

Any typographical or other corrections about these notes are welcome.

## 1 Review of the dot product

The **dot product** on  $\mathbb{R}^n$  is an operation that takes two vectors and returns a number. It is defined by

$$\vec{u} \cdot \vec{v} = \sum_{i=1}^n u_i v_i$$

where  $u_1, \dots, u_n$  are the coordinates of  $\vec{u}$ , and  $v_1, \dots, v_n$  are the coordinates of  $\vec{v}$ .

The dot product can be used to give a convenient formula for the length of a vector. We use the notation  $\|\vec{v}\|$  to denote the length of a vector (also frequently called its “norm”). It is given by the formula

$$\|\vec{v}\| = \sqrt{\vec{v} \cdot \vec{v}},$$

which is equivalent to the formula

$$\|\vec{v}\| = \sqrt{v_1^2 + v_2^2 + \dots + v_n^2}.$$

The distance *between* two vectors can be expressed by  $\|\vec{u} - \vec{v}\|$ . One of the convenient aspects of the dot product is that it satisfied certain algebraic properties, such as commutativity and distributivity. This makes it possible to “expand” the expression for the distance between two vectors, for example:

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

This type of expansion will be crucial in the analysis that follows.

## 2 Statement of the least-squares problem, and examples

These notes are concerned with a process for solving the following problem.

### The least-squares problem

Given a list of vectors  $\{\vec{v}_1, \dots, \vec{v}_n\}$  and an additional vector  $\vec{b}$ , what linear combination of  $\{\vec{v}_1, \dots, \vec{v}_n\}$  is closest to  $\vec{b}$ ? In other words, which values  $c_1, \dots, c_n$  minimize the quantity

$$\left\| \sum_{i=1}^n c_i \vec{v}_i - \vec{b} \right\|?$$

The reason that this is called “least-squares problem” is that minimizing the quantity above is the same as minimizing its square, which is the sum of the differences in the individual coordinates.

**Example: linear regression.**

A typical problem in statistics is: if we are given a list of  $n$  data points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  in the plane, how can we fit a linear model to them? That is, how can we find a line  $y = c_1x + c_2$  that passes as close as possible to these data points? To solve such a problem, it is first necessary to decide how we judge which lines are better fits than others; this depends on the application. In many applications, the line will be used to predict a  $y$  coordinate, given an  $x$  coordinate. In such applications, what is crucial is to make sure that the approximation errors  $(c_1x_i + c_2 - y_i)$  tend to be as small as possible. One very common way to measure how small these errors tend to be is to take the sum of their squares. This quantity then functions as a sort of “score” for the quality of the fit (in this game, lower scores are better; a score of 0 means a perfect fit).

$$\text{score} = \sum_{i=1}^n (c_1x_i + c_2 - y_i)^2$$

The linear regression problem is: what values of  $c_1, c_2$  should we choose to get the best score (the best sum-of-squared-error)? It turns out that we can view this as an example of the least-squares problem. Indeed, define the following three vectors.

$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \vec{1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad \vec{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

Now, observe that the “score” above is precisely the squared length of  $c_1\vec{x} + c_2\vec{1} - \vec{y}$ .

$$\text{score} = \|c_1\vec{x} + c_2\vec{1} - \vec{y}\|^2$$

Hence a solution to least-squares problem will allow us to find the best coefficients for a line of best fit (according to this choice of “score.” What we are asking for, in essence, is: what linear combination of  $\vec{x}$  and  $\vec{1}$  is closest to  $\vec{y}$ ?

**Example: projection**

An example we considered earlier in class is projecting one vector onto another. Suppose that we are given vectors  $\vec{u}$  and  $\vec{v}$ , and wish to know: what scalar multiple of  $\vec{u}$  is closest to  $\vec{v}$ ? We found, by a direct method, that the answer is given by the following projection formula:

$$\text{proj}_{\vec{u}}(\vec{v}) = \frac{\vec{u} \cdot \vec{v}}{\vec{u} \cdot \vec{u}} \vec{u}$$

This is a special case of the least squares problem: we want to know which linear combination of  $\vec{u}$  (by itself) is closest to  $\vec{v}$ .

### 3 Matrix formulation

A second formulation of the least-squares problem is as follows.

### The least-squares problem, matrix formulation

Given a matrix  $A$  and a vector  $\vec{b}$ , what vector  $\vec{x}$  minimizes the quantity

$$\|A\vec{x} - \vec{b}\|?$$

The reason this is equivalent to the first formulation is that the product of a matrix with a vector is precisely the same thing as a linear combination of the columns of the matrix.

$$A\vec{x} = x_1\vec{A}_1 + x_2\vec{A}_2 + \cdots + x_n\vec{A}_n$$

Here, as usual, the symbol  $\vec{A}_i$  refers to the  $i$ th column of the matrix  $A$ , regarded as a vector.

## 4 Orthogonality and the Pythagorean theorem in $n$ dimensions

One of the most powerful aspects of the dot product is that it provides an extremely computationally efficient way to tell if two vectors are orthogonal (math jargon for “perpendicular;” I will prefer the word “orthogonal” since it is the more common word in abstract settings, such as higher dimensions).

We will call two vectors  $\vec{u}, \vec{v}$  in  $\mathbb{R}^n$  **orthogonal** if  $\vec{u} \cdot \vec{v} = 0$ . We will sometimes also write  $\vec{u} \perp \vec{v}$  to mean the same thing, in a slightly more visually suggestive way. That is,

$$\vec{u} \perp \vec{v} \text{ is synonymous with } \vec{u} \cdot \vec{v} = 0.$$

The reason this corresponds to usual geometric notion of perpendicularity comes from the formula

$$\vec{u} \cdot \vec{v} = \|\vec{u}\| \cdot \|\vec{v}\| \cdot \cos \theta$$

where  $\theta$  denotes the angle between  $\vec{u}$  and  $\vec{v}$ . This expression will be zero precisely when  $\cos \theta = 0$ , i.e. when  $\theta = \frac{\pi}{2} = 90^\circ$  (there’s also one edge case: if either  $\vec{u}$  or  $\vec{v}$  is the zero vector  $\vec{0}$ , then the dot product is 0 even though there isn’t really a well-defined “angle” in sight. For convenience, we will regard the zero vector  $\vec{0}$  as being orthogonal to everything).

The algebraic properties of the dot product allow it to show the following  $n$ -dimensional version of the pythagorean theorem.

### Pythagorean theorem

If  $\vec{u}$  and  $\vec{v}$  are two orthogonal vectors, then

$$\|\vec{u} + \vec{v}\|^2 = \|\vec{u}\|^2 + \|\vec{v}\|^2.$$

*Proof.* Using distributivity and commutativity of the dot product,

$$\begin{aligned} \|\vec{u} + \vec{v}\|^2 &= (\vec{u} + \vec{v}) \cdot (\vec{u} + \vec{v}) \\ &= \vec{u} \cdot \vec{u} + \vec{u} \cdot \vec{v} + \vec{v} \cdot \vec{u} + \vec{v} \cdot \vec{v} \\ &= \vec{u} \cdot \vec{u} + \vec{v} \cdot \vec{v} \quad (\text{since } \vec{u} \cdot \vec{v} = 0) \\ &= \|\vec{u}\|^2 + \|\vec{v}\|^2. \end{aligned}$$

□

An easy-to-obtain consequence of this is surprisingly powerful. The following corollary can be interpreted geometrically as saying: the hypotenuse of a right triangle is always longer than either leg.

#### Corollary

If  $\vec{u} \perp \vec{v}$  and  $\vec{v} \neq \vec{0}$ , then  $\|\vec{u} + \vec{v}\| > \|\vec{u}\|$ .

*Proof.* These are both positive numbers, so its enough to show that  $\|\vec{u} + \vec{v}\|^2 > \|\vec{u}\|^2$ . The Pythagorean theorem shows that this is true as long as  $\|\vec{v}\|^2 > 0$ , which it is as long as  $\vec{v} \neq \vec{0}$ .  $\square$

## 5 The key theorem

The geometric insight needed to solve least-squares is the following: if we think we've found the very best linear combination  $\vec{v}$  of  $\{\vec{v}_1, \dots, \vec{v}_n\}$  (best in the sense that it is as close as possible to the target vector  $\vec{b}$ ), how could we check that it really is the best? One option is to "tweak" it a bit, by adding a small multiple of  $\vec{v}_1, \vec{v}_2, \dots$ , or  $\vec{v}_n$ . If we've located  $\vec{v}$  correctly, then any such tweak should only move us further away from  $\vec{b}$ . For example, if we start at  $\vec{v}$  and start adding small multiples of  $\vec{v}_1$ , we'll start traveling in a straight line that should take us further away from  $\vec{b}$ . And similarly, since we could subtract as well, traveling backwards on this line should also move us further away. If you try to picture this, you can convince yourself that this is only possible if the line sits at a right angle with the vector from  $\vec{v}$  to  $\vec{b}$ . So it should be the case that  $(\vec{b} - \vec{v}) \perp \vec{v}_1$ . And we could say that same thing from  $\vec{v}_2$  through  $\vec{v}_n$ .

(In a later version of these notes, I may try to produce a visual aid for the thought process above, as drawn in class.)

So this gives us something to work with: the best choice of  $\vec{v}$  should ensure these  $n$  orthogonality conditions. Lucky for us, these  $n$  conditions turn out to be *all* we need. This is expressed in the following theorem.

#### Theorem

Let  $\{\vec{v}_1, \dots, \vec{v}_n\}$  be a list of vectors, and  $\vec{b}$  another vector.

Suppose that  $\vec{v}$  is a linear combination of  $\{\vec{v}_1, \dots, \vec{v}_n\}$ , that satisfies the following  $n$  conditions.

$$\begin{aligned} (\vec{v} - \vec{b}) &\perp \vec{v}_1 \\ (\vec{v} - \vec{b}) &\perp \vec{v}_2 \\ &\dots \\ (\vec{v} - \vec{b}) &\perp \vec{v}_n \end{aligned}$$

Then for any other linear combination  $\vec{w}$  of  $\{\vec{v}_1, \dots, \vec{v}_n\}$ ,

$$\|\vec{w} - \vec{b}\| > \|\vec{v} - \vec{b}\|.$$

In other words,  $\vec{v}$  is the linear combination of  $\{\vec{v}_1, \dots, \vec{v}_n\}$  that comes closest to  $\vec{b}$ .

The proof doesn't take too much space on the page, and is completely algebraic, but it's a little hard to come up with without some geometric insight. In this case, the key insight is: once we know that  $\vec{v} - \vec{b}$  is orthogonal to all of the  $\vec{v}_i$ , it must in fact be orthogonal to any *linear combination* of them. So the difference  $\vec{w} - \vec{v}$  is orthogonal to  $\vec{v} - \vec{b}$ . But this means that they form two legs of a right triangle, so their sum is longer than either one of them. Their sum is precisely  $\vec{w} - \vec{b}$ . Here is the proof, written out more formally.

*Proof.* We have assumed that  $\vec{v}$  and  $\vec{w}$  are linear combinations of  $\{\vec{v}_1, \dots, \vec{v}_n\}$ . So there exists constants  $c_1, c_2, \dots, c_n$  and  $d_1, d_2, \dots, d_n$  such that

$$\begin{aligned}\vec{v} &= \sum_{i=1}^n c_i \vec{v}_i \\ \vec{w} &= \sum_{i=1}^m d_i \vec{v}_i.\end{aligned}$$

Now, observe that

$$\begin{aligned}(\vec{w} - \vec{v}) \cdot (\vec{v} - \vec{b}) &= \left( \sum_{i=1}^n (d_i - c_i) \vec{v}_i \right) \cdot (\vec{v} - \vec{b}) \\ &= \sum_{i=1}^n (d_i - c_i) [\vec{v}_i \cdot (\vec{v} - \vec{b})] \\ &= \sum_{i=1}^n (d_i - c_i) \cdot 0 \quad (\text{since } \vec{v}_i \perp (\vec{v} - \vec{b})) \\ &= 0\end{aligned}$$

So  $(\vec{w} - \vec{v}) \perp (\vec{v} - \vec{b})$ . By the Corollary to the Pythagorean theorem in the previous section, it follows that

$$\|(\vec{w} - \vec{v}) + (\vec{v} - \vec{b})\| > \|\vec{v} - \vec{b}\|,$$

which is the same (after canceling the  $\vec{v}$ s) as  $\|\vec{w} - \vec{b}\| > \|\vec{v} - \vec{b}\|$ , as desired.  $\square$

## 6 The normal equation

The theorem in the previous section shows that, to solve the least-square problem, it's enough to make sure that a bunch of perpendicularity statements are true. In other words, we need to make sure that a bunch of dot products are 0. Let's pin down what these dot products should be.

We're trying to find a linear combination  $\vec{v}$  of  $\{\vec{v}_1, \dots, \vec{v}_n\}$ , so in other words we're trying to find  $n$  coefficients  $c_1, \dots, c_n$ , where we will set

$$\vec{v} = c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n.$$

The dot products we need to get to be 0 are as follows.

$$\begin{aligned}\vec{v}_1 \cdot (c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n - \vec{b}) &= 0 \\ \vec{v}_2 \cdot (c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n - \vec{b}) &= 0 \\ &\dots \\ \vec{v}_n \cdot (c_1 \vec{v}_1 + c_2 \vec{v}_2 + \dots + c_n \vec{v}_n - \vec{b}) &= 0\end{aligned}$$

The wonderful thing here is that this is a *linear system of equations!* Indeed, after a few algebra steps, these equations become

The linear system that solves least squares

Given a list of vectors  $\{\vec{v}_1, \dots, \vec{v}_n\}$  and a target vector  $\vec{b}$ , the linear combination of  $\{\vec{v}_1, \dots, \vec{v}_n\}$  that comes nearest to  $\vec{b}$  is  $\vec{v} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_n\vec{v}_n$ , where  $c_1, \dots, c_n$  are obtained as any solution to the following system of linear equations.

$$\begin{aligned}(\vec{v}_1 \cdot \vec{v}_1)c_1 + (\vec{v}_1 \cdot \vec{v}_2)c_2 + \dots + (\vec{v}_1 \cdot \vec{v}_n)c_n &= \vec{v}_1 \cdot \vec{b} \\(\vec{v}_2 \cdot \vec{v}_1)c_1 + (\vec{v}_2 \cdot \vec{v}_2)c_2 + \dots + (\vec{v}_2 \cdot \vec{v}_n)c_n &= \vec{v}_2 \cdot \vec{b} \\&\vdots \\(\vec{v}_n \cdot \vec{v}_1)c_1 + (\vec{v}_n \cdot \vec{v}_2)c_2 + \dots + (\vec{v}_n \cdot \vec{v}_n)c_n &= \vec{v}_n \cdot \vec{b}\end{aligned}$$

This system of equations is the answer to our original problem. Its solutions<sup>1</sup> give the solution to the least squares problem.

There's a few other ways that write the normal equation that may be useful in different contexts. For example, here's how we'd write it as a matrix equation. When written as a single matrix equation, this system is called the *normal equation*.

The normal equation, in terms of a list of vectors

Given a list of vectors  $\{\vec{v}_1, \dots, \vec{v}_n\}$  and a target vector  $\vec{b}$ , the linear combination of  $\{\vec{v}_1, \dots, \vec{v}_n\}$  that comes nearest to  $\vec{b}$  is  $\vec{v} = c_1\vec{v}_1 + c_2\vec{v}_2 + \dots + c_n\vec{v}_n$ , where  $c_1, \dots, c_n$  are obtained as any solution to the following matrix equation.

$$\begin{pmatrix} \vec{v}_1 \cdot \vec{v}_1 & \vec{v}_1 \cdot \vec{v}_2 & \dots & \vec{v}_1 \cdot \vec{v}_n \\ \vec{v}_2 \cdot \vec{v}_1 & \vec{v}_2 \cdot \vec{v}_2 & \dots & \vec{v}_2 \cdot \vec{v}_n \\ \vdots & \vdots & & \vdots \\ \vec{v}_n \cdot \vec{v}_1 & \vec{v}_n \cdot \vec{v}_2 & \dots & \vec{v}_n \cdot \vec{v}_n \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} \vec{v}_1 \cdot \vec{b} \\ \vec{v}_2 \cdot \vec{b} \\ \vdots \\ \vec{v}_n \cdot \vec{b} \end{pmatrix}$$

Finally, we note that the normal equation has a particularly compact form when the least-squares problem is posed in matrix form.

The normal equation, in terms of a matrix

If  $A\vec{x} = \vec{b}$  is a possibly inconsistent linear system, then a *least-squares solution* is a solution to the equation

$$A^t A \vec{x} = A^t \vec{b}.$$

This equation is always consistent. Its solution(s) have the property that they achieve the minimum possible value of  $\|A\vec{x} - \vec{b}\|$ , i.e. they come as close as possible to being solutions of the original system.

<sup>1</sup>A point I have brushed under the rug in this discussion: we've proved that *if the normal equation has a solution*, then that solution solves the least-squares problem. But I haven't proved that there is a solution, i.e. that the linear system is consistent. Don't worry, it is; but I don't know of a proof that would be particularly illuminating this early in the course, before we have some of the machinery from later chapters of the book.

This form of the normal equation follows from the previous one, because the entries of  $A^t A$  are simply all possible dot products of columns of  $A$ , while the entries of  $A^t \vec{b}$  are the dot products of each column of  $A$  with  $\vec{b}$ .

## 7 Examples

### Finding a closest combination

What linear combination of  $\vec{v}_1 = \begin{pmatrix} 5 \\ 3 \\ 2 \end{pmatrix}$  and  $\vec{v}_2 = \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix}$  is closest to the vector  $\vec{b} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ ?

**Solution:** The best combination will be  $c_1 \vec{v}_1 + c_2 \vec{v}_2$ , where  $c_1, c_2$  are solutions to the normal equation, which is the following in this case.

$$\begin{pmatrix} \vec{v}_1 \cdot \vec{v}_1 & \vec{v}_1 \cdot \vec{v}_2 \\ \vec{v}_2 \cdot \vec{v}_1 & \vec{v}_2 \cdot \vec{v}_2 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} \vec{v}_1 \cdot \vec{b} \\ \vec{v}_2 \cdot \vec{b} \end{pmatrix}$$

$$\begin{pmatrix} 38 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} 10 \\ 0 \end{pmatrix}$$

Since this is a  $2 \times 2$  system, and the matrix has nonzero determinant  $38 \cdot 2 - 1 \cdot 1 = 75$ , a quick way to solve it is with the formula for the inverse of a  $2 \times 2$  matrix.

$$\begin{aligned} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} &= \begin{pmatrix} 38 & 1 \\ 1 & 2 \end{pmatrix}^{-1} \begin{pmatrix} 10 \\ 0 \end{pmatrix} \\ &= \frac{1}{38 \cdot 2 - 1 \cdot 1} \begin{pmatrix} 2 & -1 \\ -1 & 38 \end{pmatrix} \begin{pmatrix} 10 \\ 0 \end{pmatrix} \\ &= \frac{1}{75} \begin{pmatrix} 20 \\ -10 \end{pmatrix} \\ &= \begin{pmatrix} 4/15 \\ -2/15 \end{pmatrix} \end{aligned}$$

Hence the closest linear combination is

$$\frac{4}{15} \begin{pmatrix} 5 \\ 3 \\ 2 \end{pmatrix} - \frac{2}{15} \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 4/3 \\ 2/3 \\ 2/3 \end{pmatrix}.$$

## Approximately solving an inconsistent linear system

Define a matrix  $A$  and vector  $\vec{b}$  as follows.

$$A = \begin{pmatrix} 0 & 3 \\ 1 & 1 \\ -2 & 0 \end{pmatrix} \quad \vec{b} = \begin{pmatrix} 3 \\ 0 \\ 1 \end{pmatrix}$$

The system  $A\vec{x} = \vec{b}$  is inconsistent (you can check this using row-reduction). Make the best of this situation by finding the vector  $\vec{x}$  so that  $A\vec{x}$  gets as close as possible to  $\vec{b}$  (in other words,  $\vec{x}$  should minimize  $\|A\vec{x} - \vec{b}\|$ ).

**Solution:** We can solve the normal equation  $A^t A \vec{x} = A^t \vec{b}$ , as follows.

$$\begin{pmatrix} 0 & 1 & -2 \\ 3 & 1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 3 \\ 1 & 1 \\ -2 & 0 \end{pmatrix} \vec{x} = \begin{pmatrix} 0 & 1 & -2 \\ 3 & 1 & 0 \end{pmatrix} \begin{pmatrix} 3 \\ 0 \\ 1 \end{pmatrix}$$

$$\begin{pmatrix} 5 & 1 \\ 1 & 10 \end{pmatrix} \vec{x} = \begin{pmatrix} -2 \\ 9 \end{pmatrix}$$

This can be solved with a bit of row-reduction (we could also use the formula for inverting a  $2 \times 2$  matrix, as in the previous example; I'll row-reduce here for variety).

$$\left( \begin{array}{cc|c} 5 & 1 & -2 \\ 1 & 10 & 9 \end{array} \right) \rightarrow \left( \begin{array}{cc|c} 1 & 10 & 9 \\ 5 & 1 & -2 \end{array} \right) \rightarrow \left( \begin{array}{cc|c} 1 & 10 & 9 \\ 0 & -49 & -47 \end{array} \right) \rightarrow \left( \begin{array}{cc|c} 1 & 10 & 9 \\ 0 & 1 & 47/49 \end{array} \right) \rightarrow$$

$$\left( \begin{array}{cc|c} 1 & 0 & 9 - \frac{470}{49} \\ 0 & 1 & 47/49 \end{array} \right) = \left( \begin{array}{cc|c} 1 & 0 & -29/49 \\ 0 & 1 & 47/49 \end{array} \right)$$

So the best possible choice of  $\vec{x}$  is

$$\vec{x} = \begin{pmatrix} -29/49 \\ 47/49 \end{pmatrix}.$$



## Linear regression, solved

As we saw in Section 2, we can find a line of best fit  $y = c_1x + c_2$  (in the sense of the “score” in that example) for some data points  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  by looking for the values of  $c_1, c_2$  that minimize  $\|c_1\vec{x} + c_2\vec{1} - \vec{y}\|$ , where  $\vec{x}$  is all of the  $x$  values in a vector,  $\vec{1}$  is the all-1’s vector, and  $\vec{y}$  is all of the  $y$  values in a vector.

According to the normal equation, we can solve this problem by solving the matrix equation

$$\begin{pmatrix} \vec{x} \cdot \vec{x} & \vec{x} \cdot \vec{1} \\ \vec{1} \cdot \vec{x} & \vec{1} \cdot \vec{1} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} \vec{x} \cdot \vec{y} \\ \vec{1} \cdot \vec{y} \end{pmatrix}$$

This is a  $2 \times 2$  linear system, and can be solved quite efficiently by using the formula for the inverse of a  $2 \times 2$  matrix.

This same matrix equation can also be written  $A^t A \vec{c} = A^t \vec{y}$ , where

$$A = \begin{pmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix},$$

the matrix whose columns are  $\vec{x}$  and  $\vec{1}$ .

## Projection revisited

Suppose we want the scalar multiple of  $\vec{u}$  that is closest to  $\vec{v}$ . This is the simplest possible case of the normal equation: the matrix is just  $1 \times 1$ ! The best multiple is  $c\vec{u}$ , where  $c$  solves the normal equation, which in this case is:

$$(\vec{u} \cdot \vec{u})c = (\vec{u} \cdot \vec{v}).$$

We have simply recovered the formula for the projection of  $\vec{v}$  onto  $\vec{u}$ .

## 8 Appendix: deriving the normal equation from calculus

This appendix described another way to find the normal equation. I won’t mention this in class, and the treatment will leave out some details, but it may be of interest to some of you.

There are many other ways to derive the normal equation  $A^t A \vec{x} = A^t \vec{b}$  from the problem of minimizing  $\|A\vec{x} - \vec{b}\|$ . One is to apply calculus: one way to optimize a differentiable function of several variables is to set all of its partial derivatives equal to 0 (of course, setting all derivatives equal to 0 doesn’t guarantee that what is found is a *minimum* rather than a maximum, a saddle, or another type of critical point, but let’s not worry about this for now).

To begin, let’s writing out the *square* of the quantity being optimized, in terms of dot products. The reason to take the square is that it the minimum will occur for the same choice of  $\vec{x}$ , and we can avoid writing square roots in our formulas.

$$\begin{aligned} \|A\vec{x} - \vec{b}\|^2 &= (A\vec{x} - \vec{b}) \cdot (A\vec{x} - \vec{b}) \\ &= (A\vec{x}) \cdot (A\vec{x}) - 2(A\vec{x}) \cdot \vec{b} + \vec{b} \cdot \vec{b} \end{aligned}$$

Now, let's take the partial derivative  $\frac{\partial}{\partial x_i}$  of this expression (there's one partial derivative for each  $i$ , for  $n$  total). We need two facts, which I'll state without proof (it is worth thinking through why these are true):

1. The product rule is valid for dot products of vector-valued functions:

$$\frac{\partial}{\partial x_i} \left( \vec{f}(\vec{x}) \cdot \vec{g}(\vec{x}) \right) = \left( \frac{\partial}{\partial x_i} \vec{f}(\vec{x}) \right) \cdot \vec{g}(\vec{x}) + \vec{f}(\vec{x}) \cdot \left( \frac{\partial}{\partial x_i} \vec{g}(\vec{x}) \right)$$

2. The partial derivatives of  $A\vec{x}$  are just the columns of  $A$ :

$$\frac{\partial}{\partial x_i} (A\vec{x}) = \vec{A}_i$$

From these two facts, we can compute (using that  $\vec{b}$  is constant):

$$\begin{aligned} \frac{\partial}{\partial x_i} \|A\vec{x} - \vec{b}\|^2 &= \vec{A}_i \cdot (A\vec{x}) + (A\vec{x}) \cdot \vec{A}_i - 2\vec{A}_i \cdot \vec{b} - (A\vec{x}) \cdot \frac{\partial}{\partial x_i} \vec{b} + \frac{\partial}{\partial x_i} (\vec{b} \cdot \vec{b}) \\ &= 2\vec{A}_i \cdot (A\vec{x}) - 2\vec{A}_i \cdot \vec{b} \\ &= 2\vec{A}_i \cdot (A\vec{x} - \vec{b}) \end{aligned}$$

From these, we see that, in order for all the partial derivatives to vanish,  $A\vec{x} - \vec{b}$  must have dot product 0 with all the columns of  $A$ . This amounts to saying that  $A^t(A\vec{x} - \vec{b}) = \vec{0}$ , which rearranges to the normal equation  $A^t A\vec{x} = A^t \vec{b}$ .

**Comment:** A somewhat slicker way to carry out the calculus above is to work with all  $n$  partial derivatives at once, by taking the *gradient* of the expression. To do this requires laying some groundwork about gradients of vector-valued functions. One source I recommend (which is where I originally learned how to use techniques like least-squares in practice) is Andrew Ng's lecture notes on machine learning. Linear regression and the normal equation are discussed (in terms of gradients) in this set of notes:

<http://cs229.stanford.edu/notes/cs229-notes1.pdf>

Those notes also give a useful probabilistic interpretation of least-squares. There's a fair amount of other material you'd need to read a bit about to follow all the details in these notes, but it is worth the effort. I am happy to chat about them at office hours.