

MATH 51 SECTION NOTES

EVAN WARNER

1. JANUARY 6

1.1. **Administrative miscellany.** My office hours will be 3:00 to 4:30 on Tuesdays and Wednesdays in my office, 384K (fourth floor of the math department). You can reach me by email at ebwarner@math.stanford.edu. These notes will be posted at math.stanford.edu/~ebwarner. The course website is web.stanford.edu/class/math51.

Homework assigned each Thursday, starting the first week of classes, and due each Thursday at 3:05, starting the second week of classes. You may turn it in to me in section or under my office door. Do not put it in my mailbox or email it unless you have made prior arrangement with me. If you have a grading issue, see me within one week after the homework is returned; after one week, the grade will be considered final.

1.2. Helpful course advice.

1.2.1. *Section attendance.* Nothing in your grade is based on attendance; likewise, neither is anything in my estimation of you as a person. However, you will likely find it helpful, especially with respect to performance on the exams. If you choose to attend a section other than the one in which you are enrolled, just note that enrollees have seating priority if there is insufficient space.

1.2.2. *Reading.* As far as I am concerned, “doing the reading” is an absolutely essential part of the course. Week-by-week topics are listed on the course website (under “Syllabus”), along with the corresponding sections of the text.

1.2.3. *Questions.* Ask them, please! One of the most useful methods of learning mathematics is to repeat the following steps: realize you don’t understand something, work to formulate a precise question about it (this is sometimes difficult!), ask it, and then think about the response. You are all welcome to do this as often as you like in section, in office hours, or via email. But even if you can’t formulate your question precisely (or it’s a meta-mathematical question, like “why is Concept X relevant”), ask anyway!

1.2.4. *Extra resources.* There are exams (midterms and finals) dating back more than a decade on the course webpage. We will do a lot of exam questions in class, and they are also very useful to study with! Keep in mind that exam questions are often not like (some of the) homework questions - they tend to be more conceptual, you do not have the benefit of knowing “which section in the text you are in” (i.e., on a homework problem from Section 4, you know that you should use tools from Section 4 of the text. On an exam, this is much less clear).

Also available are linear algebra notes written by my friend Joseph Victor, which can be found at sumo.stanford.edu/pdfs/linearalgebranotes.pdf. They are wordier and chattier than the course text, although the risk of typos is correspondingly greater.

1.2.5. *Grading advice.* Homework is only ten percent of the course grade, and the lowest homework grade is automatically dropped. Thus from a practical point of view, any single homework assignment (and even more so any single homework problem!) will have a negligible affect on your grade. Homework is a pedagogical tool, helping you to learn concepts and keeping you honest and up to date with the material.

The flip side to this is, of course, that exams are rather consequential.

1.3. **Vectors.** For the purposes of this class, a *vector* is an element of \mathbb{R}^n , which is mathematical shorthand for the set of ordered n -tuples of real numbers. At the beginning, we admit two algebraic operations on vectors: addition and scalar multiplication. Addition of vectors is performed componentwise, and scalar multiplication multiplies each component by the scalar in question. For example:

$$\begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} + 4 \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} + \begin{pmatrix} 4 \\ 4 \\ -4 \end{pmatrix} = \begin{pmatrix} 6 \\ 3 \\ -4 \end{pmatrix}.$$

In order to understand *why* we make these definitions, we must give a geometric interpretation. Here, the interpretation is simple: addition corresponds to “putting vectors end to end” (if you draw vectors as arrows in \mathbb{R}^n) and scalar multiplication corresponds to scaling a vector by a given real quantity.

For any purely algebraic construction we do with \mathbb{R}^n , we should always keep an eye on “what’s happening geometrically.” It is important to keep both perspectives in mind: the algebraic definitions will allow us to calculate and prove things, while the geometric intuition will help us understand what to calculate and what to prove, and whether our answers make any sense.

1.4. **Span.** Here is an example concept to test our intuition on. Algebraically, the span of a finite set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is given by

$$\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_k) = \{c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k \mid c_i \in \mathbb{R} \text{ for } 1 \leq i \leq k\}.$$

One can read this definition in English as follows: the span of a set of vectors \mathbf{v}_1 through \mathbf{v}_k consists of all vectors of the form $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k$, where each coefficient c_1, c_2, \dots, c_n is a real number. In a more abbreviated form, we say that the span of a set of vectors is the set of all *linear combinations* of those vectors; i.e., all the vectors we get if we allow ourselves to add vectors and multiply by scalars, starting with our original set.

The geometric content of this definition might be a little opaque at first, but the idea is something like “the span is the smallest linear space that contains the given vectors and passes through the origin.” For example, that linear space could be a line passing through the origin, or a plane passing through the origin, or something in more than two dimensions that acts like a line or a plane passing through the origin, or even just the origin by itself. It may, of course, be all of \mathbb{R}^n , if we pick enough vectors in enough different directions.

Here are some examples. If

$$U = \text{Span} \left(\begin{pmatrix} 3 \\ 4 \end{pmatrix}, \begin{pmatrix} 3 \\ 5 \end{pmatrix} \right),$$

then I claim that U is all of \mathbb{R}^2 . Why? Geometrically, pick one of the vectors, and notice that you can get a whole line just by multiplying that vector by various scalars (real numbers). Then pick the other vector and note that by adding a small multiple c of that vector to every point on the line, you can get another line, slightly offset from the first. As we vary c over all real numbers, we get lines that sweep out the entire plane. There is a very nice picture of this in section 2 of the textbook. If we wanted to prove rigorously that $U = \mathbb{R}^2$, it would be far easier to do it algebraically, by showing that for any vector $\mathbf{v} = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$ in \mathbb{R}^2 , there exist real numbers c_1 and c_2 such that

$$\begin{pmatrix} v_1 \\ v_2 \end{pmatrix} = c_1 \begin{pmatrix} 3 \\ 4 \end{pmatrix} + c_2 \begin{pmatrix} 3 \\ 5 \end{pmatrix}.$$

Proving this amounts to an exercise in solving a system of two linear equations; it turns out that we can pick

$$c_1 = \frac{5}{3}v_1 - v_2, \quad c_2 = -\frac{4}{3}v_1 + v_2$$

(verify this for yourself!).

As a second example, consider

$$V = \text{Span} \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} -2 \\ -2 \end{pmatrix} \right).$$

In this case, we see that both vectors sweep out the same line (the line $y = x$ in the plane), and there is no way to take a linear combination of them and get something that is not on the line $y = x$. There is some sort of redundancy going on, and the two vectors span only a line, rather than a plane.

As a third example, consider

$$W = \text{Span} \left(\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} \right).$$

Algebraically, it is fairly clear that we can via linear combinations of these three vectors get any vector of the form $(0, y, z)$, where y and z are any real numbers. It is also clear that we can never get anything except a zero in the first component, because all three original vectors have first component equal to zero. Geometrically, this span is a plane (the yz plane in \mathbb{R}^3). Like in the last example, there is some redundancy: the span would be unchanged if we omitted the last vector, for instance.

2. JANUARY 8

2.1. Warm-up questions.

- (1) Given k vectors in \mathbb{R}^n , what are the possible dimensions of the set $\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_k)$ if $k \leq n$?
- (2) What if instead $k > n$?
- (3) Do the answers change if we require the k vectors to be distinct?

2.2. Answers.

- (1) The maximum possible dimension of the span of k vectors is equal to k , and as $k \leq n$ this is achievable in \mathbb{R}^n . Every dimension less than k is also achievable; given any linear subspace of dimension d less than or equal to k , we may pick d vectors that span it and then pick the remaining $k - d$ vectors arbitrarily in the subspace (for example, they could be the zero vectors). This also applies to dimension 0, because we can pick all of the vectors to be the zero vector. Thus the answer is the set $\{0, 1, 2, \dots, k\}$.
- (2) Any subspace of \mathbb{R}^n has dimension at most n , so we cannot exceed dimension n . The solution to the first problem tells us that this is the only constraint, so the answer is the set $\{0, 1, 2, \dots, n\}$. In general, combining the two parts, the possible dimensions of a set of k vectors in \mathbb{R}^n are the elements of the set $\{0, 1, 2, \dots, \min\{k, n\}\}$.
- (3) The only possibility that this eliminates is the zero-dimensional subspace consisting of the origin alone, for the only way for the span of a set of vectors to be the origin is if all of them are identically zero. Therefore the answer is now the set $\{1, 2, 3, \dots, \min\{k, n\}\}$.

2.3. Linear independence. We saw last section that sometimes k vectors fail to span a k -dimensional space. There is an algebraic criterion for this phenomenon called *linear dependence*. It will turn out (and we will prove easily at some point in the course, once we have the right definitions) that k vectors fail to span a k -dimensional space if and only if they are linearly dependent.

By definition, a set of vectors is linearly dependent if one of the vectors can be written as a linear combination of the others (that is, one of the vectors can be written in terms of the others using scalar multiplication and vector addition). This algebraic definition has an immediate geometric interpretation: a set of vectors is linearly dependent if and only if one lies in the span of the others. In other words, there is a “redundancy” in the span of the set of vectors.

An alternative definition is as follows: a set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ is linearly dependent if whenever we have an equality

$$c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_n\mathbf{v}_n = \mathbf{0}$$

with scalars c_1, c_2, \dots, c_n , all of the coefficients must be equal to zero (i.e. $c_1 = c_2 = \dots = c_n$). This may seem like a more complicated condition than the above, but it’s actually algebraically simpler: we only have to check one equation and make sure that all the coefficients must be zero, rather than checking that each vector individually cannot be written as a linear combination of the others. Proposition 3.1 in the text shows that these two definitions are equivalent; the proof is simple algebra.

2.4. Exercises in linear independence. For each of the following examples, let’s determine whether the given set of vectors is independent or not. Furthermore, let’s do it as efficiently as possible, doing a minimum of actual calculation.

- (1) $\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$. These are *linearly independent* because any two (nonzero) vectors that are not collinear are linearly independent (that is, they span a plane). Recall the argument in the previous lecture.

(2) $\left\{ \begin{pmatrix} 2 \\ -3 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right\}$. These are *linearly dependent* because the zero vector is always a linear combination of any other vector (you can multiply any vector by the scalar 0 and get the zero vector).

(3) $\left\{ \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix} \right\}$. These are *linearly dependent* because there are too many vectors: any four vectors in \mathbb{R}^3 must be linearly dependent because if they were not they would span a four-dimensional subspace of \mathbb{R}^3 , which does not exist.

(4) $\left\{ \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 0 \\ 0 \end{pmatrix} \right\}$. These are *linearly independent* because, again, they are two nonzero vectors that are not collinear.

(5) $\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \\ 0 \\ 0 \end{pmatrix} \right\}$. These are *linearly dependent* because we can see immediately that the third vector is equal to two times the sum of the first two, and this provides a linear dependence relation.

(6) $\left\{ \begin{pmatrix} 1 \\ -1 \\ 3 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 5 \\ 0 \\ 5 \end{pmatrix} \right\}$. Now we actually have to do some work. We should try to check using the equivalent, but more algebraically tractable, condition given by Proposition 3.1: if we have an equation of the form

$$c_1 \begin{pmatrix} 1 \\ -1 \\ 3 \end{pmatrix} + c_2 \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} + c_3 \begin{pmatrix} 5 \\ 0 \\ 5 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix},$$

then we have (equivalently) the system of three equations $c_1 + 5c_3 = 0$, $-c_1 = 0$, $3c_1 + c_2 + 5c_3 = 5$. The second equation implies that $c_1 = 0$, which upon plugging into the first equation yields $c_3 = 0$, and plugging them both into the third equation gives $c_2 = 0$. Therefore if we have an equation of the above form, the coefficients must be all equal to zero, so the vectors are *linearly independent*.

The solution of the last example illustrates a general phenomenon: the problem of determining linear independence of vectors is equivalent to the problem of solving systems of linear equations (with the zero vector on the right hand side; these equations are called *homogeneous*). We will soon learn an algorithm, row reduction, for solving any system of linear equations; it will turn out to be merely a formalized version of the process of adding multiples of one equation to another in order to eliminate variables.

2.5. Dot products. Here is a new algebraic operation we can perform on vectors. The *dot product* takes two vectors and returns a scalar according to the following

rule:

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = x_1y_1 + x_2y_2 + \dots + x_ny_n.$$

That is, we take the sum of the componentwise products. As a simple example, we can calculate

$$\begin{pmatrix} 1 \\ 0 \\ 3 \end{pmatrix} \cdot \begin{pmatrix} -2 \\ 7 \\ 1 \end{pmatrix} = 1 \cdot (-2) + 0 \cdot 7 + 3 \cdot 1 = 1.$$

The geometric interpretation is a little bit more subtle. We'll see that the dot product encapsulates information about both the length of the two vectors and how "far apart" they are pointing. The most important use of the dot product is to figure out when two vectors are perpendicular: in fact, two vectors have a dot product equal to zero if and only if they are perpendicular (or one of the vectors is the zero vector).

It is perhaps easiest to see a geometric interpretation in the case that the two vectors are equal. Then we get the square of the length of the vector: in fact, we can *define* the length (or the magnitude) $\|\mathbf{v}\|$ of a vector \mathbf{v} to be equal to $\sqrt{\mathbf{v} \cdot \mathbf{v}}$ (here we take the positive square root). For example,

$$\left\| \begin{pmatrix} 3 \\ 4 \end{pmatrix} \right\| = \sqrt{3^2 + 4^2} = \sqrt{25} = 5.$$

This is a good definition because we expect the Pythagorean theorem (in any number of dimensions) to hold, and this ensures that it does: if a vector has components v_1, v_2, \dots, v_n , its length is

$$\sqrt{v_1^2 + v_2^2 + \dots + v_n^2},$$

which recovers the Pythagorean theorem (try sketching this in two dimensions!).

The dot product and the norm obey some useful algebraic identities; see the course text (section 4) for details. For example, it is easy to see that the dot product "acts like" regular multiplication, in that it is commutative, commutes with scalar multiplication, and distributes over addition. The norm of a vector is zero if and only if the vector is the zero vector, and it is otherwise a positive real number. We also have less trivial observations, such as the *Cauchy-Schwartz inequality* $|\mathbf{v} \cdot \mathbf{w}| \leq \|\mathbf{v}\| \cdot \|\mathbf{w}\|$ and the *triangle inequality* $\|\mathbf{v} + \mathbf{w}\| \leq \|\mathbf{v}\| + \|\mathbf{w}\|$ (exercise: interpret this inequality geometrically. Why is it called the triangle inequality?).

3. JANUARY 13

3.1. Warm-up questions. For the first three questions, is the statement true or false? Think about these questions intuitively (as we haven't defined precisely what a linear subspace or its dimension are yet): a line through the origin is a linear subspace of dimension one, a plane is a linear subspace of dimension two, the origin by itself is a linear subspace of dimension zero, and so on.

- (1) If W and V are linear subspaces of \mathbb{R}^n , then $\dim V + \dim W \leq n$.
- (2) If $\mathbf{v} \neq \mathbf{0}$, then $\dim[\text{Span}(\mathbf{v})] = 1$.

- (3) If $\{\mathbf{u}, \mathbf{v}, \mathbf{w}\}$ is a linearly independent set of vectors, then $\dim[\text{Span}(\mathbf{u}, \mathbf{v})] = 2$.
- (4) Prove that if \mathbf{v} and \mathbf{w} are linearly independent, then so are $\mathbf{v} + \mathbf{w}$ and $\mathbf{v} - \mathbf{w}$.

3.2. Answers.

- (1) This is *false*; for example, W and V could both be equal to the whole space \mathbb{R}^n , in which case $\dim W = \dim V = n$ and $\dim W + \dim V = 2n$.
- (2) This is *true*; any nonzero vector spans a line, which has dimension one.
- (3) This is *true*; any subset of a set of linearly independent vectors is linearly independent, so in particular the set $\{\mathbf{u}, \mathbf{v}\}$ is linearly independent, so its span has dimension two.
- (4) We want to show that

$$d_1(\mathbf{v} + \mathbf{w}) + d_2(\mathbf{v} - \mathbf{w}) = 0$$

implies $d_1 = d_2 = 0$. To do this, rearrange the equation so that it reads

$$(d_1 + d_2)\mathbf{v} + (d_1 - d_2)\mathbf{w} = 0.$$

Now note that because \mathbf{v} and \mathbf{w} are linearly independent by assumption, we know that any equation of the form $c_1\mathbf{v} + c_2\mathbf{w} = 0$ holds we must have $c_1 = c_2 = 0$. Applying this to the above, with $c_1 = d_1 + d_2$ and $c_2 = d_1 - d_2$, we conclude that $d_1 + d_2 = 0$ and $d_1 - d_2 = 0$. But this is a simple system of linear equations, which we can solve to find that $d_1 = d_2 = 0$ is the only solution. This proves the proposition.

3.3. Parametric representations. The *parametric representation* of a line L in \mathbb{R}^n is given by

$$L = \{\mathbf{x}_0 + t\mathbf{v} \mid t \in \mathbb{R}\},$$

with $\mathbf{v} \neq \mathbf{0}$. Here \mathbf{x}_0 is any point lying on the line and \mathbf{v} is a vector pointing in the direction of the line; as t varies over the real numbers, the expression $\mathbf{x}_0 + t\mathbf{v}$ sweeps out the entire locus of points comprising the line. We require $\mathbf{v} \neq \mathbf{0}$ because otherwise the locus would be just a point; every value of t would yield the same point \mathbf{x}_0 . Every line in \mathbb{R}^n can be written in this way in many different ways, just by picking an arbitrary point on the line and an arbitrary vector in the correct direction. The representation is called parametric because it depends on an auxiliary parameter; in this case, the variable t .

As an example, let's say we want to find a parametric representation of the line L in \mathbb{R}^4 passing through the two points

$$\begin{pmatrix} 0 \\ 0 \\ 0 \\ 5 \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 1 \\ -3 \\ 2 \\ 0 \end{pmatrix}.$$

This is easy: we can just pick, for example,

$$\mathbf{x}_0 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 5 \end{pmatrix},$$

and to get a vector in the direction of the line we can just subtract one point from another, getting

$$\mathbf{v} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 5 \end{pmatrix} - \begin{pmatrix} 1 \\ -3 \\ 2 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ 3 \\ -2 \\ 5 \end{pmatrix}.$$

Thus

$$L = \left\{ \begin{pmatrix} 0 \\ 0 \\ 0 \\ 5 \end{pmatrix} + t \begin{pmatrix} -1 \\ 3 \\ -2 \\ 5 \end{pmatrix} \mid t \in \mathbb{R} \right\}.$$

The parametric representation of a plane in \mathbb{R}^n is similar, but requires two independent parameters:

$$P = \{\mathbf{x}_0 + s\mathbf{v}_1 + t\mathbf{v}_2 \mid s, t \in \mathbb{R}\},$$

where \mathbf{v}_1 and \mathbf{v}_2 are linearly independent. We require linear independence because otherwise the set of points would sweep out a line (or a point, if both happened to be the zero vector), not a plane.

As an example, let us find a parametric representation for the plane P in \mathbb{R}^3 passing through the three points $\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$, $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$, and $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$. Again, pick one of the points (say the first) to be our \mathbf{x}_0 , and we can find a suitable \mathbf{v}_1 and \mathbf{v}_2 by subtracting the first point from the second and the first point from the third:

$$\mathbf{v}_1 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}, \quad \mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}.$$

Therefore, noting that \mathbf{v}_1 and \mathbf{v}_2 are clearly linearly independent, we have

$$P = \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + s \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} + t \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} \mid s, t \in \mathbb{R} \right\}.$$

Of course, many other possibilities for \mathbf{x}_0 , \mathbf{v}_1 , and \mathbf{v}_2 would have worked as well.

Higher-dimensional examples (of translated linear subspaces) work in exactly the same way.

3.4. Equations of lines and planes. There is another, sometimes more useful way of writing down the locus of a line or a plane (or higher dimensional linear objects) without using an auxiliary parameter. Let us first consider a line lying in \mathbb{R}^2 and passing through the origin. The key insight is that this line consists precisely of the points orthogonal to a given vector, which we call the *normal vector*. Therefore the equation of this line is $\mathbf{x} \cdot \mathbf{n} = 0$, where \mathbf{x} is the vector of variables and \mathbf{n} is fixed. If the line does not go through the origin, this is no longer true, but we can shift the line to the origin to get a similar equation: if \mathbf{x}_0 is an arbitrary point on the line, and \mathbf{n} is a normal vector, then the equation of the line is

$$(\mathbf{x} - \mathbf{x}_0) \cdot \mathbf{n} = 0.$$

Written in coordinates (if, for example, \mathbf{n} has coordinates n_1 and n_2), this is

$$n_1(x - x_0) + n_2(y - y_0) = 0,$$

which is basically the point-slope formula for a line in the plane.

So we haven't really learned anything new here. The usefulness comes from the fact that exactly the same idea works to give an equation for a plane in \mathbb{R}^3 (or a three-space in \mathbb{R}^4 , and so on and so on). In the example of a plane, we can always find an equation of the form

$$(\mathbf{x} - \mathbf{x}_0) \cdot \mathbf{n} = 0,$$

where again \mathbf{n} is the normal vector (and \mathbf{x}_0 is any point on the plane). Written in coordinates, this looks like

$$n_1(x - x_0) + n_2(y - y_0) + n_3(z - z_0) = 0.$$

As an example, say that a plane P in \mathbb{R}^3 passes through the point $\begin{pmatrix} 4 \\ 2 \\ -1 \end{pmatrix}$ with normal vector $\begin{pmatrix} 5 \\ 1 \\ 5 \end{pmatrix}$. Then the corresponding equation can be read off directly; it is

$$5(x - 4) + 1(y - 2) + 5(z + 1) = 0,$$

which we can simplify to

$$5x + y + 5z = -17.$$

If instead we were given three points and asked for the equation, we would first have to calculate a normal vector, which we could do by finding two vectors in the plane and taking a cross product.

Note that this procedure only works for linear spaces that have one fewer dimension than the ambient space. If we wanted, for example, to find an equation for a line in \mathbb{R}^3 , we would find it impossible: in fact, we would need two equations to cut out a line. We will not discuss finding such equations, but you should be able to figure out how to do so geometrically by, for example, intersecting planes.

3.5. An example exam problem. This would be a good point to review the basic properties of the dot product and norm that can be found in section 4 of the text.

Here's an actual exam problem from several years ago: given vectors \mathbf{u} and \mathbf{v} such that $\|\mathbf{u}\| = \|\mathbf{v}\|$, show that $\mathbf{u} - \mathbf{v}$ and $\mathbf{u} + \mathbf{v}$ are orthogonal.

The answer is surprisingly straightforward, once you realize that vectors are orthogonal if and only if their dot product is zero. We have

$$\begin{aligned} (\mathbf{u} - \mathbf{v}) \cdot (\mathbf{u} + \mathbf{v}) &= \mathbf{u} \cdot \mathbf{u} - \mathbf{v} \cdot \mathbf{u} + \mathbf{u} \cdot \mathbf{v} - \mathbf{v} \cdot \mathbf{v} \\ &= \|\mathbf{u}\|^2 - \|\mathbf{v}\|^2 \\ &= 0. \end{aligned}$$

This calculation proves the proposition.

4. JANUARY 15

4.1. Warm-up question. Suppose you know that $\mathbf{x} \cdot \mathbf{y} = 3$, $\|\mathbf{x}\| = 2$, and $\|\mathbf{y}\| = 3$. Find the cosine of the angle between $\mathbf{x} + \mathbf{y}$ and $\mathbf{x} - \mathbf{y}$.

4.2. **Answer.** By the equation $|\mathbf{u} \cdot \mathbf{v}| = \|\mathbf{u}\| \cdot \|\mathbf{v}\| \cos \theta$, where θ is the angle between \mathbf{u} and \mathbf{v} . Applying this to $\mathbf{x} + \mathbf{y}$ and $\mathbf{x} - \mathbf{y}$, we can calculate by expanding repeatedly and using that $\mathbf{u} \cdot \mathbf{u} = \|\mathbf{u}\|^2$:

$$\begin{aligned} \cos \theta &= \frac{(\mathbf{x} + \mathbf{y}) \cdot (\mathbf{x} - \mathbf{y})}{\|\mathbf{x} + \mathbf{y}\| \cdot \|\mathbf{x} - \mathbf{y}\|} \\ &= \frac{\mathbf{x} \cdot \mathbf{x} - \mathbf{y} \cdot \mathbf{y}}{\sqrt{\mathbf{x} \cdot \mathbf{x} + 2\mathbf{x} \cdot \mathbf{y} + \mathbf{y} \cdot \mathbf{y}} \sqrt{\mathbf{x} \cdot \mathbf{x} - 2\mathbf{x} \cdot \mathbf{y} + \mathbf{y} \cdot \mathbf{y}}} \\ &= \frac{\|\mathbf{x}\|^2 - \|\mathbf{y}\|^2}{\sqrt{\|\mathbf{x}\|^2 + 2\mathbf{x} \cdot \mathbf{y} + \|\mathbf{y}\|^2} \sqrt{\|\mathbf{x}\|^2 - 2\mathbf{x} \cdot \mathbf{y} + \|\mathbf{y}\|^2}} \\ &= \frac{4 - 9}{\sqrt{19}\sqrt{7}} \\ &= -\frac{5}{\sqrt{133}}. \end{aligned}$$

4.3. **Gaussian elimination and reduced row echelon form.** We will describe (most of) an algorithm for solving linear systems of equations. The first step is to put the coefficients of the system in a box of numbers (this is called a *matrix*, but for now all we are using it for is as a box of numbers. Later we'll see more useful things one can do with matrices).

For example, the system of equations

$$\begin{aligned} 3x - 4y &= 2, \\ x + y &= -1 \end{aligned}$$

becomes the matrix

$$\begin{pmatrix} 3 & -4 & 2 \\ 1 & 1 & -1 \end{pmatrix}.$$

The process of manipulating this matrix is called *Gaussian elimination*. We are allowed three moves, each of which is easily checked to preserve the solution set of the original system of equations (when interpreted in that way):

- (1) We can swap two rows.
- (2) We can multiply any row by a nonzero scalar.
- (3) We can add any multiple of a row to another row.

The end goal is to manipulate our matrix into something called *reduced row echelon form*. This form satisfies the following three properties:

- (1) If any rows contain only zeroes, they must be at the bottom of the matrix.
- (2) The leading coefficient of all nonzero rows (we will call this leading coefficient the *pivot*) must be strictly to the right of the pivot of the row above. In particular, this means that only zeroes may lie underneath pivots.
- (3) Every pivot has entry exactly equal to one and it is the only nonzero entry in its column.

If a matrix obeys the first two rules, we say it is in *row echelon form*.

The main theorem is then that the process of Gaussian elimination can always be used to reduce a matrix to reduced row echelon form, and such a form is unique (though the sequence of steps taken to get there may not be). We will see that the reduced row echelon form will allow us to read off a description of the space of solutions of the original system; in particular, we will be able to write down parametric representations of such solution spaces.

Going back to our example, we can row reduce the matrix in the following way (remember, this is far from the only way to do this!). First flip the two rows using step (1), so we have a one in the upper right corner:

$$\begin{pmatrix} 1 & 1 & -1 \\ 3 & -4 & 2 \end{pmatrix}.$$

Next add -3 times the first row to the second row using step (2):

$$\begin{pmatrix} 1 & 1 & -1 \\ 0 & -7 & 5 \end{pmatrix}.$$

We've dealt with the first column, so let's move on to the next one. Multiply by $-\frac{1}{7}$ using step (3), getting

$$\begin{pmatrix} 1 & 1 & -1 \\ 0 & 1 & -\frac{5}{7} \end{pmatrix}.$$

Finally, we need to isolate the pivot in the second column by subtracting the second row from the first, using step (3) again:

$$\begin{pmatrix} 1 & 0 & -\frac{2}{7} \\ 0 & 1 & -\frac{5}{7} \end{pmatrix}.$$

Now the matrix is in reduced row echelon form, and we can read off the solution: it is $x = -\frac{2}{7}$ and $y = -\frac{5}{7}$.

Here's another example, which you can try yourself:

$$\begin{aligned} x + 3y + z &= 9, \\ x + y - z &= 1, \\ 3x + 11y + 5z &= 35. \end{aligned}$$

The corresponding matrix is

$$\begin{pmatrix} 1 & 3 & 1 & 9 \\ 1 & 1 & -1 & 1 \\ 3 & 11 & 5 & 35 \end{pmatrix}.$$

Going through the same procedure, which you can do on your own time, we arrive at the matrix (in reduced row echelon form)

$$\begin{pmatrix} 1 & 0 & -2 & -3 \\ 0 & 1 & 1 & 4 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

This is not quite as simple as the last example, but there's a good reason! Here, there isn't a unique solution; reading off the system gives

$$x - 2z = -3, \quad y - z = 4, \quad 0 = 0.$$

Clearly, the last equation is unhelpful; all it's saying is that we can pick z arbitrarily. The other two equations imply that then x and y are determined, so the space of solutions is a line. Choosing t as an auxiliary parameter, we can write down a parametric representation of the line of solutions:

$$\left\{ \begin{pmatrix} -3 \\ 4 \\ 0 \end{pmatrix} + t \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} \mid t \in \mathbb{R} \right\}.$$

Finally, let's look at a third example. Consider the system

$$\begin{aligned}x + y &= 1, \\2x + 2y &= 3\end{aligned}$$

We see immediately that such a system can have no solutions, but if we run Gaussian elimination anyway, we arrive at the matrix

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

which yields the equations $x + y = 0$ and $0 = 1$. As this latter equation is never satisfied, there are no solutions.

In general, we can expect that the space of solutions can be empty, a point, a line, a plane, or any higher-dimensional linear object. If it is empty, there are no solutions, if it is a point, there is a unique solution, and if it is anything else there are infinitely many solutions.

4.4. Matrices act on vectors. So far matrices are just boxes with numbers in them, but they are in fact much more interesting than that. First, we define an algebraic operation: an $m \times n$ matrix acts on a vector in \mathbb{R}^n via the following operation: we take the dot product of the first *row* of the matrix with the vector and place it in the first component, we take the dot product of the second row of the matrix with the vector and place it in the second component, and so on. In this way, we get a vector in \mathbb{R}^m . In other words, we have an operation that takes as input an $m \times n$ matrix and a vector in \mathbb{R}^n and outputs a vector in \mathbb{R}^m . For example,

$$\begin{pmatrix} 2 & 1 \\ 0 & 2 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 2 - 1 \\ 0 - 2 \end{pmatrix} = \begin{pmatrix} 1 \\ -2 \end{pmatrix}$$

and

$$\begin{pmatrix} 2 & 0 & 3 \\ 4 & 1 & -1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 + 0 + 3 \\ 4 + 1 - 1 \end{pmatrix} = \begin{pmatrix} 5 \\ 4 \end{pmatrix}.$$

We'll see only slowly why this makes geometric sense. Multiplying a fixed matrix by various vectors corresponds to something called a *linear transformation*, which we'll define later. The point is that all linear transformations can be described by multiplication by matrices, and under this identification multiplication of two matrices corresponds to composition of linear transformations (as maps).

5. JANUARY 20

5.1. Warm-up question. Define

$$M = \begin{pmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \\ 3 & 5 & t \end{pmatrix},$$

where t is an unspecified constant, and

$$\mathbf{b} = \begin{pmatrix} 9 \\ 10 \\ 11 \end{pmatrix}.$$

Consider the system of equations $M\mathbf{x} = \mathbf{b}$. For what values of t does this system have no solutions, one solution, and infinitely many solutions, respectively?

5.2. **Answer.** As usual, to find solutions to linear systems of equations, we row reduce. In fact, in this case it suffices to reduce to “row echelon form,” not all the way to “reduced row echelon form.” Without showing all the steps, we can reduce the augmented matrix from

$$\begin{pmatrix} 1 & 3 & 5 & 9 \\ 2 & 4 & 6 & 10 \\ 3 & 5 & t & 11 \end{pmatrix}$$

to

$$\begin{pmatrix} 1 & 3 & 5 & 9 \\ 0 & 1 & 2 & 4 \\ 0 & 0 & t-7 & 0 \end{pmatrix}.$$

Reinterpreting this as representing a system of equations, it is clear that once we solve the last equation $(t-7)x_3 = 0$ for x_3 , the variables x_2 and x_1 are determined uniquely by back-substitution. So it is this third equation that matters for determining the number of solutions to this system.

If $t-7$ is nonzero, then we can divide by it to determine that $x_3 = 0$, so there is one solution. On the other hand, if $t-7$ does equal zero, then x_3 can be anything at all. Thus the answer is that *the system has one solution if $t \neq 7$ and has infinitely many solutions if $t = 7$.*

This illustrates a general feature of linear systems of equations: they are extremely sensitive to small changes in the original parameters. If t is set to equal 7.0000001, then the system has exactly one solution, but if you change it just a little bit then the system has infinitely many (and is thus qualitatively extremely different!). These sensitivity issues are relevant to anyone who wants to implement linear algebra on a computer.

5.3. **Matrices as maps.** We have seen that we can multiply an $m \times n$ matrix and a vector in \mathbb{R}^n , getting a vector in \mathbb{R}^m . A slight change in perspective will be useful to understand the following concepts. Fix an $m \times n$ matrix A . We view A as inducing a map from \mathbb{R}^n to \mathbb{R}^m via multiplication; that is, A is the map that takes a vector $\mathbf{x} \in \mathbb{R}^n$ and returns the vector $A\mathbf{x} \in \mathbb{R}^m$. It will turn out that the maps of vector spaces induced by matrices are very special; they are called *linear transformations*.

5.4. **Column space.** By definition, the *column space* $C(A)$ of a matrix A is the span of the columns of A , thought of as vectors in their own right. For example, if

$$A = \begin{pmatrix} 2 & 3 \\ 1 & 0 \end{pmatrix},$$

then

$$C(A) = \text{Span} \left(\begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ 0 \end{pmatrix} \right),$$

which is all of \mathbb{R}^2 because the two vectors are linearly independent. As another example, if

$$A = \begin{pmatrix} 1 & -3 & -1 \\ -2 & 6 & -6 \end{pmatrix},$$

then

$$C(A) = \text{Span} \left(\begin{pmatrix} 1 \\ -2 \end{pmatrix}, \begin{pmatrix} -3 \\ 6 \end{pmatrix}, \begin{pmatrix} 3 \\ -6 \end{pmatrix} \right) = \text{Span} \left(\begin{pmatrix} 1 \\ -2 \end{pmatrix} \right),$$

which is the line $y = -2x$ lying in \mathbb{R}^2 . Here the columns are linearly dependent, so the column space is smaller than it otherwise might be.

A word of caution: there is no reason for the column space to be preserved under row reduction, and in general it will not be. To figure out the column space of a matrix A , you need to actually know A , not just $\text{rref}(A)$! Here's an example of what can go wrong: the matrix

$$\begin{pmatrix} 1 & 2 & 1 \\ 1 & 2 & 1 \end{pmatrix}$$

clearly has column space spanned by $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$; that is, it is the line $y = x$ in \mathbb{R}^2 . If we row reduce, we get the matrix

$$\begin{pmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \end{pmatrix},$$

which has column space spanned by $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$; that is, it is the line $y = 0$ in \mathbb{R}^2 . The column space changed!

The geometric interpretation of the column space is as follows. Notice that $C(A)$, for an $m \times n$ matrix A , is by definition a subspace of \mathbb{R}^m , the codomain of the map induced by A . My claim is that $C(A)$ is *exactly equal to the image of A as a map*. In other words, $C(A)$ is the set of all vectors $\mathbf{y} \in \mathbb{R}^m$ such that there exists a $\mathbf{x} \in \mathbb{R}^n$ such that $A\mathbf{x} = \mathbf{y}$.

To verify this claim, we first verify that if a vector \mathbf{v} is in $C(A)$, it is the image of A . This is easy: if we write

$$\mathbf{v} = c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_n,$$

where the \mathbf{v}_i are the columns of A and the c_i are scalars, then the vector

$$\mathbf{x} = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix}$$

obeys the equation

$$A\mathbf{x} = c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_n = \mathbf{v},$$

so we see that \mathbf{v} is in the image of A . Conversely, if a vector is in the image of A , it can be written as $A\mathbf{x}$ for some vector

$$\mathbf{x} = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix},$$

and multiplying out we see that

$$\mathbf{v} = A\mathbf{x} = c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_n,$$

so \mathbf{v} is in the column space of A .

5.5. Null space. By definition, the *null space* $N(A)$ of an $m \times n$ matrix A is the set of all vectors $\mathbf{x} \in \mathbb{R}^n$ such that $A\mathbf{x} = \mathbf{0}$. That is, the null space is the set of vectors whose image is zero under the map $\mathbb{R}^n \rightarrow \mathbb{R}^m$ induced by A . Sometimes we also call this set the *kernel*. Colloquially (if somewhat violently), we say that $N(A)$ is the set of vectors killed by the map induced by A .

Because $N(A)$ is defined to be the solution set of a particular equation (the equation $A\mathbf{x} = \mathbf{0}$), row reduction does preserve the null space. Remember, the whole point of row reduction is to simplify a system of linear equations without changing the solution space! So in order to figure out, say, a parametric representation of the null space, one can do the same thing one would do for any system of linear equations: row reduce, then use the free variables as parameters.

6. JANUARY 22

6.1. Warm-up questions. Are the following statements true or false?

- (1) There is a system of linear equations in real numbers with exactly two solutions.
- (2) If A is a square matrix that does not satisfy $\text{rref}(A) = I_n$ (the identity matrix), then $N(A)$ contains a line.
- (3) The set of all vectors orthogonal to a given vector is a linear subspace.
- (4) $C(A) = C(\text{rref}(A))$.
- (5) $N(A) = N(\text{rref}(A))$.
- (6) The set of solutions to a system of linear equations $A\mathbf{x} = \mathbf{b}$ is a linear subspace.

6.2. Answers.

- (1) This is *false*; in general, we have seen that there can be no solutions, one solution, or infinitely many solutions. More precisely, the solution set to a system of linear equations is an affine linear subspace (if it is nonempty); that is, a translate of a linear subspace. If an affine linear subspace contains two points, then it also contains a whole line of points.
- (2) This is *true*. If $\text{rref}(A)$ is not the identity, then because A is a square matrix it must have a column without a pivot. Such a column corresponds to a free variable, which corresponds to a line in the null space.
- (3) This is *true*. Let W be the set of vectors orthogonal to a vector \mathbf{v} . Let's check the axioms of a linear subspace. Certainly W contains the zero vector, as $\mathbf{0} \cdot \mathbf{v} = 0$. If W contains \mathbf{x} and \mathbf{y} , this means that $\mathbf{x} \cdot \mathbf{v} = \mathbf{y} \cdot \mathbf{v} = 0$, which means that

$$(\mathbf{x} + \mathbf{y}) \cdot \mathbf{v} = \mathbf{x} \cdot \mathbf{v} + \mathbf{y} \cdot \mathbf{v} = 0,$$

so $\mathbf{x} + \mathbf{y}$ lies in W . Thus W is closed under addition. If W contains \mathbf{x} , then $\mathbf{x} \cdot \mathbf{v} = 0$, so

$$(c\mathbf{x}) \cdot \mathbf{v} = c\mathbf{x} \cdot \mathbf{v} = c \cdot 0 = 0$$

for any scalar c , which means that $c\mathbf{x}$ lies in W and W is closed under scalar multiplication. Therefore W is a linear subspace.

- (4) This is *false*; row reduction does not in general preserve the column space. I gave an example last class.
- (5) This is *true*; the null space is a solution space and row reduction preserves solution spaces.

- (6) This is *false*, because such a solution set might be a linear translation of a linear subspace (a so-called *affine linear subspace*). A translation of a linear subspace in general satisfies none of the axioms of a linear subspace. For example, a solution set might be any line, but only lines passing through the origin are linear subspaces.

6.3. Bases. A *basis* of a linear subspace V is a set of vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ such that

- (i) $V = \text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_k)$, and
- (ii) $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is a linearly independent set.

In other words, a basis is a spanning set that is as small as possible (it contains no redundancies). Bases in general are highly non-unique.

For example, one possible basis for the line $y = x$ in \mathbb{R}^2 is given by the set $\left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}$. Another is given by the set $\left\{ \begin{pmatrix} -2 \\ -2 \end{pmatrix} \right\}$. All possible bases for this linear subspace will contain exactly one vector: we need one to span the line, and if there were more than one they would be linearly dependent.

The *standard basis* for the vector space \mathbb{R}^n is the set of n vectors

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix}, \dots, \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{pmatrix} \right\}.$$

For example, the standard basis for \mathbb{R}^2 is the set

$$\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}.$$

It is easy to see that these vectors span everything and are linearly independent. The standard basis is, of course, not the only possible basis for \mathbb{R}^n , but it is the simplest and often the easiest to work with.

A basis for the trivial subspace $\{\mathbf{0}\}$ is given by the empty set $\{\}$.

6.4. Finding a basis for the null space. The general algorithm is as follows: we row reduce the matrix, we determine the free variables, we write a parametric representation for the null space, and from the parametric representation we can read off a basis immediately.

For example, set

$$A = \begin{pmatrix} 1 & 2 & 1 & -1 \\ 2 & 3 & 0 & 0 \end{pmatrix}.$$

Row reducing this matrix, we get

$$\text{rref}(A) = \begin{pmatrix} 1 & 0 & -3 & 3 \\ 0 & 1 & 1 & -1 \end{pmatrix}.$$

Our free variables, corresponding to the columns with no pivots, are x_3 and x_4 . Therefore we can write the solution space as the set of all vectors of the form

$$\begin{pmatrix} 3x_3 - 3x_4 \\ -x_3 + x_4 \\ x_3 \\ x_4 \end{pmatrix},$$

where x_3 and x_4 are free (this is the parametric representation). We rewrite this as

$$x_3 \begin{pmatrix} 3 \\ -1 \\ 1 \\ 0 \end{pmatrix} + x_4 \begin{pmatrix} -3 \\ 1 \\ 0 \\ 1 \end{pmatrix}.$$

Because these two vectors are linearly independent, and by the parametric representation they span the null space, we find that a basis for the null space is

$$\left\{ \begin{pmatrix} 3 \\ -1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} -3 \\ 1 \\ 0 \\ 1 \end{pmatrix} \right\}.$$

Note that this is far from the only possible basis we could write down. For example, we could multiply one of the vectors by a scalar, or add a nonzero multiple of one vector to the other, and we would still have a basis.

It is a general fact that the parametric representation we get will always give us linearly independent vectors (hence a basis for the null space). This is because if there were a linear dependence relation between the vectors, the corresponding variables could not have all been free variables (exercise: see if you can make this argument precise to furnish a proof of our algorithm!).

6.5. Finding a basis for the column space. The column space of a matrix is just the subspace spanned by the column vectors. To find a basis, we have to find a subset of these vectors that still span the whole space, but are linearly independent.

One way of doing this is the following algorithm: row reduce the matrix, and determine which of the columns have pivots. Take the vectors of the original matrix corresponding to these columns. That will be a basis for the column space.

Here's an example. Take

$$A = \begin{pmatrix} 1 & 2 & -3 \\ 5 & 0 & 5 \\ 0 & 2 & -4 \end{pmatrix}.$$

This matrix row reduces to

$$\text{rref}(A) = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -2 \\ 0 & 0 & 0 \end{pmatrix}.$$

The pivot columns are the first and second, so we take the first and second columns of the original matrix A and conclude that one basis for $C(A)$ is

$$\left\{ \begin{pmatrix} 1 \\ 5 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 0 \\ 2 \end{pmatrix} \right\}.$$

Yet again, this is far from the only possible basis for $C(A)$.

Why does this work? Essentially, the reason is that although row reduction does not preserve $C(A)$ itself, it *does* preserve the linear dependence relationships between the columns (this is equivalent to the statement that row reduction preserves the null space). For example, we see that in $\text{rref}(A)$ as above, the third column is equal to the first column plus -2 times the second column (this is clear precisely because $\text{rref}(A)$ is in reduced row echelon form). We can check that the same is

true in A : the first column, plus -2 times the second column, is equal to the third column.

Because row reduction preserves linear dependence relations, and the columns in a row-reduced matrix are all generated by the pivot columns which are themselves linearly independent (check this!), the columns in the original matrix corresponding to the pivot columns also generate all the columns while being linearly independent. In other words, they form a basis.

6.6. Nonempty solution spaces are translations of the null space. Moving away from bases for the moment, I want to discuss briefly an important geometric statement about solution spaces to systems $A\mathbf{x} = \mathbf{b}$. The statement is essentially Proposition 8.2 in the textbook, but I want to emphasize the geometric meaning:

Proposition 6.1. *Let V be the solution set to $A\mathbf{x} = \mathbf{b}$ and let $N = N(A)$. If V is not the empty set, then V is a translation of N .*

Proof. By assumption, V is nonempty, so we can find some $\mathbf{x}_p \in V$ (the p stands for “particular,” as in “particular solution to $A\mathbf{x} = \mathbf{b}$ ”). I claim that V is a translation of N by the vector \mathbf{x}_p . To show this, we need to show two things: first, we need to show that if $\mathbf{n} \in N$, then $\mathbf{n} + \mathbf{x}_p$ is in V . Second, we need to show that if $\mathbf{v} \in V$, then $\mathbf{v} - \mathbf{x}_p$ is in N .

To show the first, we calculate that

$$A(\mathbf{n} + \mathbf{x}_p) = A\mathbf{n} + A\mathbf{x}_p = \mathbf{0} + \mathbf{b} = \mathbf{b},$$

so by the definition of V we know that $\mathbf{n} + \mathbf{x}_p$ is in V .

To show the second, we calculate that

$$A(\mathbf{v} - \mathbf{x}_p) = A\mathbf{v} - A\mathbf{x}_p = \mathbf{b} - \mathbf{b} = \mathbf{0},$$

so by the definition of null space we know that $\mathbf{v} - \mathbf{x}_p$ is in N . □

The proposition is false without the assumption that V is nonempty. For example, let

$$A = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}.$$

Then $N(A) = \mathbb{R}^2$ because the matrix sends *every* vector to $\mathbf{0}$, but if $\mathbf{b} \neq \mathbf{0}$ then there are no solutions to $A\mathbf{x} = \mathbf{b}$ (again, precisely because the matrix sends every vector to $\mathbf{0}$). So the null space N is very large, but the solution space V is often empty. In general, the larger the null space, the easier it is for V to be totally empty.

This tells us a lot about the geometry of solutions to systems of linear equations. For example, if we know that the null space of a given matrix is a line in a certain direction, then the solution space to a system of equations involving that matrix must be a line in the same direction, or empty. Less trivially, it is the easiest way to rigorously prove (as you were asked to do on the last homework assignment) that the intersection of two planes in \mathbb{R}^3 cannot be a point. One first shows that the intersection of two planes is described by the solution set V of a system of three equations in four unknowns. The corresponding matrix has three rows and four columns, so its row reduction must have a free variable (there are not enough possible pivots for all the columns), so its null space must contain a line. By the above proposition, V must either be empty or a translation of a linear subspace that contains a line. In particular, V cannot be a point.

7. JANUARY 27: EXAM 1 REVIEW

7.1. Some things you need to know.

- Definitions are important to know, both because you may be asked to give one directly and more importantly because in order to prove anything one has to start with a definition. For example, we have learned definitions for the terms *linearly independent*, *linearly dependent*, *subspace*, *basis*, *dimension*...
- How to calculate reduced row echelon form (rref), including how to recognize if a matrix is in rref. Almost all of the algorithms in this class use row reduction in some capacity; for example, solving linear systems of equations, finding a basis for the null space of a matrix, finding a basis for the column space of a matrix, finding equations for a subspace (e.g. the column space of a matrix or a given parametric representation)...
- Properties of the dot product. Most important is the statement that two vectors are orthogonal if and only if their dot product is zero. Also recall the law of cosines and various properties of the cross product in \mathbb{R}^3 .
- Finding parametric equations for lines and planes passing through specified points, and writing down the equation of a line in \mathbb{R}^2 or a plane in \mathbb{R}^3 by finding a normal vector and a point.
- Anything from the homework is fair game; for example, the projection operator that was defined in the first homework can be used to quickly find the point on a plane of shortest distance from the origin.
- And so on...

7.2. **Some problems.** These are all problems (some slightly modified) that have appeared on recent past exams. Most are on the more difficult side.

- (1) Consider the equation

$$A \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \\ -2 \end{pmatrix}.$$

Write down matrices A such that the above equation has infinitely many solutions, exactly one solution, and zero solutions, or show that no such matrix can exist.

- (2) Assume that

$$\text{rref}(A) = \begin{pmatrix} 1 & 0 & -2 & 0 \\ 0 & 1 & 5 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Let \mathbf{a}_1 , \mathbf{a}_2 , \mathbf{a}_3 , and \mathbf{a}_4 be the columns of A . Write down all of the subsets of the set $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4\}$ that are bases for $C(A)$.

- (3) Let P be the plane in \mathbb{R}^3 containing the points

$$\begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix}, \quad \text{and} \quad \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}.$$

Find a parametric representation *and* an equation defining P .

(4) Let

$$A = \begin{pmatrix} 1 & 3 & 2 \\ a & 6 & 2 \\ 0 & 9 & 5 \end{pmatrix}.$$

Calculate $\text{rank}(A)$ for all possible values of a .

(5) Let

$$A = \begin{pmatrix} 1 & 3 & 0 \\ 2 & 7 & 0 \\ 1 & 2 & 0 \end{pmatrix}.$$

Find all possible \mathbf{b} such that $A\mathbf{x} = \mathbf{b}$ has a solution (that is, find equations in the components of \mathbf{b}).

(6) Suppose \mathbf{u} , \mathbf{v} , and \mathbf{w} are unit vectors (i.e., they have length 1) and are mutually orthogonal. Find a scalar c such that $\mathbf{u} + 2\mathbf{v} + 3\mathbf{w}$ and $5\mathbf{u} + \mathbf{v} + c\mathbf{w}$ are orthogonal.

7.3. Some solutions.

(1) Clearly, for all of the three possibilities A must be a 3×2 matrix (otherwise the multiplication does not make sense).

In order for there to be infinitely many solutions to the given equation,

the vector $\mathbf{b} = \begin{pmatrix} 2 \\ 1 \\ -2 \end{pmatrix}$ must be in the column space of A and the null

space of A must contain a line (that is, we need $\text{nullity}(A) \geq 1$). This is certainly possible; for example, we can make one of the columns \mathbf{b} itself and the other column linearly dependent (i.e., a multiple of \mathbf{b}). One possibility is

$$A = \begin{pmatrix} 2 & 0 \\ 1 & 0 \\ -2 & 0 \end{pmatrix}.$$

In order for there to be exactly one solution, we need \mathbf{b} to be in $C(A)$ but $\text{nullity}(A) = 0$; i.e., we can work as we did above but instead of making the other column linearly dependent of \mathbf{b} we make it linearly independent. Therefore one other possibility is

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 0 \\ -2 & 0 \end{pmatrix}.$$

As an alternative argument, we could have made the first two equations easily and uniquely solvable (for example, $x_1 = 2$ and $x_2 = 1$) and added a third row of the matrix that yields a true but unnecessary equation, like $-x_1 = -2$. This would yield the matrix

$$A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \end{pmatrix},$$

for instance.

In order for there to be no solutions, we need to arrange that \mathbf{b} does not lie in $C(A)$. This is easy; for example, the matrix A with all entries zero will do the trick. In particular, notice that if we set the last row to be all

zeroes, then we get the inconsistent equation $0 = -2$, so any such matrix will work. In fact, if you write down a “random” 3×2 matrix, you will likely get no solutions: it is the generic case.

One can visualize this problem by viewing A as a map from \mathbb{R}^2 to \mathbb{R}^3 ; that is, we are mapping a plane into a 3-dimensional space. In order for $A\mathbf{x} = \mathbf{b}$ to have no solutions, it suffices that the image miss \mathbf{b} (that is, $\mathbf{b} \notin C(A)$), which is easy to arrange. In order for $A\mathbf{x} = \mathbf{b}$ to have one solution, it suffices that the image hit \mathbf{b} (that is, $\mathbf{b} \in C(A)$) and that there be no “collapsing” of the map, which translates into the criterion that the nullity should be zero. In order for $A\mathbf{x} = \mathbf{b}$ to have infinitely many solutions, we must have $\mathbf{b} \in C(A)$ and there must be nontrivial “collapsing” of the map (so the nullity is nonzero); that is, the image of the plane should be a line passing through \mathbf{b} .

- (2) From the algorithm we know for finding a basis for $C(A)$, the set $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_4\}$ is a basis (picking the vectors corresponding to the pivot columns). Given any subspace such as $C(A)$, any basis must have the same size, so we know any possible basis drawn from the set $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4\}$ must have three elements. But we can read off the linear dependence relationships among the columns from the reduced row echelon form. In particular, we have $\mathbf{a}_3 = -2\mathbf{a}_1 + 5\mathbf{a}_2$ but no other relationships. Therefore the other possible bases with elements drawn from the columns of A are $\{\mathbf{a}_1, \mathbf{a}_3, \mathbf{a}_4\}$ and $\{\mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4\}$. The last possibility, $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3\}$, is a linearly dependent set and therefore not a basis.
- (3) To find a parametric representation we need one point \mathbf{x}_0 on the plane and two linearly independent vectors \mathbf{u} and \mathbf{v} in the direction of the plane. In this case, we might as well pick

$$\mathbf{x}_0 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \quad \mathbf{u} = \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix}, \quad \text{and} \quad \mathbf{v} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 2 \\ 2 \end{pmatrix}.$$

A parametric representation is then

$$P = \left\{ \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + s \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} + r \begin{pmatrix} 0 \\ 2 \\ 2 \end{pmatrix} \mid s, t \in \mathbb{R} \right\}.$$

There are many other possibilities.

To find an equation, we need to find a normal vector \mathbf{n} and a point in the plane \mathbf{x}_0 . Using the same point \mathbf{x}_0 as before, we can find a normal vector in several ways: we could set $\mathbf{n} \cdot \mathbf{u} = \mathbf{n} \cdot \mathbf{v} = 0$ with \mathbf{u}, \mathbf{v} as above and solve for the components of \mathbf{n} (we will have one free parameter, because the length of \mathbf{n} is irrelevant), or, perhaps more easily, we can take the cross product of \mathbf{u} and \mathbf{v} . We calculate

$$\mathbf{n} = \mathbf{u} \times \mathbf{v} = \begin{pmatrix} -1 \\ 1 \\ -2 \end{pmatrix} \times \begin{pmatrix} 0 \\ 2 \\ 2 \end{pmatrix} = \begin{pmatrix} 6 \\ 2 \\ -2 \end{pmatrix},$$

so the equation of the plane, $\mathbf{n} \cdot (\mathbf{x} - \mathbf{x}_0) = 0$, simplifies to

$$6(x - 1) + 2(y - 0) - 2(z - 1) = 0,$$

which simplifies further to

$$3x + y - z = 2.$$

- (4) We row reduce as usual, carrying around the baggage of the unknown a :

$$\begin{aligned} \begin{pmatrix} 1 & 3 & 2 \\ a & 6 & 2 \\ 0 & 9 & 5 \end{pmatrix} &\rightsquigarrow \begin{pmatrix} 1 & 3 & 2 \\ 0 & 6-3a & 2-2a \\ 0 & 9 & 5 \end{pmatrix} \rightsquigarrow \begin{pmatrix} 1 & 3 & 2 \\ 0 & 9 & 5 \\ 0 & 6-3a & 2-2a \end{pmatrix} \rightsquigarrow \\ &\begin{pmatrix} 1 & 3 & 2 \\ 0 & 1 & 5/9 \\ 0 & 6-3a & 2-2a \end{pmatrix} \rightsquigarrow \begin{pmatrix} 1 & 0 & 1/3 \\ 0 & 1 & 5/9 \\ 0 & 0 & -4/3 - a/3 \end{pmatrix}. \end{aligned}$$

I switched rows in the second step for convenience, but it isn't absolutely necessary. Now notice that if the lower right corner is zero, then the matrix is already row reduced and there will be one free variable, so the rank will be 2. If the lower right corner is nonzero, we will be able to row reduce all the way to the identity, so the rank will be 3. The lower right corner is zero if and only if $a = -4$, so the answer is that the rank is 3 if $a \neq -4$ and 2 if $a = -4$.

This accords with common sense: if we pick three "random" vectors in \mathbb{R}^3 , they will likely not all lie in a plane, so they will form a basis; correspondingly, a "random" 3×3 matrix will have rank 3, and will only have lower rank in exceptional circumstances (in this case, only when $a = -4$).

- (5) We write down the augmented matrix

$$\begin{pmatrix} 1 & 3 & 0 & b_1 \\ 2 & 7 & 0 & b_2 \\ 1 & 2 & 0 & b_3 \end{pmatrix}$$

and do row reduction. I won't write out all the steps, but at the end of the day we get two rows with pivots and one row that looks like $(0, 0, 0, -3b_1 + b_2 + b_3)$. The rows with pivots do not constrain the components b_i , they just let you know what some of the variables x_i will have to be in terms of the b_i . The last row, however, expresses the constraint $-3b_1 + b_2 + b_3 = 0$, which is what we are after.

Incidentally, we can use this procedure (which will find equations for the column space of a matrix in general) to find a matrix B such that $N(B) = C(A)$. Since we know here that $C(A)$ is the locus of points where $-3b_1 + b_2 + b_3 = 0$, we can translate this into vector language by writing

$$(-3 \ 1 \ 1) \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} = 0.$$

Therefore setting

$$B = (-3 \ 1 \ 1)$$

gives us a matrix with $N(B) = C(A)$.

- (6) This is a fairly simple calculation: we want to find a c such that $(\mathbf{u} + 2\mathbf{v} + 3\mathbf{w}) \cdot (5\mathbf{u} + \mathbf{v} + c\mathbf{w}) = 0$. Expand out everything and notice that the dot

products of all distinct elements of the set $\{\mathbf{u}, \mathbf{v}, \mathbf{w}\}$ are all zero because the vectors are mutually orthogonal. Then we get the equation

$$5\|\mathbf{u}\|^2 + 2\|\mathbf{v}\|^2 + 3c\|\mathbf{w}\|^2 = 0.$$

By definition, unit vectors have norm one, so this simplifies further to

$$5 + 2 + 3c = 0.$$

We calculate $c = -7/3$.

8. JANUARY 29

8.1. Warm-up questions. Recall that, by definition, a *linear transformation* \mathbf{T} is a function from some \mathbb{R}^n to some \mathbb{R}^m that obeys the rules $\mathbf{T}(\mathbf{x} + \mathbf{y}) = \mathbf{T}(\mathbf{x}) + \mathbf{T}(\mathbf{y})$ for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\mathbf{T}(c\mathbf{x}) = c\mathbf{T}(\mathbf{x})$ for all $c \in \mathbb{R}$, $\mathbf{x} \in \mathbb{R}^n$.

Which of the following functions are linear transformations?

- (1) $\mathbf{f}(x_1, x_2, x_3) = (0, 0)$.
- (2) $\mathbf{f}(x_1, x_2) = (x_1 + 7x_2, x_1)$.
- (3) $\mathbf{f}(x_1, x_2) = (x_1, x_1)$.
- (4) $\mathbf{f}(x_1, x_2) = (2, x_1)$.
- (5) $\mathbf{f}(x_1) = (x_1 + 2, 0)$.
- (6) $\mathbf{f}(x_1, x_2) = (x_1^2 - x_2, 2x_2)$.

8.2. Answers and elaborations. We will see that it is possible to pick out linear transformations at a glance: they are exactly those functions that are linear functions of the inputs with no constant term in each variable.

- (1) This map sends each vector to the zero vector. Checking the two rules is easy: both sides of both equation are always zero, no matter what, so the equalities hold. This is a linear transformation.
- (2) Let's check the two rules. Picking any (x_1, x_2) and any (y_1, y_2) , we calculate

$$\begin{aligned} \mathbf{f}(x_1 + y_1, x_2 + y_2) &= (x_1 + y_1 + 7(x_2 + y_2), x_1 + y_1) \\ &= (x_1 + 7x_2, x_1) + (y_1 + 7y_2, y_1) \\ &= \mathbf{f}(x_1, x_2) + \mathbf{f}(y_1, y_2), \end{aligned}$$

so the first rule holds. Picking any (x_1, x_2) and any constant c , we calculate

$$\mathbf{f}(cx_1, cy_1) = (cx_1 + 7(cx_2), cx_1) = c(x_1 + 7x_2, x_1) = c\mathbf{f}(x_1, x_2),$$

so the second rule holds. Therefore \mathbf{f} is a linear transformation.

- (3) We can run the same check as in the previous part to show that \mathbf{f} is a linear transformation.
- (4) This time, if we try to run the same check, we have a problem. For example, with the specific point $(x_1, x_2) = (1, 0)$ and the specific constant $c = 2$, the second rule tells us that we should have $2 \cdot \mathbf{f}(1, 0) = \mathbf{f}(2 \cdot (1, 0))$. But this is not the case: the left hand side here is $(4, 2)$ and the right hand side is $(2, 2)$ (one can also see that the other rule fails). Therefore \mathbf{f} is not a linear transformation. This illustrates a general fact: just as lines and planes that do not pass through the origin are not considered linear subspaces, linear functions with nonzero constant terms are not considered linear transformations.

- (5) Again, the constant 2 is messing things up. We can plug in the specific example of $x_1 = 1$ and $c = 2$ and get that $2\mathbf{f}(1) = (6, 0)$ and $\mathbf{f}(2) = (4, 0)$, which are not equal.
- (6) This time there are more obvious problems with the map: there is a quadratic term! Check for yourself using the same point and constant as in (4) above that this function is not a linear transformation.

Each linear transformation can be represented by a matrix. More precisely, this means that every linear transformation $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ can be written as $\mathbf{f}(\mathbf{x}) = A\mathbf{x}$ for some $m \times n$ matrix A ; that is, every linear transformation “is really” matrix multiplication with some matrix. Conversely, it is easy to see that multiplication by a matrix is a linear transformation. Thus linear transformations and matrices for us are really the same thing.¹

It is easy to find the matrix corresponding to a linear transformation that is given as a function: we apply the function to the standard basis vectors and concatenate the answers we get into a matrix. For example, in the first example, we calculate

$$\mathbf{f}\left(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}\right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{f}\left(\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}\right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad \mathbf{f}\left(\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}\right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

so the matrix corresponding to \mathbf{f} is

$$A = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

Admittedly, this is not a very interesting matrix. For the second example, we can calculate

$$\mathbf{f}\left(\begin{pmatrix} 1 \\ 0 \end{pmatrix}\right) = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \mathbf{f}\left(\begin{pmatrix} 0 \\ 1 \end{pmatrix}\right) = \begin{pmatrix} 7 \\ 0 \end{pmatrix},$$

so the corresponding matrix is

$$A = \begin{pmatrix} 1 & 7 \\ 1 & 0 \end{pmatrix}.$$

In the third example, the only standard basis vector of \mathbb{R}^1 is the vector (1) , which evaluates to $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$, so the corresponding matrix is

$$A = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

8.3. Examples of linear transformations. From now on we will freely move between the concepts of a linear transformation as a function and as a matrix.

¹A word of caution that you should feel free to ignore: this is only true because we are only considering vector spaces that come with a preferred basis, the standard basis of \mathbb{R}^n . In abstract linear algebra, this is no longer quite true: a matrix is then “the same” as a linear transformation only once a basis has been chosen.

8.3.1. *The identity transformation.* The $n \times n$ matrix

$$I_n = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix},$$

which has 1 in each diagonal entry and zero in every other entry, is called the *identity matrix*. It corresponds to the function that does nothing to each vector, because one can calculate

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} x_1 + 0 + 0 + \dots \\ 0 + x_2 + 0 + \dots \\ \vdots \\ \dots + 0 + 0 + x_n \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}.$$

8.3.2. *Scalar transformations.* Let $\alpha \in \mathbb{R}$, and consider the matrix

$$\alpha I_n = \begin{pmatrix} \alpha & 0 & \dots & 0 \\ 0 & \alpha & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & \alpha \end{pmatrix}.$$

We can calculate just as above that this is the linear transformation that multiplies each component by α :

$$\begin{pmatrix} \alpha & 0 & \dots & 0 \\ 0 & \alpha & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & \alpha \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \alpha x_1 \\ \alpha x_2 \\ \vdots \\ \alpha x_n \end{pmatrix}.$$

Geometrically, this linear transformation stretches each vector by a factor of α ; that is, it “scales” the whole vector space by α .

8.3.3. *Diagonal transformations.* Let’s generalize a bit. Take $\alpha_1, \alpha_2, \dots, \alpha_n \in \mathbb{R}$ and consider the matrix

$$\begin{pmatrix} \alpha_1 & 0 & \dots & 0 \\ 0 & \alpha_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & \alpha_n \end{pmatrix}.$$

Clearly, this matrix multiplies the i th entry in a vector by the corresponding α_i ; i.e., it stretches the vector space by a factor of α_i in the i th direction, for each i .

8.3.4. *Rotations.* Now we’ll specialize to $m = n = 2$ for concreteness, although it is possible to write down the general form of a $n \times n$ rotation matrix. The linear transformation \mathbf{R}_θ that rotates \mathbb{R}^2 around the origin by an angle θ in the counterclockwise direction is a linear transformation: see pages 90 and 91 in the textbook for a “proof by picture.” We can calculate the corresponding matrix in the usual way, by seeing what happens to the standard basis elements $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and

$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$. See the picture halfway down page 91 in the textbook for an illustration. Simple trigonometry tells us that

$$\mathbf{R}_\theta \left(\begin{pmatrix} 1 \\ 0 \end{pmatrix} \right) = \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix}, \quad \mathbf{R}_\theta \left(\begin{pmatrix} 0 \\ 1 \end{pmatrix} \right) = \begin{pmatrix} -\sin \theta \\ \cos \theta \end{pmatrix},$$

so the corresponding matrix is

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}.$$

It is useful to keep this picture in mind, because sometimes it's hard to remember this formula!

8.3.5. *Projections.* We saw the formula for the projection of \mathbf{x} onto \mathbf{v} on the first homework:

$$\text{Proj}_{\mathbf{v}}(\mathbf{x}) = \left(\frac{\mathbf{x} \cdot \mathbf{v}}{\mathbf{v} \cdot \mathbf{v}} \right) \mathbf{v}.$$

It is a simple algebra exercise to show that this gives a linear transformation. To find the corresponding matrix, we do the same procedure: calculate what happens to the standard basis vectors. We'll work in \mathbb{R}^2 again for definitiveness, although similar calculations work for general n . If $\mathbf{v} = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$, then

$$\text{Proj}_{\mathbf{v}} \left(\begin{pmatrix} 1 \\ 0 \end{pmatrix} \right) = \begin{pmatrix} \frac{v_1^2}{v_1^2 + v_2^2} \\ \frac{v_1 v_2}{v_1^2 + v_2^2} \end{pmatrix}, \quad \text{Proj}_{\mathbf{v}} \left(\begin{pmatrix} 0 \\ 1 \end{pmatrix} \right) = \begin{pmatrix} \frac{v_1 v_2}{v_1^2 + v_2^2} \\ \frac{v_2^2}{v_1^2 + v_2^2} \end{pmatrix},$$

so the corresponding matrix is

$$\frac{1}{v_1^2 + v_2^2} \begin{pmatrix} v_1^2 & v_1 v_2 \\ v_1 v_2 & v_2^2 \end{pmatrix}.$$

This formula simplifies considerably if \mathbf{v} is a unit vector, because then the factor in front is identically 1.

8.3.6. *Reflections.* Let L be a line through the origin in \mathbb{R}^n spanned by a vector \mathbf{v} . Reflecting points about L is easily seen to be a linear transformation; in fact, using the picture on top of page 95 in the textbook, we can see that the reflection of \mathbf{x} about L is equal to two times its projection onto \mathbf{v} minus \mathbf{x} itself. In matrices,

$$\text{Ref}_L(\mathbf{x}) = 2 \text{Proj}_{\mathbf{v}}(\mathbf{x}) - I_n.$$

9. FEBRUARY 3

9.1. **Warm-up questions.** Write down matrices corresponding to the following linear transformations:

- (1) Rotation counterclockwise by $\pi/3$ radians in \mathbb{R}^2 .
- (2) Scale everything by $1/2$ in \mathbb{R}^3 .
- (3) Reflect everything about the line $y = x$ in \mathbb{R}^2 .
- (4) Project onto the x -axis in \mathbb{R}^2 , then rotate by $\pi/2$ radians counterclockwise.

9.2. Warm-up answers.

- (1) The general form for a rotation matrix in \mathbb{R}^2 is

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix},$$

so we can just plug in $\theta = \pi/3$ to get

$$\begin{pmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{pmatrix}.$$

Alternatively, we could have worked from scratch and figured out where the standard basis vectors were sent using a little simple trigonometry. We would have found that the linear transformation in question sends

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} \mapsto \begin{pmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} 0 \\ 1 \end{pmatrix} \mapsto \begin{pmatrix} -\frac{\sqrt{3}}{2} \\ \frac{1}{2} \end{pmatrix}.$$

Concatenating these two vectors gives the same matrix.

- (2) This is straightforward; we merely take $\frac{1}{2}$ multiplied by the identity matrix I_3 and get

$$\begin{pmatrix} \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} \end{pmatrix}.$$

Again, we could have figured this out by noting where the three standard basis vectors go and concatenating the result.

- (3) Here the best way forward is definitely to figure out where the standard basis vectors go, rather than try anything using the formulas for reflection matrices in the textbook. We see (perhaps by drawing a quick picture) that the basis vector $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$, which lies on the x -axis, is reflected up to the point $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ on the y -axis. Symmetrically, the second basis vector $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ gets sent to $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$. Therefore the matrix is

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

- (4) Here we have a composition of two maps, which we can figure out separately and then multiply together. Projection onto the x -axis does not change the first basis vector $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and moves the second basis vector $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ to the origin $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$, so the matrix is

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

Rotation by $\pi/2$ radians counterclockwise is represented by plugging in $\theta = \pi/2$ to the standard rotation matrix, yielding

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

Now we just multiply these two matrices together. Because we are projecting *first*, then rotating, we put the projection matrix *on the right*. Convince yourself that this makes sense, because it's counterintuitive! Function composition goes right-to-left, which is the opposite of most conventions in mathematics. In all,

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

is the desired matrix.

9.3. Matrix multiplication practice. We went through an example like

$$\begin{pmatrix} 1 & 2 \\ 3 & 0 \\ 1 & 1 \end{pmatrix} \cdot \begin{pmatrix} 5 & -1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 7 & 1 \\ 15 & -3 \\ 6 & 0 \end{pmatrix}.$$

Remember that the product of an $m \times n$ matrix and a $p \times q$ matrix is defined if and only if $n = p$, and the result will be a $m \times q$ matrix. Also remember that if AB and BA are both defined (which means they're both square matrices of the same size!), there is no reason for them to be equal! We say that matrix multiplication is not commutative. It does obey all the other rules you would expect: it distributes over matrix addition (i.e. $A(B + C) = AB + AC$) and commutes with scalar multiplication (i.e. $A(cB) = c(AB)$). See the textbook for an exhaustive list of properties that matrix multiplication satisfies.

9.4. Inverses. If you are not familiar with the concepts of one-to-one, onto, and invertible maps of sets (also sometimes called injective, surjective, and bijective respectively), you should read the first part of Section 16 of the textbook carefully! We will use the key fact, which is "obvious" once you think about it for long enough but needs to be carefully parsed, that a map of sets is invertible if and only if it is one-to-one and onto.

Let's derive conditions for the invertibility of a matrix. We'd like to make the same statement for linear transformations as for maps of sets: i.e., we'd like to state that a matrix is invertible if and only if its corresponding linear transformation is one-to-one and onto. In order to assert this statement, we need the second key fact: the inverse of a linear transformation as a map of sets is itself always a linear transformation. This is stated, but not proved, in the textbook (it's left as an exercise). Therefore we can speak interchangeably about invertibility as a map of sets and invertibility as a linear transformation.

The linear transformation associated to an $m \times n$ matrix A is onto if and only if $C(A)$ is everything; i.e., all of \mathbb{R}^m . In other words, $\text{rank}(A) = m$.

The linear transformation associated to A is one-to-one if and only if for each \mathbf{b} in the image (which is $C(A)$), the solution space of the system of equations $A\mathbf{x} = \mathbf{b}$ consists of exactly one point. By definition, since \mathbf{b} is in the image, there is at least one solution \mathbf{x}_p . Therefore by Proposition 8.2 in the text, which states that a nonzero solution set is the translate of the null space, the linear transformation associated to A is one-to-one if and only if the null space consists of a point; i.e., $N(A) = \{\mathbf{0}\}$. In other words, $\text{nullity}(A) = 0$, or equivalently by the rank-nullity theorem, $\text{rank}(A) = n$.

Putting these two paragraphs together with the statement that invertibility is the same as onto and one-to-one, we see that invertibility of an $m \times n$ matrix is

equivalent to $m = n = \text{rank}(A)$. That is, an invertible matrix is square and has full rank, and conversely a square matrix of full rank is invertible. It is easy to see from this that a matrix is invertible if and only if it row reduces to the identity matrix I_n .

When we think about what the above discussion means, a lot of it makes intuitive sense. For example, if $m > n$ there should not be a linear transformation from \mathbb{R}^n to \mathbb{R}^m that hits every point of \mathbb{R}^m . The latter space is simply “too big.” This intuition is supported by the statement above that a linear transformation is onto if and only if $\text{rank}(A) = m$, which in particular implies that there are at least m columns, so $n \geq m$.

To calculate inverse matrices, there is a simple algorithm using row reduction. If A is invertible, then the augmented matrix $[A|I_n]$ will row reduce to $[I_n|A^{-1}]$, and if A is not invertible then it will row reduce to something not of that form (because A will not row reduce to I_n).

For example, let

$$A = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}.$$

Augmenting with the identity matrix, we get

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

which row reduces in two steps (check!) to

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

Therefore A is invertible and

$$A^{-1} = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}.$$

Whenever we calculate a matrix inverse, we can check easily to see that we have the right answer:

$$AA^{-1} = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \cdot \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I_2.$$

10. FEBRUARY 5

10.1. **Warm-up questions.** Calculate A^{-1} , where

(1)

$$A = \begin{pmatrix} 1 & 3 \\ 2 & -5 \end{pmatrix}.$$

(2)

$$A = \begin{pmatrix} 1 & 1 & -1 \\ 1 & 2 & 3 \\ 2 & 1 & 1 \end{pmatrix}.$$

10.2. Warm-up answers. This is merely an exercise in row reduction, following the algorithm explained in the last section. I will show all the steps for the first one and merely the answer for the second.

(1) We row reduce as follows:

$$\begin{pmatrix} 1 & 3 & 1 & 0 \\ 2 & -5 & 0 & 1 \end{pmatrix} \rightsquigarrow \begin{pmatrix} 1 & 3 & 1 & 0 \\ 0 & -11 & -2 & 1 \end{pmatrix} \rightsquigarrow \begin{pmatrix} 1 & 3 & 1 & 0 \\ 0 & 1 & \frac{2}{11} & -\frac{1}{11} \end{pmatrix} \\ \rightsquigarrow \begin{pmatrix} 1 & 0 & \frac{5}{11} & \frac{3}{11} \\ 0 & 1 & \frac{2}{11} & -\frac{1}{11} \end{pmatrix}.$$

Therefore

$$A^{-1} = \begin{pmatrix} \frac{5}{11} & \frac{3}{11} \\ \frac{2}{11} & -\frac{1}{11} \end{pmatrix}.$$

(2) After augmenting by the identity matrix I_3 and row reducing, we find that

$$A^{-1} = \begin{pmatrix} -\frac{1}{7} & -\frac{2}{7} & \frac{5}{7} \\ \frac{5}{7} & \frac{3}{7} & -\frac{4}{7} \\ -\frac{3}{7} & \frac{1}{7} & \frac{1}{7} \end{pmatrix}.$$

After finding both answers, we should do a quick sanity check that AA^{-1} really is the identity matrix.

10.3. Determinants. The *determinant* of a matrix is given by a rather strange formula, and may at first seem like a bizarre concept. Recursively, we define $\det((c)) = c$ for a 1×1 matrix (c) , and let

$$\det(A) = \sum_{j=1}^n (-1)^{1+j} a_{1j} \det(A_{1j}),$$

where A_{1j} is the matrix you get from A when you delete the first row and the j th column. This is called *expanding in minors*². The $(-1)^{1+j}$ factor just adds a minus sign on every other summand. We notate the determinant by large vertical bars around a matrix. For example, we calculate

$$\begin{vmatrix} 1 & 3 \\ 2 & -5 \end{vmatrix} = 1 \cdot \det((-5)) - 3 \cdot \det((2)) = -5 - 6 = -11.$$

In general, the determinant of a 2×2 matrix is the product of the diagonal entries minus the product of the anti-diagonal entries.

Keep in mind that this is not the only formula one can use to express the determinant; there is a similar formula where we expand in minors along any row (not just the first row) or even any column. For example, the formula for the determinant expanded along the i th row is

$$\det(A) = \sum_{j=1}^n (-1)^{i+j} a_{ij} \det(A_{ij}).$$

These formulas are often useful because if we see a row (or column) that has many zero entries, the formula becomes particularly simple: any summand corresponding to a zero entry a_{ij} vanishes entirely. In particular, if there is an entire row or column consisting of all zeroes, the determinant vanishes and no further calculation is necessary.

²the matrices A_{1j} are examples of *minors* of a matrix, which are just the matrices you get when you delete one row and one column.

Although this is a reasonable formula for small matrices, it quickly becomes unwieldy for large ones, as the number of calculations increases exponentially in the size of the matrix. A better algorithm, at least if one were to implement it on a computer, involves row reduction, combined with the observation that upper (or lower) triangular matrices have particularly easy-to-calculate determinants. Let's see how to implement it.

Row reduction involves three steps, and it is easy to check (see the textbook for details!) how each of them affects the determinant. Adding a multiple of one row to another leaves the determinant entirely unchanged. Multiplying a row by a scalar c multiplies the corresponding determinant by c . Swapping two rows multiplies the corresponding determinant by -1 (i.e., the determinant switches sign). So if, as we row reduce, we keep track of these numbers to be multiplied together, the problem reduces to figuring out the determinant of a row reduced matrix.

In fact, we don't have to row reduce all the way to reduced row echelon form; an upper triangular matrix will do. For definitiveness, let's consider the example

$$A = \begin{pmatrix} 1 & 5 & 10^{100} \\ 1 & 2 & 5 \\ 0 & 0 & -3 \end{pmatrix}.$$

A priori this determinant might be a very nasty thing; look at that upper right hand corner! But in fact it's very simple. Let's expand along the third row, because it has a lot of zeroes. The only entry that contributes is the lower left:

$$\det(A) = 0 - 0 + (-3) \begin{vmatrix} 1 & 5 \\ 0 & 2 \end{vmatrix}.$$

And if we expand this 2×2 matrix along the lower row, we see that only the 2 contributes; i.e.

$$\det(A) = (-3)(-0 + 2 \cdot 1) = -3 \cdot 2 \cdot 1 = -6.$$

We find that the determinant is just the product of the diagonal entries! This is true in general, for any triangular matrix, and this gives us a very good algorithm for finding determinants.

As an example, let's calculate

$$\begin{vmatrix} 1 & 1 & -1 \\ 1 & 2 & 3 \\ 2 & 1 & 1 \end{vmatrix}$$

in two different ways: by expanding in minors and by row reduction. Expanding in minors along the first row, we have

$$\begin{aligned} \det(A) &= 1 \cdot \begin{vmatrix} 2 & 3 \\ 1 & 1 \end{vmatrix} - 1 \cdot \begin{vmatrix} 1 & 3 \\ 2 & 1 \end{vmatrix} + (-1) \cdot \begin{vmatrix} 1 & 2 \\ 2 & 1 \end{vmatrix} \\ &= 1 \cdot (2 - 3) - 1 \cdot (1 - 6) - 1 \cdot (1 - 4) \\ &= -1 + 5 + 3 \\ &= 7. \end{aligned}$$

If instead we row reduce, we make the steps

$$\begin{pmatrix} 1 & 1 & -1 \\ 1 & 2 & 3 \\ 2 & 1 & 1 \end{pmatrix} \rightsquigarrow \begin{pmatrix} 1 & 1 & -1 \\ 0 & 1 & 4 \\ 0 & -1 & 3 \end{pmatrix} \rightsquigarrow \begin{pmatrix} 1 & 1 & -1 \\ 0 & 1 & 4 \\ 0 & 0 & 7 \end{pmatrix}$$

by only adding a multiple of one row to another row, which we know does not affect the determinant at all. We've reduced to where the matrix is upper-triangular, with diagonal entries that multiply to 7, so we find that this method also gives us $\det(A) = 7$.

10.4. What's the point of determinants? Here's two answers to that question, one of which will be genuinely useful in this course and one of which I am just mentioning for completeness' sake. The useful statement is that *the absolute value of the determinant is the factor by which a linear transformation changes the area of a region*. More specifically, if $R \subset \mathbb{R}^2$ is any region to which we can assign a reasonable notion of area, and $\mathbf{T} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is a linear transformation with associated matrix A , then we have the formula

$$\text{area}(\mathbf{T}(R)) = |\det(A)| \cdot \text{area}(R).$$

Similarly, if S is a region in \mathbb{R}^3 , and $\mathbf{T} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$, then

$$\text{vol}(\mathbf{T}(R)) = |\det(A)| \cdot \text{vol}(R).$$

Similar formulas hold in higher dimensions. This is actually very useful: it tells us that in order to figure out how the area of a region changes under a linear transformation, we don't need to know anything about the geometry of that region. All we need to do is calculate a determinant.

For example, let's say we are asked to calculate the area of $\mathbf{T}(R)$, where R is the unit disc (the locus of points $(x, y) \in \mathbb{R}^2$ satisfying $x^2 + y^2 \leq 1$) and \mathbf{T} is the linear transformation

$$\mathbf{T} \left(\begin{pmatrix} x \\ y \end{pmatrix} \right) = \begin{pmatrix} x + 3y \\ 2x - 5y \end{pmatrix}.$$

The region $\mathbf{T}(R)$ is going to be an ellipse whose axes are not parallel to the axes of the plane; we could find its area directly but it will be difficult. Much easier: we notice that the matrix associated to \mathbf{T} is

$$A = \begin{pmatrix} 1 & 3 \\ 2 & -5 \end{pmatrix},$$

which we have previously calculated to have determinant -11 . The area of the unit disc is π . Therefore

$$\text{area}(\mathbf{T}(R)) = |-11| \cdot \pi = 11\pi.$$

Adopting a bigger-picture viewpoint for the moment, the other "point" of determinants is the following. In the textbook, it is stated that the determinant is *alternating* and *multilinear* as a function of the rows. In fact, the determinant is the unique operation on square matrices that is alternating, multilinear, and such that $\det(I_n) = 1$. This, more or less, is how one defines the determinant in a more abstract setting, and at least suggests why the determinant might be useful mathematically.

10.5. Systems of coordinates. For the moment, let's consider \mathbb{R}^3 for concreteness. Recall that we have defined the standard basis vectors

$$\mathbf{e}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \quad \mathbf{e}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \quad \text{and} \quad \mathbf{e}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}.$$

When we write a vector in coordinates, such as

$$\mathbf{w} = \begin{pmatrix} c_1 \\ c_2 \\ c_3 \end{pmatrix},$$

this is clearly equivalent to writing

$$\mathbf{w} = c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 + c_3 \mathbf{e}_3.$$

On the general principle that anything we do with one basis ought to be possible with another, let's pick another basis $\mathcal{B} = \{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\}$. Since this is a basis, we can uniquely write

$$\mathbf{w} = d_1 \mathbf{v}_1 + d_2 \mathbf{v}_2 + d_3 \mathbf{v}_3.$$

We suggestively write this as

$$[\mathbf{w}]_{\mathcal{B}} = \begin{pmatrix} d_1 \\ d_2 \\ d_3 \end{pmatrix}.$$

What's the point here? Sometimes it is useful to use a basis other than the standard basis, and when we do we find that essentially all of the linear algebra we have developed adapts with minimal changes to the new setting of "vectors (and matrices) written with respect to \mathcal{B} ." For example, we can talk about spans and linear subspaces and bases, do row reduction, calculate inverses, and so on.³ To go back and forth between these two worlds, write a matrix $C_{\mathcal{B}}$ whose columns are \mathbf{v}_1 , \mathbf{v}_2 , and \mathbf{v}_3 . Then it is easy to check that

$$\mathbf{w} = C_{\mathcal{B}}[\mathbf{w}]_{\mathcal{B}}.$$

Equivalently, as $C_{\mathcal{B}}$ is always invertible (why?) we can write

$$[\mathbf{w}]_{\mathcal{B}} = C_{\mathcal{B}}^{-1} \mathbf{w}.$$

These are called the *change of basis formulas*.

Of course, this theory is applicable in any number of dimensions.

11. FEBRUARY 10

11.1. **Warm-up question.** Let \mathcal{B} be the basis $\left\{ \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 2 \end{pmatrix} \right\}$, and let $\mathbf{w} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$. Find $[\mathbf{w}]_{\mathcal{B}}$.

11.2. **Warm-up answer.** By definition of the notation $[\mathbf{w}]_{\mathcal{B}}$, we are looking for the constants c_1 and c_2 such that

$$\begin{pmatrix} 2 \\ 2 \end{pmatrix} = c_1 \begin{pmatrix} 2 \\ 1 \end{pmatrix} + c_2 \begin{pmatrix} 1 \\ 2 \end{pmatrix}.$$

Solving this system of equations yields $c_1 = c_2 = \frac{2}{3}$, so

$$[\mathbf{w}]_{\mathcal{B}} = \begin{pmatrix} \frac{2}{3} \\ \frac{2}{3} \end{pmatrix}.$$

Alternatively (but equivalently), we can use the formula

$$\mathbf{w} = C_{\mathcal{B}}[\mathbf{w}]_{\mathcal{B}},$$

³It is important to note that there is a very big caveat here: the definition of the dot product is rather special to the standard basis, or more generally to any so-called *orthogonal basis*.

where $C_{\mathcal{B}}$ is the change-of-basis matrix whose columns are the elements of \mathcal{B} ; that is,

$$C_{\mathcal{B}} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}.$$

Since we wish to solve for $[\mathbf{w}]_{\mathcal{B}}$, we multiply both sides by $C_{\mathcal{B}}^{-1}$ to get

$$C_{\mathcal{B}}^{-1}\mathbf{w} = C_{\mathcal{B}}^{-1}C_{\mathcal{B}}[\mathbf{w}]_{\mathcal{B}} = [\mathbf{w}]_{\mathcal{B}}.$$

So finding the inverse of $C_{\mathcal{B}}$ (by row reduction, say; I will skip the details) and applying it to \mathbf{w} will give us the same answer:

$$[\mathbf{w}]_{\mathcal{B}} = \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} \end{pmatrix} \begin{pmatrix} 2 \\ 2 \end{pmatrix} = \begin{pmatrix} \frac{2}{3} \\ \frac{2}{3} \end{pmatrix}.$$

11.3. Eigenvalues and eigenvectors. To motivate, and geometrically explain, the concepts of eigenvalues and eigenvectors, consider the following special case. Recall that diagonal matrices D , which are matrices with arbitrary elements (say $(\lambda_1, \lambda_2, \dots, \lambda_n)$) along the diagonal and zeroes everywhere else, have a very simple associated linear transformation: along the x_i coordinate axis, the matrix scales by the factor λ_i . In symbols, recalling that we write \mathbf{e}_i for the i th standard basis vector, we have

$$D\mathbf{e}_i = \lambda_i\mathbf{e}_i$$

for each coordinate axis, labelled by i . The matrix D thus acts particularly simply on the standard basis vectors: it just multiplies them by various constants.

This example is special to the standard basis, but we can generalize it as follows: for an $n \times n$ matrix A , we say that a scalar λ is an *eigenvalue* and a nonzero vector \mathbf{v} is an *eigenvector* if the equation

$$A\mathbf{v} = \lambda\mathbf{v}$$

holds. In words, \mathbf{v} is a special direction along which A acts very simply, just by multiplying by the constant λ . If we can find lots of eigenvectors (enough to form a basis of \mathbb{R}^n), then we will get a similar geometric description of A as we did for the diagonal matrix D above - the only difference is that instead of the standard basis, we will have a new basis (called an *eigenbasis*). We will see that, unfortunately, eigenbases do not always exist. When they do, we call the matrix *diagonalizable*.

11.4. How to calculate eigenvalues and eigenvectors. The above discussion does not tell us much about how to actually find eigenvalues and eigenvectors in practice. To get an algorithm, let's first rearrange the defining equation $A\mathbf{v} = \lambda\mathbf{v}$. Notice that applying the scalar λ to the vector \mathbf{v} is the same as applying the multiple of the identity matrix λI_n to the vector. Therefore it is equivalent to write $A\mathbf{v} = \lambda I_n\mathbf{v}$. Putting everything on one side and pulling out \mathbf{v} , we are left with the equivalent equation

$$(A - \lambda I_n)\mathbf{v} = \mathbf{0}.$$

This is more promising! We are looking for nonzero \mathbf{v} satisfying this equation, which is the same as saying that we are looking for nonzero \mathbf{v} in the null space $N(A - \lambda I_n)$.

This is well and good, but doing row reduction (which is our general plan of attack for finding null spaces) is a pain in the rear if we have to carry around the unknown quantity λ with us. Looking for λ such that the null space $N(A - \lambda I_n)$ is nontrivial is the same as looking for λ such that $A - \lambda I_n$ is not invertible (by

the rank-nullity theorem, say). And we know that a matrix is not invertible if and only if its determinant is zero. So to find all the eigenvalues λ it suffices to solve the equation

$$\det(A - \lambda I_n) = 0$$

for λ .

The function $\det(A - \lambda I_n)$ is always a polynomial in λ , and it is important enough to get its own name: the *characteristic polynomial*. We will call it $p_A(\lambda)$. Its roots are the eigenvalues. Once we have found the eigenvalues by finding the roots of this polynomial, we can plug them back in one by one to the equation $(A - \lambda I_n)\mathbf{v} = \mathbf{0}$ to find the possible eigenvectors associated to each eigenvalue. The set of eigenvectors associated to an eigenvalue is a null space minus the origin (because we do not allow $\mathbf{0}$ to be an eigenvector), so it is a linear subspace (minus the origin). This linear subspace is called the *eigenspace* associated to an eigenvalue. It will usually be a line (i.e., the eigenvector will be unique up to a constant multiple), but it sometimes has larger dimension.

11.5. Examples. Let's find the eigenvalues and eigenvectors associated to the matrix

$$A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}.$$

First we have to calculate and factor $p_A(\lambda) = \det(A - \lambda I_2)$, the characteristic polynomial. We have

$$p_A(\lambda) = \begin{vmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{vmatrix} = (2 - \lambda)^2 - 1 = \lambda^2 - 4\lambda + 3 = (\lambda - 1)(\lambda - 3).$$

Therefore the eigenvalues are 1 and 3. For $\lambda = 1$, we have

$$A - \lambda I_2 = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix},$$

which has null space spanned by the vector $\mathbf{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$, say. Thus the eigenvectors associated to the eigenvalue 1 are all nonzero multiples of \mathbf{v}_1 . For $\lambda = 3$, we have

$$A - \lambda I_2 = \begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix},$$

which has null space spanned by the vector $\mathbf{v}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$. Therefore the eigenvectors associated to the eigenvalue 3 are all nonzero multiples of \mathbf{v}_2 .

Note that we have found an eigenbasis for the matrix A : we can take

$$\left\{ \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\},$$

for instance. In fact, eigenvectors associated to distinct eigenvalues are always linearly independent (see Proposition 23.3 in the text for a proof), so if we have n distinct eigenvalues for an $n \times n$ matrix we can always find an eigenbasis (that is, such a matrix is always diagonalizable). The converse is false: if we don't have n distinct eigenvalues, we may still be able to find an eigenbasis. See below.

As a second example, let's make sure this machinery works with our original example of a diagonal matrix. Let A be the matrix with $(4, 5, -1)$ along the diagonal.

Then $A - \lambda I_3$ is a diagonal matrix too, so its determinant is easy to calculate: it's just the product of the diagonal elements. Therefore

$$p_A(\lambda) = (4 - \lambda)(5 - \lambda)(-1 - \lambda).$$

The roots of this polynomial are easy to read off; they're 4, 5, and -1 . Thus we find that the eigenvalues of a diagonal matrix are precisely the diagonal entries (this also works for an upper or lower triangular matrix!).

To find the eigenvectors, let's take $\lambda = 4$ first. Then

$$A - \lambda I_3 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & -1 \end{pmatrix}.$$

This matrix has one free variable, so its null space is a line; a little thought or scratch work shows that the line in question is the one spanned by \mathbf{e}_1 , the first standard basis vector. Similarly, we find that the eigenspace associated to $\lambda = 5$ is spanned by \mathbf{e}_2 , and the eigenspace associated to $\lambda = -1$ is spanned by \mathbf{e}_3 .

So we see that our original motivation checks out: for our example matrix (and for diagonal matrices in general, as it is not hard to show), the eigenvalues are the diagonal entries and the eigenvectors are along the coordinate axes. That is, the coordinate axes get stretched by a factor of the corresponding diagonal entry. Because the coordinate axes form a basis, we find that there always exists an eigenbasis for a diagonal matrix: the standard basis! In particular, diagonal matrices are always diagonalizable (as one might expect from the name!). This is true even though some diagonal matrices, such as the identity matrix, may have "repeated" eigenvalues, so we cannot apply Proposition 23.3.

11.6. What can go wrong? Here are two "problems" with the concepts of eigenvalues and eigenvectors. First, consider the matrix

$$A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}.$$

If we calculate its characteristic polynomial, we find $p_A(\lambda) = \lambda^2 + 1$, which does not have any real roots. Hence A does not have any real eigenvalues. Geometrically this makes sense: A is a rotation matrix, and no nonzero vector in \mathbb{R}^2 is scaled by a rotation: a rotation (unless it is a rotation by an integer multiple of π radians) always changes the direction of a vector unless it is the zero vector. To deal with this issue, we could simply allow complex eigenvalues and complex eigenvectors, which effectively causes us to replace the real numbers \mathbb{R} with the complex numbers \mathbb{C} in all our vector fields. This seems drastic, but it's the usual way of dealing with the problem: we really like it when we find eigenbases, even if they are complex! Alternatively, we could refuse to do this, and simply say that such a matrix does not have real eigenvalues and hence is not (real) diagonalizable. This is the approach we will take in this course.

The second issue is actually much more mathematically interesting: there are matrices which are not diagonalizable even if you allow complex numbers. So far we have not seen an example of this, but it's not hard to cook one up. Consider the matrix

$$A = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}.$$

It has only one eigenvalue, 1, with multiplicity two. The corresponding eigenspace is the null space of the matrix

$$A - I_2 = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix},$$

which is one-dimensional, spanned by the standard basis vector \mathbf{e}_1 . So there is no way we can ever find an eigenbasis of A : the only eigenvectors lie along the x -axis, and thus can never span all of \mathbb{R}^2 .

There are ways of dealing with this second issue; if you're curious, the key words are *Jordan normal form* and *generalized eigenvector*. A more in-depth linear algebra class would cover these concepts.

12. FEBRUARY 12

Instead of talking about new material (quadratic forms) today, by popular request we went over some more conceptual questions related to the material of the last couple of weeks.

12.1. Some conceptual questions.

- (1) Let A be an $n \times n$ matrix such that A^2 is the zero matrix. What can we conclude: that $I_n - A$ is always invertible, that $I_n + A$ is never invertible, or that either case is possible?
- (2) Which of the following linear transformations have diagonalizable matrices over \mathbb{R} ?
 - (a) The identity.
 - (b) Projection onto the line $y = x$ in \mathbb{R}^2 .
 - (c) Rotation clockwise by $\pi/2$ radians.
 - (d) $\mathbf{f}(x_1, x_2) = (x_1 + x_2, x_2)$.
 - (e) $\mathbf{f}(x_1, x_2) = (2x_1 - 3x_2, -3x_1 + x_2)$.
- (3) Suppose $P^2 = P$ (geometrically, this means that P is a projection matrix. Think about why projection matrices always obey this relation). True or false?
 - (a) -1 is a possible eigenvalue of P .
 - (b) P is not invertible.
 - (c) P is diagonalizable.

12.2. Answers.

- (1) There are at least two ways of doing this problem: one using a bit of cleverness, and one using eigenvalues. By using a bit of cleverness and the fact that A^2 is the zero matrix, we see that

$$(I_n - A)(I_n + A) = I_n + I_n A - I_n A - A^2 = I_n,$$

so by definition the matrix $I_n + A$ is an inverse to $I_n - A$. Thus $I_n - A$ is *always invertible in this situation*.

If we didn't see this trick, what could we do? First, let's see what we can say about the eigenvalues of A . We saw on a homework assignment that the eigenvalues of a power of a matrix are just the corresponding powers of the eigenvalues. So if $\lambda_1, \lambda_2, \dots, \lambda_n$ are the eigenvalues of A , then $\lambda_1^2, \lambda_2^2, \dots, \lambda_n^2$ are the eigenvalues of A^2 . But A^2 is the zero matrix, which has the eigenvalue zero (with multiplicity n), so $\lambda_1 = \lambda_2 = \dots = \lambda_n = 0$.

That is, we can conclude that the eigenvalues of A are also zero. Note that this does not imply that A is the zero matrix (just as A^2 being the zero matrix does not imply that A is the zero matrix). For example, A could be the matrix

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

Moving on, knowing that A has all zero eigenvalues, let's see what we can say about the eigenvalues of $I_n - A$. If λ' is an eigenvalue of $I_n - A$, with corresponding eigenvector \mathbf{v} , then we have

$$(I_n - A)\mathbf{v} = \lambda'\mathbf{v}.$$

Upon rearranging and noticing that $I_n\mathbf{v} = \mathbf{v}$,

$$A\mathbf{v} = (1 - \lambda')\mathbf{v}.$$

But we know that A has all eigenvalues equal to zero, so $1 - \lambda' = 0$, so $\lambda' = 1$. Therefore all eigenvalues of $I_n - A$ are equal to one. A matrix is invertible if and only if all its eigenvalues are nonzero, so we conclude that $I_n - A$ is always invertible.

Notice that on our last step here we calculated that because I_n has all eigenvalues equal to 1 and A has all eigenvalues equal to 0, the matrix $I_n - A$ has all eigenvalues equal to $1 - 0 = 1$. This logic only works when the matrices share eigenvectors! For example, in our case all of the eigenvectors of A , whatever they are, are also eigenvectors of I_n (because all nonzero vectors are eigenvectors of the identity matrix); this allowed our calculation to go through. It is absolutely not true in general that the eigenvalues of $A + B$ are the sums of the eigenvalues of A and the eigenvalues of B .

One final remark: the statement remains true if you replace A^2 with A^N for any positive integer N . The second proof is largely the same (using the same homework problem about eigenvalues of powers of matrices), and the first proof method also works if you notice that

$$(I_n - A)(I_n + A + A^2 + \dots + A^{N-1}) = I_n$$

if A^N is the zero matrix (you get a lot of cancellation).

- (2) Most of these can be dealt with by “pure thought” (i.e., no calculation).
- (a) The identity is a diagonal matrix, so it is automatically diagonalizable.
 - (b) The projection onto the line $y = x$ is given by the matrix

$$\begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

(just see where the two standard basis vectors go). We can calculate that the eigenvectors of this matrix are 0 and 1, which are distinct, so by Proposition 23.3 the matrix is diagonalizable. Alternatively, by the answer to the third question, projection matrices are always diagonalizable.

(c) We could write down the matrix, find the characteristic polynomial, and so on, but it is easiest to notice that a nontrivial rotation in the plane will *never* have real eigenvalues (unless it is a rotation by an integer multiple of π), because every nonzero vector will be carried to a vector pointing in a different direction. Therefore this matrix is not (real) diagonalizable.

If we wanted to work with complex numbers, we would find that it *is* diagonalizable over the complex numbers.

(d) The corresponding matrix is

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

which has eigenvalue 1 with multiplicity two (just look at the diagonal entries because it's a triangular matrix!). The corresponding eigenspace is the nullspace of the matrix

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

which is only one-dimensional. Therefore there can be no basis of \mathbb{R}^2 consisting only of eigenvectors of this matrix, and it is therefore not diagonalizable.

(e) The corresponding matrix is

$$\begin{pmatrix} 2 & -3 \\ -3 & 1 \end{pmatrix},$$

which is symmetric. By the spectral theorem (Proposition 25.2 in the text), it is therefore diagonalizable.

(3) Part (c) of this problem is probably the hardest of all of these.

(a) We can again use the fact, proven on the homework, that the eigenvalues of a power of a matrix are the powers of the eigenvalues. Let $\lambda_1, \dots, \lambda_n$ be the eigenvalues of P . Then $\lambda_1^2, \dots, \lambda_n^2$ are the eigenvalues of P^2 , corresponding to the same eigenvectors as previously. But since $P^2 = P$, we must have $\lambda_1 = \lambda_1^2, \dots, \lambda_n = \lambda_n^2$. Each of these equations is solved only by the values 0 and 1, so each eigenvalue must be 0 or 1. Therefore, it is *true* that -1 is not a possible eigenvalue of P .

(b) The identity matrix I_n satisfies $I_n^2 = I_n$, and the identity matrix is invertible, so it is *false* that P is always not invertible. (It turns out that this is the only invertible example of such a P .)

(c) There are a number of ways of thinking about this question; here is a particularly slick way. We have already shown that the eigenvalues of P are all 0 or 1. First note that the eigenspace associated to the eigenvalue 0 is precisely the kernel of P (i.e., the null space of the corresponding matrix), as it is the set of \mathbf{v} such that $(P - 0 \cdot I_n)\mathbf{v} = P\mathbf{v} = \mathbf{0}$.

The eigenspace associated to the eigenvalue 1 is the set of \mathbf{v} such that $(P - I_n)\mathbf{v} = \mathbf{0}$ (or $P\mathbf{v} = \mathbf{v}$); that is, the set of vectors that are unchanged by the action of P . I claim that this set is exactly the image of P (i.e., the column space of the corresponding matrix). To show this in one direction, assume that $P\mathbf{v} = \mathbf{v}$. Then \mathbf{v} is certainly in the image of P . Conversely, if \mathbf{v} is in the image of P , then there is a \mathbf{w} such that $P\mathbf{w} = \mathbf{v}$. Applying P to both sides, we find that $P^2\mathbf{w} = P\mathbf{v}$, so as $P^2 = P$ we have $P\mathbf{w} = P\mathbf{v}$. Therefore if such a \mathbf{w} exists, then $P\mathbf{v} = \mathbf{v}$ as well. This verifies the claim.

By the rank-nullity theorem, the dimension of the kernel of P plus the dimension of the image of P is equal to the total dimension n . Therefore by picking any basis of the kernel of P and any basis of the image of P and

putting them together, we have an eigenbasis. Therefore the claim that P is diagonalizable is *true*.

13. FEBRUARY 17

13.1. **Warm-up question.** Determine the “definiteness” of the quadratic forms associated to the following matrices:

$$(a) \begin{pmatrix} 1 & 2 \\ 2 & -1 \end{pmatrix}.$$

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

See if you can avoid calculating any eigenvalues!

13.2. **Warm-up answer.** If all else fails, we can calculate the eigenvalues of these matrices, which will then tell us the definiteness of the corresponding quadratic form (for example, all positive eigenvalues implies positive definite, some positive and some negative implies indefinite, and so on).

I claim there’s a fast way in this case, though: imagine plugging in the standard basis vectors into the quadratic form. For part (a), we have

$$Q(\mathbf{e}_1) = (1 \ 0) \begin{pmatrix} 1 & 2 \\ 2 & -1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = (1 \ 0) \begin{pmatrix} 1 \\ 2 \end{pmatrix} = 1$$

and

$$Q(\mathbf{e}_2) = (0 \ 1) \begin{pmatrix} 1 & 2 \\ 2 & -1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = (0 \ 1) \begin{pmatrix} 2 \\ -1 \end{pmatrix} = -1.$$

We have found two nonzero vectors that evaluate to different signs, so the quadratic form in (a) is indefinite. Notice what we are calculating: the quadratic form evaluated at \mathbf{e}_1 is the first row of the first column of the matrix, and the quadratic form evaluated at \mathbf{e}_2 is the second row of the second column. This clearly generalizes to any quadratic form. We conclude that *if the diagonal of a symmetric matrix contains both positive and negative elements, the corresponding quadratic form must be indefinite*. Thus the quadratic form in (b) is indefinite as well.

If this is not clear, try writing out the quadratic form explicitly. For part (a), we have

$$Q(x_1, x_2) = x_1^2 + 4x_1x_2 - x_2^2.$$

If we plug in a vector with x_2 set to zero, then the only surviving term is the first term. Thus evaluating on (say) the standard vector $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ will only involve the term x_1^2 , which is always positive except at $x_1 = 0$. Likewise, evaluating on a vector with $x_1 = 0$ will only involve the term $-x_2^2$, which is always negative except at $x_2 = 0$.

A word of caution: there are indefinite matrices that are not of this form! For example, the form associated to the matrix

$$\begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}$$

is indefinite, because if we plug in the vector $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ we get the value 6 while if we plug in the vector $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$ we get the value -2 . To figure out the definiteness of this matrix, the easiest method is almost certainly just to find the eigenvalues.

13.3. One last problem in linear algebra. Here's a nice hard exam-style question: prove that if A_1 and A_2 are $n \times n$ symmetric matrices with only positive eigenvalues, then $A_1 + A_2$ is as well.

A solution is as follows: we have to show two things; first, that $A_1 + A_2$ is symmetric, and second that it has only positive eigenvalues. The first task is easier: we have

$$(A_1 + A_2)^T = A_1^T + A_2^T = A_1 + A_2,$$

where we used that A_1 and A_2 are both symmetric and therefore equal to their transposes. For the second task, we recall what we know about symmetric matrices with only positive eigenvalues: they are associated to positive-definite quadratic forms. So we have two quadratic positive-definite forms Q_{A_1} and Q_{A_2} . By the very definition of positive-definite, these functions are positive except at zero, so their sum $Q_{A_1} + Q_{A_2}$ is certainly positive except at zero as well. But this sum is the quadratic form associated to the symmetric matrix $A_1 + A_2$:

$$(Q_{A_1} + Q_{A_2})(\mathbf{x}) = \mathbf{x}^T A_1 \mathbf{x} + \mathbf{x}^T A_2 \mathbf{x} = \mathbf{x}^T (A_1 + A_2) \mathbf{x}.$$

So we conclude that $A_1 + A_2$ is associated to a positive-definite quadratic form, hence has only positive eigenvalues, which is what we wanted to show.

13.4. Some vocabulary. So far we have discussed a very special type of function from \mathbb{R}^n to \mathbb{R}^m , the linear transformations. From here on out, we will be discussing much more general functions $\mathbf{f} : D \rightarrow \mathbb{R}^m$, where D is a subset of \mathbb{R}^n (we want to allow functions that might not be defined everywhere on \mathbb{R}^n). The basic, essential vocabulary is as follows: the *domain* of such a function is the set D , the *codomain* is \mathbb{R}^m , and the *range* (or image) is the subset of the codomain that actually gets hit by the function. More precisely, the range of \mathbf{f} is the set

$$\{\mathbf{y} \in \mathbb{R}^m : \text{there exists } \mathbf{x} \in D \text{ such that } \mathbf{f}(\mathbf{x}) = \mathbf{y}\}.$$

As an example, let D be the real numbers minus the origin and $f : D \rightarrow \mathbb{R}$ be the function $f(x) = \frac{1}{x^2}$. Then the domain is D , the codomain is \mathbb{R} , and the range is the set of all positive numbers, because these are the numbers that can be written in the form $\frac{1}{x^2}$ for some x in D .

13.5. Graphs of functions. If we have a function $\mathbf{f} : D \rightarrow \mathbb{R}^m$, the *graph* of the function is the subset of the product $D \times \mathbb{R}^m$ (the set of all ordered pairs of points consisting of one point in D and one in \mathbb{R}^m) consisting of all of the ordered pairs $(\mathbf{x}, \mathbf{f}(\mathbf{x}))$. In set notation,

$$\text{graph}(f) = \{(\mathbf{x}, \mathbf{f}(\mathbf{x})) : \mathbf{x} \in D\} \subset D \times \mathbb{R}^m.$$

If f happens to be a function from $D \rightarrow \mathbb{R}$, where D is a subset of the real numbers, then we recover the usual notion of the graph of a real-valued function on the real line. We can also visualize graphs of real-valued functions of two real variables and \mathbb{R}^2 -valued functions of one real variable, because these graphs will lie in \mathbb{R}^3 . For example, the graph of the function $f(x, y) = e^{-x^2 - y^2}$ looks like a very symmetric

mountain centered at the origin, and the graph of the function $\mathbf{f}(x) = (\cos x, \sin x)$ looks like a helix.

As an illustration of these definitions, let's prove the following proposition: the graph of a linear transformation is a linear subspace. To this end, let $\mathbf{T} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a linear transformation. We need to check that $\text{graph}(\mathbf{T})$ contains the zero vector, is closed under addition, and is closed under scalar multiplication.

We have

$$\text{graph}(\mathbf{T}) = \{(\mathbf{x}, \mathbf{T}(\mathbf{x})) : \mathbf{x} \in \mathbb{R}^n\}.$$

It is clear that $(\mathbf{0}, \mathbf{0})$ lies in $\text{graph}(\mathbf{T})$ because $\mathbf{T}(\mathbf{0}) = \mathbf{0}$. To check addition, let \mathbf{a} and \mathbf{b} be two points in $\text{graph}(\mathbf{T})$. They are of the form $(\mathbf{x}_1, \mathbf{T}(\mathbf{x}_1))$ and $(\mathbf{x}_2, \mathbf{T}(\mathbf{x}_2))$ for some points \mathbf{x}_1 and \mathbf{x}_2 in \mathbb{R}^n . Their sum is

$$\mathbf{a} + \mathbf{b} = (\mathbf{x}_1 + \mathbf{x}_2, \mathbf{T}(\mathbf{x}_1) + \mathbf{T}(\mathbf{x}_2)).$$

But \mathbf{T} is a linear transformation, so $\mathbf{T}(\mathbf{x}_1) + \mathbf{T}(\mathbf{x}_2) = \mathbf{T}(\mathbf{x}_1 + \mathbf{x}_2)$, and therefore $\mathbf{a} + \mathbf{b}$ is in the graph of \mathbf{T} , equal as it is to the point $(\mathbf{x}_1 + \mathbf{x}_2, \mathbf{T}(\mathbf{x}_1 + \mathbf{x}_2))$. Closure under scalar multiplication is proven in the same manner: if we have some \mathbf{a} in the graph of \mathbf{T} , of the form $(\mathbf{x}_1, \mathbf{T}(\mathbf{x}_1))$, then $c\mathbf{a}$ is also in the graph of \mathbf{T} because it is equal to $(c\mathbf{x}_1, \mathbf{T}(c\mathbf{x}_1))$, using the fact that $c\mathbf{T}(\mathbf{x}_1) = \mathbf{T}(c\mathbf{x}_1)$.

In fact, this whole argument is reversible (check!): if the graph of a function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear subspace, then \mathbf{f} is a linear transformation. So we could have *defined* a linear transformation to be a function $\mathbb{R}^n \rightarrow \mathbb{R}^m$ whose graph is a linear subspace (though it would have been a little perverse to do so).

14. FEBRUARY 19

14.1. Warm-up question. Here is an open-ended question: what can we say about the contour map of a positive-definite quadratic form in two variables?

14.2. Warm-up answer. It's an open-ended question, so here's an open-ended answer. Let's start with the simplest possible positive-definite quadratic form, just to get our bearings: if $Q(x, y) = x^2 + y^2$, then the contours are the loci defined by

$$x^2 + y^2 = c,$$

where c varies over the real numbers. These are clearly all circles centered at the origin, so the contour plot of Q consists of concentric circles around the origin.

Now let's generalize a little bit and consider a diagonal positive-definite quadratic form; i.e., one of the form

$$Q'(x, y) = ax^2 + by^2$$

where a and b are both positive numbers. Then the contours are of the form

$$ax^2 + by^2 = c,$$

which are all ellipses centered at the origin with major and minor axes aligned with the coordinate axes.

Now we can generalize all the way! Since every quadratic form is diagonalizable by the spectral theorem, there exists a basis of \mathbb{R}^2 in which the quadratic form is just the above example. In the original basis, our ellipses are now tilted and stretched (recall your homework problems about ellipses with respect to various bases). So the contour plot of a general positive-definite quadratic form in two variables consists of concentric ellipses around the origin.

If you're happy with this example, think about two possible extensions: what would the contour plot of an indefinite quadratic form in two variables look like? And what would the "contour plot" of a positive-definite quadratic form in three variables look like (the contours are now surfaces in \mathbb{R}^3)?

14.3. Contour plot example. Let's do one more example of a contour plot. How would we draw the contour plot of the function $f(x, y) = y^2 - x$? As usual, we know that the contours are given by the curves

$$y^2 - x = c,$$

where c varies over all real numbers. So we just have to try a few values of c and plot the resulting curves.

Let's first try $c = 0$, because it seems like the simplest. Then the contour is $x = y^2$, which is a parabola opening up to the right. If we didn't know what this looked like, we could plug in a few points: the locus of points where $x = y^2$ passes through the points $(4, -2)$, $(1, -1)$, $(0, 0)$, $(1, 1)$, and $(4, 2)$, for instance, and it is clear how to interpolate and extrapolate from here.

Now let's try $c = 1$. The contour is now $y^2 - x = 1$, which we can write as $x = y^2 - 1$. This is still a parabola opening to the right; it is the same as the above parabola except "shifted by one." We can see exactly how by trying some points: the curve passes through $(3, -2)$, $(0, -1)$, $(-1, 1)$, $(0, 1)$, and $(3, 2)$. So this is curve shifted one unit to the left from the above.

It is clear that all of these contours will be translates of the original parabola, shifted to the right and left. In drawing the contour plot, it is helpful to label the heights of the contours; for example, $x = y^2$ might be labelled by " $f = 0$ " or something similar, because it is the contour at height 0. It is a worthwhile exercise to visualize the graph of f : it slopes down to the right with parabolas as level curves, perhaps sort of like a very large slide.

14.4. Parametrized curves. A *parametrized curve* is nothing more or less than a (let's say continuous for simplicity) function $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^m$. We can visualize a parametrized curve with a physical metaphor: let the coordinate of the domain be t , and for each t think of $f(t)$ as the position of a particle at time t . Thus a parametrized curve describes the time evolution of a particle. The image of f will be some curve in \mathbb{R}^m , but keep in mind that a curve in \mathbb{R}^m will have many possible parametrizations! The particle may be moving along the curve at any speed at any point in time, or it may even be doubling back and so on. When we say parametrized curve, we mean the data of the function f , not just the data of the image curve.

As a first basic example, recall that we already know how to work certain parametrizations of lines: we just pick a point \mathbf{x}_0 on the line, a vector \mathbf{v} pointing in the direction of the line, and note that

$$L = \{\mathbf{x}_0 + t\mathbf{v} : t \in \mathbb{R}\}.$$

This parametrization corresponds to the function

$$\mathbf{f}(t) = \mathbf{x}_0 + t\mathbf{v}.$$

Of course, there are many possible parametrizations of a line of this form (and even more that are not of this form, corresponding to particles that are not traveling at a constant speed).

As a second example, consider the parametrized curve defined by

$$\mathbf{f}(t) = (\cos t, \sin t).$$

As t varies, this describes the unit circle in \mathbb{R}^2 , traversed at a constant rate in the counterclockwise direction.

We can analyze the geometry of parametrized curves by using calculus. Let's make the definition that $\mathbf{f}'(t)$ denotes the function $\mathbb{R} \rightarrow \mathbb{R}^m$ that you get when you differentiate each component separately in t . (This definition will later be superseded when we define the derivative of an arbitrary differentiable function $\mathbb{R}^n \rightarrow \mathbb{R}^m$.) In our particle metaphor, $\mathbf{f}'(t)$ corresponds to the *velocity* of the particle at the point t , and is therefore sometimes written $\mathbf{v}(t)$. In our second example,

$$\mathbf{f}'(t) = \left(\frac{d}{dt} \cos t, \frac{d}{dt} \sin t \right) = (-\sin t, \cos t).$$

Geometrically, if we pick a time t and place the vector $\mathbf{f}'(t)$ starting at the point $\mathbf{f}(t)$, we will get a *tangent vector* to the curve. The length of this vector represents how fast the particle is going: it is called the *speed*. In the same example, we can calculate the speed as a function of time:

$$\|\mathbf{f}'(t)\| = \sqrt{(-\sin t)^2 + (\cos t)^2} = 1,$$

where we have used the all-important trigonometric identity $(\sin t)^2 + (\cos t)^2 = 1$. We can interpret this as telling us that our particle is traveling at a constant speed for all time. Note that the velocity is not constant, even though its length is, because it is continuously changing direction. In general, the vector

$$\frac{\mathbf{f}'(t)}{\|\mathbf{f}'(t)\|}$$

is called the *unit tangent vector*. It is the vector tangent to the curve with unit length. In our example, the velocity vector is already a unit tangent vector.

We can carry this further, if we like. The *acceleration* is defined to be the componentwise derivative (in t) of the velocity. We can write it as $\mathbf{f}''(t)$ or $\mathbf{a}(t)$. In our example, we calculate

$$\mathbf{f}''(t) = \left(\frac{d}{dt}(-\sin t), \frac{d}{dt} \cos t \right) = (-\cos t, -\sin t).$$

Again, this has a geometric interpretation: if we place the vector $\mathbf{f}''(t)$ starting at the point $\mathbf{f}(t)$, then it points in the direction the curve is turning, and its length (essentially) represents how quickly the curve is turning. In our example, the acceleration vector always points straight towards the center of the circle, showing us that the curve is constantly bending in that direction.

One side note: although $\mathbf{f}'(t)$ and $\mathbf{f}''(t)$ can be viewed as maps $\mathbb{R} \rightarrow \mathbb{R}^m$, and therefore as parametrized curves in their own right, this is not a particularly useful way of viewing them. Rather, we should view them the way we have been: as a vector that is somehow “attached” to each corresponding point of the original curve. When we generalize below, we will see that the best way of viewing a derivative is as a linear transformation $\mathbb{R}^n \rightarrow \mathbb{R}^m$ attached to each point in the domain; the case of parametrized curves corresponds to the case where the domain is just (a subset of) \mathbb{R} , so the linear transformation at each point can be represented by a $n \times 1$ matrix, which can be thought of as a vector.

15. FEBRUARY 24

15.1. Warm-up question. The position of a particle at time t is $(\cos t, t^2, \sin t)$. Find the velocity of the particle at time t , the acceleration at time t , the speed at time t , and find a parametric representation of the tangent line to the path of the particle at the point $(1, 0, 0)$.

15.2. Warm-up answer. To find the velocity, we just differentiate in each coordinate, getting

$$(-\sin t, 2t, \cos t).$$

To find the speed, we take the magnitude of this vector, getting

$$\sqrt{(-\sin t)^2 + (2t)^2 + (\cos t)^2} = \sqrt{4t^2 + 1}.$$

To find the acceleration, we differentiate the velocity in each coordinate, getting

$$(-\cos t, 2, -\sin t).$$

Finally, to find a parametric representation of a line all we need is a point in the line and a vector in the direction of the line. The point here is $(1, 0, 0)$. To find the vector in the direction of the line, we just need the velocity vector at that point. To find it, we need to find out which value of t corresponds to $(1, 0, 0)$ so we can plug it into our formula. Looking at the second coordinate, we have $t^2 = 0$, so $t = 0$ is the time we need. Plugging in $t = 0$ we get $(0, 0, 1)$ as the requisite velocity vector, so the line in question is parametrized as

$$\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + s \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} : s \in \mathbb{R} \right\}.$$

Of course, this parametrization is far from unique.

15.3. The derivative. For an ordinary function $f : \mathbb{R} \rightarrow \mathbb{R}$, the derivative at a point, if it exists, should be thought of as the slope of the tangent line of the graph of f at that point. The slope is a single number, an element of \mathbb{R} . To generalize this to functions $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we can draw the graph, but now the tangent space will no longer be a line: it will be some higher-dimensional affine linear subspace. (You can visualize this for a function $\mathbb{R}^2 \rightarrow \mathbb{R}$, which if sufficiently nice will have a tangent plane at each point of its graph.) How should we represent the data of the “slope” of a pretty arbitrary affine linear subspace?

Certainly we don’t want the derivative to possess information about the value of the function at the point in question; i.e., it shouldn’t change if we change the function by a constant, just like in the one-dimensional case. So without loss of generality we can shift the affine linear subspace to the origin, so it is an honest linear subspace. Now here is the key insight: this linear subspace is the graph of some linear transformation $\mathbb{R}^n \rightarrow \mathbb{R}^m$. *We define the derivative of \mathbf{f} to be this linear transformation.*

In the case of a function $f : \mathbb{R} \rightarrow \mathbb{R}$, this is a conceptual shift but it doesn’t really change anything: a linear transformation $\mathbb{R} \rightarrow \mathbb{R}$ can be represented by a 1×1 matrix; that is, a scalar. This scalar is the derivative as we usually imagine it. Concretely, the “new derivative” at a point a is the linear transformation

$$h \mapsto \frac{df}{dx}(a) \cdot h,$$

where $\frac{df}{dx}(a)$ is the “old derivative” at a . Notice that we get a different linear transformation at each point, so the “derivative of a function” is really a collection of many linear transformations, one for each point.

In the general case, it is not difficult to trace through this reasoning keeping the standard basis in mind to find a formula for the derivative $(D\mathbf{f})(\mathbf{a})$ of a function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ at a point $\mathbf{a} \in \mathbb{R}^n$: if

$$\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})),$$

then $(D\mathbf{f})(\mathbf{a})$ is the linear transformation associated to the matrix

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{a}) & \frac{\partial f_1}{\partial x_2}(\mathbf{a}) & \dots & \frac{\partial f_1}{\partial x_n}(\mathbf{a}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{a}) & \frac{\partial f_2}{\partial x_2}(\mathbf{a}) & \dots & \frac{\partial f_2}{\partial x_n}(\mathbf{a}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{a}) & \frac{\partial f_m}{\partial x_2}(\mathbf{a}) & \dots & \frac{\partial f_m}{\partial x_n}(\mathbf{a}) \end{pmatrix}.$$

So the calculation of the derivative of a general function reduces to the calculation of a bunch of partial derivatives, which in any given instance are easy enough to calculate. To remember this formula (which way do the indices go?), it helps to remember that if \mathbf{f} is a function from \mathbb{R}^n to \mathbb{R}^m , then so is its derivative at any point, so the derivative is a $m \times n$ matrix. Therefore the function indices have to correspond to the rows, and the coordinate indices to the columns.

As an example, let's calculate $(D\mathbf{f})(\mathbf{a})$ where

$$\mathbf{f}(x, y) = (3x^2y, \cos x \sin y), \quad \mathbf{a} = (0, 0).$$

To do so, we first calculate all the first partial derivatives and stick them in a matrix as above:

$$(D\mathbf{f})(x, y) = \begin{pmatrix} \frac{\partial}{\partial x}(3x^2y) & \frac{\partial}{\partial y}(3x^2y) \\ \frac{\partial}{\partial x}(\cos x \sin y) & \frac{\partial}{\partial y}(\cos x \sin y) \end{pmatrix} = \begin{pmatrix} 6xy & 3x^2 \\ -\sin x \sin y & \cos x \cos y \end{pmatrix}.$$

Then we evaluate at $(x, y) = (0, 0)$, getting

$$(D\mathbf{f})(0, 0) = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

Geometrically, since $\mathbf{f}(0, 0) = \mathbf{0}$, the graph of $(D\mathbf{f})(0, 0)$ is precisely the tangent plane of \mathbf{f} at the origin (we don't have to shift at all). In other words, \mathbf{f} “looks like” the map $(D\mathbf{f})(0, 0)$ near the origin; we call it the *linearization* of \mathbf{f} at $\mathbf{0}$. Next time we'll calculate some linearizations in the general case, when we do have to shift the graph of the derivative.

Here's a second example: consider a linear transformation $\mathbf{T}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, given by $\mathbf{T}(\mathbf{x}) = A\mathbf{x}$ where A is an $m \times n$ matrix. What is $(D\mathbf{T})(\mathbf{a})$ for any point \mathbf{a} ?

The answer is: we have $(D\mathbf{T})(\mathbf{a}) = \mathbf{T}$! That is, the derivative of a linear transformation at any point is the linear transformation itself. Proving this is a problem on this week's homework. Let's think about why this makes sense: we know that the graph of a linear transformation is a linear subspace, so the “best linear approximation” to \mathbf{T} (at any point!) should be \mathbf{T} itself!

We will soon see more justification for defining the derivative at a point as a linear transformation: it makes the chain rule very easy to state, and it makes multivariable Taylor expansions easy to write down.

15.4. Tangent spaces of graphs. If we are interested in finding equations that cut out the tangent space to a point \mathbf{a} in the graph of a differentiable function $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, we can use the following fact, which is easy to derive from knowledge about linearizations (to be discussed next class): the tangent space is defined by the system of equations

$$\mathbf{z} = \mathbf{f}(\mathbf{a}) + (D\mathbf{f})(\mathbf{a})(\mathbf{x} - \mathbf{a}),$$

where $(\mathbf{x}, \mathbf{z}) \in \mathbb{R}^n \times \mathbb{R}^m$ are the coordinates. This is a system of m equations in $m + n$ unknowns, and hence should in general have a solution space of dimension n , which is what we want.

Specializing to the case where $n = 2$ and $m = 1$, we see that the tangent space (which is a plane in this case) is given by the single equation

$$z = f(\mathbf{a}) + (Df)(\mathbf{a})(\mathbf{x} - \mathbf{a}),$$

which simplifies to (if we let $\mathbf{x} = (x, y)$ and $\mathbf{a} = (a_1, a_2)$) the equation

$$z = f(a_1, a_2) + \frac{\partial f}{\partial x}(a_1, a_2)(x - a_1) + \frac{\partial f}{\partial y}(a_1, a_2)(y - a_2).$$

As an example, consider the function $f(x, y) = e^{x(y+1)}$ at the point $\mathbf{a} = (0, 0)$. We have

$$\frac{\partial f}{\partial x}(x, y) = (y + 1)e^{x(y+1)}, \quad \frac{\partial f}{\partial y}(x, y) = xe^{x(y+1)},$$

and evaluated at $(0, 0)$ we get

$$\frac{\partial f}{\partial x}(0, 0) = 1, \quad \frac{\partial f}{\partial y}(0, 0) = 0.$$

Furthermore we have $f(0, 0) = 1$. Therefore, after plugging in and simplifying, the equation for the tangent plane at $(0, 0)$ is given by the equation

$$z = 1 + x.$$

Of course, one can consider tangent spaces of sets that are more general than just graphs of functions. For example, the textbook contains a discussion about how to find the tangent plane to a level set of a function from \mathbb{R}^3 to \mathbb{R} . We will not cover this in the course, but it is worth reading over once.

16. FEBRUARY 26

16.1. Warm-up question. Does there exist a differentiable function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ such that $\frac{\partial f}{\partial x} = x + y$ and $\frac{\partial f}{\partial y} = 2y$? Explain.

16.2. Warm-up answer. If f were such a function, then clearly its two derivatives are themselves differentiable, and we can calculate

$$\frac{\partial^2 f}{\partial y \partial x} = 1, \quad \frac{\partial^2 f}{\partial x \partial y} = 0.$$

These are both continuous functions, so by Clairaut's theorem (mixed partials are equal if a function is twice continuously differentiable), they should be equal. Contradiction! Therefore no such function exists.

It is also solve this problem in a slightly more hands-on way (by writing down the general form that such an f could have, given the constraints on its derivatives, and deriving a contradiction that way), but citing Clairaut's theorem is much easier. In this warm-up question, we are secretly solving a system of differential equations (or

rather, showing that they have no solution); Clairaut's theorem gives a nontrivial condition on the solvability of such equations. This idea will come back if you take Math 53.

16.3. Some remarks on Clairaut's theorem. I wish there were a proof of Clairaut's theorem that I could present in a few minutes on the board, but there really isn't: there is one proof in the appendix to your textbook, which uses the mean value theorem, and there a somewhat simpler proof using basic theorems of multivariable integral calculus that are not part of Math 51. However, I think it might be useful to investigate why the proof is not so simple (and, in the process, why we need to assume that the second derivatives are continuous in order for the theorem to hold).

Let's make an attempt at proving Clairaut's theorem by writing down the definition of the second derivatives in question. Let f be a real-valued function of two variables x and y which possesses all second derivatives. Then, by definition,

$$\begin{aligned} \frac{\partial^2 f}{\partial y \partial x}(x, y) &= \lim_{k \rightarrow 0} \left[\frac{\frac{\partial f}{\partial x}(x, y+k) - \frac{\partial f}{\partial x}(x, y)}{k} \right] \\ &= \lim_{k \rightarrow 0} \left[\frac{\lim_{h \rightarrow 0} \frac{f(x+h, y+k) - f(x, y+k)}{h} - \lim_{h \rightarrow 0} \frac{f(x+h, y) - f(x, y)}{h}}{k} \right] \\ &= \lim_{k \rightarrow 0} \lim_{h \rightarrow 0} \left[\frac{f(x+h, y+k) - f(x, y+k) - f(x+h, y) + f(x, y)}{hk} \right]. \end{aligned}$$

Similarly, we can expand to calculate

$$\frac{\partial^2 f}{\partial x \partial y}(x, y) = \lim_{h \rightarrow 0} \lim_{k \rightarrow 0} \left[\frac{f(x+h, y+k) - f(x, y+k) - f(x+h, y) + f(x, y)}{hk} \right].$$

So *if we could exchange the limits, we would find that the mixed partials are equal.*

Unfortunately, we can't! It simply isn't true in this level of generality that limits can be rearranged: consider a function g of two variables, continuous except at the origin, such that when you approach the origin along the x -axis, you get one value (say, zero) and when you approach the origin along the y -axis, you get another (say, one). Then

$$\lim_{x \rightarrow 0} \lim_{y \rightarrow 0} g(x, y) = 0 \neq 1 = \lim_{y \rightarrow 0} \lim_{x \rightarrow 0} g(x, y).$$

If one is patient, one can use this idea to write down a function that is twice-differentiable (but not twice continuously differentiable!) and such that its mixed partials at the origin are not equal. The standard example of this is the function $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$f(x, y) = \begin{cases} xy \frac{x^2 - y^2}{x^2 + y^2} & \text{if } (x, y) \neq (0, 0), \\ 0 & \text{if } (x, y) = (0, 0). \end{cases}$$

Then (exercise!), f is twice differentiable everywhere, but

$$\frac{\partial^2 f}{\partial x \partial y}(0, 0) = 1 \neq -1 = \frac{\partial^2 f}{\partial y \partial x}.$$

To actually prove Clairaut's theorem, we somehow need to use that the second partials are themselves continuous, and that's why things get a bit trickier. In practice, one rarely has to worry about these subtleties: usually, a function that is given to you is going to be obviously twice continuously differentiable. Even if it

isn't at a given point or points, one can still apply Clairaut's theorem away from those points.⁴

16.4. Linearization. We've discussed how to find the tangent space of the graph of a differentiable function \mathbf{f} at a point \mathbf{a} . This tangent space is itself the graph of a function $\mathbf{L} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, which is called the *linearization* of \mathbf{f} at \mathbf{a} . The linearization is given by the formula

$$\mathbf{L}(\mathbf{x}) = \mathbf{f}(\mathbf{a}) + (D\mathbf{f})(\mathbf{a})(\mathbf{x} - \mathbf{a}).$$

One should think of this as a first-order Taylor approximation. This formula is not at all difficult to derive: we simply take the derivative function and translate it to the point $(\mathbf{a}, \mathbf{f}(\mathbf{a}))$.

16.5. Quadratic approximation for real-valued functions. One can ask about better approximations: what if we want, say, a second-order Taylor approximation, assuming our function is twice differentiable? We could certainly do this! Notationally, however, things get quite complicated, so in this class we will only consider second-order expansions of real-valued functions (that is, functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$). In this case, the derivative at a point is given by the $1 \times n$ "horizontal" matrix

$$\left(\frac{\partial f}{\partial x_1}(\mathbf{a}) \quad \frac{\partial f}{\partial x_2}(\mathbf{a}) \quad \dots \quad \frac{\partial f}{\partial x_n}(\mathbf{a}) \right).$$

We define the Hessian matrix $(Hf)(\mathbf{a})$ to be the following $n \times n$ matrix:

$$\begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2}(\mathbf{a}) & \frac{\partial^2 f}{\partial x_1 \partial x_2}(\mathbf{a}) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n}(\mathbf{a}) \\ \frac{\partial^2 f}{\partial x_2 \partial x_1}(\mathbf{a}) & \frac{\partial^2 f}{\partial x_2^2}(\mathbf{a}) & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n}(\mathbf{a}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1}(\mathbf{a}) & \frac{\partial^2 f}{\partial x_n \partial x_2}(\mathbf{a}) & \dots & \frac{\partial^2 f}{\partial x_n \partial x_n}(\mathbf{a}) \end{pmatrix}.$$

This matrix is easier to remember than the derivative matrix, because if f is twice continuously differentiable then by Clairaut's theorem it is symmetric. We will assume this in the following.

Define the *quadratic approximation* (or *second-order approximation*) to f at \mathbf{a} to be the function $T_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ defined by

$$T_2(\mathbf{x}) = f(\mathbf{a}) + (Df)(\mathbf{a})(\mathbf{x} - \mathbf{a}) + \frac{1}{2}(\mathbf{x} - \mathbf{a})^T (Hf)(\mathbf{a})(\mathbf{x} - \mathbf{a}).$$

Notice that the last term is half the quadratic form associated to the matrix $(Hf)(\mathbf{a})$, shifted by \mathbf{a} . Thus we see quadratic forms appearing "in the wild" (so to speak).

As an example, let's calculate the quadratic approximation to the function

$$f(u, v) = u^2 \cos v + v \sin u$$

at $(u, v) = (0, 0)$. We calculate $f(0, 0) = 0$,

$$(Df)(u, v) = (2u \cos v + v \cos u \quad -u^2 \sin v + \sin u)$$

so $(Df)(0, 0) = (0 \quad 0)$, and

$$(Hf)(u, v) = \begin{pmatrix} 2 \cos v - v \sin u & -2u \sin v + \cos u \\ -2u \sin v + \cos u & -u^2 \cos v \end{pmatrix},$$

⁴There's a small technical issue here involving topology that might come up for *really* nasty functions, which I'm going to ignore completely.

so

$$(Hf)(0,0) = \begin{pmatrix} 2 & 1 \\ 1 & 0 \end{pmatrix}.$$

Therefore

$$\begin{aligned} T_2(u,v) &= 0 + (0 \ 0) \begin{pmatrix} u \\ v \end{pmatrix} + \frac{1}{2}(u \ v) \begin{pmatrix} 2 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix} \\ &= u^2 + uv. \end{aligned}$$

With a little knowledge about the Taylor series expansions of the sine and cosine functions around the origin, we could have saved ourselves some work: we have

$$\cos v = 1 - \frac{v^2}{2} + \dots, \quad \sin u = u - \frac{u^3}{6} + \dots,$$

so plugging in we get

$$f(u,v) = u^2 \left(1 - \frac{v^2}{2} + \dots\right) + v \left(u - \frac{u^3}{6} + \dots\right),$$

which is the full Taylor series expansion of f (in both variables), at least when we simplify it. We are only interested in terms up to quadratic, so we are left with

$$T_2(u,v) = u^2 + uv,$$

as before.

17. MARCH 3

17.1. Warm-up question. Find the second-order Taylor approximation at the origin to the function $f(x,y) = y + xe^{-3y}$. Bonus: do this in two different ways.

17.2. Warm-up answer. First, the standard method, which invokes the formula

$$T_2(x,y) = f(0,0) + (Df)(0,0) \begin{pmatrix} x \\ y \end{pmatrix} + \frac{1}{2}(x \ y) (Hf)(0,0) \begin{pmatrix} x \\ y \end{pmatrix}.$$

So we calculate: $f(0,0) = 0$, and taking two partial derivatives,

$$(Df)(x,y) = (e^{-3y} \ 1 - 3xe^{-3y}),$$

so $(Df)(0,0) = (1 \ 1)$; likewise, taking more partial derivatives,

$$(Hf)(x,y) = \begin{pmatrix} 0 & -3e^{-3y} \\ 3e^{-3y} & 9xe^{-3y} \end{pmatrix},$$

so

$$(Hf)(0,0) = \begin{pmatrix} 0 & -3 \\ -3 & 0 \end{pmatrix}.$$

Putting this all together, we get

$$T_2(x,y) = x + y - 3xy.$$

As usual, the constant term (here zero) depends on the value of the function at the given point, the linear terms depend on the derivative of the function at the given point, and the quadratic terms depend on the Hessian of the function at the given point.

Alternatively, without calculating any derivatives, we can proceed as follows. We recall that the Taylor series of the exponential function at the origin is easy to remember and everywhere convergent. We have

$$e^z = 1 + z + \frac{z^2}{2} + \dots$$

for all z . Plugging in $z = -3y$, we get

$$e^{-3y} = 1 - 3y + \frac{9}{2}y^2 - \dots,$$

so

$$f(x, y) = y + x \left(1 - 3y + \frac{9}{2}y^2 - \dots \right).$$

As we are interested in the quadratic approximation, we can just get rid of any terms that are of third order and higher, leaving us with

$$T_2(x, y) = y + x(1 - 3y) = y + x - 3xy,$$

as before.

17.3. The chain rule. Let's recall the one-variable chain rule: if $f : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable at a and $g : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable at $f(a)$, then

$$(g \circ f)'(a) = g'(f(a))f'(a).$$

The multivariable chain rule is remarkably similar, thanks to our definition of the derivative. It reads as follows: if $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is differentiable at \mathbf{a} and $\mathbf{g} : \mathbb{R}^m \rightarrow \mathbb{R}^p$ is differentiable at $\mathbf{f}(\mathbf{a})$, then

$$(D(\mathbf{g} \circ \mathbf{f}))(\mathbf{a}) = (D\mathbf{g})(\mathbf{f}(\mathbf{a}))(D\mathbf{f})(\mathbf{a}).$$

In words: if everything is differentiable, the derivative of the composition of two functions at a point is equal to the derivative of the second function, evaluated at the first function evaluated at the point, multiplied by the derivative of the first function, evaluated at the point. Let's check that this makes sense: the left hand side is the derivative at a point of a function from \mathbb{R}^n to \mathbb{R}^p , so is therefore a $p \times n$ matrix. The right hand side is a $p \times m$ matrix (the derivative at a point of a function from \mathbb{R}^m to \mathbb{R}^p) times a $m \times n$ matrix (the derivatives at a point of a function from \mathbb{R}^n to \mathbb{R}^m). So we see that the dimensions check out.

In this formula, order matters in a way that it doesn't in the one-dimensional chain rule! If we switch the two matrices on the right hand side, they will not have compatible dimensions in general, and even if they do we will not get the right answer in general, as $AB \neq BA$ for most square matrices A and B . This does not come up in the one-dimensional case, because one-dimensional matrix multiplication is just multiplication of real numbers, which does commute.

There is a sense in which the multivariable chain rule clarifies even the one-dimensional chain rule via our interpretation of derivatives as linear transformations. If we consider the matrices above as linear transformations, then matrix multiplication is composition of linear transformations, and the chain rule just tells us the following: the derivative of $\mathbf{g} \circ \mathbf{f}$ at \mathbf{a} is equal to the derivative of \mathbf{g} at $\mathbf{f}(\mathbf{a})$ composed with the derivative of \mathbf{f} at \mathbf{a} . Even more succinctly: *the derivative of the composition of two functions is the composition of the derivatives* (evaluated at the obvious points)! When we rephrase the chain rule in this way, how could it be otherwise?

The proof is annoying but not particularly difficult; the difficulty is simply in adapting the one-dimensional case appropriately to deal with matrices.

Let's do an example. Suppose $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ and $\mathbf{g} : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ are two functions such that

$$(D\mathbf{f})(0, 0) = \begin{pmatrix} 0 & 2 \\ 1 & 3 \\ 0 & 1 \end{pmatrix}$$

and

$$\mathbf{f}(0, 0) = \begin{pmatrix} -2 \\ 0 \\ 2 \end{pmatrix},$$

whereas

$$\mathbf{g}(x, y, z) = \begin{pmatrix} x^2 + yz \\ y \end{pmatrix}.$$

What is $(D(\mathbf{g} \circ \mathbf{f}))(0, 0)$?

Well, by the chain rule, it's equal to $(D\mathbf{g})(\mathbf{f}(0, 0))(D\mathbf{f})(0, 0)$. We can calculate partial derivatives to get

$$(D\mathbf{g})(x, y, z) = \begin{pmatrix} 2x & z & y \\ 0 & 1 & 0 \end{pmatrix},$$

and since we're interested in this matrix evaluated at $\mathbf{f}(0, 0)$, we can plug in $x = -2$, $y = 0$, $z = 2$ to get

$$(D\mathbf{g})(\mathbf{f}(0, 0)) = \begin{pmatrix} -4 & 2 & 0 \\ 0 & 1 & 0 \end{pmatrix}.$$

We can then multiply the two matrices to get the answer:

$$(D\mathbf{g})(\mathbf{f}(0, 0))(D\mathbf{f})(0, 0) = \begin{pmatrix} -4 & 2 & 0 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 2 \\ 1 & 3 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 2 & -2 \\ 1 & 3 \end{pmatrix}.$$

This example demonstrates that in order to find the derivative of a composition of two functions at a point, we do not need to know formulas for both of the functions! All we need to know is the information the chain rule asks for. If we are given full formulas for both functions, we don't even need to use the multivariable chain rule (although sometimes it may be helpful): we can just plug in one function into the other to get a formula for their composition, and take a derivative as usual.

If we are only interested in one element of the matrix of derivatives of a composition - that is, we only care about one partial derivative of one of the components - then we are often better off by ditching the matrix notation. Let g_i be the coordinate functions of \mathbf{g} , f_i the coordinate functions of \mathbf{f} , and h_i the coordinate functions of $\mathbf{g} \circ \mathbf{f}$. Let y_i and x_i be the independent variables that \mathbf{g} and \mathbf{f} depend on, respectively. The chain rule, expanding the matrix multiplication to find one entry of the left hand side, gives us

$$\frac{\partial h_i}{\partial x_j}(\mathbf{a}) = \left(\frac{\partial g_i}{\partial y_1}(\mathbf{f}(\mathbf{a})) \quad \dots \quad \frac{\partial g_i}{\partial y_m}(\mathbf{f}(\mathbf{a})) \right) \begin{pmatrix} \frac{\partial f_1}{\partial x_j}(\mathbf{a}) \\ \vdots \\ \frac{\partial f_m}{\partial x_j}(\mathbf{a}) \end{pmatrix} = \sum_{k=1}^m \frac{\partial g_i}{\partial y_k}(\mathbf{f}(\mathbf{a})) \frac{\partial f_k}{\partial x_j}(\mathbf{a}).$$

This seems more complicated than it actually is. What it is saying is that in order to find this partial derivative, we need to add up products of derivatives over all possible ways of getting from h_i to x_j via dependence of variables. Here, h_i

depends on each of f_1 through f_m , all of whom depend on x_j . So there are m terms in the sum: one which corresponds to h_i depending on f_1 depending on x_j , one which corresponds to h_i depending on f_2 depending on x_j , and so on. A good way to visualize this is via a dependency (or “tree”) graph, of which there are examples in the textbook on pages 76 through 78.

In fact, this whole idea works for compositions of more than two functions, because we are still just multiplying derivative matrices at the end of the day. For example, say we have the following functions that depend on the given variables: $f(u, v)$, $u(x, y, z)$ and $v(x, y, z)$, and $x(s, t)$, $y(t)$, and $z(s, t)$. We might be interested in the partial derivative of f with respect to s . To find it, we have to find all of the dependency paths linking f with s , of which there are four: f depends on u depends on x depends on s , f depends on u depends on z depends on s , f depends on v depends on x depends on s , and f depends on v depends on z depends on s . So we have four terms, and we shorthandedly write this as

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial u} \frac{\partial u}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial f}{\partial u} \frac{\partial u}{\partial z} \frac{\partial z}{\partial s} + \frac{\partial f}{\partial v} \frac{\partial v}{\partial x} \frac{\partial x}{\partial s} + \frac{\partial f}{\partial v} \frac{\partial v}{\partial z} \frac{\partial z}{\partial s},$$

even though strictly speaking we have to evaluate each of these partial derivatives at the correct point, so this is sort of meaningless as written.

17.4. Directional derivatives. Consider a real-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. We have defined partial derivatives with respect to the n coordinate directions of f , but geometrically there is nothing particularly special about these directions. For example, if f depends on two real variables, we can easily visualize its graph (or its contour map); the two partial derivatives correspond to the slope of the graph in the two directions parallel to the x and y axes. But nothing is stopping us from considering the slope in other directions as well! To this end, we define the *directional derivative* of f at \mathbf{a} in the direction \mathbf{v} as

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

The directional derivative is a scalar and it has a geometric interpretation as a slope if \mathbf{v} is a unit vector (in general, it is something like “the total amount the height of the graph changes over the length of \mathbf{v} if the graph happened to be linear”). In particular, the directional derivative evaluated at a standard basis vector is just the partial derivative:

$$(D_{\mathbf{e}_i}f)(\mathbf{a}) = \frac{\partial f}{\partial x_i}(\mathbf{a}).$$

How do we actually calculate this? Fortunately, we have a convenient theorem:

$$(D_{\mathbf{v}}f)(\mathbf{a}) = (Df)(\mathbf{a})\mathbf{v}.$$

The right hand side is the matrix multiplication of the derivative matrix of f at \mathbf{a} and the vector \mathbf{v} , which yields a scalar (a 1×1 matrix). This is often written in terms of the *gradient vector* $(\nabla f)(\mathbf{a})$, which is defined as the transpose of the derivative; thus we also have

$$(D_{\mathbf{v}}f)(\mathbf{a}) = (\nabla f)(\mathbf{a}) \cdot \mathbf{v},$$

where \cdot is the dot product.

Let's do an example. Let $f(x, y, z) = 3x + y \cos z$, and let $\mathbf{v} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$. Then we have

$$(D_{\mathbf{v}}f)(0, 0, 0) = (\nabla f)(0, 0, 0) \cdot \mathbf{v}.$$

Calculating three partial derivatives,

$$(\nabla f)(x, y, z) = \begin{pmatrix} 3 \\ \cos z \\ -y \sin z \end{pmatrix} \implies (\nabla f)(0, 0, 0) = \begin{pmatrix} 3 \\ 1 \\ 0 \end{pmatrix}.$$

Therefore

$$(D_{\mathbf{v}}f)(0, 0, 0) = \begin{pmatrix} 3 \\ 1 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = 3.$$

The proof of this theorem is easy using the chain rule, and it tells us something interesting about the gradient vector. Consider all possible unit vector directions \mathbf{v} , and ask for which is the directional derivative maximal. By the law of cosines, the dot product is equal to the magnitude of the vectors multiplied by the cosine of the angle between them. Because we are only considering unit vectors, the maximum therefore occurs when the cosine of the angle between them is at a maximum; that is, precisely at $\theta = 0$, when the two vectors are pointing in the same direction. This tells us that the direction of the gradient vector itself is the direction of greatest increase of the function, and similarly we can conclude that the negative of the gradient vector points in the direction of greatest decrease of the function. If we visualize the function via a contour map, the gradient always points perpendicularly to the level sets, in the direction that the function is increasing.

18. MARCH 5

18.1. Warm-up question. Suppose $S \subset \mathbb{R}^3$ is the surface defined by the equation $xy + yz + zx = 1$. Find the equation of the tangent plane to S at $(-1, 2, 3)$.

18.2. Warm-up answer. I wanted to go over one example of a problem like this, because we have not explicitly described how to solve these and there is a convenient trick. View the surface S as a level set of the function

$$F(x, y, z) = xy + yz + zx.$$

(Then we have $S = F^{-1}(1)$.) By our previous discussion of the gradient vector, we know that the gradient vector at a point always points perpendicularly to the level set; in particular, it will be a normal vector to the tangent plane of S . Since we know how to find the equation of a plane in \mathbb{R}^3 given a normal vector and a point, we are good to go.

Calculating, we get

$$(\nabla F)(x, y, z) = \begin{pmatrix} y + z \\ x + z \\ x + y \end{pmatrix},$$

so

$$(\nabla F)(-1, 2, 3) = \begin{pmatrix} 5 \\ 2 \\ 1 \end{pmatrix}.$$

This is our normal vector, and we already know that the plane passes through the point $(-1, 2, 3)$, so the equation of the plane is

$$5(x + 1) + 2(y - 2) + (z - 3) = 0,$$

which we can simplify to

$$5x + 2y + z = 2.$$

18.3. Local extrema. In one-variable calculus, we have the “first derivative test” and the “second derivative test,” which are really quite different things but both have to do with finding local extrema of functions. Briefly, the first derivative test tells us that all local extrema are *critical points* (which we define as points where the function is either not differentiable or has derivative zero). Note that it is certainly possible to have critical points that are not local extrema! The second derivative test tells us that if everything is twice continuously differentiable (the same criterion we needed to apply Clairaut’s theorem), then we can usually determine whether we have a local minimum or a local maximum by looking at the second derivative.

Both of these tests have straightforward extensions to multivariable real-valued functions. We define a critical point of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ to be a point where either f is not differentiable or the derivative matrix is the zero matrix. Then the first derivative test says that if f has a local extremum at \mathbf{a} , then \mathbf{a} is a critical point of f . Now assume that f has continuous second derivatives at \mathbf{a} and let $Q_{\mathbf{a}}$ be the quadratic form associated to the Hessian matrix $(Hf)(\mathbf{a})$. The second derivative test tells us that $Q_{\mathbf{a}}$ is positive definite if and only if f has a local minimum at \mathbf{a} , it is negative definite if and only if f has a local maximum at \mathbf{a} , and it is indefinite if and only if f has a saddle point at \mathbf{a} . If the quadratic form is semidefinite without being definite, we cannot conclude anything from the second derivative test alone. As an example, consider the functions $f_1(x) = x^3$, $f_2(x) = x^4$, and $f_3(x) = -x^4$. All three functions have zero first and second derivatives at the origin, but f_3 has a local maximum there, f_2 has a local minimum there, and f_1 has neither a maximum nor a minimum there.

Note that a *saddle point* is defined to be a point where the function is differentiable but has neither a maximum nor a minimum. A saddle point has the property that any neighborhood will contain a point that gives a larger value and a point that gives a smaller value.

Here’s an example. Let’s classify (i.e., determine which are maxima, minima, etc.) the critical points of the function $f(x, y) = ye^x + x - 2y$. This is clearly always differentiable as many times as we would like, so the critical points are just the points where the derivative is zero. We calculate

$$(Df)(x, y) = (ye^x + 1 \quad e^x - 2),$$

and solve the two equations $ye^x + 1 = 0$, $e^x - 2 = 0$. These give us the single solution $x = \ln 2$, $y = -\frac{1}{2}$, so there is only one critical point. We can apply the second derivative test to it. The Hessian matrix is

$$(Hf)(x, y) = \begin{pmatrix} ye^x & e^x \\ e^x & 0 \end{pmatrix},$$

so at the critical point in question we have

$$(Hf)\left(\ln 2, -\frac{1}{2}\right) = \begin{pmatrix} -1 & 2 \\ 2 & 0 \end{pmatrix}.$$

We want to determine the definiteness of the corresponding quadratic form, which involves calculating eigenvalues. The characteristic polynomial is

$$p(\lambda) = (-1 - \lambda)(-\lambda) - 4 = \lambda^2 + \lambda - 4,$$

which by the quadratic formula has roots

$$\frac{-1 \pm \sqrt{17}}{2}.$$

Since $\sqrt{17}$ is clearly greater than 1, one of these roots is positive and one is negative, which means the quadratic form associated to the Hessian is indefinite, which means that our critical point is a saddle point.

18.4. Global extrema - an example. Consider a set $D \subset \mathbb{R}^n$. We call D *closed* if it contains each point of its boundary, and we call D *bounded* if it is contained in some ball centered at the origin (that is, it does not “go off to infinity” in any direction). See the textbook for more precise definitions of these concepts. In any event, we have a beautiful theorem, the *extreme value theorem*, which tells us that any continuous function defined on a closed and bounded subset of \mathbb{R}^n is guaranteed to possess a global minimum and a global maximum. The proof is rather difficult, as it requires a close study of basic properties of the real numbers; as such, it is beyond the scope of this course.

If we are given a function defined on a closed and bounded domain and asked to find the global minimum and global maximum, we can certainly still apply the first and second derivative test to the interior of the domain. The boundary, however, presents problems, because we cannot apply these tests directly and the global minimum and maximum may very well lie there. The following example shows one method of dealing with this issue.

Let’s find the global minimum and maximum of the function $f(x, y) = y^2 + (x + 2)^2$ on the domain defined by $x^2 + y^2 \leq 1$; that is, the closed unit disc. First we look for critical points in the interior, because a maximum or minimum might be found there. Calculating,

$$(Df)(x, y) = (2(x + 2) \quad 2y),$$

so solving $2(x + 2) = 0$ and $2y = 0$ tells us that the only critical point of this function is $x = -2$, $y = 0$, which lies outside the given domain and is therefore irrelevant. So our minimum and maximum must lie on the boundary somewhere.

In order to study the boundary, we will write it as the image of a parametrized curve. Here the boundary is the unit circle, so it is the image of the curve

$$\mathbf{g}(t) = \begin{pmatrix} \cos t \\ \sin t \end{pmatrix}.$$

So as t varies over all real numbers (or even over all numbers in $[0, 2\pi)$), the quantity $f(\mathbf{g}(t))$ varies over all values of the function f on the boundary. To find the maximum and minimum, we are therefore reduced to a one-variable problem: we have to find the critical points of $(f \circ \mathbf{g})$ and plug them in to see which is the biggest and which is the smallest.

To do this, we can use the chain rule or just plug in. I will do the latter, getting

$$f(\mathbf{g}(t)) = (\sin t)^2 + (\cos t + 2)^2.$$

The derivative is

$$\frac{d}{dt}f(\mathbf{g}(t)) = 2 \cos t \sin t - 2 \sin t(\cos t + 2) = -4 \sin t.$$

Setting this equal to zero, we get $-4 \sin t = 0 \implies \sin t = 0$, which has solutions $t = \dots - 2\pi, -\pi, 0, \pi, 2\pi, \dots$. But adding 2π to t does nothing to the value of f because we just go around the circle once, so we only have to check, say, $t = 0$ and $t = \pi$. We have

$$f(\mathbf{g}(0)) = 0^2 + 3^2 = 9, \quad f(\mathbf{g}(\pi)) = 0^2 + (-1 + 2)^2 = 1,$$

so the former is the global maximum and the latter is the global minimum. The global maximum is reached at the point $\mathbf{g}(0) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and the global minimum is reached at the point $\mathbf{g}(\pi) = \begin{pmatrix} -1 \\ 0 \end{pmatrix}$.

At this point, it would be a good exercise to visualize what is going on in this example (by thinking about the graph of f , for example) to see if you could have guessed this answer beforehand.

19. MARCH 9

19.1. Warm-up question. Find the minimum value of the function

$$f(x, y) = 3x^2 + 5y^2 - xy + y$$

restricted to the half-plane $S = \{(x, y) : x \leq -1\} \subset \mathbb{R}^2$. You may assume that a minimum exists.

19.2. Warm-up answer. We can do this with the basic method from last section: first we check for critical points in the interior of S using the first derivative test, and second we check for critical points on the boundary via a parametrization.

We calculate

$$(\nabla f)(x, y) = \begin{pmatrix} 6x - y \\ 10y - x + 1 \end{pmatrix}.$$

Clearly f is differentiable everywhere, so the critical points will just be the points at which the gradient is the zero vector. This yields the two equations $6x - y = 0$ and $10y - x + 1 = 0$. From the first equation, $y = 6x$, so we can plug this back in to the second equation to get $60x - x + 1 = 0$, or $x = -\frac{1}{59}$. And here we can stop: we see that any critical point of f has an x -value that excludes it from being in the set S , because $-\frac{1}{59} > -1$. Therefore there are no critical points in the interior of S .

Now we parametrize the boundary, which is easy to do: the boundary is the line $x = -1$, which we can parametrize with the curve

$$\mathbf{g}(t) = \begin{pmatrix} -1 \\ t \end{pmatrix}.$$

Therefore

$$(f \circ \mathbf{g})(t) = 3 + 5t^2 + 2t.$$

Note that this is just what we get when we let $x = -1$ and call the y coordinate t . To find the critical points here, again it suffices to set the derivative equal to zero because this function is differentiable everywhere. So we calculate

$$(f \circ \mathbf{g})'(t) = 10t + 2,$$

and setting this equal to zero we get $10t + 2 = 0$, which implies $t = -\frac{1}{5}$. The image of this value of t under \mathbf{g} is

$$\mathbf{g}\left(-\frac{1}{5}\right) = \begin{pmatrix} -1 \\ -\frac{1}{5} \end{pmatrix},$$

so this is our only critical point.

Because we assumed that a minimum of f on S exists, it must exist here. Its value is

$$f\left(-1, -\frac{1}{5}\right) = 3 + \frac{5}{25} - \frac{2}{5} = \frac{14}{5}.$$

19.3. Lagrange multipliers: theory. A more robust way of dealing with extrema of function on boundaries, or subject to constraints more generally, is given by the theory of Lagrange multipliers. Here is the setup: suppose we have a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that we want to minimize or maximize on a set $S \subset \mathbb{R}^n$, and say that we can write $S = g^{-1}(c)$ for some $c \in \mathbb{R}$ and some $g : \mathbb{R}^n \rightarrow \mathbb{R}$. That is, we assume that S is a level set of a function. Also assume that f and g are continuously differentiable.

In this case, we can draw (at least in two dimensions) the set S and the level sets of the function f , and we can notice that a minimal or maximal point of f on S cannot lie at a point at which the level set of f crosses S – that is; if \mathbf{a} is the intersection of S with a level set of f , and these two sets are not tangent at \mathbf{a} , then \mathbf{a} cannot be an extremum (because otherwise there would be a point very close to \mathbf{a} that has a smaller value, and one that has a larger value, just by picking a point on one side and then on the other side of the given level set of f). Note that it is possible for S and the level set of f to be parallel at a point without there being an extremum there, which is essentially equivalent to the statement that a point at which a derivative is zero need not be a minimum or maximum (it could be a saddle point). This can all be proven rigorously, as is done in the textbook.

How do we algebraically characterize the level set of f and S being parallel? Well, first notice that the gradient of f always points orthogonally to a tangent space, and as S is a level set of g we can apply this to g as well. We find that at an extremum the gradient of f and the gradient of g must point in the same direction; i.e., they are proportional. If, say $(\nabla g)(\mathbf{a}) \neq \mathbf{0}$, then we have

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

for some $\lambda \in \mathbb{R}$. If $(\nabla g)(\mathbf{a})$ does happen to be zero, then that's not such a big problem either; we can either check such points separately when we hunt for possible extrema or if $(\nabla f)(\mathbf{a}) \neq \mathbf{0}$ we can write the similar equation

$$\mu(\nabla f)(\mathbf{a}) = (\nabla g)(\mathbf{a})$$

for some $\mu \in \mathbb{R}$.

Summarizing, we have the following theorem: if $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ are continuously differentiable functions, $c \in \mathbb{R}$, \mathbf{a} is an extreme value of f when restricted to $S = g^{-1}(c)$, and $(\nabla g)(\mathbf{a}) \neq \mathbf{0}$, then

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

for some $\lambda \in \mathbb{R}$.

This gives us an algorithm for solving minimization or maximization problems restricted to a subset S : we express S as $g^{-1}(c)$ for some c and g , calculate (∇f)

and (∇g) , and then the Lagrange multiplier equation gives us a system of n real-variable equations. We have $n + 1$ unknowns (the coordinates of the critical point and λ), so we expect we need one more equation: this is true, and it is given by the original constraint $g(\mathbf{x}) = c$. Solving these equations gives us the critical points, which we can test individually to see which is the largest and which is the smallest.

19.4. Lagrange multipliers: examples. Let's redo part of the warm-up problem with Lagrange multipliers. In the process, we'll see that sometimes it might be easier *not* to use Lagrange multipliers if there's an easy way to parametrize something.

We're interested in finding the local extrema of the function $f(x, y) = 3x^2 + 5y^2 - xy + y$ on the line L given by $x = -1$ (the boundary of the set in the warm-up). There are many functions for which $x = -1$ is a level set, but perhaps the most obvious is the function $g(x, y) = x$; then $L = g^{-1}(-1)$. The gradient of g is

$$(\nabla g)(x, y) = \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

in particular, it is never zero so we can apply the Lagrange multiplier theorem, concluding that

$$(\nabla f)(x, y) = \lambda(\nabla g)(x, y)$$

for some $\lambda \in \mathbb{R}$ at any local extremum (x, y) . Using the calculation of $(\nabla f)(x, y)$ from above as well as the constraint equation $g(x, y) = -1$, we get the equations

$$\begin{aligned} 6x - y &= \lambda, \\ 10y - x + 1 &= 0, \\ x &= -1. \end{aligned}$$

These are easy to solve; we have $x = -1$, so plugging into the second equation gives us $y = -\frac{1}{5}$, and this is the only possible critical point. We recover our previous result, but possibly doing more work in the process.

Here's an example of a question for which Lagrange multipliers really are helpful. Let's find the point or points on the surface $x^2 - yz = 4$ closest to the origin. One of the homework problems proves that such points always exist, so we don't have to worry about that.

Rephrased, we are trying to find points on the surface $S = \{(x, y, z) : x^2 - yz = 4\}$ that minimize the distance to the origin, which is given by the function $\sqrt{x^2 + y^2 + z^2}$. Actually, because it is much easier to differentiate, we will minimize the function $f(x, y, z) = x^2 + y^2 + z^2$ instead; clearly, any minimum of the former will be a minimum of the latter and vice versa. In order to use Lagrange multipliers, we write $S = g^{-1}(4)$, where $g(x, y, z) = x^2 - yz$ (in general, if we are given our constraints in terms of equations, it is very easy to write them as level sets).

We calculate

$$(\nabla f)(x, y, z) = \begin{pmatrix} 2x \\ 2y \\ 2z \end{pmatrix}, \quad (\nabla g)(x, y, z) = \begin{pmatrix} 2x \\ -z \\ -y \end{pmatrix}.$$

We see that ∇g is equal to zero only at the origin $(0, 0, 0)$, which is not on the surface S , so we can apply the Lagrange multiplier theorem. Together with the

constraint equation, we get the four equations

$$\begin{aligned}2x &= \lambda 2x, \\2y &= -\lambda z, \\2z &= -\lambda y, \\x^2 - yz &= 4.\end{aligned}$$

To solve, we split into cases. The first equation tells us that either $\lambda = 1$ or $x = 0$. In the first case, we have $2y = -z$ and $2z = -y$, which together imply that $y = z = 0$ (just by plugging one into the other). Then the last equation gives $x^2 = 4$, so $x = \pm 2$. Therefore we get two critical points, $(\pm 2, 0, 0)$. In the second case, where $x = 0$, we plug in the equation $y = -\frac{\lambda}{2}z$ into the third equation to get

$$2z = -\lambda \left(-\frac{\lambda}{2}z\right) \implies \left(2 - \frac{\lambda^2}{2}\right)z = 0.$$

Therefore we have two cases; either $z = 0$ or $\lambda^2 = 4$. In the first case, the fourth equation gives a contradiction: we cannot have $0 \cdot z = 4$. So that is not possible. In the second case, $\lambda = \pm 2$, so $y = \pm z$. If $y = z$ then the fourth equation gives $-y^2 = 4$, which is not possible (a square must be nonnegative). So we must have $y = -z$, and $y^2 = 4$ gives $y = \pm 2$. Therefore we get two additional critical points, $(0, 2, -2)$ and $(0, -2, 2)$.

It is easy to check that the first two critical points, $(\pm 2, 0, 0)$, are two units away from the origin, while the second two critical points are $2\sqrt{2}$ units away from the origin. Therefore the solution to the problem of finding the closest points to zero is the two points $(\pm 2, 0, 0)$.

19.5. A generalization. There is a generalization of the Lagrange multiplier theorem to the situation where one can write S as an intersection of several level curves (for example, maybe S is given to us as the solution set of several equations). It reads as follows: assume that there exist finitely many functions $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ and constants $c_i \in \mathbb{R}$ such that S is the intersection of the sets $g_i^{-1}(c_i)$. Assume that all of the g_i and a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ are continuously differentiable. Then if \mathbf{a} is a local extremum of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and the vectors $(\nabla g_i)(\mathbf{a})$ are linearly independent, there exist constants $\lambda_i \in \mathbb{R}$ such that

$$(\nabla f)(\mathbf{a}) = \sum_i \lambda_i (\nabla g_i)(\mathbf{a}).$$

Two things to notice: first, the requirement that the $(\nabla g_i)(\mathbf{a})$ are linearly independent is the natural generalization of the requirement that $(\nabla g)(\mathbf{a})$ is nonzero (as a vector is a linearly independent set if and only if it is nonzero). Second, if we apply this theorem in a situation where $1 \leq i \leq r$, then we will have n equations from the theorem and r from the constraint equations, for a total of $n+r$ equations, and we will have n unknowns from the coordinates of the vector and r unknowns from the λ_i , for a total of $n+r$ variables. So we still have as many variables as unknowns. In this class, you will almost surely never have to apply this theorem for $r > 2$.

20. MARCH 11

Today was all review, from the whole course. Some questions follow, in approximate increasing order of difficulty, followed by answers.

20.1. Questions that would be very reasonable on an exam.

- (1) True or false (and why?): if A , B , and C are matrices, and $AB = AC$, then $B = C$.
- (2) What is the definition of the dimension of a subspace?
- (3) True or false (and why?): if $f(x, y) = \sin(\cos(\sin(x^2 - y) + \cos(x)))$, then $f_{xyxyx} = f_{yyxxx}$.
- (4) True or false (and why?): if A is an $m \times n$ matrix such that $A\mathbf{x} = \mathbf{0}$ has infinitely many solutions, then $A\mathbf{x} = \mathbf{b}$ has infinitely many solutions for all $\mathbf{b} \in \mathbb{R}^m$.
- (5) True or false (and why?): there exists an $n \times n$ matrix A such that $A\mathbf{x} = \mathbf{b}$ has infinitely many solutions for every $\mathbf{b} \in \mathbb{R}^n$.
- (6) True or false (and why?): the function $f(x, y) = -x^{1000} - y^{1000}$ has a local minimum at the origin.
- (7) Let V be a subspace of \mathbb{R}^n . Suppose that every $\mathbf{x} \in V$ can be written uniquely as a linear combination of the vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$. Prove that $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ is a basis for V .
- (8) Suppose you are hiking on a mountain whose shape is given by the graph of $f(x, y) = x^2 + y^2 - xy - 4x - 4y + 100$. What is the slope of the mountain at the point $(2, 3)$?
- (9) Suppose we know that

$$A = \begin{pmatrix} 1 & ? & 4 \\ 2 & ? & 5 \\ 3 & ? & 6 \end{pmatrix}, \quad \text{rref}(A) = \begin{pmatrix} 1 & 0 & 2 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}.$$

What is the second column of A ?

- (10) Find the quadratic (second order) approximation to the function $f(x, y) = x + 3y + x^2 - 2xy + xy^5$ at the origin.
- (11) Let A be any 2×3 matrix. Prove that $f(\mathbf{x}) = \|A\mathbf{x}\|^2$ is a positive semidefinite quadratic form that is not positive definite.

20.2. Brief answers.

- (1) False; for example, let A be the zero matrix. There are other examples as well!
- (2) The dimension of a subspace is the number of vectors in any basis of the subspace, or the maximal number of linearly independent vectors in the subspace.
- (3) True by Clairaut's theorem, used repeatedly.
- (4) False; for example, let A be the zero matrix of any size (though again, there are many other counterexamples).
- (5) False; if $A\mathbf{x} = \mathbf{0}$ has infinitely many solutions then the nullity of A is at least 1, and if $A\mathbf{x} = \mathbf{b}$ has a solution for every \mathbf{b} then the rank is n (because the column space is the image of the linear transformation that corresponds to a matrix). But by the rank-nullity theorem the rank plus the nullity equals n , so we have a contradiction.
- (6) False; it has a local maximum at the origin. We can't detect this with the second derivative test, because the Hessian at the origin will be all zeroes; we actually have to think about what the graph looks like.
- (7) We have to prove that $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ spans V and is linearly independent. For the first, if every vector in V can be written as a unique linear combination

of these vectors, then in particular every vector in V is in their span. For the second, let $c_1 \mathbf{v}_1 + \dots + c_k \mathbf{v}_k = \mathbf{0}$ be a linear dependence relation. But we also know that $0 \cdot \mathbf{v}_1 + \dots + 0 \cdot \mathbf{v}_k = \mathbf{0}$, so by the assumed uniqueness of linear combinations we have $c_1 = 0, c_2 = 0, \dots, c_k = 0$. This proves linear independence.

- (8) The slope at a point is just the magnitude of the gradient vector. In this case, the gradient at this point is $\begin{pmatrix} -3 \\ 0 \end{pmatrix}$ so the slope is 3.
- (9) Row reduction preserves column dependence relations, so we know that two times the first row plus the second row equals the third row. Some arithmetic yields that the second column is $\begin{pmatrix} 2 \\ 1 \\ 0 \end{pmatrix}$.
- (10) We could solve this in the usual way by calculating the value of the function, its derivative, and its Hessian matrix, but it is easier to notice that the Taylor series associated to this function at the origin is just the function itself (as for any polynomial), and to get the second order approximation we just need to keep only the terms of second order and below. Thus we can simply get rid of the last term, getting $f(x, y) = x + 3y + x^2 - 2xy$. Note that this trick only works like this if we are expanding about the origin.
- (11) We have

$$\|A\mathbf{x}\|^2 = (A\mathbf{x})^t \cdot (A\mathbf{x}) = \mathbf{x}^t A^t A \mathbf{x},$$

so this is a quadratic form with associated matrix $A^t A$ once we check that $A^t A$ is symmetric, which is easy to do: $(A^t A)^t = A^t (A^t)^t = A^t A$, so it's invariant under transpose, hence symmetric. It is at least positive semidefinite because the magnitude of a vector is always positive. It is not positive definite because if we pick any nonzero vector \mathbf{x} in the null space of A , which must be nontrivial by its dimensions (we'll always get a free variable upon reducing), then we find that $f(\mathbf{x}) = \mathbf{0}$. Therefore the function is not positive definite.

20.3. Challenge questions.

- (1) Write down a square matrix A such that A^{10} is the zero matrix but none of A, A^2, \dots, A^9 are the zero matrix.
- (2) True or false (and why?): there is a continuously differentiable function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ such that for every line L through the origin, the function f restricted to L has a local minimum at the origin, but the function f does not have a local minimum at the origin.