbf{F} that satisfies (9.6) to define a one-parameter family Φ_t of continuous one-to-one transformations of \mathbb{R}^2 onto \mathbb{R}^2 :

Definition 93 (Flow transformations). *Let \mathbf{F} be a vector field on \mathbb{R}^n that satisfies (9.6). Then for each $t \in \mathbb{R}$, the functions Φ_t from \mathbb{R}^n to \mathbb{R}^n is defined as follows: For any $\mathbf{x}_0 \in \mathbb{R}^n$,*

$$\Phi_t(\mathbf{x}_0) = \mathbf{x}(t)$$

where $\mathbf{x}(t)$ is the point at time t along the unique flow curve through \mathbf{x}_0 . The transformations $\{ \Phi_t : t \in \mathbb{R} \}$ are called the flow transformations generated by \mathbf{F} . Note that Φ_0 is the identity transformation.

Example 133 (Rotational flow). *Let us once more consider our rotational vector field*

$$\mathbf{F}(x, y) = (-y, x) .$$

Then as we saw in Example 130, the solution of $\mathbf{x}'(t) = \mathbf{F}(\mathbf{x}(t))$ is

$$\mathbf{x}(t) = [R(t)]\mathbf{x}_0 \quad \text{where} \quad [R(t)] := \begin{bmatrix} \cos t & -\sin t \\ \sin t & \cos t \end{bmatrix} .$$

Hence, in this case, $\Phi_t(\mathbf{x})$ is the linear transformation

$$\Phi_t(\mathbf{x}) = [R(t)]\mathbf{x} ,$$

which is rotation of the plane \mathbb{R}^2 counterclockwise through the angle t .

Now notice that for any t_1 and t_2 , a simple calculation using the angle addition formulas shows that

$$[R(t_1)][R(t_2)] = [R(t_1 + t_2)] .$$

Therefore,

$$\Phi_{t_1} \circ \Phi_{t_2} = \Phi_{t_1 + t_2} .$$

Furthermore, since rotations are invertible, and in fact $[R(t)]^{-1} = [R(-t)]$, each Φ_t is invertible, and $\Phi_t^{-1} = \Phi_{-t}$. That is, the inverse transformation is obtained by running the flow backwards in time.

Here is one more simple, but important, example:

Example 134 (Constant flow). *Consider the constant vector field $\mathbf{F}(x, y) = (a, b)$ from Example 129. As we have seen there, The the flow curves are given by $\mathbf{x}(t) = \mathbf{x}_0 + t(a, b)$. Thus,*

$$\Phi_t(\mathbf{x}_0) = \mathbf{x}_0 + t(a, b) = (x_0 + ta, y_0 + tb) ,$$

Since this is true for any \mathbf{x}_0 , we may now drop the subscript 0, and simply write

$$\Phi_t(\mathbf{x}) = (x + ta, y + tb) .$$

It is easy to check from this formula that for any t_1 and t_2 , $\Phi_{t_1} \circ \Phi_{t_2} = \Phi_{t_1+t_2}$ and that each Φ_t is invertible with $\Phi_t^{-1} = \Phi_{-t}$. That is, as in the previous example, the inverse transformation is obtained by running the flow backwards in time.

The reason the last example is important is that for *any* differentiable, and hence continuous, vector field \mathbf{F} defined on a neighborhood of some point \mathbf{x}_0 , as long as $\mathbf{F}(\mathbf{x}_0) \neq \mathbf{0}$, if you “zoom in” on the vector field in a small neighborhood of \mathbf{x}_0 , you will have

$$\mathbf{F}(\mathbf{x}) \approx \mathbf{F}(\mathbf{x}_0) ,$$

and so locally, the vector field will look like a constant vector field, and thus, locally, the flow lines will be nearly straight parallel lines.

This is not necessarily true as a point \mathbf{x}_0 where $\mathbf{F}(\mathbf{x}_0) = \mathbf{0}$. Consider the rotational vector field $\mathbf{F}(x, y) = (-y, x)$. Then $\mathbf{F}(0, 0) = (0, 0)$, but the vector field takes on *every* direction in any small neighborhood of $(0, 0)$, and the flow lines are concentric circles, not straight lines.

The next theorem tells us that some of what we saw in the last two examples always happens for flows generated by nice vector fields \mathbf{F} . We shall prove all of it here except the last part, which we postpone until later.

Theorem 87 (Fundamental properties of the flow transformations). *Let \mathbf{F} be a continuously differentiable vector field on \mathbb{R}^n that satisfies (9.6), and let $\{ \Phi_t : t \in \mathbb{R} \}$ be the set of flow transformations generated by \mathbf{f} . Then:*

- (1) *For each t , Φ_t is a continuous, one-to-one transformation of \mathbb{R}^n onto \mathbb{R}^n , and therefore invertible.*
- (2) *For each t_1, t_2 ,*

$$\Phi_{t_1} \circ \Phi_{t_2} = \Phi_{t_1+t_2} .$$

- (3) *For each t , the inverse of Φ_t is Φ_{-t} .*
- (4) *For each t , Φ_t is is not only continuous; it is continuously differentiable.*

Proof: For (1), fix any t . Since by definition $\|\Phi_t(\mathbf{x}_0) - \Phi_t(\mathbf{y}_0)\| = \|\mathbf{x}(t) - \mathbf{y}(t)\|$, (9.8) says that for any $\epsilon > 0$, with $\delta(\epsilon) := e^{|t|L}\epsilon$,

$$\|\mathbf{y}_0 - \mathbf{x}_0\| < \delta(\epsilon) \quad \Rightarrow \quad \|\Phi_t(\mathbf{x}_0) - \Phi_t(\mathbf{y}_0)\| < \epsilon .$$

This shows that Φ_t is continuous at \mathbf{x}_0 , and since \mathbf{x}_0 was an arbitrary point in \mathbb{R}^n , Φ_t is continuous on \mathbb{R}^n .

To see that Φ_t transforms \mathbb{R}^n onto \mathbb{R}^n , consider any $\mathbf{x} \in \mathbb{R}^n$, and let $\mathbf{x}(s)$ be the flow curve passing through \mathbf{x} at time $s = 0$. (We used s as the parameter now since t is already in use and fixed.) consider the curve $\tilde{\mathbf{x}}(s)$ defined by

$$\tilde{\mathbf{x}}(s) = \tilde{\mathbf{x}}(s - t) .$$

Differentiating, we find $\frac{d}{ds}\tilde{\mathbf{x}}(s) = \mathbf{x}'(s - t) = \mathbf{F}(\mathbf{x}(s - t)) = \mathbf{F}(\tilde{\mathbf{x}}(s))$. Thus, $\tilde{\mathbf{x}}(s)$ is a flow curve. Define

$$\tilde{\mathbf{x}}_0 := \tilde{\mathbf{x}}(0) = \mathbf{x}(-t) .$$

Then, by the definition of Φ_t ,

$$\Phi_t(\tilde{\mathbf{x}}_0) = \tilde{\mathbf{x}}(t) = \mathbf{x}(t - t) = \mathbf{x}(0) = \mathbf{x} .$$

Thus Φ_t transforms $\tilde{\mathbf{x}}_0$ onto \mathbf{x} , and since \mathbf{x} is an arbitrary point in \mathbb{R}^2 , we have shown that Φ_t transforms \mathbb{R}^n onto \mathbb{R}^n . Finally, if Φ_t were not one-to-one, there would be some $\mathbf{x}_0, \mathbf{y}_0 \in \mathbb{R}^n$ with $\mathbf{x}_0 \neq \mathbf{y}_0$ but with $\Phi_t(\mathbf{x}_0) = \Phi_t(\mathbf{y}_0)$. But this is impossible, since flow lines never cross when \mathbf{F} satisfies (9.6), as we have explained after Theorem 85. Altogether, we have proved that Φ_t is a continuous invertible transformation from \mathbb{R}^n onto \mathbb{R}^n .

We next prove (2), Fix any $\mathbf{x}_0 \in \mathbb{R}^n$. Let $\mathbf{x}(t)$ be the flow curve of \mathbf{F} with $\mathbf{x}(0) = \mathbf{x}_0$, Let $\tilde{\mathbf{x}}(t) = \mathbf{x}(t + t_1)$. Differentiating as above, we see that $\tilde{\mathbf{x}}$ is a flow curve of \mathbf{F} , and $\tilde{\mathbf{x}}(0) = \mathbf{x}(t_1)$.

By definition, $\Phi_{t_1}(\mathbf{x}_0) = \mathbf{x}(t_1) = \tilde{\mathbf{x}}(0)$, and so

$$\Phi_{t_2} \circ \Phi_{t_1}(\mathbf{x}_0) = \Phi_{t_2}(\Phi_{t_1}(\mathbf{x}_0)) = \Phi_{t_2}(\tilde{\mathbf{x}}(0)) = \tilde{\mathbf{x}}(t_2) = \mathbf{x}(t_1 + t_2) = \Phi_{t_1+t_2}(\mathbf{x}_0) .$$

Finally since Φ_0 is the identity transformation, (3) follows directly from (2). \square

The theorem we have just proved says that the set $\{\Phi_t : t \in \mathbb{R}\}$ of transformations of \mathbb{R}^n generated by \mathbf{F} is a transformation group. As we work with these transformations, we shall see why the group property is important.

Let us close this subsection by going back and considering our two “badly behaved” vector fields from Examples 131 and 132, starting with the latter.

When $\mathbf{F}(\mathbf{x})$ is given by $\mathbf{F}(x, y) = (x^{1/3}, 1)$, we compute that for $x \neq 0$,

$$[D_{\mathbf{F}}(x, y)] = \begin{bmatrix} |x|^{-2/3}/3 & 0 \\ 0 & 0 \end{bmatrix} .$$

Hence

$$\|[D_{\mathbf{F}}(x, y)]\|_{\mathbf{F}} = \frac{1}{3}|x|^{-2/3}$$

which tends to infinity as x tends to zero. Thus (9.6) cannot be satisfied for any finite L , and even worse, \mathbf{F} is not even differentiable at any point $(0, y)$. Vector fields with such irregular behavior evidently can be badly behaved, and must be treated with care.

Things are somewhat better with the vector field in Example 131. This vector field $\mathbf{F}(x, y) = (x^2, 1)$ is clearly continuously differentiable, and

$$[D_{\mathbf{F}}(x, y)] = \begin{bmatrix} 2x & 0 \\ 0 & 0 \end{bmatrix}.$$

Hence $\|[D_{\mathbf{F}}(x, y)]\|_F = 4x^2$. Note (9.6) cannot be satisfied for any finite L , but only because of what happens for very large values of $\|\mathbf{x}\|$.

Very often in what follows, we will be interested only in what the flow transformation is doing to points in some bounded set for a short interval of time. In this case, the following *localization procedure* is often useful.

We begin with some preliminary calculations. Note that the function $(x-1)(x-2)$ is negative for $1 < x < 2$ and equals zero for $x = 1$ and $x = 2$. Hence the function $\phi(y)$ defined by

$$\phi(y) = \int_1^y (x-1)(x-2)dx = \frac{1}{6}(2y^3 - 9y^2 + 12y - 5)$$

is monotone decreasing on the interval $(1, 2)$ with $\phi(1) = \phi'(1) = \phi'(2) = 0$.

Now define the function $\chi(r)$ on $[0, \infty)$ by

$$\chi(r) = \begin{cases} 1 & r \leq 1 \\ 1 + \phi(r)/\phi(2) & 1 < r < 2 \\ 0 & r > 2 \end{cases}.$$

By what we have proved above about ϕ , this function is continuously differentiable. Hence by the Chain Rule, for any $R > 0$, the function $\chi(\|\mathbf{x}\|/R)$ is continuously differentiable on \mathbb{R}^n . and equals 1 for $\|\mathbf{x}\| \leq R$, and equals 0 for $\|\mathbf{x}\| \geq 2R$, and

$$\nabla\chi(\|\mathbf{x}\|/R) = \chi'(R\|\mathbf{x}\|) \frac{\mathbf{x}}{R\|\mathbf{x}\|}.$$

A simple calculation shows that $|\chi'(r)| \leq |\chi'(3/2)| = 3/2$, and so

$$\|\nabla\chi(\|\mathbf{x}\|/R)\| \leq \frac{3}{2R}.$$

Now, given any continuously differentiable vector field \mathbf{F} on \mathbb{R}^2 , suppose we are only interested in what the flow is doing inside $B_R(\mathbf{0})$, the ball of radius R centered at the origin. Then let us define the *localized vector field* $\tilde{\mathbf{F}}$ by

$$\tilde{\mathbf{F}}(\mathbf{x}) = \chi(\|\mathbf{x}\|/R)\mathbf{F}(\mathbf{x}).$$

Then, by construction, $\tilde{\mathbf{F}}(\mathbf{x}) = \mathbf{F}(\mathbf{x})$ everywhere in $B_R(\mathbf{0})$. Hence the flow curves of $\tilde{\mathbf{F}}$ and \mathbf{F} coincide as long as they remain inside $B_R(\mathbf{0})$. Next, using the product rule, one sees that the Jacobian of $\tilde{\mathbf{F}}(\mathbf{x})$ is continuous and every entry is identically zero outside of $B_{2R}(\mathbf{0})$. Hence $\|[D_{\tilde{\mathbf{F}}}(\mathbf{x})]\|_F$ is a continuous function of \mathbf{x} , and is zero outside of $B_{2R}(\mathbf{0})$. Let L be its maximum value in $B_{2R}(\mathbf{0})$. Then we have

$$\|[D_{\tilde{\mathbf{F}}}(\mathbf{x})]\|_F \leq L$$

for all $\mathbf{x} \in \mathbb{R}^n$. Hence the localized vector field $\tilde{\mathbf{F}}$ is a nice vector field satisfying the hypotheses of all of the theorems of this subsection.

Example 135 (Localizing a flow with explosion). Consider the vector field $\mathbf{F}(x, y) = (x^2, 1)$. from Example 131.

Then as we saw in Example 131, the flow curve of \mathbf{F} through (x_0, y_0) is given by

$$\mathbf{x}(t) = (x_0/(1 - tx_0), y_0 + t) \quad \text{for } tx_0 < 1 .$$

Now fix any $R > 0$. For all $\mathbf{x} \in B_R(\mathbf{0})$, $\|\mathbf{F}(\mathbf{x})\| \leq \sqrt{R^4 + 1}$, and hence if $\mathbf{x}(t)$ is any flow curve of \mathbf{F} or of $\tilde{\mathbf{F}}$, as long as this curve is in $B_R(\mathbf{0})$, its speed is bounded above by $\sqrt{R^4 + 1}$. Evidently any flow curve starting in $B_{R/2}(\mathbf{0})$ at time $t = 0$ must travel a distance of at least $R/2$ to exit $B_R(\mathbf{0})$, and since its speed is limited by $\sqrt{R^4 + 1}$,

$$|t| < \frac{R}{2\sqrt{R^4 + 1}} \quad \Rightarrow \quad \mathbf{x}(t) \in B_R(\mathbf{0}) .$$

Hence, if we only look at times t with $|t| < R/\sqrt{R^4 + 1}$ and points \mathbf{c} with $\|\mathbf{x}\| < R$, we do not see the difference between the “explosive” vector field $\tilde{\mathbf{F}}$ and the “nice” vector field \mathbf{F} . In particular, for such t and \mathbf{x} , $\Phi_t(\mathbf{x})$ is the same for $\tilde{\mathbf{F}}$ and for \mathbf{F} , and may be computed using the flow curves we have found in Example 131:

$$\Phi_t(x, y) = (x/(1 - tx), y + t) .$$

However, for larger times explosion may have occurred for some points in $B_{R/2}(\mathbf{0})$, and so for such t , the transformation cannot be defined on $B_R(\mathbf{0})$.

9.2 Flux integrals in \mathbb{R}^2

Let \mathbf{F} be a continuously differentiable vector field defined on \mathbb{R}^2 , and suppose that for some $a > 0$, the corresponding flow transformation Φ_t is well defined for $|t| < a$. Let C be a parameterized curve in \mathbb{R}^2 given by $\mathbf{x}(u)$ for $0 \leq u \leq b$, some $b > 0$. The curve C may be open or closed.

For $0 < t < a$, consider the set of points A_t that get “swept across C ” by time t . To be concrete, let us take \mathbf{F} to be the rotational vector field

$$\mathbf{F}(x, y) = (-y, x)$$

from Example 130, and let us take C to be the line segment running from $(0, 1)$ to $(0, 3)$. We may parameterize C as

$$\mathbf{x}(u) = (0, 1 + 2u) \quad \text{for } 0 \leq u \leq 1 .$$

Now it is easy to see that $\Phi_t(\mathbf{x}) = \begin{bmatrix} \cos t & -\sin t \\ \sin t & \cos t \end{bmatrix} \mathbf{x}$ crosses C by time t if and only if \mathbf{x} belongs to the keystone-shaped region given in polar coordinates by

$$1 \leq r \leq 3 \quad \text{and} \quad \frac{\pi}{2} - t \leq \theta \leq \pi . \quad (9.11)$$

Thus, the region A_t is the part of the plane given by (9.11). We then compute

$$\text{area}(A_t) = \frac{1}{2}(3^2 - 1^2)t = 4t .$$

Thus, area is being swept across the segment C , from right to left, at a steady rate of 4 units of area per unit of time. We express this by saying that *the flux across C , from right to left, produced by \mathbf{F} is 4.*

In this example, all of the area that got swept across the curve at by time t got swept across from right to left. For other curves, this need not be the case. For example, replace C by the line segment running from $(0, -1)$ to $(0, 1)$. Then there is a region in the upper half plane that gets swept right to left across the curve by time t , but there is a region of equal area in the lower half plane the swept left to right across the curve by time t . In this case, there is no net area flowing across the curve from right to left (or from left to right), and so in this case we say the flux is zero. *Flux refers to the net rate of flow of area from one side of a curve to the other.*

Now that we have in mind a picture of what we mean by flux, let us proceed to precise mathematical definitions.

Definition 94 (Oriented curve). *Consider a differentiable curve C , and let curve $\mathbf{x}(s)$, $0 \leq s \leq s_*$ be its arc length parameterization. The curve is simple in case for $s_1 < s_2$, $\mathbf{x}(s_1) \neq \mathbf{x}(s_2)$, except possibly if $s_1 = 0$ and $s_2 = s_*$, in which case it is a simple closed curve. At each point $\mathbf{x}(s)$, $0 \leq s \leq s_*$, there are two unit vectors that are orthogonal to $\mathbf{T}(s)$, namely $\pm \mathbf{T}(s)^\perp$. An orientation of such a curve C is a specification of either $\mathbf{T}(s)^\perp$ or $-\mathbf{T}(s)^\perp$, making the same sign choice for all s , as the preferred normal $\mathbf{N}(s)$. We think of $\mathbf{N}(s)$ as pointing to the “positive side” of C . We call such a curve with an orientation an oriented curve.*

For a simple closed curve in \mathbb{R}^2 , there are always an inside and an outside. To orient a simple closed is to specify whether we regard the inside as the positive side, in which case we choose $\mathbf{N}(s)$ to point inward everywhere along the curve, or whether we regard the outside as the positive side, in which case we choose $\mathbf{N}(s)$ to point outward everywhere along the curve.

Now consider an oriented curve C . Let us “zoom in” and look at what the flow is doing near a segment of the curve between $\mathbf{x}(s_0)$ and $\mathbf{x}(s_1)$ with $s_1 - s_0$ small but positive. Since for very small values of t ,

If the segment is very short, then the approximation

$$\mathbf{F}(x) \approx \mathbf{F}(\mathbf{x}(s_0))$$

will be a good one. We know the flow associated to the constant vector field

$$\tilde{\mathbf{F}}(\mathbf{x}) := \mathbf{F}(\mathbf{x}(s_0)) ;$$

it is simply

$$\tilde{\Phi}_t(\mathbf{x}) = \mathbf{x} + t\mathbf{F}(\mathbf{x}(s_0)) .$$

For small values of t and \mathbf{x} near the segment, this should be a good approximation to $\Phi_t(\mathbf{x})$, and in the limit in which we take t to zero, and the length of our segment to zero, this approximation will become exact.

In particular, a point \mathbf{x} gets carried across the segment between $\mathbf{x}(s_0)$ and $\mathbf{x}(s_1)$ by time t if and only for some s with $s_0 \leq s \leq s_1$ and some u with $0 \leq u \leq t$,

$$\mathbf{x} = \mathbf{x}(s) - u\mathbf{F}(\mathbf{x}(s_0)) ,$$

because in this case, and only this case,

$$\tilde{\Phi}_u(\mathbf{x}) = [\mathbf{x}(s) - u\mathbf{F}(\mathbf{x}(s_0))] + u\mathbf{F}(\mathbf{x}(s_0)) = \mathbf{x}(s)$$

which is on the segment.

The region of points that get swept across the segment by the flow is an approximate parallelogram with vertices

$$\mathbf{x}(s_0) , \quad \mathbf{x}(s_1) , \quad \mathbf{x}(s_0) - t\mathbf{F}(\mathbf{x}(s_0)) , \quad \text{and} \quad \mathbf{x}(s_1) - t\mathbf{F}(\mathbf{x}(s_0)) .$$

Let $\Delta s = s_1 - s_0$ be the arc length of our segment. Then the sides of the parallelogram that meet at $\mathbf{x}(s_0)$ run along the vectors

$$\Delta s \mathbf{T}(s_0) \quad \text{and} \quad t\mathbf{F}(\mathbf{x}(s_0)) .$$

If we regard the first vector as running along the base of the parallelogram, the height of the parallelogram is $|\mathbf{N}(s_0) \cdot t\mathbf{F}(\mathbf{x}(s_0))|$. Thus the area swept across this segment by the flow in time t is approximately

$$|\mathbf{N}(s_0) \cdot t\mathbf{F}(\mathbf{x}(s_0))| \Delta s .$$

Next, notice that

$$\mathbf{N}(s_0) \cdot \mathbf{F}(\mathbf{x}(s_0)) > 0$$

if and only if $\mathbf{F}(\mathbf{x}(s_0))$ points to the positive side of C . Thus, the rate at which area is flowing across this segment of C , from the negative side to the positive side, is

$$t\mathbf{F}(\mathbf{x}(s_0)) \cdot \mathbf{N}(s_0) \Delta s ,$$

where we have reordered terms in what will turn out to be a convenient way. Dividing by t to get the *rate* at which area is flowing, this gives us the *flux element* for this segment of the curve.

Now adding up the flux elements all along the curve, and taking the limit as the length of the segments goes to zero, the sum becomes an integral and our approximations become exact, and we find that the flux is given by

$$\int_0^{s_*} \mathbf{F}(\mathbf{x}(s)) \cdot \mathbf{N}(s) ds .$$

We shall also write this simply as

$$\int_C \mathbf{F} \cdot \mathbf{N} ds ,$$

and in case C is a simple closed curve, as

$$\oint_C \mathbf{F} \cdot \mathbf{N} ds ,$$

where the special integration symbol simply emphasizes that the curve is closed.

Definition 95 (Flux integral along an oriented curve). *Let C be an oriented curve, and let \mathbf{F} be a differentiable vector field defined in a neighborhood of C . The flux integral, representing the rate of flow of area across C , from the negative side to the positive side, under the flow generated by \mathbf{F} , is*

$$\int_C \mathbf{F} \cdot \mathbf{N} ds := \int_0^{s_*} \mathbf{F}(\mathbf{x}(s)) \cdot \mathbf{N}(s) ds ,$$

where $\mathbf{x}(s)$, $0 \leq s \leq s_*$ is an arclength parameterization of C (there are two of them), and $\mathbf{N}(s)$ is the preferred unit normal at $\mathbf{x}(s)$.

9.2.1 Computing flux integrals in \mathbb{R}^2

Based on the discussion in the previous section, you might think that computing a flux integral will involve finding an arc length parameterization of a curve. Fortunately, things are much easier than that.

Let $\mathbf{x}(t) = (x(t), y(t))$, $b \leq t \leq c$, be *any* parameterization of a curve C . Then the element of arc length is

$$ds = \sqrt{(x'(t))^2 + (y'(t))^2} dt .$$

The unit tangent vector $\mathbf{T}(t)$ is

$$\mathbf{T}(t) = \frac{1}{\sqrt{(x'(t))^2 + (y'(t))^2}} (x'(t), y'(t)) ,$$

and so the preferred unit normal $\mathbf{N}(t)$ is

$$\mathbf{N}(t) = \pm \frac{1}{\sqrt{(x'(t))^2 + (y'(t))^2}} (-y'(t), x'(t)) ,$$

with the sign depending on our preference.

Putting things together,

$$\begin{aligned} \mathbf{N}(t) ds &= \pm \frac{1}{\sqrt{(x'(t))^2 + (y'(t))^2}} (-y'(t), x'(t)) \sqrt{(x'(t))^2 + (y'(t))^2} dt \\ &= \pm (-y'(t), x'(t)) dt . \end{aligned}$$

The good news is that the factors involving square roots have cancelled out. Let us define

$$d\mathbf{x}^\perp(t) := (-y'(t), x'(t)) dt .$$

Then we have, for any parameterization $\mathbf{x}(t)$, with $s = s(t)$,

$$\mathbf{N}(s) ds = \pm d\mathbf{x}^\perp(t) ,$$

so that

$$\mathbf{F}(\mathbf{x}(s)) \cdot \mathbf{N}(s) ds = \pm \mathbf{F}(\mathbf{x}(t)) \cdot d\mathbf{x}^\perp(t) , \quad (9.12)$$

and hence

$$\int_C \mathbf{F} \cdot \mathbf{N} ds = \pm \int_b^c \mathbf{F}(\mathbf{x}(t)) \cdot d\mathbf{x}^\perp(t) . \quad (9.13)$$

Example 136 (Computing a flux integral in \mathbb{R}^2). *Let*

$$\mathbf{F}(x, y) = (-y, x)$$

from Example 130, and let us take C to be the line segment running from $(0, 1)$ to $(0, 3)$, parameterized as

$$\mathbf{x}(t) = (0, 1 + 2t) \quad \text{for} \quad 0 \leq t \leq 1 .$$

We orient C so the positive side is to the left.

The we compute

$$\mathbf{F}(\mathbf{x}(t)) = (-1 - 2t, 0) \quad \text{and} \quad \mathbf{N}(t) ds = \pm (-2, 0) dt .$$

Choosing the plus sign, the first component of \mathbf{N} is negative, so \mathbf{N} points to the positive side. Thus, this is the correct choice of the sign.

Putting things together, we have from (9.13) that

$$\int_C \mathbf{F} \cdot \mathbf{N} ds = \int_0^1 (-1 - 2t, 0) \cdot (-2, 0) dt = \int_0^1 (2 + 4t) dt = 4 ,$$

which is what we found earlier by computing areas.

Example 137 (Computing a flux integral in \mathbb{R}^2). Let

$$\mathbf{F}(x, y) = (xy, x^2 - y^2) ,$$

and let us take C to be the circle of unit radius centered on $(1, 1)$, oriented so the outside is the positive side.

We first parameterize C . The standard way is

$$\mathbf{x}(t) = (1 + \cos t, 1 + \sin t) ,$$

for $0 \leq t \leq 2\pi$.

Then we compute

$$\mathbf{F}(\mathbf{x}(t)) = (1 + \cos t + \sin t + \cos t \sin t , \cos^2 t - \sin^2 t + 2(\cos t - \sin t))$$

and

$$\mathbf{N}(t) ds = \pm (\cos t, \sin t) dt$$

Choosing the plus sign, the first component of \mathbf{N} is positive at $t = 0$, so \mathbf{N} points to the positive side. Thus, this is the correct choice of the sign. Therefore with $s = s(t)$,

$$\mathbf{F}(\mathbf{x}(s)) \cdot \mathbf{N}(s) ds = [\cos t(1 + \cos t + \sin t + \cos t \sin t) + \sin t(\cos^2 t - \sin^2 t + 2(\cos t - \sin t))] dt .$$

Putting things together, we have from (9.13) that

$$\begin{aligned} \int_C \mathbf{F} \cdot \mathbf{N} ds = \\ \int_0^{2\pi} [\cos t(1 + \cos t + \sin t + \cos t \sin t) + \sin t(\cos^2 t - \sin^2 t + 2(\cos t - \sin t))] dt . \end{aligned} \quad (9.14)$$

There are many terms, but most integrate to zero. Discarding all such terms, we are left with

$$\int_C \mathbf{F} \cdot \mathbf{N} ds = \int_0^{2\pi} [\cos^2 t - 2\sin^2 t] dt = -\pi .$$

In this example, there is more area being swept into the disc bounded by C than there is being swept out.

9.2.2 The divergence and flux density

Definition 96 (Divergence of a vector field). *Let $\mathbf{F} = (f_1, \dots, f_n)$ be a differentiable vector field defined on an open set $U \subset \mathbb{R}^n$. Then the divergence of \mathbf{F} is the real valued function $\operatorname{div}(\mathbf{F})$ defined by*

$$\operatorname{div}(\mathbf{F})(\mathbf{x}) = \sum_{j=1}^n \frac{\partial}{\partial x_j} f_j(\mathbf{x}) .$$

Example 138 (Computing a divergence). *Let*

$$\mathbf{F}(x, y) = (xy, x^2 - y^2) ,$$

Then

$$\operatorname{div}(\mathbf{F})(x, y) = \frac{\partial}{\partial x} xy + \frac{\partial}{\partial y} (x^2 - y^2) = y - 2y = -y .$$

We are now ready to state the Divergence Theorem for flux integrals in \mathbb{R}^2 :

Theorem 88 (The Divergence Theorem for flux integrals in \mathbb{R}^2). *Let C be a simple closed curve in \mathbb{R}^2 , and let D denote the region bounded by C . Orient C so that the outside is the positive side. Let \mathbf{F} be any continuously differentiable vector field defined on a neighborhood of D . Then*

$$\oint_C \mathbf{F} \cdot \mathbf{N} ds = \int_D \operatorname{div}(\mathbf{F}) dA . \quad (9.15)$$

Example 139 (Using the Divergence Theorem to compute flux). *Let*

$$\mathbf{F}(x, y) = (xy, x^2 - y^2) .$$

Let C be the circle of unit radius centered on $(1, 1)$, oriented so the outside is the positive side. Then since, as computed in Example 138, $\operatorname{div}(\mathbf{F})(x, y) = -y$,

$$\oint_C \mathbf{F} \cdot \mathbf{N} ds = - \int_D y dA ,$$

where D is the disk of unit radius centered on $(1, 1)$. By symmetry, the average value of y in D is 1. Hence

$$\frac{\int_D y dA}{\int_D 1 dA} = 1 .$$

Since

$$\begin{aligned} \int_D 1 dA &= \operatorname{area}(D) = \pi , \\ - \int_D y dA &= -\pi . \end{aligned}$$

This is what we found in Example 137, but here the computation is simpler.

The Divergence Theorem is also useful for computing the flux across a curve that is not closed.

Example 140 (The Divergence Theorem and flux across open curves). *Let C_1 be the part of the parabola $y = 4 - x^2$ lying above the x -axis oriented so the upward side is the positive side. Let $\mathbf{F}(x, y) = (x^3 y - y^2 + x, x^2 y - 3x + 5y)$.*

The endpoints of C_1 are $(-2, 0)$ and $(2, 0)$. Let C_2 be the straight line segment from $(-2, 0)$ and $(2, 0)$. Finally, let C be the simple closed curve that runs from $(-2, 0)$ to $(2, 0)$ along C_2 , and then from $(2, 0)$ and $(-2, 0)$ along C_1 .

Notice that orienting C so the outside is the positive side coincides with the original orientation on C_1 , and induces the orientation on C_2 in which the positive side is the downward side.

Because integrals are limits of sums,

$$\oint_C \mathbf{F} \cdot \mathbf{N} ds = \int_{C_1} \mathbf{F} \cdot \mathbf{N} ds + \int_{C_2} \mathbf{F} \cdot \mathbf{N} ds .$$

But by the Divergence Theorem,

$$\oint_C \mathbf{F} \cdot \mathbf{N} ds = \int_D \operatorname{div}(\mathbf{F}) dA$$

where D is the region bounded by C . Therefore,

$$\int_{C_1} \mathbf{F} \cdot \mathbf{N} ds = \int_D \operatorname{div}(\mathbf{F}) dA - \int_{C_2} \mathbf{F} \cdot \mathbf{N} ds . \quad (9.16)$$

We now compute the two integrals on the right, each of which is much easier than the integral on the left.

Indeed, we can parameterize C_2 by

$$\mathbf{x}(u) = (-2 + 4s, 0)$$

for $0 \leq s \leq 4$, and this is an arc length parameterization. (Arc length parameterizations are easy for straight line segments!) Then since $y = 0$ all along C_2 , it is easy to compute $\mathbf{F}(\mathbf{x}(s))$:

$$\mathbf{F}(\mathbf{x}(u)) := (-2 + s, 6 - 3s) .$$

Also, since \mathbf{N} is the downward pointing unit vector, $\mathbf{N}(s) ds = (0, -1) ds$. Putting it all together,

$$\int_{C_2} \mathbf{F} \cdot \mathbf{N} ds = \int_0^4 (3s - 6) ds = \frac{3}{2} 16 - 24 = 0 . \quad (9.17)$$

Next, let us compute $\int_D \operatorname{div}(\mathbf{F}) dA$. The first step is to compute

$$\operatorname{div}(\mathbf{F})(x, y) = 3x^2 y + (6 + x^2) .$$

The region D is given by

$$0 \leq y \leq 4 - x^2 \quad \text{and} \quad -2 \leq x \leq 2 .$$

Thus,

$$\begin{aligned} \int_D \operatorname{div}(\mathbf{F}) dA &= \int_{-2}^2 \left(\int_0^{4-x^2} [3x^2 y + (6 + x^2)] dy \right) dx \\ &= \int_{-2}^2 \left[\frac{3}{2} x^2 (4 - x^2)^2 + (6 + x^2)(4 - x^2) \right] dx \\ &= \int_{-2}^2 \left[48 - 14x^2 + \frac{1}{2} x^4 \right] dx = \frac{1856}{15} . \end{aligned}$$

Combining this result with (9.16) and (9.17), we have

$$\int_{C_2} \mathbf{F} \cdot \mathbf{N} ds = \frac{1856}{15} .$$

To really appreciate this example, you should carry out the direct computation of $\int_{C_2} \mathbf{F} \cdot \mathbf{N} ds$, which you will find to be somewhat messy.

The point of the last example is that the Divergence Theorem specifies the price that must be paid to “trade” a curve C_1 connecting two point \mathbf{x}_0 and \mathbf{x}_1 in on another, simpler curve C_2 connecting the same points. When C_1 does not intersect C_2 except at the endpoints \mathbf{x}_0 and \mathbf{x}_1 , so that following C_1 from \mathbf{x}_0 to \mathbf{x}_1 , and then C_2 from \mathbf{x}_1 back to \mathbf{x}_0 produces a simple closed curve, the “trade in” is done just as in the previous example.

9.2.3 Proof and interpretation of the Divergence Theorem

Now that we have seen some examples of how to use the Divergence Theorem, we ask:

- Why is the Divergence Theorem true?
- Why is the divergence related to flux?

The two questions are closely related. To answer the first question, we go back to the notion of flux as a rate of flow of area.

Proof of the Divergence Theorem in \mathbb{R}^2 : Let C be a simple closed curve bounding the region $D \subset \mathbb{R}^2$, and orient C so that the positive side of C is the outside of D .

Let \mathbf{F} be a differentiable vector field defined on a neighborhood of D , and let Φ_t denote the flow transformation generated by \mathbf{F} at time t .

Define

$$D_t = \{ \Phi_t(\mathbf{x}) : \mathbf{x} \in D \} .$$

Then notice that D_{-t} is precisely the set of points that are in D after running the flow for a time t . Therefore, the net area swept out of D by the flow in time t is

$$\text{area}(D) - \text{area}(D_{-t}) .$$

(Whatever was in D that is not replaced by what comes in from D_t has gone out.)

To compute $\text{area}(D_{-t})$ we use the change of variables given by the flow transformation:

$$(u, v) := \Phi_{-t}(x, y) .$$

Note that by part (4) of Theorem 87, Φ_{-t} is continuously differentiable, and so it has a Jacobian matrix $[D_{\Phi_{-t}}(\mathbf{x})]$.

Then by the change of variables formula,

$$\begin{aligned} \text{area}(D_{-t}) &= \int_{D_{-t}} 1 d^2 \mathbf{u} = \int_D 1 |\det (D_{\Phi_{-t}}(\mathbf{x}))| d^2 \mathbf{x} . \\ &= \int_D \det (D_{\Phi_{-t}}(\mathbf{x})) d^2 \mathbf{x} , \end{aligned}$$

where in the second line we have dropped the factor of 1 and the absolute value signs. These are not needed since for $t = 0$, Φ_{-1} is the identity transformation, and so $\det(D_{\Phi_0}(\mathbf{x})) = 1$. By continuity, $\det(D_{\Phi_{-t}}(\mathbf{x}))$ is positive for small t which is what concerns us.

Therefore, the flux out of D is given by

$$\begin{aligned} \text{flux out of } D &= \lim_{t \rightarrow 0} \frac{\text{area}(D) - \text{area}(D_{-t})}{t} \\ &= \lim_{t \rightarrow 0} \frac{1}{t} \int_D [1 - \det(D_{\Phi_{-t}}(\mathbf{x}))] d^2\mathbf{x} \\ &= - \int_D \frac{\partial}{\partial t} \det(D_{\Phi_{-t}}(\mathbf{x})) \Big|_{t=0} d^2\mathbf{x} . \end{aligned}$$

To complete the proof, we only need to show that

$$\frac{\partial}{\partial t} \det(D_{\Phi_{-t}}(\mathbf{x})) \Big|_{t=0} = -\text{div}(\mathbf{F}(\mathbf{x})) . \quad (9.18)$$

Here is one way to do this. We compute an approximation to $\det(D_{\Phi_{-t}}(\mathbf{x}))$ that is accurate to the leading order in t . For small values of t ,

$$\Phi_{-t}(\mathbf{x}) \approx \mathbf{x} - t\mathbf{F}(\mathbf{x}) .$$

That is, with $\mathbf{F}(x, y) = (f(x, y), g(x, y))$, we have

$$u(x, y) = x - tf(x, y) \quad \text{and} \quad v(x, y) = y - tg(x, y) .$$

It follows that

$$D_{\Phi_{-t}}(\mathbf{x}) \approx \begin{bmatrix} 1 - t \frac{\partial}{\partial x} f(x, y) & -t \frac{\partial}{\partial y} f(x, y) \\ -t \frac{\partial}{\partial x} g(x, y) & 1 - t \frac{\partial}{\partial y} g(x, y) \end{bmatrix} .$$

Therefore, to leading order in t , which is all that concerns us in the derivative (9.18) that we are aiming to compute,

$$\begin{aligned} \det(D_{\Phi_{-t}}(\mathbf{x})) &\approx 1 - t \left(\frac{\partial}{\partial x} f(x, y) + \frac{\partial}{\partial y} g(x, y) \right) \\ &= 1 - t \text{div}(\mathbf{F}(\mathbf{x})) , \end{aligned}$$

□

The key to the proof we have just given the identity (9.18). We now give a second proof of this identity that is valid in \mathbb{R}^n for all n , which we shall use later to prove the divergence theorem in \mathbb{R}^3 . We shall use the following definition and theorem:

The following definition and theorem complete out work:

Definition 97 (Trace of an $n \times n$ matrix). *Let A be an $n \times n$ matrix. The trace of A , $\text{tr}(A)$ is defined by*

$$\text{tr}(A) = \sum_{i=1}^n A_{i,i} .$$

That is, $\text{tr}(A)$ is the sum of the diagonal elements of A .

For example, if \mathbf{F} is any vector field on \mathbb{R}^n , then since $[D_{\mathbf{F}}]_{i,j} = \frac{\partial}{\partial x_j} f_i(\mathbf{x})$,

$$\operatorname{tr}(D_{\mathbf{F}}(\mathbf{x})) = \sum_{i=1}^n \frac{\partial}{\partial x_i} f_i(\mathbf{x}) = \operatorname{div}(\mathbf{F}(\mathbf{x})) . \quad (9.19)$$

• In other words, the divergence of \mathbf{F} is the trace of the Jacobian $[D_{\mathbf{F}}]$.

The following important theorem relates the trace and the determinant:

Theorem 89 (The trace, determinant and derivatives). *Let $A(t)$ be an $n \times n$ matrix valued function of $t \in \mathbb{R}$. Suppose that for each $1 \leq i, j \leq n$, $A_{i,j}(t)$ is differentiable at $t = 0$, with*

$$B_{i,j} = \frac{d}{dt} A_{i,j}(t) \Big|_{t=0} .$$

Suppose also that $A(0) = I_{n \times n}$. Then

$$\frac{d}{dt} \det(A(t)) \Big|_{t=0} = \operatorname{tr}(B) . \quad (9.20)$$

Proof By the determinant formula, $\det(A(t)) = \sum_{\sigma \in S_n} \chi(\sigma) \prod_{i=1}^n A_{i,\sigma(i)}(t)$. By the product rule and the definition of B ,

$$\frac{d}{dt} \prod_{i=1}^n A_{i,\sigma(i)}(t) \Big|_{t=0} = \sum_{j=1}^n \left(B_{j,\sigma(j)} \prod_{i=1, i \neq j}^n A_{i,\sigma(i)}(0) \right) .$$

since $A(0) = I_{n \times n}$,

$$\prod_{i=1, i \neq j}^n A_{i,\sigma(i)}(0) = 0$$

unless $\sigma(i) = i$ for each $i \neq j$ from 1 to n . But then since σ is one-to-one, it must also be the case that $\sigma(j) = j$. That is, $\prod_{i=1, i \neq j}^n A_{i,\sigma(i)}(0) = 0$ unless σ is the identity permutation, in which case $\prod_{i=1, i \neq j}^n A_{i,\sigma(i)}(0) = 1$. Therefore,

$$\frac{d}{dt} \left(\sum_{\sigma \in S_n} \chi(\sigma) \prod_{i=1}^n A_{i,\sigma(i)}(t) \right) = \sum_{j=1}^n B_{j,j} = \operatorname{tr}(B) .$$

□

We are now ready to prove the n -dimensional version of (9.18).

Theorem 90 (The divergence and flows). *Let \mathbf{F} be a continuously differentiable vector field on \mathbb{R}^n , and suppose that for some finite L $\|D_{\mathbf{F}}(\mathbf{x})\|_{\mathbf{F}} \leq L$ for all $\mathbf{x} \in \mathbb{R}^n$ so that the group of flow transformations Φ_t generated by \mathbf{F} is well defined, and each Φ_t is a differentiable one-to-one transformation of \mathbb{R}^n onto \mathbb{R}^n . Then*

$$\frac{\partial}{\partial t} \det(D_{\Phi_{-t}}(\mathbf{x})) \Big|_{t=0} = -\operatorname{div}(\mathbf{F}(\mathbf{x})) . \quad (9.21)$$

Proof: Let $(\Phi_{-t}(\mathbf{x}))_i$ denote the i th component of the vector $\Phi_{-t}(\mathbf{x})$. Then, by definition,

$$[D_{\Phi_{-t}}(\mathbf{x})]_{i,j} = \frac{\partial}{\partial x_j} (\Phi_{-t}(\mathbf{x}))_i ,$$

and so, by Clairault's Theorem,

$$\begin{aligned} \frac{\partial}{\partial t} [D_{\Phi_{-t}}(\mathbf{x})] &= \frac{\partial}{\partial t} \frac{\partial}{\partial x_j} (\Phi_{-t}(\mathbf{x}))_i \\ &= \frac{\partial}{\partial x_j} \left(\frac{\partial}{\partial t} (\Phi_{-t}(\mathbf{x}))_i \right) . \end{aligned}$$

Next, since $t \mapsto \Phi_{-t}(\mathbf{x})$ is the flow line *backwards* along the vector field $\mathbf{F} := (f_1, \dots, f_n)$, the derivative of this curve at any time t is $-\mathbf{F}$ at the point $\Phi_{-t}(\mathbf{x})$. That is,

$$\frac{\partial}{\partial t} (\Phi_{-t}(\mathbf{x})) = -\mathbf{F}(\Phi_{-t}(\mathbf{x})) .$$

Therefore,

$$\frac{\partial}{\partial t} [D_{\Phi_{-t}}(\mathbf{x})]_{i,j} \Big|_{t=0} = -\frac{\partial}{\partial x_j} f_i(\mathbf{x}) = -[D_{\mathbf{F}}(\mathbf{x})]_{i,j} . \quad (9.22)$$

Next, since Φ_0 is the identity transformation; i.e., $\Phi_0(\mathbf{x}) = \mathbf{x}$,

$$[D_{\Phi_0}(\mathbf{x})] = I_{n \times n} , \quad (9.23)$$

the $n \times n$ identity matrix. Therefore, combining (9.22), (9.23) and Theorem 89, we conclude

$$\frac{\partial}{\partial t} \det(D_{\Phi_{-t}}(\mathbf{x})) \Big|_{t=0} = \text{tr} \left(-[D_{\mathbf{F}}(\mathbf{x})]_{i,j} \right) ,$$

which is equivalent to (9.21). □

9.3 Flux integrals in \mathbb{R}^3

9.3.1 The flux out of a region $\mathcal{V} \subset \mathbb{R}^3$

Let \mathcal{V} be a bounded region in \mathbb{R}^3 bounded by a simple (non-self intersecting) closed surface \S . Let \mathbf{F} be a differentiable vector field defined on a neighborhood of \mathcal{V} . Then the flux across \S generated by \mathbf{F} , from inside to outside is the rate at which volume is swept out of \mathcal{V} by the flow Φ_t generated by \mathbf{F} .

To compute the flux, we reason exactly as in the last section: After running the flow a short time t , the set of points that are inside \mathcal{V} are the points in

$$\mathcal{V}_{-t} := \{ \Phi_{-t}(\mathbf{x}) : \mathbf{x} \in \mathcal{V} \} .$$

Therefore, the net flux out of \mathcal{V} is given by

$$\text{net flux out of } \mathcal{V} = \lim_{t \rightarrow 0} \frac{1}{t} [\text{vol}(\mathcal{V}) - \text{vol}(\mathcal{V}_{-t})] .$$

Also, just as in the last section, the change of variables formula gives

$$\text{vol}(\mathcal{V}_{-t}) = \int_{\mathcal{V}} \det(D_{\Phi_{-t}}(\mathbf{x})) \, d^3 \mathbf{x} ,$$

so that

$$\lim_{t \rightarrow 0} \frac{1}{t} [\text{vol}(\mathcal{V}) - \text{vol}(\mathcal{V}_{-t})] = - \int_{\mathcal{V}} \frac{\partial}{\partial t} \det([D_{\Phi_{-t}}(\mathbf{x})]) \Big|_{t=0} d^3 \mathbf{x} .$$

Then by Theorem 90 for $n = 3$, we have that

$$\text{net flux out of } \mathcal{V} = \int_{\mathcal{V}} \text{div}(\mathbf{F}(\mathbf{x})) d^3 \mathbf{x} . \quad (9.24)$$

9.3.2 Flux across an oriented surface \S

There is also a direct way to calculate the flux across a surface \S in terms of a surface integral. In this case, the surface need not be closed, but it must be orientable.

Definition 98 (Orientable surface). *Let \S be a differentiable parameterized surface in \mathbb{R}^3 . At each point of the surface, there are two sides to the tangent plane to the surface at that point, and hence two unit normal vectors to the surface at each point. The surface is orientable if it is possible to specify a preferred unit normal vector $\mathbf{N}(\mathbf{x})$ at each point of the surface so that $\mathbf{N}(\mathbf{x})$ is a continuous function of the point \mathbf{x} on the surface. Such a specification of a preferred unit normal, if one exists, is called an orientation of \S .*

If \S is simple and closed, then it bounds a region \mathcal{V} , and it clearly has two sides: an inside and an outside. We can choose $\mathbf{N}(\mathbf{x})$ at each point $\mathbf{x} \in \S$ to either point outward from \mathcal{V} or inward into \mathcal{V} . In the first case we say \mathbf{N} is the *outward unit normal vector*, and in the second case we say \mathbf{N} is the *inward unit normal vector*.

Surfaces \S that are not closed may or may not be orientable. The Möbius band is an example of a non-orientable surface \mathbb{R}^3 .

Now let \S be an oriented surface with preferred unit normal \mathbf{N} . Let \mathbf{F} be a differentiable vector field defined in a neighborhood of \S . Reasoning exactly as in the last section, the rate at which the flow associated to \mathbf{F} sweeps volume across the surface \S , from the negative side to the positive side, is

$$\int_{\S} \mathbf{F} \cdot \mathbf{N} dS .$$

Also, just as in the previous section, it is easy to work out the integral in a concrete parameterization of \S .

Let $\mathbf{x}(u, v)$ with $(u, v) \in U \subset \mathbb{R}^2$ be a parameterization of \S . Then with

$$\mathbf{T}_u = \frac{\partial}{\partial u} \mathbf{x}(u, v) \quad \text{and} \quad \mathbf{T}_v = \frac{\partial}{\partial v} \mathbf{x}(u, v) ,$$

we know that

$$\mathbf{N}(u, v) = \pm \frac{1}{\|\mathbf{T}_u \times \mathbf{T}_v(u, v)\|} \mathbf{T}_u \times \mathbf{T}_v(u, v)$$

and

$$dS = \|\mathbf{T}_u \times \mathbf{T}_v(u, v)\| du dv .$$

Therefore the *flux element* for an infinitesimal tile on the surface is

$$\mathbf{F} \cdot \mathbf{N} dS = \mathbf{F}(\mathbf{x}(u, v)) \cdot \mathbf{T}_u \times \mathbf{T}_v(u, v) du dv ,$$

and the flux integral is given, in ready-to-be-computed form as

$$\int_{\S} \mathbf{F} \cdot \mathbf{N} dS = \int_U \mathbf{F}(\mathbf{x}(u, v)) \cdot \mathbf{T}_u \times \mathbf{T}_v(u, v) du dv . \quad (9.25)$$

Example 141 (Computing the flux across a surface). *Let \mathbf{F} be the vector field*

$$\mathbf{F} = (2xyz - y^2, x^2z - 2xy, x^2y) .$$

Let \S be the part of the paraboloid $z = 1 - x^2 - y^2$ that lies above the x, y plane, oriented so its preferred normal points upward. We will now compute the flux

$$\int_{\S} \mathbf{F} \cdot \mathbf{N} dS$$

using (9.25).

The first step is to parameterize the surface. Let us use cylindrical coordinates. Then the equation defining \S is $z = 1 - r^2$ and $z = 0$, and $z \geq 0$ becomes $r \leq 1$. So

$$\mathbf{x}(r, \theta) = (r \cos \theta, r \sin \theta, 1 - r^2) ,$$

and the parameter domain U is given by

$$0 \leq r \leq 1 \quad \text{and} \quad 0 \leq \theta \leq 2\pi .$$

Differentiating, we find

$$\mathbf{T}_r(r, \theta) = (\cos \theta, \sin \theta, -2r) ,$$

and

$$\mathbf{T}_\theta(r, \theta) = (-r \sin \theta, r \cos \theta, 0) .$$

We then compute

$$\mathbf{T}_r \times \mathbf{T}_\theta(r, \theta) = (2r^2 \cos \theta, 2r^2 \sin \theta, r) .$$

notice that the third component is positive, so this vector points upwards. Thus,

$$\mathbf{N} dS = \mathbf{T}_r \times \mathbf{T}_\theta dr d\theta .$$

We then compute

$$\mathbf{F}(\mathbf{x}(r, \theta)) = (2r^2(1 - r^2) \cos \theta \sin \theta - r^2 \sin^2 \theta, r^2(1 - r^2) \cos^2 \theta - 2r^2 \cos \theta \sin \theta, r^3 \cos^2 \theta \sin \theta) .$$

Therefore, the flux element is

$$\begin{aligned} \mathbf{F}(\mathbf{x}(r, \theta)) \cdot \mathbf{N}(r, \theta) dS &= [2r^2(1 - r^2) \cos \theta \sin \theta - r^2 \sin^2 \theta][2r^2 \cos \theta] dr d\theta \\ &+ [r^2(1 - r^2) \cos^2 \theta - 2r^2 \cos \theta \sin \theta][2r^2 \sin \theta] dr d\theta \\ &+ [r^3 \cos^2 \theta \sin \theta][r] dr d\theta . \end{aligned}$$

We now integrate over U . But since

$$\int_0^{2\pi} \sin^2 \theta \cos \theta d\theta = 0 \quad \text{and} \quad \int_0^{2\pi} \cos^2 \theta \sin \theta d\theta = 0 .$$

all of the integral give zero. Hence there is no net flux across \S .

9.3.3 The Divergence Theorem in \mathbb{R}^3

For surfaces \S that bound a connected region $\mathcal{V} \subset \mathbb{R}^3$, which are necessarily simple closed surfaces, we usually write

$$\oint_{\S} \mathbf{F} \cdot \mathbf{N} dS$$

to denote the flux out across \S . That is, we canonically take the outward unit normal as the orientation of \S in this case; this is indicated by the special symbol for the integral. Combining our two ways of computing flux in this case, we have:

Theorem 91 (The Divergence Theorem in \mathbb{R}^3). *Let \mathcal{V} be a bounded connected region in \mathbb{R}^3 , and let \S be its boundary, equipped with the outward unit normal orientation. Then*

$$\oint_{\S} \mathbf{F} \cdot \mathbf{N} dS = \int_{\mathcal{V}} \operatorname{div}(\mathbf{F}(\mathbf{x})) d^3\mathbf{x} .$$

The Divergence Theorem is useful for computing flux even for surfaces that are not closed: *It tells you how much you will change the flux by changing the surface \S into something simpler.* Here is an example of this:

Example 142 (Trading in one surface on another). *There is a better way to compute the flux integral in the Example 141. Notice that if we let \S_2 denote the unit disk in the x, y plane, then together \S and \S_2 bound the region \mathcal{V} consisting of points (x, y, z) in \mathbb{R}^3 with*

$$0 \leq z \leq 1 - x^2 - y^2 .$$

Also, the outward unit normal on the boundary of \mathcal{V} coincides with the preferred unit normal on \S . Thus the net flux outward across the boundary of \mathcal{V} is given by

$$\int_{\S} \mathbf{F} \cdot \mathbf{N} dS + \int_{\S_2} \mathbf{F} \cdot \mathbf{N} dS ,$$

where on \S_2 we take the downward unit normal, since this is the outward unit normal.

We can now use (9.24) to compute the net flux outward across the boundary of \mathcal{V} , finding

$$\int_{\mathcal{V}} \operatorname{div}(\mathbf{F}(\mathbf{x})) d^3\mathbf{x} .$$

That is,

$$\int_{\S} \mathbf{F} \cdot \mathbf{N} dS = \int_{\S_2} \mathbf{F} \cdot \mathbf{N} dS - \int_{\mathcal{V}} \operatorname{div}(\mathbf{F}(\mathbf{x})) d^3\mathbf{x} . \quad (9.26)$$

It turn out that both of the integrals on the right are very easy to compute. First, since \S_2 is simply the unit disk in the x, y plane, the downward unit normal on \S_2 is simply $-\mathbf{e}_3 = (0, 0, -1)$. Also, $\mathbf{F}(x, y, 0) = (-y^2, -2xy, x^2y)$, and the area element in the x, y plane is simply $dS = dx dy$. Thus,

$$\mathbf{F} \cdot \mathbf{N} dS = -x^2 y dx dy .$$

Since the integrand is odd under reflection in y , and the region of integration is even, it is then clear that

$$\int_{\S_2} \mathbf{F} \cdot \mathbf{N} dS = 0 .$$

Next, we compute

$$\operatorname{div}(\mathbf{F}(\mathbf{x})) = 2yz - 2x .$$

Notice how much simpler this is than \mathbf{F} itself! To compute the volume integral, we use cylindrical coordinates. The limits of integration are given by

$$0 \leq r \leq 1 \quad 0 \leq \theta \leq 2\pi \quad \text{and} \quad 0 \leq z \leq 1 - r^2 .$$

Thus,

$$\int_{\mathcal{V}} \operatorname{div}(\mathbf{F}(\mathbf{x})) d^3\mathbf{x} = \int_0^1 \left(\int_0^{1-r^2} \left(\int_0^{2\pi} [2r \sin \theta z - 2r \cos \theta] d\theta \right) dz \right) r dr = 0 ,$$

since $\int_0^{2\pi} \sin \theta d\theta = \int_0^{2\pi} \cos \theta d\theta = 0$.

Going back to (9.26), we learn that $\int_{\mathcal{S}} \mathbf{F} \cdot \mathbf{N} dS = 0$, with much less computation than in Example 141.

9.4 Line integrals and circulation

9.4.1 Line integrals, force fields and work

Let \mathbf{F} be a continuous vector field on \mathbb{R}^n . In this section, we think of \mathbf{F} as representing a force field; that is \mathbf{F} gives the force that acts on a point particle located at \mathbf{x} .

For instance, if some electric charges are distributed in \mathbb{R}^3 , they will produce an electric field $\mathbf{E}(x)$, and then any point particle with an electrical charge q will be acted upon by a force $\mathbf{F}(\mathbf{x}) = q\mathbf{E}(\mathbf{x})$.

Let $\mathbf{x}(t)$, $a \leq t \leq b$, be a differentiable parameterized curve in \mathbb{R}^n . Suppose we move the point particle along the path $\mathbf{x}(t)$. We ask: How much work is done *on* the point particle as it moves along the curve from $\mathbf{x}_0 := \mathbf{x}(a)$ to $\mathbf{x}_1 := \mathbf{x}(b)$?

Let $h > 0$ be a small time step. As the particle moves from $\mathbf{x}(t)$ to $\mathbf{x}(t+h)$, the work $\Delta W(t)$ done is approximately given by the dot product of the displacement of the particle and the force acting time t :

$$\Delta W(t) \approx \mathbf{F}(\mathbf{x}(t)) \cdot (\mathbf{x}(t+h) - \mathbf{x}(t)) .$$

This is not exact since the force \mathbf{F} is not constant, but if the segment is very short, the variation in the force is a small percentage of the force itself. In this same small step limit, there is one more useful approximation to make:

$$\mathbf{F}(\mathbf{x}(t)) \cdot (\mathbf{x}(t+h) - \mathbf{x}(t)) = \mathbf{F}(\mathbf{x}(t)) \cdot \frac{\mathbf{x}(t+h) - \mathbf{x}(t)}{h} h \approx \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t) h .$$

Thus, if we divide the path into many such small segments, and then add up all of the contributions from all of the segments, and take the limit as the length of the segments tends to zero, we obtain an integral giving the exact value of the work that gets done: This is

$$\int_a^b \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t) dt . \tag{9.27}$$

Example 143 (Computing a line integral). Let $\mathbf{F}(x, y, z) = (z, x, y)$. Let $\mathbf{x}(t) = (\cos t, \sin t, t)$ for $0 \leq t \leq 2\pi$. Let us compute (9.27) in this case:

$$\mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t) = (t, \cos t, \sin t) \cdot (-\sin t, \cos t, 1) = -t \sin t + \cos^2(t) + \sin(t) .$$

Thus the total work is

$$\int_0^{2\pi} [-t \sin t + \cos^2(t) + \sin(t)] dt = 3\pi .$$

For computational purposes, it is best to represent the line integral in terms of some explicit parameterization of the curve. But the work done on the particle as it moves along the curve C , with a specified direction of travel, has a well defined meaning that is independent of any particular parameterization.

To write the line integral in such a way, introduce the unit tangent vector $\mathbf{T}(t)$ along the curve that points in the specified direction of motion. At each point along a differentiable curve, there are two unit vectors tangent to the curve. If $\mathbf{x}(t)$ is any parameterization of the curve, these are

$$\pm \frac{1}{\|\mathbf{x}'(t)\|} \mathbf{x}'(t) .$$

Taking whichever choice agrees with the specified direction of motion, we have

$$\mathbf{T}(t) = \pm \frac{1}{\|\mathbf{x}'(t)\|} \mathbf{x}'(t) \quad \text{and} \quad ds = \|\mathbf{x}'(t)\| dt .$$

Thus,

$$\mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{T}(t) ds = \pm \mathbf{F}(\mathbf{x}(t)) \cdot \mathbf{x}'(t) dt = . \quad (9.28)$$

This gives us the geometric form of the line integral of \mathbf{F} along the curve C :

$$\int_C \mathbf{F} \cdot \mathbf{T} ds .$$

We do not need to be given a parameterization of the curve C , we only need to be given the curve and, what is crucial, the direction in which the curve is traversed. If we are given the force field \mathbf{F} and are told that the path of the particle is the part of the parabola in the x, y plane given by $y = 1 - x^2$, $y \geq 0$, this is *not* enough to determine the work integral. We must also be given the direction of motion along the parabola. At each point along the parabola, there are two unit vectors that are tangent to the parabola. The direction of motion singles one of them out to be used as \mathbf{T} . *This specification of the direction of motion is called the orientation of the curve.*

Example 144 (Computing another line integral). Let $\mathbf{F}(x, y, z) = (z, x, y)$ be a given force field, and suppose a particle moves from $(1, 0, 0)$ to $(-1, 0, 0)$ along the parabola $y = 1 - x^2$ in the x, y plane. How much work is done on the particle?

First, we parameterize the path as $\mathbf{x}(t) = (t, 1 - t^2, 0)$ with $-1 \leq t \leq 1$. This traces out the parabola in question, but does so backwards, starting at $(-1, 0, 0)$ and ending at $(1, 0, 0)$. Thus, the correct unit tangent vector is the opposite of the one associated to this parameterization. Therefore, we choose the $-$ sign in (9.28) and have

$$\int_C \mathbf{F} \cdot \mathbf{T} ds = - \int_{-1}^1 (0, t, 1 - t^2) \cdot (1, -2t, 0) dt = - \int_{-1}^1 2t^2 dt = -\frac{4}{3} .$$

9.4.2 Conservative vector fields

There is a particularly nice kind of line integral: One in which the vector field $\mathbf{F}(\mathbf{x})$ is the gradient of some function $\varphi(\mathbf{x})$. Indeed, by the chain rule of Chapter 3,

$$\frac{d}{dt}\varphi(\mathbf{x}(t)) = \nabla(\mathbf{x}(t)) \cdot \mathbf{x}'(t) .$$

Therefore, if C is the path running along $\mathbf{x}(t)$ for, say, $a \leq t \leq b$, the fundamental Theorem of Calculus gives us

$$\varphi(\mathbf{x}(b)) - \varphi(\mathbf{x}(a)) = \int_a^b \nabla\varphi(\mathbf{x}(t)) \cdot \mathbf{x}'(t)dt = \int_C \nabla\varphi \cdot \mathbf{T}ds .$$

Notice that the left hand side depends only on the initial and final points along the curve C . Therefore,

$$\int_C \nabla\varphi \cdot \mathbf{T}ds$$

only depends on the curve C through its starting points and endpoint.

Definition 99 (Conservative vector field). *Let \mathbf{F} be a differentiable vector field defined on an open set $U \subset \mathbb{R}^n$. Then \mathbf{F} is a conservative vector field in case whenever C_1 and C_2 are any two piecewise differentiable curves with the same initial point and the same final point, both of which stay inside the set U where \mathbf{F} is defined,*

$$\int_{C_1} \mathbf{F} \cdot \mathbf{T}ds = \int_{C_2} \mathbf{F} \cdot \mathbf{T}ds .$$

What we have seen just above gives us one class of conservative vector fields – gradients vector fields: The value of the line integral of $\nabla\varphi$ along C is given by the difference in values of φ at the endpoints. *In particular, the line integral is zero whenever C is an oriented closed curve.* In fact, this is true for *any* conservative vector field:

Theorem 92 (Closed curves and conservative vector fields). *Let \mathbf{F} be a continuous vector field defined on $U \subset \mathbb{R}^3$. Then \mathbf{F} is conservative if and only if for all be any closed oriented piecewise differentiable curves C in U ,*

$$\oint_C \mathbf{F} \cdot \mathbf{T}ds = 0 . \tag{9.29}$$

Note: In (9.29), we have used the special integral symbol to emphasize that we are integrating over a closed curve.

Proof of Theorem 92: Suppose first that \mathbf{F} is conservative. Let C be any closed oriented piecewise differentiable curve C in U . Pick two distinct points \mathbf{x}_0 and \mathbf{x}_1 on C . Define C_1 to be the curve obtained by following C from \mathbf{x}_0 to \mathbf{x}_1 , following the given orientation. Let C_2 be the curve obtained by continuing onwards from \mathbf{x}_1 back to \mathbf{x}_0 , still following the given orientation. Let $-C_2$ denote the reversal of the curve C_2 : This is the curve on the same path, but with the orientation reversed.

Then C_1 and $-C_2$ are two curves in U running from \mathbf{x}_0 to \mathbf{x}_1 . Since \mathbf{F} is a conservative vector field,

$$\int_{C_1} \mathbf{F} \cdot \mathbf{T}ds = \int_{-C_2} \mathbf{F} \cdot \mathbf{T}ds .$$

Since changing the orientation changes the sign of \mathbf{T} ,

$$\int_{-C_2} \mathbf{F} \cdot \mathbf{T} ds = - \int_{C_2} \mathbf{F} \cdot \mathbf{T} ds .$$

But by the additivity property of integrals,

$$\int_C \mathbf{F} \cdot \mathbf{T} ds = \int_{C_1} \mathbf{F} \cdot \mathbf{T} ds + \int_{C_2} \mathbf{F} \cdot \mathbf{T} ds .$$

Combining the last three identities, we have

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = 0 .$$

The argument reverses: Now suppose that \mathbf{F} is such that (9.29) is true for any closed oriented piecewise differentiable curve C in U . If C_1 and C_2 are two piecewise differentiable curves in U running from \mathbf{x}_0 to \mathbf{x}_1 , define the simple closed curve C by following C_1 from \mathbf{x}_0 to \mathbf{x}_1 , then return from \mathbf{x}_1 to \mathbf{x}_0 along $-C_2$. Then

$$\begin{aligned} \oint_C \mathbf{F} \cdot \mathbf{T} ds &= \int_{C_1} \mathbf{F} \cdot \mathbf{T} ds + \int_{-C_2} \mathbf{F} \cdot \mathbf{T} ds \\ &= \int_{C_1} \mathbf{F} \cdot \mathbf{T} ds - \int_{C_2} \mathbf{F} \cdot \mathbf{T} ds = 0 , \end{aligned}$$

Thus, \mathbf{F} is conservative. □

Theorem 93 (Potential functions). *Let U be a pathwise connected open set in \mathbb{R}^3 . Then a continuous vector field \mathbf{F} defined on U is conservative if and only if there is continuously differentiable function φ such that $\mathbf{F}(\mathbf{x}) = \nabla\varphi(\mathbf{x})$ for all \mathbf{x} in U .*

Proof: We have already seen that all gradient vector fields are conservative. Now suppose that \mathbf{F} is some conservative vector field. Pick any point $\mathbf{x}_0 \in U$ and then for any $\mathbf{x} \in U$, define

$$\varphi(\mathbf{x}) = \int_{C_{\mathbf{x}_0, \mathbf{x}}} \mathbf{F} \cdot \mathbf{T} ds$$

where $C_{\mathbf{x}_0, \mathbf{x}}$ is any piecewise differentiable curve starting at \mathbf{x}_0 and ending at \mathbf{x} . This is a well defined function since \mathbf{F} is conservative.

Now fix any \mathbf{x} in U . Since U is open, $\mathbf{x} + h\mathbf{e}_1 \in U$ for all sufficiently small values of $|h|$. Pick any piecewise differentiable curve $C_{\mathbf{x}_0, \mathbf{x}}$ starting at \mathbf{x}_0 and ending at \mathbf{x} . Let $C_{\mathbf{x}_0, \mathbf{x} + h\mathbf{e}_1}$ be the curve obtained by continuing $C_{\mathbf{x}_0, \mathbf{x}}$ by moving along the straight line segment from \mathbf{x} to $\mathbf{x} + h\mathbf{e}_1$. Then

$$\varphi(\mathbf{x} + h\mathbf{e}_1) - \varphi(\mathbf{x}) = \int_{C_{\mathbf{x}_0, \mathbf{x} + h\mathbf{e}_1}} \mathbf{F} \cdot \mathbf{T} ds - \int_{C_{\mathbf{x}_0, \mathbf{x}}} \mathbf{F} \cdot \mathbf{T} ds = \int_0^1 \mathbf{F}(\mathbf{x} + t h \mathbf{e}_1) \cdot h \mathbf{e}_1 dt .$$

Therefore,

$$\nabla\varphi(\mathbf{x}) \cdot \mathbf{e}_1 = \lim_{h \rightarrow 0} \frac{1}{h} (\varphi(\mathbf{x} + h\mathbf{e}_1) - \varphi(\mathbf{x})) = \mathbf{F}(\mathbf{x}) \cdot \mathbf{e}_1 .$$

The same argument may be repeated for the other entries, and we obtain the result. □

9.4.3 Circulation

Theorem 92 motivates the the following definition:

Definition 100. Let C be a closed oriented piecewise differentiable curve in \mathbb{R}^3 . Let \mathbf{F} be a continuous vector field defined on a neighborhood of the curve C . Then the circulation of \mathbf{F} around C is the quantity

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds .$$

We may restate Theorem 92 by saying a continuous vector field \mathbf{F} defined in $U \subset \mathbb{R}^3$ is conservative if and only if for every closed oriented piecewise differentiable curve C in U , the circulation of \mathbf{F} around C is zero.

The computation of circulation for gradient vector fields is trivial: By what we have seen above, the circulation is always zero.

A circulation integral is therefore nothing other than a special case of a line integral – it is the case in which the curve C is a simple closed curve. When it is given parametrically, we have a beginning point and an end point which are the same, and letting the parameter increase specifies a direction of motion. However, if the curve is specified in purely geometric terms, say, as the circle of unit radius centered on $(0, 0, 1)$ in the plane $x + y + z = 1$, then we need additional information to specify the orientation.

Circulation is very easy to compute for a conservative vector field.

Example 145 (Computation of a circulation integral). Let C be the circle of unit radius centered on $(0, 0, 1)$ in the plane $x + y + z = 1$ oriented so that the direction of motion is counter-clockwise when viewed from above. Let $\mathbf{F} = (xy, 1, xy)$. Let us compute the circulation $\oint_C \mathbf{F} \cdot \mathbf{T} ds$.

First, we need to parameterize the circle. The normal vector to the plane is $(1, 1, 1)$. If $\{\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3\}$ is an orthonormal basis of \mathbb{R}^3 with \mathbf{u}_3 parallel to $(1, 1, 1)$, then $\{\mathbf{u}_1, \mathbf{u}_2\}$ is an orthonormal basis for the plane in question, and then

$$\mathbf{x}(t) = (0, 0, 1) + \cos t \mathbf{u}_1 + \sin t \mathbf{u}_2 \quad 0 \leq t \leq 2\pi ,$$

is a parameterization of the circle. To find the basis explicitly, note that $(1, -1, 0)$ is orthogonal to $(1, 1, 1)$. We therefore take

$$\mathbf{u}_3 = \frac{1}{\sqrt{3}}(1, 1, 1) \quad , \quad \mathbf{u}_1 = \frac{1}{\sqrt{2}}(1, -1, 0) \quad \text{and} \quad \mathbf{u}_2 := \mathbf{u}_3 \times \mathbf{u}_1 = \frac{1}{\sqrt{6}}(1, 1, -2) .$$

Thus,

$$\mathbf{x}(t) = (2^{-1/2} \cos t + 6^{-1/2} \sin t , -2^{-1/2} \cos t + 6^{-1/2} \sin t , 1 - 2^{1/2} 3^{-1/2} \sin t) .$$

Differentiating, we find

$$\mathbf{x}'(t) = (-2^{-1/2} \sin t + 6^{-1/2} \cos t , 2^{-1/2} \sin t + 6^{-1/2} \cos t , -2^{1/2} 3^{-1/2} \cos t) .$$

In particular, $\mathbf{x}(0) = (2^{-1/2}, -2^{-1/2}, 1)$ and $\mathbf{x}'(0) = (6^{-1/2}, 6^{-1/2}, 0)$

At $t = 0$, the x -coordinate is positive and increasing, and the y coordinate is negative. This means that when viewed from above, we see counter-clockwise motion. Thus, with this parameterization,

$$\mathbf{T}ds = +(-2^{-1/2}\sin t + 6^{-1/2}\cos t, 2^{-1/2}\sin t + 6^{-1/2}\cos t, -2^{1/2}3^{-1/2}\cos t)dt ;$$

i.e., we choose the $+$ sign in (9.28). We then have

$$\oint_C \mathbf{F} \cdot \mathbf{T}ds = \int_0^{2\pi} \left(\frac{1}{6}\sin^2 t - \frac{1}{2}\cos^2 t - 1 \right) \left(-\frac{1}{\sqrt{2}}\sin t - \frac{1}{\sqrt{6}}\cos t \right) dt = 0 .$$

9.4.4 The curl of a vector field on \mathbb{R}^3

Let $\mathbf{F} = (f_1, f_2, f_3)$ and $\mathbf{G} = (g_1, g_2, g_3)$ be two differentiable vector fields defined on an open set $U \subset \mathbb{R}^3$. We can build a new vector field out of \mathbf{F} and \mathbf{G} by taking their cross product, giving us the vector field

$$\mathbf{F} \times \mathbf{G} = (f_2g_3 - f_3g_2, f_3g_1 - f_1g_3, f_1g_2 - f_2g_1) .$$

Let us compute the divergence of $\mathbf{F} \times \mathbf{G}$. We find:

$$\begin{aligned} \operatorname{div}(\mathbf{F} \times \mathbf{G}) &= \left(\frac{\partial f_3}{\partial y} - \frac{\partial f_2}{\partial z} \right) g_1 + \left(\frac{\partial f_1}{\partial z} - \frac{\partial f_3}{\partial x} \right) g_2 + \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) g_3 \\ &\quad - \left(\frac{\partial g_3}{\partial y} - \frac{\partial g_2}{\partial z} \right) f_1 - \left(\frac{\partial g_1}{\partial z} - \frac{\partial g_3}{\partial x} \right) f_2 - \left(\frac{\partial g_2}{\partial x} - \frac{\partial g_1}{\partial y} \right) f_3 . \end{aligned} \quad (9.30)$$

We have grouped the terms so that we can write this in terms of dot products. We now make the following definition:

Definition 101 (The curl of a vector field on \mathbb{R}^3). *Let $\mathbf{F} = (f_1, f_2, f_3)$ be a differentiable vector field defined on an open set $U \subset \mathbb{R}^3$. Then the curl of \mathbf{F} , $\operatorname{curl}(\mathbf{F})$ is the vector field on U defined by*

$$\operatorname{curl}(\mathbf{F}) = \left(\left(\frac{\partial f_3}{\partial y} - \frac{\partial f_2}{\partial z} \right), \left(\frac{\partial f_1}{\partial z} - \frac{\partial f_3}{\partial x} \right), \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) \right) . \quad (9.31)$$

Example 146 (Computing a curl). *Let $\mathbf{F} = (xy, 1, xy)$ is in Example 145 . Then we find*

$$\operatorname{curl}(\mathbf{F}) = (x, -y, -x) .$$

Going back to our computation (9.30), we immediately deduce:

Theorem 94 (Curl and the divergence of a cross product). *Let $\mathbf{F} = (f_1, f_2, f_3)$ and $\mathbf{G} = (g_1, g_2, g_3)$ be two differentiable vector fields defined on an open set $U \subset \mathbb{R}^3$. Then*

$$\operatorname{div}(\mathbf{F} \times \mathbf{G}) = \operatorname{curl}(\mathbf{F}) \cdot \mathbf{G} - \operatorname{curl}(\mathbf{G}) \cdot \mathbf{F} .$$

This is the first of several important identities relating gradients, divergences and curls. Here is another:

Theorem 95 (The curl a gradient is zero). *Let φ be a twice differentiable function on an open set $U \subset \mathbb{R}^3$ so that $\nabla\varphi$ is a vector field on U . Then*

$$\operatorname{curl}(\nabla\varphi(\mathbf{x})) = \mathbf{0}$$

for all \mathbf{x} in U .

Proof: From the definition (9.31),

$$\operatorname{curl}(\nabla\varphi) = \left(\left(\frac{\partial^2\varphi}{\partial y\partial z} - \frac{\partial^2\varphi}{\partial z\partial y} \right), \left(\frac{\partial^2\varphi}{\partial z\partial x} - \frac{\partial^2\varphi}{\partial x\partial z} \right), \left(\frac{\partial^2\varphi}{\partial x\partial y} - \frac{\partial^2\varphi}{\partial y\partial x} \right) \right),$$

and each entry on the right is zero by Clairault's Theorem. \square

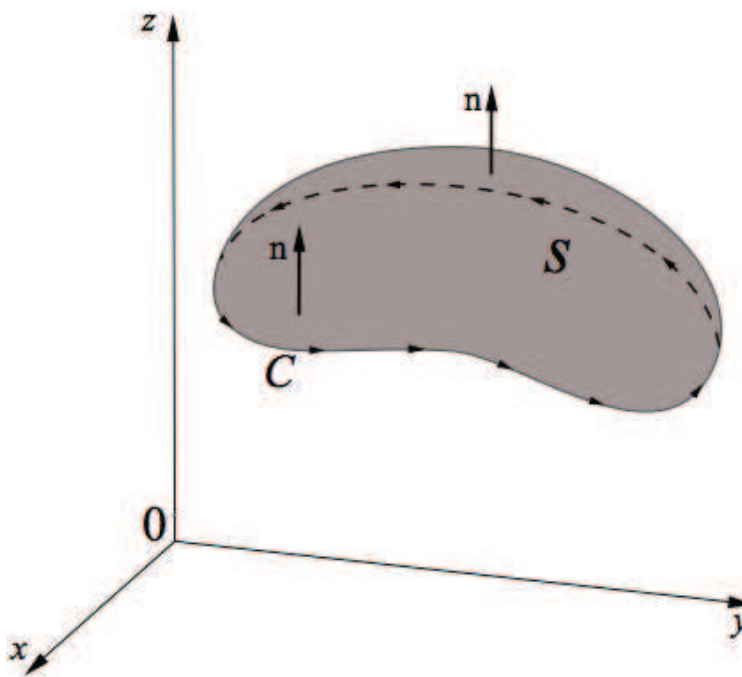
9.4.5 Stokes' Theorem

The curl of a vector field can be thought of as a *circulation density*, giving the circulation per unit area, in much the same way that the divergence can be thought of as a *flux density*, giving the flux per unit volume. The theorem that is the basis of this statement is Stokes' Theorem, which we state next:

Theorem 96 (Stokes' Theorem). *Let \mathcal{S} be an differentiable oriented surface in \mathbb{R}^3 , with unit normal \mathbf{N} , and suppose that \mathcal{S} is bounded by a differentiable simple closed curve C . Orient C so that the unit tangent vector \mathbf{T} has the property that at any point on the boundary, $\mathbf{T} \times \mathbf{N}$ points outward from the surface. Then*

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = \int_{\mathcal{S}} \operatorname{curl} \mathbf{F} \cdot \mathbf{N} dS.$$

Here is a picture showing how the orientations of the surface and its boundary “match up”.



Example 147 (Verification of Stokes' Theorem in an example). *Let \mathbf{F} and C be the vector field and curve from Example 145. That is, $\mathbf{F} = (xy, 1, xy)$, and C is the circle of unit radius centered on $(0, 0, 1)$ in the plane $x + y + z = 1$ oriented so that the direction of motion is counter-clockwise when viewed from above. We have already computed that*

$$\int_C \mathbf{F} \cdot \mathbf{T} ds = 0.$$

We now take \S to be the disk in the plane $x + y + z = 1$ that is bounded by C . To orient \S consistently with C , as in Stokes' Theorem, we must take \mathbf{N} to be the upward unit normal on \S . Thus, at each point of \S ,

$$\mathbf{N} = \frac{1}{\sqrt{3}}(1, 1, 1) .$$

In Example 146, we have already computed that $\text{curl}(\mathbf{F}) = (x, -y, -x)$. Thus,

$$\text{curl}(\mathbf{F}) \cdot \mathbf{N} = -y ,$$

and so

$$\int_{\S} \text{curl}(\mathbf{F}) \cdot \mathbf{N} dS = - \int_{\S} y dS .$$

Since \S is symmetric under the transformation $y \mapsto -y$, it is clear that $\int_{\S} y dS = 0$. Thus,

$$\int_{\S} \text{curl}(\mathbf{F}) \cdot \mathbf{N} dS = 0 ,$$

which based on our results in Example 145 is consistent with Stokes' Theorem.

Example 148 (Computing circulation using Stokes' Theorem). Let C be the contour that runs from $(1, 0, 0)$ to $(0, 1, 0)$, and from there to $(0, 0, 1)$, and from there back to $(1, 0, 0)$. Let $\mathbf{G} = (y + z^2, x + z^2, 2x + 2y)$. Compute the total circulation

$$\oint_C \mathbf{G} \cdot \mathbf{T} ds .$$

When asked to compute a circulation, or more generally, a work integral, unless the answer is totally obvious, the first step is to compute the curl of the vector field. We find:

$$\text{curl}(\mathbf{G}) = (2 - 2z, 2z - 2, 0) .$$

This is pretty simple, so it will be good to use Stokes' Theorem, which says,

$$\oint_C \mathbf{G} \cdot \mathbf{T} ds = \int_S \text{curl}(\mathbf{G}) \cdot \mathbf{N} dS ,$$

where S is the triangle with the specified vertices.

The triangle S lies in the plane given by $x + y + z = 1$, and for this plane the unit normal is

$$\mathbf{N} = \pm \frac{1}{\sqrt{3}}(1, 1, 1) .$$

Therefore, $\text{curl}(\mathbf{G}) \cdot \mathbf{N} = 0$, and so

$$\oint_C \mathbf{G} \cdot \mathbf{T} ds = 0 .$$

Stokes' Theorem may be applied to compute that change in the value of a line integral $\int_C \mathbf{F} \cdot \mathbf{T} ds$ that is induced by a change in the curve C . To see how to do this, let C_1 and C_2 be two differentiable curves running from \mathbf{x}_0 to \mathbf{x}_1 . Let $\mathbf{x}_1(t)$ and $\mathbf{x}_2(t)$, both for $0 \leq t \leq 1$ be parameterization of C_1 and C_2 respectively.

Define a parameterized surface

$$\mathbf{x}(s, t) \quad 0 \leq s, t \leq 1$$

by

$$\mathbf{x}(s, t) := (1 - s)\mathbf{x}_1(t) + s\mathbf{x}_2(t) .$$

This stretches out a “sheet” between C_1 and C_2 . Now let C denote $C_1 - C_2$; that is, the curve obtained by following C_1 from \mathbf{x}_0 to \mathbf{x}_1 , and then following C_2 backwards from \mathbf{x}_1 to \mathbf{x}_0 . Then C is the boundary of \S , and we orient \S consistently with C . (If C_1 and C_2 intersect, \S will consist of several pieces, each of which should be oriented separately. We will explain this further in examples; perhaps for now it is best to think of C_1 and C_2 as non-intersecting.)

Now suppose \mathbf{F} is a continuously differentiable vector field that is defined everywhere on a neighborhood of \S . Then by Stokes’ Theorem,

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = \int_{\S} \text{curl}(\mathbf{F}) \cdot \mathbf{N} dS .$$

However,

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = \int_{C_1} \mathbf{F} \cdot \mathbf{T} ds - \int_{C_2} \mathbf{F} \cdot \mathbf{T} ds ,$$

and so

$$\int_{C_1} \mathbf{F} \cdot \mathbf{T} ds = \int_{C_2} \mathbf{F} \cdot \mathbf{T} ds + \int_{\S} \text{curl}(\mathbf{F}) \cdot \mathbf{N} dS .$$

In particular, if $\text{curl}(\mathbf{F}) = \mathbf{0}$ everywhere, at least everywhere on \S , we have

$$\int_{C_1} \mathbf{F} \cdot \mathbf{T} ds = \int_{C_2} \mathbf{F} \cdot \mathbf{T} ds .$$

Thus, any continuously differentiable vector field \mathbf{F} that is defined on all of \mathbb{R}^3 and satisfies $\text{curl}(\mathbf{F}) = \mathbf{0}$ everywhere on \mathbb{R}^3 is a conservative vector field.

If the vector field is not defined on all of \mathbb{R}^3 , then this need not be the case.

Example 149 (Zero curl, but not conservative). *Consider open set*

$$U := \{(x, y, z) : x^2 + y^2 > 0\} .$$

That is, U is \mathbb{R}^3 with the z -axis removed. Consider the vector field \mathbf{F} defined on given by

$$\mathbf{F}(x, y, z) = \frac{1}{x^2 + y^2} (-y, x, 0) ,$$

which is well-defined everywhere on U . Then direct calculation yields

$$\text{curl}(\mathbf{F})(x, y, z) = \mathbf{0} .$$

However, if C is the unit circle in the x, y plane, oriented to run counter-clockwise as usual, then

$$\mathbf{x}(t) = (\cos t, \sin t, 0) , \quad 0 \leq t \leq 2\pi ,$$

is a parameterization of C , and we compute

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = \int_0^{2\pi} (-\sin t, \cos t, 0) \cdot (-\sin t, \cos t, 0) dt = 2\pi .$$

Thus, \mathbf{F} is not conservative. The problem is that one cannot find any surface in U of which C is the boundary: Any such surface must cross the z -axis somewhere, and \mathbf{F} is not defined on the z -axis.

The previous example brings us to a definition:

Definition 102 (Simply connected region). *An open set $U \subset \mathbb{R}^3$ is simply connected in case whenever C is a simple closed curve in U , there is an oriented surface \S in U of which C is the boundary. An open set $O \subset \mathbb{R}^2$ is simply connected if the cylinder $U := \{(x, y, z) : (x, y) \in O\}$ is simply connected in \mathbb{R}^3 .*

The reason for the second line in the definition is that, as seen in the last example, we can consider any vector field $\mathbf{F}(x, y)$ defined on an open set $O \subset \mathbb{R}^2$ as a vector field $\tilde{\mathbf{F}}(x, y, z)$ defined on the cylinder $U := \{(x, y, z) : (x, y) \in O\}$ by the simple device

$$\tilde{\mathbf{F}}(x, y, z) := \mathbf{F}(x, y) .$$

In this way, we may apply Stokes' Theorem to vector fields on \mathbb{R}^2 . We summarize our results in a Theorem:

Theorem 97 (Zero curl and conservation). *Let U be a simply connected open set in \mathbb{R}^3 . Let \mathbf{F} be a continuously differentiable vector field defined on U . Then \mathbf{F} is conservative if and only if $\text{curl}(\mathbf{F}) = \mathbf{0}$ everywhere on U . Likewise, let O be a simply connected open set in \mathbb{R}^2 . Let $\mathbf{F} = (f, g)$ be a continuously differentiable vector field defined on O . Then \mathbf{F} is conservative if and only if*

$$\frac{\partial}{\partial x}g(x, y) - \frac{\partial}{\partial y}f(x, y) = 0$$

everywhere on O .

Proof: By the remarks made above, it suffices to note that

$$\text{curl}((f, g, 0)) = \left(0, 0, \frac{\partial}{\partial x}g(x, y) - \frac{\partial}{\partial y}f(x, y)\right) .$$

□

We now have the means to determine whether a vector field, defined in a simply connected open set U is the gradient of some potential function: Compute the curl. In case the curl is zero, so that the vector field is the gradient of some potential function, we can even use the method of proof of Theorem 93 to compute such a potential function. (The potential function φ is only defined up to an additive constant: Adding a constant to φ does not change its gradient, and if φ and ψ are two potential functions for \mathbf{F} , $\nabla(\varphi - \psi) = \mathbf{0}$, so φ and ψ differ by a constant.)

Example 150 (Finding a potential function). *Consider the two vector fields*

$$\mathbf{F} = (y + z^2, x + z^2, 2zx + 2zy) \quad \text{and} \quad \mathbf{G} = (y + z^2, x + z^2, 2x + 2y) .$$

One of the vector fields \mathbf{F} and \mathbf{G} is equal to $\nabla\varphi$ for some potential function φ . Which one is it? Find such a potential function.

To do this, we compute

$$\text{curl}(\mathbf{F}) = \mathbf{0} \quad \text{and} \quad \text{curl}(\mathbf{G}) = (2 - 2z, 2z - 2, 0) .$$

A vector field on \mathbb{R}^3 is a gradient if and only if its curl is zero at every point in \mathbb{R}^3 . Hence \mathbf{F} is the gradient of some potential function φ .

To find φ , we compute line integrals. Pick $\mathbf{x}_0 = \mathbf{0}$ as our base point. Then for any point along the z axis, we find

$$\varphi(0, 0, z) = \int_0^1 \mathbf{F}(0, 0, tz) \cdot \mathbf{e}_3 z dt = 0 .$$

We next compute

$$\varphi(0, y, z) = \varphi(0, 0, z) = \int_0^1 \mathbf{F}(0, ty, z) \cdot \mathbf{e}_2 y dt = z^2 y .$$

We finally compute

$$\varphi(x, y, z) = \varphi(0, y, z) = \int_0^1 \mathbf{F}(tx, y, z) \cdot \mathbf{e}_1 x dt = (y + z^2)x .$$

Altogether,

$$\varphi(x, y, z) = xy + z^2(x + y) ,$$

and you can now easily verify that $\nabla\varphi = \mathbf{F}$.

9.4.6 Proof of Stokes Theorem

The key to proving Stoke's Theorem in general it to prove it when C is an oriented triangle in \mathbb{R}^3 . Let \mathbf{p}_1 , \mathbf{p}_2 and \mathbf{p}_3 be the non-colinear points in \mathbb{R}^3 , so that they are the vertices of a non-degenerate triangle in \mathbb{R}^3 . Let C be the oriented curve that traverses the boundary of the triangle starting at \mathbf{p}_1 , going next to \mathbf{p}_2 , then on to \mathbf{p}_3 , and finally returning to \mathbf{p}_1 . Let $\mathbf{x}(t)$, $0 \leq t \leq 1$ be a parameterization of C that is consistent with the orientation. In particular $\mathbf{x}(0) = \mathbf{p}_1$.

Now let us suppose that \mathbf{p}_1 , \mathbf{p}_2 and \mathbf{p}_3 are all very close together, so that the triangle is very small. We will eventually be concerned with what happens in the limit as these side-lengths go to zero.

When the distances are very small, the linear approximation

$$\mathbf{F}(\mathbf{x}(t)) \approx \mathbf{F}(\mathbf{x}(0)) + [J_{\mathbf{F}}(\mathbf{x}_0)](\mathbf{x}(t) - \mathbf{x}(0))$$

will be a good approximation, with the errors vanishing percentage-wise in the limit as the side-lengths go to zero.

Lemma 22. *Using the notation established above,*

$$\oint_C \mathbf{F}(\mathbf{x}(t)) \cdot d\mathbf{x}(t) \approx \alpha \text{curl}(\mathbf{F}(\mathbf{x}(0))) \cdot \mathbf{N}$$

where α is the area of the triangle, and where \mathbf{N} is its unit normal consistent with the specified orientation. The errors in this approximation go to zero as a percentage of the right hand side as the maximum side length of the triangle goes to zero, so that this approximation becomes exact in this limit.

Proof: Using the linear approximation $\mathbf{F}(\mathbf{x}) = \mathbf{x}(0) + [J_{\mathbf{F}}(\mathbf{x}_0)](\mathbf{x}(t) - \mathbf{x}(0))$,

$$\oint_C \mathbf{F}(\mathbf{x}(t)) \cdot d\mathbf{x}(t) \approx \oint_C \mathbf{F}(\mathbf{x}(0)) \cdot d\mathbf{x}(t) + \oint_C [J_{\mathbf{F}}(\mathbf{x}_0)](\mathbf{x}(t) - \mathbf{x}(0)) \cdot d\mathbf{x}(t) .$$

Since $\mathbf{F}(\mathbf{x}(0))$ is independent of t , and since $\mathbf{x}(1) = \mathbf{x}(0) = \mathbf{p}_1$,

$$\oint_C \mathbf{F}(\mathbf{x}(0)) \cdot d\mathbf{x}(t) = \mathbf{F}(\mathbf{x}(0)) \cdot \oint_C d\mathbf{x}(t) = \mathbf{F}(\mathbf{x}(0)) \cdot \int_0^1 \mathbf{x}'(t) dt = 0$$

Next, define $\mathbf{z}(t) = \mathbf{x}(t) - \mathbf{x}(0)$. Then

$$\oint_C [J_{\mathbf{F}}(\mathbf{x}_0)](\mathbf{x}(t) - \mathbf{x}(0)) \cdot d\mathbf{x}(t) = \int_0^1 ([J_{\mathbf{F}}(\mathbf{x}_0)]\mathbf{z}(t)) \cdot \mathbf{z}'(t) dt .$$

Now define matrices A and B by

$$A := \frac{1}{2} ([J_{\mathbf{F}}(\mathbf{x}_0)] - [J_{\mathbf{F}}(\mathbf{x}_0)]^T) \quad \text{and} \quad B := \frac{1}{2} ([J_{\mathbf{F}}(\mathbf{x}_0)] + [J_{\mathbf{F}}(\mathbf{x}_0)]^T) . \quad (9.32)$$

Notice that

$$A = -A^T, \quad B = B^T \quad \text{and} \quad [J_{\mathbf{F}}(\mathbf{x}_0)] = A + B .$$

The matrix A is called the *antisymmetric part of* $[J_{\mathbf{F}}(\mathbf{x}_0)]$, and the matrix B is called the *symmetric part of* $[J_{\mathbf{F}}(\mathbf{x}_0)]$.

Since B is symmetric,

$$\frac{d}{dt} \mathbf{z}(t) \cdot B\mathbf{z}(t) = 2(B\mathbf{z}(t)) \cdot \mathbf{z}'(t) .$$

Therefore, since $\mathbf{z}(0) = \mathbf{z}(1)$,

$$\int_0^1 (B\mathbf{z}(t)) \cdot \mathbf{z}'(t) dt = 0 .$$

Thus the symmetric part of $[J_{\mathbf{F}}(\mathbf{x}_0)]$ plays no role in our circulation computation, and we have that

$$\oint_C [J_{\mathbf{F}}(\mathbf{x}_0)](\mathbf{x}(t) - \mathbf{x}(0)) \cdot d\mathbf{x}(t) = \int_0^1 (A\mathbf{z}(t)) \cdot \mathbf{z}'(t) dt .$$

Since A is antisymmetric, it has the form

$$A = \begin{bmatrix} 0 & -c & b \\ c & 0 & -a \\ -b & a & 0 \end{bmatrix} . \quad (9.33)$$

Define the vector $\mathbf{a} := (a, b, c)$. Then as we have seen in Example 69,

$$A\mathbf{z} = \mathbf{b} \times \mathbf{z} .$$

Therefore, by the triple product identity,

$$(A\mathbf{z}(t)) \cdot \mathbf{z}'(t) = (\mathbf{a} \times \mathbf{z}(t)) \cdot \mathbf{z}'(t) = \mathbf{a} \cdot (\mathbf{z}(t) \times \mathbf{z}'(t)) .$$

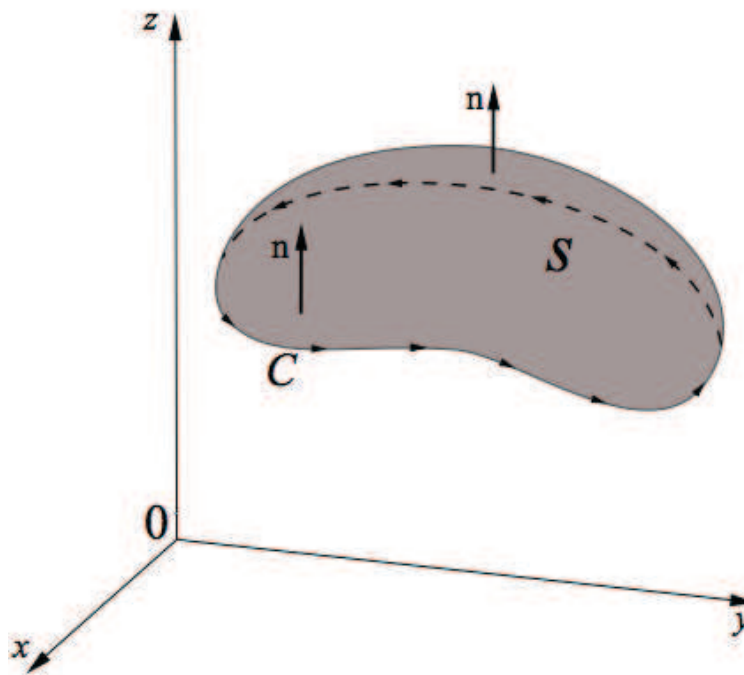
Therefore,

$$\oint_C [J_{\mathbf{F}}(\mathbf{x}_0)](\mathbf{x}(t) - \mathbf{x}(0)) \cdot d\mathbf{x}(t) = \mathbf{a} \cdot \left(\int_0^1 \mathbf{z}(t) \times \mathbf{z}'(t) dt \right) .$$

From familiar calculations, we recognize $\int_0^1 \mathbf{z}(t) \times \mathbf{z}'(t) dt$ as $\alpha \mathbf{N}$ where α is the area of our triangle, and \mathbf{N} is the unit normal pointing to the positive side according to the orientation induced by C .

The final step in our computation is to recognize \mathbf{a} as $\text{curl}(\mathbf{F}(\mathbf{x}(0)))$: This follows directly from (9.32) and (9.33). \square

Form here, the proof of Stoke's Theorem is easy. Consider any nice surface such as the one shown in



“Chop” the surface up into small triangular tiles. Each edge of any of these triangles in the interior of \S is traversed twice, because it is part of the boundary of two triangular tiles. But it is traversed in opposite directions, so that all of the contributions to the circulation from the interior triangles is zero: Adding up the circulation around of all of the triangular tiles, everything except the contribution coming from the boundary of \S cancels out.

Thus we have that the circulation about C , the boundary of \S is the sum of the circulations about each of the triangular tiles. Taking the limit as the maximum side length of these triangles goes to zero, and using Lemma 22,

$$\begin{aligned} \oint_C \mathbf{F} \cdot d\mathbf{x} &= \lim_{\text{side length to } 0} \sum_{\text{triangular tiles in } \S} (\text{circulation about tile}) \\ &= \lim_{\text{side length to } 0} \sum_{\text{triangular tiles in } \S} (\text{area of tile}) \times (\mathbf{N} \text{ in tile}) \times (\text{curl}(\mathbf{F}) \text{ in tile}) \\ &= \int_{\S} \text{curl}(\mathbf{F}(\mathbf{x})) \cdot \mathbf{N}(\mathbf{x}) dS . \end{aligned}$$

This completes the proof of Stoke's Theorem. \square .

9.4.7 Vector Potentials

Lemma 23 (The divergence of a curl is zero). *Let \mathbf{A} be a twice continuously differentiable vector field. Then*

$$\operatorname{div}(\operatorname{curl} \mathbf{A}) = 0 .$$

Proof: Write $\mathbf{A}(\mathbf{x}) = (f(\mathbf{x}), g(\mathbf{x}), h(\mathbf{x}))$. Then since

$$\begin{aligned} \operatorname{curl}(\mathbf{A}) &= \left(\left(\frac{\partial h}{\partial y} - \frac{\partial g}{\partial z} \right), \left(\frac{\partial f}{\partial z} - \frac{\partial h}{\partial x} \right), \left(\frac{\partial g}{\partial x} - \frac{\partial f}{\partial y} \right) \right) , \\ \operatorname{div}(\operatorname{curl} \mathbf{A}) &= \frac{\partial}{\partial x} \left(\frac{\partial h}{\partial y} - \frac{\partial g}{\partial z} \right) + \frac{\partial}{\partial y} \left(\frac{\partial f}{\partial z} - \frac{\partial h}{\partial x} \right) + \frac{\partial}{\partial z} \left(\frac{\partial g}{\partial x} - \frac{\partial f}{\partial y} \right) \\ &= \left(\frac{\partial^2 f}{\partial y \partial z} - \frac{\partial^2 f}{\partial z \partial y} \right) + \left(\frac{\partial^2 g}{\partial z \partial x} - \frac{\partial^2 g}{\partial x \partial z} \right) + \left(\frac{\partial^2 h}{\partial x \partial y} - \frac{\partial^2 h}{\partial y \partial x} \right) , \end{aligned}$$

and each of the last three terms are zero by Clairault's Theorem. \square

Lemma 23 gives us a necessary condition for a vector field \mathbf{F} to be the curl of some other vector field \mathbf{A} : It must be the case that $\operatorname{div}(\mathbf{F}) = 0$.

This condition turns out to be necessary as well. To see why this is true, let us consider a continuously differentiable vector field $\mathbf{F}(\mathbf{x}) = (P(\mathbf{x}), Q(\mathbf{x}), R(\mathbf{x}))$ defined on all of \mathbb{R}^3 such that $\operatorname{div}(\mathbf{F}(\mathbf{x})) = 0$ for all $\mathbf{x} \in \mathbb{R}^3$.

Let us first observe that If $\mathbf{F} = \operatorname{curl}(\mathbf{A})$ for some other vector field \mathbf{A} , then \mathbf{A} is far from unique: Since $\operatorname{curl}(\nabla \varphi) = \mathbf{0}$, and since the operation of taking a curl is linear,

$$\operatorname{curl}(\mathbf{A} + \nabla \varphi) = \operatorname{curl}(\mathbf{A}) + \operatorname{curl}(\nabla \varphi) = \mathbf{F} + \mathbf{0} = \mathbf{F} .$$

Hence one might hope that among the possible choices of \mathbf{A} , there are some that are particularly simple. This turns out to be the case:

Consider a vector field \mathbf{A} of the form

$$\mathbf{A}(\mathbf{x}) = (f(\mathbf{x}), 0, h(\mathbf{x}))$$

for twice continuously differentiable real valued functions f and h on \mathbb{R}^3 . Then by the formula for $\operatorname{curl}(\mathbf{A})$,

$$\operatorname{curl}(\mathbf{A}) = \left(\frac{\partial h}{\partial y}, \frac{\partial f}{\partial z} - \frac{\partial h}{\partial x}, -\frac{\partial f}{\partial y} \right) . \quad (9.34)$$

Therefore, if $\operatorname{curl}(\mathbf{A}) = \mathbf{F} = (P, Q, R)$, then by the Fundamental Theorem of Calculus, we must have

$$\begin{aligned} f(x, y, z) &= - \int_0^y R(x, t, z) dt + \alpha(x, z) \\ h(x, y, z) &= \int_0^y P(x, t, z) dt + \beta(x, z) \end{aligned} \quad (9.35)$$

for some functions $\alpha(x, z)$ and $\beta(x, z)$, since this is equivalent to

$$-\frac{\partial f}{\partial y} = R \quad \text{and} \quad \frac{\partial h}{\partial y} = P . \quad (9.36)$$

Then with f and h defined by (9.35), we compute the middle component of $\text{curl}(\mathbf{A})$:

$$\frac{\partial f}{\partial z} - \frac{\partial h}{\partial x} = - \int_0^y \left[\frac{\partial R}{\partial z}(x, t, z) + \frac{\partial P}{\partial x}(x, t, z) \right] dt + \frac{\partial \alpha}{\partial z}(x, z) - \frac{\partial \beta}{\partial x}(x, z) . \quad (9.37)$$

However, since $\text{div}(\mathbf{F}) = 0$,

$$\left[\frac{\partial R}{\partial z}(x, t, z) + \frac{\partial P}{\partial x}(x, t, z) \right] = - \frac{\partial Q}{\partial y}(x, t, z) .$$

Using this and the Fundamental Theorem of Calculus, (9.37) becomes

$$\frac{\partial f}{\partial z} - \frac{\partial h}{\partial x} = Q(x, y, z) - Q(x, 0, z) + \frac{\partial \alpha}{\partial z}(x, z) - \frac{\partial \beta}{\partial x}(x, z) . \quad (9.38)$$

We must choose $\alpha(x, z)$ and $\beta(x, z)$ so that the right hand side reduces to $Q(x, y, z)$, since then by (9.34) and (9.36) we will have $\text{curl}(\mathbf{A}) = (P, Q, R) = \mathbf{F}$. Note that if we chose

$$\alpha(x, z) = \int_0^z Q(x, 0, t) dt \quad \text{and} \quad \beta(x, z) = 0 ,$$

we do indeed obtain

$$\frac{\partial f}{\partial z} - \frac{\partial g}{\partial x} = Q(x, y, z) ,$$

and hence $\text{curl}(\mathbf{A}) = (P, Q, R) = \mathbf{F}$. Thus, whenever $\text{div}(\mathbf{F}) = 0$ on all of \mathbb{R}^3 , there is a vector field \mathbf{A} so that $\mathbf{F} = \text{curl}(\mathbf{A})$ everywhere on \mathbb{R}^3 . such a vector field \mathbf{A} is called a *vector potential* for \mathbf{F} .

We have proved:

Theorem 98 (Vector potentials). *Let \mathbf{F} be a continuously differentiable vector field on \mathbb{R}^3 such that $\text{div}(\mathbf{F}(\mathbf{x})) = 0$ for all $\mathbf{x} \in \mathbb{R}^3$. Then there is a continuously differentiable vector field $\mathbf{A}(\mathbf{x})$ on \mathbb{R}^3 such that $\text{curl}(\mathbf{A}(\mathbf{x})) = \mathbf{F}(\mathbf{x})$ for all \mathbf{x} in \mathbb{R}^3 . If $\mathbf{F} = (P, Q, R)$, then one such vector potential \mathbf{A} is given by $\mathbf{A} = (f, 0, h)$ where*

$$\begin{aligned} f(x, y, z) &= - \int_0^y R(x, t, z) dt + \int_0^z Q(x, 0, t) dt \\ h(x, y, z) &= \int_0^y P(x, t, z) dt . \end{aligned} \quad (9.39)$$

Example 151 (Computing a vector potential). *Let $\mathbf{F} = (-y(2+x), x, yz)$. We readily compute*

$$\text{div}(\mathbf{F}(x, y, z)) = -y + 0 + y = 0 .$$

Hence, \mathbf{F} has a vector potential, and the recipe in Theorem 98 provides one. Since $P(x, t, z) = -t(2+x)$, we have

$$h(x, y, z) = -(2+x) \int_0^y t dt = -\frac{1}{2}(2+x)y^2 .$$

Next, since $Q(x, 0, t) = x$ and $R(x, t, z) = tz$,

$$f(x, y, z) = -\frac{1}{2}zy^2 + xz .$$

Altogether,

$$\mathbf{A}(x, y, z) = \left(-\frac{1}{2}zy^2 + xz , 0 , -\frac{1}{2}(2+x)y^2 \right) .$$

As one readily checks, we do indeed have $\text{curl}(\mathbf{A}) = \mathbf{F}$.

9.5 Exercises

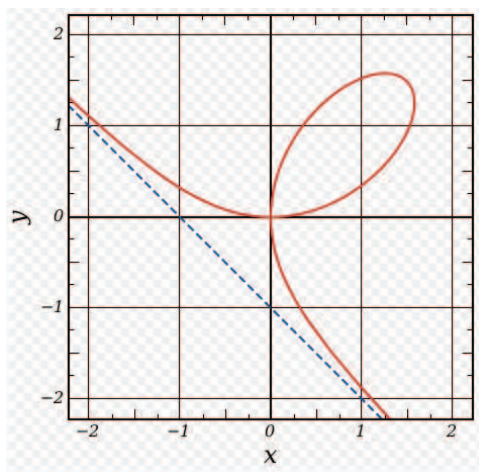
1. (a) Let $\mathbf{F}(\mathbf{x})$ be a vector field on \mathbb{R}^n of the form $\mathbf{F}(\mathbf{x}) = -\nabla V(\mathbf{x})$ for some twice continuously differentiable function V on \mathbb{R}^n . Show that along any flow curve $\mathbf{x}(t)$ of \mathbf{F} , $V(\mathbf{x}(t))$ is a non-increasing function of t .

(b) Let $\mathbf{F}(\mathbf{x})$ be a vector field on \mathbb{R}^n of the form $\mathbf{F}(\mathbf{x}) = A\nabla V(\mathbf{x})$ for some twice continuously differentiable function V on \mathbb{R}^n and some antisymmetric $n \times n$ matrix A . Show that along any flow curve $\mathbf{x}(t)$ of \mathbf{F} , $V(\mathbf{x}(t))$ is a constant function of t .

2. Let $\mathbf{F}(\mathbf{x})$ be the vector field on \mathbb{R}^2 given by $\mathbf{F}(x, y) = (y, x)$. Proceeding as in Example 130, find an explicit formula for the flow curve through the general point \mathbf{x}_0 in \mathbb{R}^2 . Also, for each t , find an explicit formula for the flow transformation $\Phi_t(\mathbf{x})$. Your answers will involve the hyperbolic trigonometric functions.

3. Let C be the path consisting of straight line segments running from $(0, 0)$ to $(3, 3)$ and from there to $(4, 5)$, and from there to $(0, 7)$. Let $\mathbf{F}(x, y) = (\sin(x) + y, 3x + y)$. Compute the flux integral $\int_C \mathbf{F} \cdot \mathbf{N} ds$ and the circulation integral $\int_C \mathbf{F} \cdot \mathbf{T} ds$, using the orientation induced by the parameterization of C .

4. The curve in the plane given by the equation $x^3 + y^3 = 3xy$ is known as the folium of Descartes. Here is a plot of the part of the curve that we shall consider in this exercise:



(Folium is Latin for leaf, as in the English word foliage.) Consider the parameterized curve $\mathbf{x}(t)$ given by

$$\mathbf{x}(t) = \left(\frac{3t}{1+t^3}, \frac{3t^2}{1+t^3} \right)$$

for $0 \leq t < \infty$.

(a) Show that each point on $\mathbf{x}(t)$ lies on the folium of Descartes, and that

$$(0, 0) = \mathbf{x}(0) = \lim_{t \rightarrow \infty} \mathbf{x}(t),$$

so that the parameterized curve $\mathbf{x}(t)$ is a closed loop. (It is in fact the “leaf” of the folium, which lies in the upper right quadrant.)

(b) Show that

$$\mathbf{x}(t) \cdot d\mathbf{x}^\perp(t) = \frac{9t^2}{(1+t^3)^2} dt ,$$

and compute the area enclosed by the curve. (That is, compute the area of the leaf.)

(c) Let $\mathbf{F}(x, y) = (x^2 - y^2, 2xy)$. Compute the flux integral $\oint_C \mathbf{F} \cdot \mathbf{N} ds$ where C is the leaf in the folium of Descartes, and \mathbf{N} is its outward unit normal.

5. Let C be path along the unit circle that is in the upper right quadrant in \mathbb{R}^2 , starting at $(1, 0)$ and ending at $(0, 1)$. Let \mathbf{F} be a vector field of the form

$$\mathbf{F}(x, y) = \mathbf{G}(x, y) + \nabla\varphi(x, y)$$

where

$$\varphi(1, 0) = 1 \quad \text{and} \quad \varphi(0, 1) = 2 .$$

Suppose also that with $\mathbf{G}(x, y) = (P(x, y), Q(x, y))$,

$$\frac{\partial Q}{\partial x} - \frac{\partial P}{\partial y} = 1$$

for all x and y , and finally, $P(x, 0) = 0$ for all x and $Q(0, y) = 0$ for all y .

Using this information, compute $\int_C \mathbf{F} \cdot \mathbf{T} ds$.

6. Verify Stokes' Theorem by calculating both sides of

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = \int_S \text{curl}(\mathbf{F}) \cdot \mathbf{N} dS$$

where

$$\mathbf{F}(x, y, z) = (y^2, x, z^2) ,$$

S is the part of the paraboloid $z = x^2 + y^2$ lying below the plane $z = 1$ with \mathbf{N} being the upward unit normal to S , and C the boundary of S with the orientation that is consistent with the choice of \mathbf{N} .

7. Consider the following vector fields:

$$\begin{aligned} \mathbf{F}(x, y, z) &= (x, xy + z, y^2) \\ \mathbf{G}(x, y, z) &= (1, xy - z, y^2 - xz) \\ \mathbf{H}(x, y, z) &= (x, x^2 + y + z, y^2) \end{aligned}$$

Which of these vector fields are curls, and which are not? That is, does there exist a vector field $\mathbf{A}(x, y, z)$ such that $\text{curl}(\mathbf{A}) = \mathbf{F}$, and likewise for \mathbf{G} and \mathbf{H} ? Justify your answers.

8. Consider the following vector fields:

$$\begin{aligned} \mathbf{F}(x, y, z) &= (y^3 + z^2y, 3y^2x + z^2x, 2xyz) \\ \mathbf{G}(x, y, z) &= (xyz, xy, z) \\ \mathbf{H}(x, y, z) &= (yz, x, x) \end{aligned}$$

Which of these vector fields are gradients, and which are not? That is, does there exist a potential function $\varphi(x, y, z)$ such that $\nabla\varphi = \mathbf{F}$, and likewise for \mathbf{G} and \mathbf{H} ? Justify your answers, and in case a potential function exists, explicitly find one.

9. Let C be any simple closed curve in \mathbb{R}^3 that lies on the surface of the cone $z = \sqrt{x^2 + y^2}$. Let $\mathbf{F}(x, y, z) = (x^2, y^2, z(x^2 + y^2))$. Show that

$$\oint_C \mathbf{F} \cdot \mathbf{T} ds = 0 .$$

10. Let \mathcal{V} be the rectangular box $[0, 1] \times [2, 4] \times [1, 5]$. Let $\mathbf{F}(x, y, z) = (xy, yz, x^2z + z^2)$. Let \mathcal{S} be the boundary of \mathcal{V} . Compute the outward flux

$$\int_{\mathcal{S}} \mathbf{F} \cdot \mathbf{N} dS .$$

11: Let \mathcal{S} be the part of the paraboloid $z = 1 - x^2 - y^2$ that lies above the plane $x + z = 1$. Let \mathbf{F} be the vector field $\mathbf{F}(x, y, z) = (xy, yz, zx)$. Compute the flux integral

$$\int_{\mathcal{S}} \mathbf{F} \cdot \mathbf{N} dS$$

where \mathbf{N} is the downward unit normal to the surface. That is, compute the flux across the surface from top to bottom.

12: Let \mathcal{S} be the part of the ellipsoid $4x^2 + 9y^2 + z^2 = 36$ that lies above the plane $z = 3$. Let \mathbf{F} be the vector field $\mathbf{F}(x, y, z) = (x, 0, z)$. Compute the flux integral

$$\int_{\mathcal{S}} \mathbf{F} \cdot \mathbf{N} dS$$

where \mathbf{N} is the downward unit normal to the surface. That is, compute the flux across the surface from top to bottom.

13: Let \mathcal{S} be the boundary of the region \mathcal{V} that is above the sphere $x^2 + y^2 + z^2 = 6$ and below the paraboloid $z = 4 - x^2 - y^2$. Let $\mathbf{F}(x, y, z)$ be the vector field $\mathbf{F}(x, y, z) = (z, y, x)$. Compute the flux integral $\int_{\mathcal{S}} \mathbf{F} \cdot \mathbf{N} dS$ for the flux em out of the region D .

14: Let \mathcal{V} be the region in \mathbb{R}^3 that is inside ellipsoid $4x^2 + 9y^2 + z^2 = 36$, and above the plane $z = 3$. Let \mathcal{S} be the boundary of \mathcal{V} . Let \mathbf{F} be the vector field $\mathbf{F}(x, y, z) = (x, x, z)$.

(a) Compute the flux integral $\int_{\mathcal{S}} \mathbf{F} \cdot \mathbf{N} dS$ where \mathbf{N} is the outward unit normal to the surface.

(b) Let C be the curve at which the plane $z = 3$ intersects the ellipsoid $4x^2 + 9y^2 + z^2 = 36$, oriented to run counterclockwise when viewed from above. Compute $\int_C \mathbf{F} \cdot \mathbf{T} ds$.

15: Consider the two vector fields

$$\mathbf{F}(x, y, z) := (2x(y - z) + 2, x^2 - 2yz, -x^2 - y^2 - 3)$$

and

$$\mathbf{G}(x, y, z) := (2x(y + z) + 2, x^2 - 2yz, x^2 + y^2 + 3) .$$

(a) Compute $\text{curl}(\mathbf{F})$ and $\text{curl}(\mathbf{G})$.

(b) One of the vector fields is the gradient of some function $\varphi(x, y, z)$, and the other is not. Which one is, and how do you know? For the one that is the gradient of some function φ , find such a function φ .

16: Consider the two vector fields

$$\mathbf{F} = (y + z^2, x + z^2, 2zx + 2zy) \quad \text{and} \quad \mathbf{G} = (y + z^2, x + z^2, 2x + 2y) .$$

(a) Compute the divergence and curl of \mathbf{F} and \mathbf{G} .

(b) Let S be the unit sphere, and \mathbf{N} its outward normal. Compute **either**

$$\int_S \mathbf{F} \cdot \mathbf{N} dS \quad \text{or} \quad \int_S \mathbf{G} \cdot \mathbf{N} dS .$$

The choice is yours. Do whichever one you find easier, and justify your answer to receive credit.

(c) One of the vector fields \mathbf{F} and \mathbf{G} is equal to $\nabla\varphi$ for some potential function φ . Which one is it? Find such a potential function.

(d) Let C be the curve that is given by

$$x^2 + y^2 + z^2 = 4 \quad \text{and} \quad x + y + z = 1 .$$

Orient C so that it is traversed in the counter-clockwise direction when viewed from above. Compute **either**

$$\int_C \mathbf{F} \cdot \mathbf{T} ds \quad \text{or} \quad \int_C \mathbf{G} \cdot \mathbf{T} ds .$$

The choice is yours. Do whichever one you find easier, and justify your answer to receive credit.

17: Let \mathcal{V} be the region in \mathbb{R}^3 that lies inside the sphere $x^2 + y^2 + z^2 = 4$, and above the graph of $z = 1/\sqrt{x^2 + y^2}$, as in problem 8. Let $\mathbf{F} = (y + z^2, x + z^2, 2z(x + y))$ and let \mathbf{N} be the outward normal to S , the boundary of \mathcal{V} . Compute the total flux

$$\int_S \mathbf{F} \cdot \mathbf{N} dS .$$

18: Let C be the contour that runs from $(1, 0, 0)$ to $(0, 1, 0)$, and from there to $(0, 0, 1)$, and from there back to $(1, 0, 0)$. Let $\mathbf{G} = (y + z^2, x + z^2, 2x + 2y)$. Compute the total circulation

$$\oint_C \mathbf{G} \cdot \mathbf{T} ds .$$

19: Consider the two vector fields

$$\mathbf{F} = (2xyz - y^2, x^2z - 2xy, x^2y) \quad \text{and} \quad \mathbf{G} = (2yz - y^2, x^2z - 2x, x^2y) .$$

(a) Compute the divergence and curl of \mathbf{F} and \mathbf{G} .

(b) Let S be the part of the paraboloid $z = 1 - x^2 - y^2$ that lies above the x, y plane. Compute

$$\int_S \mathbf{G} \cdot \mathbf{N} dS .$$

(c) One of the vector fields \mathbf{F} and \mathbf{G} is equal to $\nabla\varphi$ for some potential function φ . Which one is it? Find such a potential function.

(d) Let C be the curve that is parametrized by

$$\mathbf{x}(t) = (t^3 - 2t^2, t - 3t^2, t + t^3) \quad \text{for} \quad 0 \leq t \leq 1 .$$

Compute

$$\int_C \mathbf{F} \cdot \mathbf{T} ds .$$

20: Let S be the part of the surface in \mathbb{R}^3 given by $\sqrt{x^2 + y^2} = 8 - z$ that lies outside the cylinder $x^2 + y^2 = 4$. With $\mathbf{G} = (2yz - y^2, x^2z - 2x, x^2y)$, compute the flux

$$\int_S \mathbf{G} \cdot \mathbf{N} dS ,$$

where \mathbf{N} is taken to point outward from the z -axis.

21: Let C be the contour that runs from $(0, 0, 0)$ to $(0, 1, 2)$, and from there to $(2, 2, 2)$, and from there back to $(0, 0, 0)$. Let $\mathbf{G} = (z, x, y)$. Compute the total circulation

$$\oint_C \mathbf{G} \cdot \mathbf{T} ds .$$

22: Consider the two vector fields

$$\mathbf{F}(\mathbf{x}) := (-2y + z, 2x + 4yz, x - 2y^2) \quad \text{and} \quad \mathbf{G}(\mathbf{x}) := (2y + z, 2x - 4yz, x - 2y^2)$$

both defined everywhere on \mathbb{R}^3 .

(a) Compute $\text{curl}(\mathbf{F})$ and $\text{curl}(\mathbf{G})$

(b) One of \mathbf{F} and \mathbf{G} is a gradient vector field and the other is not. Which one is the gradient of some potential function $\varphi(\mathbf{x})$, and how do you know? For the one that is, find such a potential function $\varphi(\mathbf{x})$.

(c) For the vector field that is not a gradient vector field, compute its circulation around the unit circle in the x, y plane, given the usual counter-clockwise orientation.

23: Let \mathcal{V} be the region in \mathbb{R}^3 specified in problem 7. Let $\mathbf{F} = (y + z^2, x + z^2, 2z(x + y))$ and let \mathbf{N} be the outward normal to S , the boundary of \mathcal{V} . Compute the total flux

$$\int_S \mathbf{F} \cdot \mathbf{N} dS .$$

24: Let C be the contour that runs from $(1, 0, 0)$ to $(0, 1, 0)$, and from there to $(0, 0, 1)$, and from there back to $(1, 0, 0)$. Let $\mathbf{G} = (y + z^2, x + z^2, 2x + 2y)$. Compute the total circulation

$$\oint_C \mathbf{G} \cdot \mathbf{T} ds .$$

25: Let \mathbf{F} be the vector field

$$\mathbf{F} = (2xyz - y^2, x^2z - 2xy, x^2y) .$$

Let \S be the part of the paraboloid $z = 1 - x^2 - y^2$ that lies above the x, y plane, oriented so its preferred normal points upward. Compute the flux

$$\int_{\S} \mathbf{F} \cdot \mathbf{N} dS$$

26: Let C be the curve given by the intersection of the surfaces $z = x^2$ and $z = 4 - y^2$. Orient C to run counterclockwise when viewed from above. Let $\mathbf{F}(x, y, z) = (x, x, z)$. Compute

$$\int_C \mathbf{F} \cdot \mathbf{T} ds .$$

27: Let \mathcal{V} be the region in \mathbb{R}^3 that lies inside the sphere $x^2 + y^2 + z^2 = 4$, and above the graph of $z = 1/\sqrt{x^2 + y^2}$. Compute the total surface area of its boundary \mathcal{S} . (There are two pieces to the boundary.)

(a) Compute the surface area of \S .

(b) Let $\mathbf{F}(x, y, z) = (x, x, z)$. Compute $\int_{\S} \mathbf{F} \cdot \mathbf{N} dV$ where \mathbf{N} is the outward unit normal vector.

28: Let C be the contour that runs from $(0, 0, 0)$ to $(1, 0, 0)$, and from there to $(1, 0, 1)$, and from there to $(0, 0, 1)$. Let $\mathbf{F} = (x, x, z)$. Compute the line integral

$$\int_C \mathbf{F} \cdot \mathbf{T} ds .$$