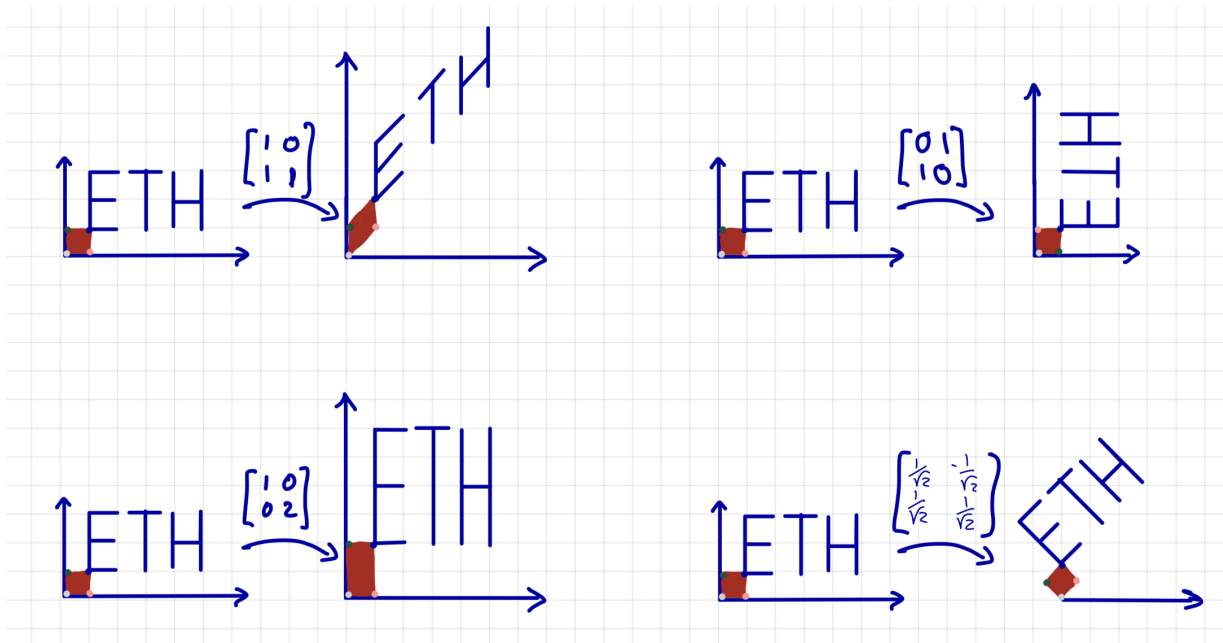


Linear Algebra

Fall 2023

(ETHZ 401-0131-00L)

Lecture Notes Part II



Afonso S. Bandeira
ETH Zürich

Last update on December 17, 2023

“READ ME” FOR PART II

My webpage, with contact information, is: <https://people.math.ethz.ch/~abandeira>

These lecture notes serve as a **continuation**¹ of **Part I**, taught by Prof. Bernd Gärtner, available at https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_I.pdf. Please read the Preface there. Please note there may be some changes in notation. Furthermore, we will try to stay close to the notation in [Str23], but there also be some differences.

[Str23] Gilbert Strang. Introduction to Linear Algebra. Wellesley - Cambridge Press, 6th ed., 2023.

The **course page** has relevant information for the course: <https://ti.inf.ethz.ch/ew/courses/LA23>.

I offer **office hours** (in HG G23.1) almost weekly, feel free to stop by, to chat about the course, Mathematics, Computer Science, or University in general. Most office hours visitors stop by to learn more about research in Mathematics and Computer or Data Science. You can see the schedule on my webpage, on the calendar applet on the left.

There are countless excellent Linear Algebra books with the material covered in this course. For Part II we will roughly continue to follow, in structure and content, [Str23], with some small deviations. I will try to keep the numbering of Chapters/Sections and Sections/Subsections consistent with [Str23] (as far as the deviations allow). See Appendix A for some important preliminaries and some remarks on notation.

Throughout the notes, and the lectures, I will try to motivate some of the material with **Guiding Questions**. For students who would like to explore the topic further, I will include some **Exploratory Challenges** and **Further Reading**, these often will include difficult problems or topics. I will also take some opportunities to share some active **Research Questions** related to the topics we covered (we are still discovering new phenomena in Linear Algebra today and for many years to come!).

After deriving a result, we will often do some **Sanity Checks**, and some things I will leave as a **Challenge**: these should be accessible and of difficulty comparable to homework questions, a \star indicates a harder problem (but still within the scope). On the other hand, **Exploratory Challenges** are generally outside the scope of the course or of substantial higher difficulty.

¹If you are reading these notes and did not follow Part I, please read Appendix A.

Linear Algebra is a beautiful topic, connecting Algebra with Geometry², and has countless applications making it a key component of almost all quantitative pursuits. I sincerely hope you will enjoy the course as much as I enjoy teaching this subject!

MISCELLANEOUS THOUGHTS

I believe the Questions, Sanity Checks, Challenges, etc are very useful to learn the material, but **when you want to review the material**, or do a last read before the exam, you can focus on the Definitions, Propositions, Theorems, etc (and **focus less on the blue parts**).

In many of my side comments (usually in blue), and in some of the [CS Lens Lectures](#), I do not include specific citations, and sometimes use technical terms that you might not have seen before. My goal is, quoting my collaborator Dustin Mixon, “to provide enough breadcrumbs for the interested reader to find more information online”.³ If you would like specific references, tell me a bit more about your interests and I would be happy to point you to some references (different references are better for different takes/interest on each of the topics). While [CS Lens Lectures](#) are not covered in the lecture notes, slides can be accessed in the course website: <https://ti.inf.ethz.ch/ew/courses/LA23/index.html>

As your mathematical level matures over the semester, the notes will have less illustrations and more definitions and mathematical statements. My recommendation is to read the notes with pen & paper next to you and to draw the picture yourself, this “translation” you will be doing — from mathematical statement to picture — will (I believe) help you greatly in the learning of Mathematics!

There are also countless high-quality videos and other content online about Linear Algebra, for example there is also an excellent series of videos by Gil Strang filmed ~15 years ago: <https://www.youtube.com/playlist?list=PLE7DDD91010BC51F8>.

Strang actually retired just a few months ago, at almost 90 years of age! You can see his last lecture online: <https://www.youtube.com/watch?v=1UUte2o2Sn8>

²and Analysis, as you will likely see later in your academic life. For example, when Joseph Fourier invented Fourier Series to developed a theory of heat transfer he was essentially finding good orthonormal bases for functions.

³This itself also provides enough breadcrumbs for you to find the lecture notes I am quoting; they are excellent Linear Algebra lecture notes! (the order and content is somewhat different from our course)

Moreover, there are many excellent animations online giving lots of great intuition on several Linear Algebra topics and phenomena. While it is a great idea to take advantage of this, I would recommend first trying yourself to develop an intuition of the concept/phenomenon (e.g. by drawing a picture) and using these tools only after — use them to improve your intuition, not to create it!

As these Lecture Notes are being continuously updated, and sometimes the discussion in lectures leads us into proving an extra result, or suggests a remark, etc, I will try to add them and not change the numbering of things downstream, I do this by numbering them with +1000.

After each lecture, we post the handwritten notes from lecture on the course website <https://ti.inf.ethz.ch/ew/courses/LA23/index.html>. My suggestion would be to use the Lecture Notes to review the material, not the handwritten notes (which are mainly meant to support my oral exposition).

CONTENTS

“Read me” for Part II	2
Miscellaneous Thoughts	3
4. Orthogonality, Projections, and Least Squares	6
4.2. Projections	6
4.3. Least Squares Approximation	10
4.4. Orthonormal Bases and Gram Schmidt	14
4.5. The Pseudoinverse, also known as Moore–Penrose Inverse	19
5. Linear Transformations and Determinants	23
5.1. The Determinant	28
6. Eigenvalues and Eigenvectors	35
6.0. Complex Numbers	36
6.1. Introduction to Eigenvalues and Eigenvectors	39
6.2. Diagonalizing a Matrix and Change of Basis of a Linear Transformation	49

	5
6.3. Symmetric Matrices and the Spectral Theorem	51
7. Singular Value Decomposition; and some open questions in Linear Algebra	57
7.1. The Singular Value Decomposition	57
7.2. Vector and Matrix Norms	60
7.3. Low-Rank Modelling, Images, Data, and Principal Component Analysis	61
7.10. Some Mathematical Open Problems	62
Acknowledgements	64
Appendix A. Some Important Preliminaries and Remarks on Notation	64
Appendix B. Weekly Schedule	65
Appendix C. CS Lens Lectures (not part of Chapter 7)	65
Appendix D. A “Simple proof” of Fundamental Theorem of Algebra	65

4. ORTHOGONALITY, PROJECTIONS, AND LEAST SQUARES

Guiding Question 1. If we have a system of linear equations that has no solution, how do we find the “solution” that has the smallest error? This question is central in countless applications⁴.

Before diving into systems of equations, we will study Projections of vectors in a subspace.

4.2. Projections.

Definition 4.2.1 (Projection of a vector onto a subspace). *The projection of a vector $b \in \mathbb{R}^m$ on a subspace S (of \mathbb{R}^m) is the point in S that is closest to b . In other words*

$$(1) \quad \text{proj}_S(b) = \underset{p \in S}{\text{argmin}} \|b - p\|.$$

Sanity Check 2. This is only a proper definition if the minimum exists and is unique. Can you show it exists and is unique? (perhaps at the end of the lecture?)

Let us build us some intuition by starting with projections to a line. Let S be the subspace corresponding to the line that goes through the vector a , i.e. $S = \text{Span}(a)$.

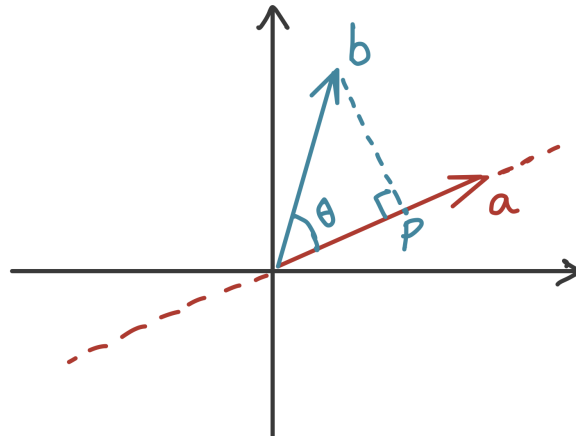


FIGURE 1. Projection on a line.

⁴as you will see later on, it is in a sense what Machine Learning is all about.

The projection p is the vector in the subspace S such that the “error vector” $e = b - p$ is perpendicular to a (i.e. $b - p \perp a$). Since $p \in S$ we have $p = \hat{x}a$, for some $\hat{x} \in \mathbb{R}$. Since $b - p \perp a$ we have $a^\top(b - p) = 0$. Substituting gives

$$a^\top(b - \hat{x}a) = 0 \iff \hat{x} = \frac{a^\top b}{a^\top a} \iff p = \frac{a^\top b}{a^\top a}a \iff p = \frac{aa^\top}{a^\top a}b.$$

Indeed, we have the following Proposition.

Proposition 4.2.2. *Let $a \in \mathbb{R}^m$ be a non-zero vector. The projection of a vector $b \in \mathbb{R}^m$ on $S = \text{Span}(a)$ the span of a , is given by*

$$\text{proj}_S(b) = \frac{aa^\top}{a^\top a}b.$$

Sanity Check 3. The projection of a vector that is already a multiple of a should be the identity operation. Check that this is the case! (and do it later, again, for general subspaces).

For general subspaces the idea is precisely the same. Let S be a subspace in \mathbb{R}^m with dimension n . Let a_1, \dots, a_n be a basis of S , meaning that $S = \text{Span}(a_1, \dots, a_n)$ and $S = C(A)$ where

$$A = \begin{bmatrix} | & | & & | \\ a_1 & a_2 & \cdots & a_n \\ | & | & & | \end{bmatrix}.$$

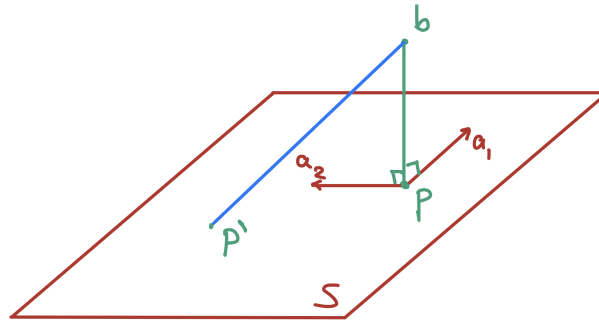


FIGURE 2. Projection on a subspace.

Similarly to the case of a line, it is easy to see (see Figure 2) that the projection p of a vector b on S is such that the error vector $e = b - p$ is perpendicular to each of the a_k 's. To prove this fact rigorously we start by showing existence of such a vector p : take the orthogonal complement S^\perp of S and write b as a sum of $e \in S^\perp$ and $p \in S$, then $e = b - p$ is orthogonal to the subspace S . Now, let us assume that there exists another point p' (as in Figure 2) and note that since $p' - p \in S$ we have that $b - p \perp p' - p$, and so, by Pythagoras' Theorem we have $\|p' - b\|^2 =$

$\|p - p'\|^2 + \|p - b\|^2$, which implies that $\|p' - b\|^2 \geq \|p - b\|^2$ (with equality holding only when $p = p'$).^{5 6}

We just showed that $a_k^\top(b - p) = 0$ for $k = 1, \dots, n$. In matrix-vector notation

$$\begin{bmatrix} | & | & & | \\ a_1 & a_2 & \cdots & a_n \\ | & | & & | \end{bmatrix}^\top (b - p) = 0 \iff A^\top(b - p) = 0.$$

Since $p \in C(A)$ we have $p = A\hat{x}$ for some $\hat{x} \in \mathbb{R}^n$. This means that

$$A^\top(b - A\hat{x}) = 0 \iff A^\top A\hat{x} = A^\top b.$$

We just proved the following Proposition.

Proposition 4.2.3. *The projection p of a vector $b \in \mathbb{R}^m$ on a subspace S with a basis a_1, \dots, a_n can be written as $p = A\hat{x}$ where $\hat{x} \in \mathbb{R}^n$ satisfies the normal equations*

$$A^\top A\hat{x} = A^\top b,$$

where $A = \begin{bmatrix} | & | & & | \\ a_1 & a_2 & \cdots & a_n \\ | & | & & | \end{bmatrix}$ is the matrix whose columns are a basis of S .

If we can show that $A^\top A$ is invertible then we would have $p = A\hat{x} = A(A^\top A)^{-1}A^\top b$. Let's make a detour to show that it is indeed invertible.

Proposition 4.2.4. *$A^\top A$ is invertible if and only if A has linearly independent columns.*

Proof. We show this by showing that $A^\top A$ and A have the same nullspace. This is enough because, since $A^\top A$ is a square matrix it is invertible if and only if its nullspace only has the 0 vector, and A has linearly independent columns if and only if its nullspace only has the 0 vector.⁷

If $x \in N(A)$ then $Ax = 0$ and so $A^\top Ax = 0$, thus $x \in N(A^\top A)$. The other implication is more interesting.

If $x \in N(A^\top A)$ then $A^\top Ax = 0$. This implies that $x^\top A^\top Ax = x^\top 0 = 0$. But $x^\top A^\top Ax = (Ax)^\top (Ax) = \|Ax\|^2$ so Ax must be a vector with norm 0 which implies that $Ax = 0$ and so $x \in N(A)$.

□

⁵We have also, as a byproduct, answered the question in Sanity Check 2.

⁶Sometimes the projection is simply defined as the point on the subspace such that the error vector is orthogonal to the subspace, here we showed the two possible definitions are equivalent.

⁷We usually call a nullspace with only the zero vector, a trivial nullspace.

Corollary 4.2.5. *If A has linearly independent columns then $A^\top A$ is a square matrix, it is invertible and symmetric.*⁸

Now back to deriving a formula for projections: Since the columns of A are a basis they are linearly independent and so $A^\top A$ is indeed invertible. We just proved the following.

Theorem 4.2.6. *Let S be a subspace in \mathbb{R}^m and A a matrix whose columns are a basis of S . The projection of $b \in \mathbb{R}^m$ to S is given by*

$$\text{proj}_S(b) = Pb,$$

where $P = A(A^\top A)^{-1}A^\top$ is the projection matrix.

The matrix $P = A(A^\top A)^{-1}A^\top$ is known as a Projection Matrix, it maps a vector b to its projection Pb on a subspace S . For the case of lines, P was given by $P = \frac{aa^\top}{a^\top a} = a \frac{1}{a^\top a} a^\top$.

Caution! 4. The matrix A (and A^\top) are not necessarily square, and so they don't have inverses. The expression $A(A^\top A)^{-1}A^\top$ **cannot** be simplified by expanding $(A^\top A)^{-1}$ (which would yield $I = P$, this would only make sense if S was all of \mathbb{R}^m and note that, unsurprisingly, this would correspond exactly to the case when A is invertible).

Just as with the "sanity check" above, we should have $P^2 = P$, because if we project a point twice, the second time should not do anything as the point is already in S and indeed

$$P^2 = \left(A(A^\top A)^{-1}A^\top \right)^2 = A(A^\top A)^{-1}A^\top A(A^\top A)^{-1}A^\top = A(A^\top A)^{-1}A^\top = P.$$

Challenge 5. Is $I - P$ a projection? If so, which projection does it correspond to?

Challenge 6. How does the rank of P depend on properties of the subspace S ?

Exploratory Challenge 7. We derived all of the formulas for projections using geometry. If you have taken Analysis/Calculus (I know many of you haven't, but you will in a few months) you can try to re-derive everything using the fact that derivatives at the minimum should be zero. You will see that you will get exactly the same answers.

In lecture, when discussing Figure 2 we explicitly proved the following proposition.

Proposition 4.2.1006. *Let S be a subspace in \mathbb{R}^m with a basis a_1, \dots, a_n . For $v \in \mathbb{R}^m$, v being orthogonal to all vectors in S is equivalent to being orthogonal to a_k for all $1 \leq k \leq n$.*

⁸Corollary is like a Theorem or a Proposition but one that follows directly from another one, this one follows directly from the Proposition above.

Proof. Since a_1, \dots, a_n are in S , if v is perpendicular to all vectors in S , it is in particular perpendicular to a_1, \dots, a_n . On the other hand, any $w \in S$ can be written as $w = \alpha_1 a_1 + \dots + \alpha_n a_n$ and $w^\top v = \alpha_1 a_1^\top v + \dots + \alpha_n a_n^\top v = 0$. \square

Linear Algebra — A. Bandeira (ETHZ) — Week 8 - 2023.11.10 & 2023.11.15

Please find most up to date notes at: https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_II.pdf

4.3. Least Squares Approximation. We go back to the guiding question of what to do when we want to “solve” a linear system that does not have an exact solution. More precisely let us suppose we have a linear system

$$Ax = b,$$

for which no solution x exists (for example, with too many equations, which would happen if $A \in \mathbb{R}^{m \times n}$ and $m > n$). A natural approach is to try to find x for which Ax is as close as possible to b

$$(2) \quad \min_{\hat{x} \in \mathbb{R}^n} \|A\hat{x} - b\|^2.$$

Further Remark 8. This seemingly simple observation is key to countless technologies. Measurement systems often have errors and so it is impossible to find the target object/signal x that satisfies them all exactly, and we look for the one that satisfies them the best possible. In Data Science and Learning Theory we often want to find a predictor that best describes a set of *training data*, but usually no predictor described the data exactly, so we look for the best possible, etc. We’ll see a couple of applications later.

We can solve this problem using the ideas we developed above. What we are looking for is a vector \hat{x} for which the error $e = b - A\hat{x}$ is as small as possible. Since the set of possible vectors $y = A\hat{x}$ is exactly $C(A)$, $A\hat{x}$ is precisely the projection of b on $C(A)$. As we saw in the Section above, this means that

$$A^\top (b - A\hat{x}) = 0.$$

These are known as the *normal equations* and can be rewritten as

$$(3) \quad A^\top A\hat{x} = A^\top b.$$

Remark 4.3.1. For this to make sense it must be that (3) always has a solution. If we think geometrically, it is relatively easy to see that it must, because of how we constructed the normal equations. Can you give a rigorous algebraic proof of this fact? Note that essentially what you are proving is the Proposition below.

Proposition 4.3.2. For any matrix A , $C(A^\top) = C(A^\top A)$.

Challenge 9. Try to prove this Proposition. This can be done in a few different ways. I suggest starting by trying to show that $\text{rank}(A) = \text{rank}(A^\top) = \text{rank}(A^\top A) = \text{rank}(AA^\top)$.

We know that if A has linearly independent columns, then $A^\top A$ is invertible and so we can write $\hat{x} = (A^\top A)^{-1} A^\top b$. We will address the case in which A has dependent columns shortly.

Fact 4.3.3. A minimizer of (2) is also a solution of (3). When A has independent columns the unique minimizer \hat{x} of (2) is given by

$$(4) \quad \hat{x} = (A^\top A)^{-1} A^\top b$$

Exploratory Challenge 10. Similarly to the projections derivation, this derivation can also be done by differentiating (2). Try it.

4.3.2. *Linear Regression — fitting a line to data points.* One of the most common tasks in data analysis is linear regression, to fit a line through data points. Let us consider data points

$$(t_1, b_1), (t_2, b_2), \dots, (t_m, b_m),$$

perhaps representing some attribute b over time t . If the relation between t and b is (at least partly) explained by a linear relationship then it makes sense to search for constants $\alpha_0 \in \mathbb{R}$ and $\alpha_1 \in \mathbb{R}$ such that

$$b_k \approx \alpha_0 + \alpha_1 t_k.$$

See Figure 3. In particular, it is natural to search for α_0 and α_1 that minimize the sum of squares of the errors (“least squares”),

$$\min_{\alpha_0, \alpha_1} \sum_{k=1}^m (b_k - [\alpha_0 + \alpha_1 t_k])^2.$$

In matrix-vector notation

$$(5) \quad \min_{\alpha_0, \alpha_1} \left\| b - A \begin{bmatrix} \alpha_0 \\ \alpha_1 \end{bmatrix} \right\|^2,$$

where

$$b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{m-1} \\ b_m \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_{m-1} \\ 1 & t_m \end{bmatrix}.$$

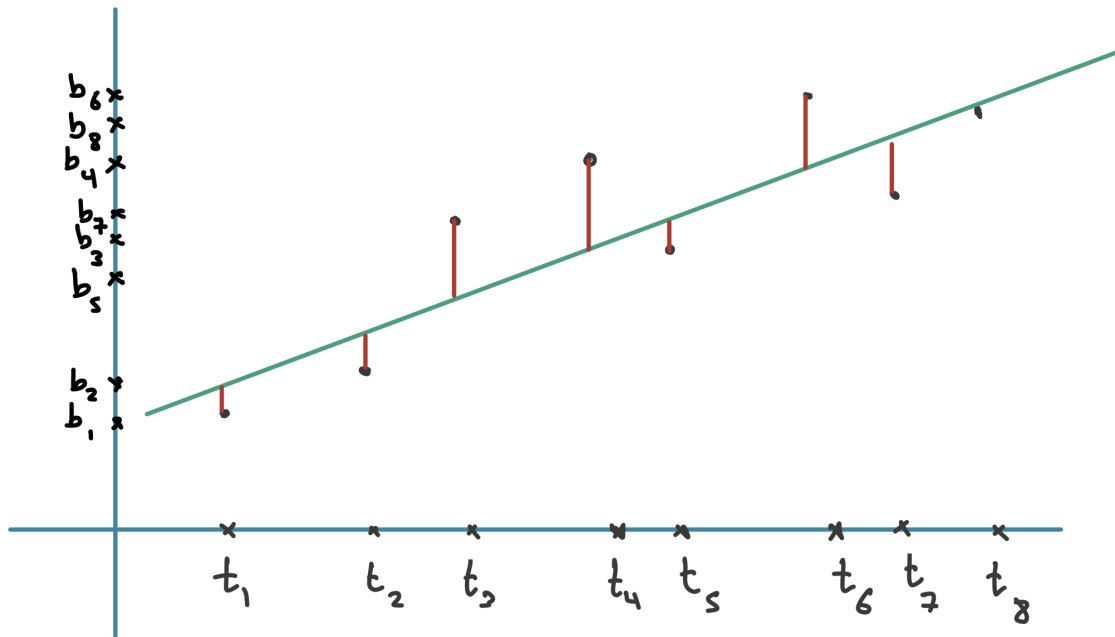


FIGURE 3. Fitting a line to points

As long as A has independent columns (see Remark 4.3.4) the solution to (5) is given by

$$\begin{bmatrix} \alpha_0 \\ \alpha_1 \end{bmatrix} = (A^\top A)^{-1} A^\top b = \begin{bmatrix} m & \sum_{k=1}^m t_k \\ \sum_{k=1}^m t_k & \sum_{k=1}^m t_k^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{k=1}^m b_k \\ \sum_{k=1}^m t_k b_k \end{bmatrix}$$

Remark 4.3.4. *It is worth working out what it means for the columns of A , in this example, to be linearly dependent. It essentially corresponds to all points t_k being the same, which is clearly a degenerate case of linear regression.*

Remark 4.3.5. *If the columns of A are pairwise orthogonal, then $A^\top A$ is a diagonal matrix, which is easy to invert. In this example, the columns of A being orthogonal corresponds to $\sum_{k=1}^m t_k = 0$. We could simply do a change of variables to a new time $t_k^{\text{new}} = t_k - \frac{1}{m} \sum_{i=1}^m t_i$ to achieve this. If indeed $\sum_{k=1}^m t_k = 0$ then the equation above could be easily simplified:*

$$\begin{aligned} \begin{bmatrix} \alpha_0 \\ \alpha_1 \end{bmatrix} &= \begin{bmatrix} m & 0 \\ 0 & \sum_{k=1}^m t_k^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_{k=1}^m b_k \\ \sum_{k=1}^m t_k b_k \end{bmatrix} = \begin{bmatrix} \frac{1}{m} & 0 \\ 0 & \frac{1}{\sum_{k=1}^m t_k^2} \end{bmatrix} \begin{bmatrix} \sum_{k=1}^m b_k \\ \sum_{k=1}^m t_k b_k \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{m} \sum_{k=1}^m b_k \\ (\sum_{k=1}^m t_k b_k) / (\sum_{k=1}^m t_k^2) \end{bmatrix}, \end{aligned}$$

this is an instance where having orthogonal vectors is beneficial. In this next Section we will see how to build orthonormal basis for subspaces, and some of the many benefits they have.

Challenge 11. Try to work out the actual change of variables that makes the t_k 's add up to zero and derive a formula for fitting a line to points without the assumption in Remark 4.3.5

Example 4.3.6 (Fitting a Parabola). We can use Linear Algebra to do fits of many other curves (or functions), not just lines. If we believe the relationship between t_k and b_k is quadratic we could attempt to fit a Parabola:

$$b_k \approx \alpha_0 + \alpha_1 t_k + \alpha_2 t_k^2.$$

While this isn't a linear function in t_k , this is still a linear function on the coefficients α_0 , α_1 , and α_2 , and this is what is important. Similarly as with linear regression, it is natural to attempt to minimize

$$(6) \quad \min_{\alpha_0, \alpha_1, \alpha_2} \left\| b - A \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \alpha_2 \end{bmatrix} \right\|^2,$$

where

$$b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{m-1} \\ b_m \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} 1 & t_1 & t_1^2 \\ 1 & t_2 & t_2^2 \\ \vdots & \vdots & \vdots \\ 1 & t_{m-1} & t_{m-1}^2 \\ 1 & t_m & t_m^2 \end{bmatrix},$$

and we can use the technology we developed in this section to solve this problem as well.

Challenge 12. Try to work out the example of fitting a parabola further. What is $A^T A$? When is $A^T A$ diagonal?

Further Reading 13. There is a whole (beautiful) area of Mathematics related to studying so-called *Orthogonal Polynomials*. The basic idea can be already hinted at from these examples: In the example of the parabola we wrote a function of t as a linear combination of the polynomials 1 , t , and t^2 . But we could have picked other polynomials, we could have e.g. written something like $b \approx \alpha'_0 + \alpha'_1(t - 2023) + \alpha'_2(t^2 + t)$, and a particularly good choice (that would depend on the distribution of the points t_k) might have resulted in a diagonal matrix $A^T A$... search "orthogonal polynomials" online to learn more.

Further Reading 14. A lot of Machine Learning includes Linear Regression as a key component. The idea is to create, find, or *learn* features of the data points. Given n data points t_1, \dots, t_n (which now can be perhaps pixel images, rather than just timepoints) we might want to do classification (for example, in the case of images, maybe we want a function that is large when the picture has a dog in it and small when it has a cat in it). It is hard to imagine that this can be done with

a linear fit, but if we build good feature vectors $\varphi(t_k) \in \mathbb{R}^p$ for very large p then the function can depend on all coordinates of $\varphi(t_k)$ (the p features) and this is incredible powerful. There are several ways to construct features, a bit over a decade ago they were sometimes handmade, now they are often learned (this is in a sense what Deep Learning does). Another important way to build (or compute with) features are the so-called Kernel Methods, you can see more in the CS Lens Lecture (Appendix C).

4.4. Orthonormal Bases and Gram Schmidt. When we think of (or draw) a basis of a subspace, we tend to think of (or draw) vectors that are orthogonal (have an angle of 90°) and that have the same length (length 1). Indeed, these bases have many advantages, this section is about these bases, some of their advantages, and how to find them.

Definition 4.4.1 (Orthonormal vectors). *We say n vectors $q_1, \dots, q_n \in \mathbb{R}^m$ are orthonormal if they are orthogonal and have norm 1. In other words, for all $i, j = 1, \dots, n$*

$$q_i^T q_j = \delta_{ij},$$

where δ_{ij} is the Kronecker delta

$$(7) \quad \delta_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j. \end{cases}$$

If Q is the matrix whose columns are the vectors q_i 's, then the condition that the vectors are orthonormal can be rewritten as $Q^T Q = I$.

Caution! 15. Q may not be a square matrix, and so it is not necessarily the case that $QQ^T = I$.

Example 4.4.2. *A classical example of an orthonormal set of vectors is the canonical basis, $e_1, \dots, e_n \in \mathbb{R}^n$ where e_i is the vector with a 1 in the i -th entry and 0 in all other entries: $(e_i)_j = \delta_{ij}$.*

When Q is a square matrix then $Q^T Q = I$ implies also that $QQ^T = I$ and so $Q^{-1} = Q^T$. We call such matrices *orthogonal matrices*. This corresponds to the case when the q_i 's are an orthonormal basis for all of \mathbb{R}^n .

Definition 4.4.3 (Orthogonal Matrix). *A square matrix $Q \in \mathbb{R}^{n \times n}$ is an Orthogonal Matrix when $Q^T Q = I$. In this case, $QQ^T = I$, $Q^{-1} = Q^T$, and the columns of Q form an orthonormal basis for \mathbb{R}^n .*

Remark 4.4.4. *It is often useful to think of an $m \times n$ matrix A as a function from \mathbb{R}^n to \mathbb{R}^m , that takes $x \in \mathbb{R}^n$ to $Ax \in \mathbb{R}^m$.*

$$\begin{aligned} A : \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ x &\rightarrow Ax \end{aligned}$$

Later in the course, when we discuss Linear Transformations, we will, among other things, discuss which functions can be described by a matrix this way (and some properties of these functions/transformations). For now, let us just keep in mind that a matrix can be thought of as a function. It is also worth noting that this explains why in some Linear Algebra books the Nullspace is called the Kernel (it is the set of vectors x that are mapped to 0 by this function) and the Column Space is called Image, or Range, as it is the set of vectors in \mathbb{R}^m that is the image of this function.

Example 4.4.5. The 2×2 matrix Q that corresponds to rotating, counterclockwise, the plane by θ ,

$$R_\theta = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

is an orthogonal matrix.

Challenge 16. Prove that the rotation matrices R_θ are orthogonal matrices. Are there other 2×2 orthogonal matrices? If so, can you describe them all?

Example 4.4.6. Permutation matrices are another example of orthogonal matrices.

Challenge 17. Show that indeed permutation matrices are orthogonal matrices

Exploratory Challenge 18. One of the most important structures in Algebra is that of a group. The set of Permutations of n elements is an example of a group, two permutations can be composed to form another permutation and for every permutation there is one corresponding to undoing it (called the inverse). Permutation matrices represent the permutations, composing corresponds to matrix multiplication and the inverse permutation corresponds to the matrix inverse of the permutation matrix. There is a whole field of Mathematics, called Representation Theory, that studies matrix representations of groups (and in many important cases the matrices involved are orthogonal). Can you come up with a matrix representation of addition modulo 2? What about addition modulo 5?

Challenge 19 (*). Show that for every permutation matrix P there exists a positive integer k such that $P^k = I$.

Proposition 4.4.7. Orthogonal matrices preserve norm and inner product of vectors. In other words, if $Q \in \mathbb{R}^{n \times n}$ is orthogonal then, for all $x, y \in \mathbb{R}^n$

$$\|Qx\| = \|x\| \text{ and } (Qx)^\top (Qy) = x^\top y$$

Proof. To show the second inequality note that, for $x, y \in \mathbb{R}^n$ we have that $(Qx)^\top(Qy) = x^\top Q^\top Qy = x^\top Iy = x^\top y$. To show the first equality note that, since for $x \in \mathbb{R}^n$ we have that $\|Qx\| \geq 0$ and $\|x\| \geq 0$, it suffices to show that the squares are equal and indeed $\|Qx\|^2 = (Qx)^\top(Qx) = x^\top x = \|x\|^2$. \square

4.4.1. *Projections with Orthonormal Basis.* One of the advantages of orthonormal basis is that projections become much simpler. The reason is simple: when we discussed projections and least squares, many of the expressions we derived included $A^\top A$, but in the case when A has orthonormal columns, these all simplify as $A^\top A = I$. We collect these observations in the following proposition.

Proposition 4.4.8. *Let S be a subspace of \mathbb{R}^m and q_1, \dots, q_n be an orthonormal basis for S . Let Q be the $m \times n$ matrix whose columns are the q_i 's; $Q = \begin{bmatrix} q_1 & \cdots & q_n \end{bmatrix}$. Then the Projection Matrix that projects to S is given by QQ^\top and the Least Squares solution to $Qx = b$ is given by $\hat{x} = Q^\top b$.*

Remark 4.4.9. *When Q is a square matrix then the projection QQ^\top is simply the identity (corresponding to projecting to the entire ambient space \mathbb{R}^n). Even in this seemingly trivial instance, it is useful to look closer at what this operation does: For a vector $x \in \mathbb{R}^n$ it gives*

$$x = q_1 \left(q_1^\top x \right) + q_2 \left(q_2^\top x \right) + \cdots + q_n \left(q_n^\top x \right).$$

It is writing x as a linear combination of the orthonormal basis $\{q_i\}_{i=1}^n$ (as we will see later this is sometimes referred to as a change of basis).⁹

4.4.2. *Gram-Schmidt Process.* Hopefully by now you are convinced that orthonormal basis are useful, now we discuss how to construct them. Fortunately, there is a relatively simple process to construct orthonormal bases, that will also suggest a new matrix factorization.

The idea is simple: If we have 2 linearly independent vectors a_1 and a_2 which span a subspace S , it is straightforward to transform them into an orthonormal basis of S : we first normalize a_1 : $q_1 = \frac{a_1}{\|a_1\|}$, then subtract from a_2 a multiple of q_1 so that it becomes orthogonal to q_1 , followed by a normalization step:

$$q_2 = \frac{a_2 - (a_2^\top q_1)q_1}{\|a_2 - (a_2^\top q_1)q_1\|}.$$

⁹There are countless instances in which doing this operation is beneficial, for example one of the most important algorithms, the *Fast Fourier Transform*, is an instance of this operation.

Let us check that indeed these vectors are orthonormal: By construction they have unit norm, and

$$q_1^\top q_2 = q_1^\top \frac{a_2 - (a_2^\top q_1)q_1}{\|a_2 - (a_2^\top q_1)q_1\|} = \frac{q_1^\top a_2 - (a_2^\top q_1)q_1^\top q_1}{\|a_2 - (a_2^\top q_1)q_1\|} = \frac{0}{\|a_2 - (a_2^\top q_1)q_1\|} = 0.$$

Note that the denominator is not zero because a_1 and a_2 are linearly independent; and that, since q_1 has unit norm, $(a_2^\top q_1)q_1 = \text{proj}_{\text{Span}(q_1)}(a_2)$.

For more vectors, the idea is to this process recursively, by removing from a vector a_{k+1} the projection of it on the subspace spanned by the k vectors before it. More formally:

Algorithm 4.4.10. [Gram-Schmidt Process] Given n linearly independent vectors a_1, \dots, a_n that span a subspace S , the Gram-Schmidt process constructs q_1, \dots, q_n the following way:

- $q_1 = \frac{a_1}{\|a_1\|}$.
- For $k = 2, \dots, n$ do

$$q'_k = a_k - \sum_{i=1}^{k-1} (a_k^\top q_i)q_i$$

$$q_k = \frac{q'_k}{\|q'_k\|}.$$

Theorem 4.4.11 (Correctness of Gram-Schmidt). Given n linearly independent vectors a_1, \dots, a_n , the Gram-Schmidt process outputs an orthonormal basis for the span of a_1, \dots, a_n .

Proof. ¹⁰ We prove this by induction. ¹¹ Let S_k be the subspace spanned by a_1, \dots, a_k . Then $S = S_n$. We will show, by induction, that q_1, \dots, q_k are an orthonormal basis for S_k . It is enough to show that they are orthonormal and are in S_k since orthonormality implies linearly independence and S_k has dimension k .

For the base case, note that $\|q_1\| = 1$ and q_1 is a multiple of a_1 and so $q_1 \in S_1$.

Now we assume the hypothesis for $i = 1, \dots, k-1$ and prove it for k . By the hypothesis q_1, \dots, q_{k-1} are orthonormal, so we have to show that $\|q_k\| = 1$ and that $q_i^\top q_k = 0$ for all $1 \leq i \leq k-1$.

- Since a_k is linearly independent from the other original vectors it is not in S_{k-1} and so $q'_k \neq 0$. Thus $\|q_k\| = 1$.
- By construction $a_k \in S_k$ and so $q_k \in S_k$.
- Let $1 \leq j \leq k-1$. Since q_1, \dots, q_{k-1} are orthonormal, we have

$$q_j^\top \left(a_k - \sum_{i=1}^{k-1} (a_k^\top q_i)q_i \right) = q_j^\top a_k - \sum_{i=1}^{k-1} (a_k^\top q_i)q_j^\top q_i = q_j^\top a_k - (a_k^\top q_j) = 0,$$

$$\text{and } q_j^\top q_k = \frac{1}{\|q'_k\|} q_j^\top q'_k = 0.$$

¹⁰This is a good Theorem to try to prove yourself before reading the proof.

¹¹Since this is our first proof by Induction, we will do it slowly.

Challenge 20. Try to do the Gram-Schmidt process for the columns of

$$\begin{bmatrix} 1 & 2 & 3 & 0 \\ 0 & 4 & 5 & 6 \\ 0 & 0 & 7 & 8 \\ 0 & 0 & 0 & 9 \end{bmatrix}.$$

Is it the case that the Gram-Schmidt process of the columns of an upper triangular matrix (with non-zero diagonal elements) is always a subset of the canonical basis? Can you come up with an example of a set of vectors for which Gram-Schmidt does not output elements of the canonical basis?

Linear Algebra — A. Bandeira (ETHZ) — Week 9 - 2023.11.17 & 2023.11.22

Please find most up to date notes at: https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_II.pdf

Gram-Schmidt actually provides us with a new matrix factorization. Let A be an $m \times n$ matrix with linearly independent columns a_1, \dots, a_n and Q the $m \times n$ matrix whose columns are q_1, \dots, q_n as outputted by Algorithm 4.4.10. Let $R = Q^\top A$, since each q_k is orthogonal to every a_i for $i < k$ we have that R is upper triangular. Q is not necessarily a square matrix, and so not necessarily invertible. But QQ^\top is the projection on the span of the q_i 's and thus also on the a_i 's, this means that $QQ^\top A = A$ and so we have that $QR = QQ^\top A = A$. We call $A = QR$ the QR decomposition.

Definition 4.4.12 (QR decomposition). *Let A be an $m \times n$ matrix with linearly independent columns the QR decomposition is given by*

$$A = QR,$$

where Q is an $m \times n$ matrix with orthonormal columns (they are the output of Gram Schmidt, Algorithm 4.4.10, on the columns of A) and R is an upper triangular matrix given by $R = Q^\top A$.

Remark 4.4.1012. *Note that R is a square matrix ($n \times n$), and since the columns of A are linearly independent we have $N(A) = \{0\}$ and so, since $A = QR$, we have also $N(R) = \{0\}$ and so both R and R^\top are invertible.*

Fact 4.4.13. *The QR decomposition greatly simplifies calculations involving Projections and Least Squares.*

- *Since the $C(A) = C(Q)$ then projections on $C(A)$ can be done with Q which means they are given by $\text{proj}_{C(A)}(b) = QQ^\top b$.*

- The least squares solution to $Ax = b$ is \hat{x} solution of the normal equations (recall (3))

$$A^\top A \hat{x} = A^\top b.$$

Furthermore, $A^\top A = (QR)^\top (QR) = R^\top Q^\top QR = R^\top R$, and so we can write

$$(8) \quad R^\top R \hat{x} = R^\top Q^\top b.$$

Since R has independent columns (is full column rank) then $N(R) = \{0\}$ and so we can simplify (8) to

$$(9) \quad R \hat{x} = Q^\top b,$$

which can be efficiently solved by back-substitution since R is a triangular matrix.

4.5. The Pseudoinverse, also known as Moore–Penrose Inverse. The goal of this Section is to construct an analogue to the inverse of a matrix A for matrices that have no inverse, this is called the Pseudoinverse, or the Moore-Penrose Inverse, and we will denote it by A^\dagger . It is also commonly denoted by A^+ .

Guiding Question 21. While not all matrices are A invertible, we saw that we can still aim to find the (or a) vector x such that Ax is as close as possible to a target vector b . Can we develop this idea to define a “pseudoinverse” for any matrix A , a matrix that is, in a sense, closest to being an inverse for A ? What should “closest to being an inverse” even mean?

There are (at least) three issues we need to overcome to try to define a *pseudoinverse* for a non-invertible matrix A : (i) For some vectors b there might not be a vector x such that $Ax = b$, (ii) For some vectors b there may be more than one x such that $Ax = b$ and we would have to pick one, and (iii) even if we make such choices, it is not clear that such operation will correspond to multiplying by a matrix A^\dagger .

Let $A \in \mathbb{R}^{m \times n}$ be an $m \times n$ matrix. There are a couple of different ways we could try to define a *pseudoinverse* A^\dagger for a non-invertible matrix A . Let us start by building on what we discussed on Section 4.3 (Least Squares Approximations), if the columns of A are linearly independent that it would make sense to build A^\dagger such that $A^\dagger b$ is the Least Squares Solution $\hat{x} = (A^\top A)^{-1} A^\top b$ (the vector \hat{x} such that $A\hat{x}$ is as close as possible to b), and so for matrices A with independent columns we will define $A^\dagger = (A^\top A)^{-1} A^\top$. This motivates the following definition.

Definition 4.5.1 (Pseudoinverse for matrices with full column rank). *For $A \in \mathbb{R}^{m \times n}$ with $\text{rank}(A) = n$ we define the pseudo-inverse $A^\dagger \in \mathbb{R}^{n \times m}$ of A as*

$$A^\dagger = (A^\top A)^{-1} A^\top.$$

Proposition 4.5.2. For $A \in \mathbb{R}^{m \times n}$ with $\text{rank}(A) = n$, the pseudoinverse A^\dagger is a left inverse of A , meaning that $A^\dagger A = I$.

Proof. Since $\text{rank}(A) = n$, $A^\top A$ is invertible. Furthermore, $A^\dagger A = (A^\top A)^{-1} A^\top A = I$. \square

Let us now consider the case for which the rows are linearly independent (in other words, $A \in \mathbb{R}^{m \times n}$ is full row rank; or equivalently $\text{rank}(A) = m$). One natural way to define pseudoinverse is by noting that A^\top is full column rank and to define A^\dagger as

$$\left((A^\top)^\dagger \right)^\top = \left(\left((A^\top)^\top (A^\top) \right)^{-1} (A^\top)^\top \right)^\top = \left((AA^\top)^{-1} A \right)^\top = A^\top (AA^\top)^{-1}.$$

Definition 4.5.3 (Pseudoinverse for matrices with full row rank). For $A \in \mathbb{R}^{m \times n}$ with $\text{rank}(A) = m$ we define the pseudo-inverse $A^\dagger \in \mathbb{R}^{n \times m}$ of A as

$$A^\dagger = A^\top (AA^\top)^{-1}.$$

Proposition 4.5.4. For $A \in \mathbb{R}^{m \times n}$ with $\text{rank}(A) = m$, the pseudoinverse A^\dagger is a right inverse of A , meaning that $AA^\dagger = I$.

Proof. Since $\text{rank}(A) = m$, AA^\top is invertible. Furthermore, $AA^\dagger = AA^\top (AA^\top)^{-1} = I$. \square

Let us try to understand what A^\dagger is achieving for full row rank matrices A . Since A is full row rank, for all $b \in \mathbb{R}^m$, there exists $x \in \mathbb{R}^n$ such that $Ax = b$. The issue is that there are potentially many such vectors. A natural strategy in this case is to pick, among all such vectors, the one with smallest norm.¹² In other words to solve

$$(10) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & \|x\|^2 \\ \text{s.t.} \quad & Ax = b, \end{aligned}$$

where s.t. stands for “subject to” or “such that”. If x_1 and x_2 are vectors such that $Ax_1 = Ax_2 = b$ then $x_1 - x_2 \in \mathcal{N}(A)$, and conversely, if $Ax = b$ and $y \in \mathcal{N}(A)$ then $A(x + y) = b$. Thus, given one vector x_1 such that $Ax_1 = b$ the set of solutions to $Ax = b$ are all vectors of the form $x_1 + y$ where $y \in \mathcal{N}(A)$. So we would like to find the minimum $\|x_1 + y\|$ among all vectors $y \in \mathcal{N}(A)$. Let us write $x_1 = \left(x_1 - \text{proj}_{\mathcal{N}(A)}(x_1) \right) + \text{proj}_{\mathcal{N}(A)}(x_1)$. Since $y \in \mathcal{N}(A)$ we have that

¹²This idea, of picking the smallest (or simplest) solution among many possibilities goes far beyond Linear Algebra and is known as “regularization” in Statistics, Machine Learning, Signal Processing, and Image Processing, etc. It can be viewed as a mathematical version of the famous “Occam’s razor” principle in Philosophy.

$(x_1 - \text{proj}_{N(A)}(x_1)) \perp (y + \text{proj}_{N(A)}(x_1))$ and so, by Pythagoras,

$$\begin{aligned} \|x_1 + y\|^2 &= \left\| (x_1 - \text{proj}_{N(A)}(x_1)) + \text{proj}_{N(A)}(x_1) + y \right\|^2 \\ &= \left\| x_1 - \text{proj}_{N(A)}(x_1) \right\|^2 + \left\| \text{proj}_{N(A)}(x_1) + y \right\|^2, \end{aligned}$$

and so picking $y = -\text{proj}_{N(A)}(x_1)$ yields the smallest norm solution. Since the vectors orthogonal to $N(A)$ are precisely the vectors that are in the row space of A , $C(A^\top)$. We just proved:

Proposition 4.5.5. *For a full row rank matrix A , the (unique) solution to (10) is given by the vector $\hat{x} \in C(A^\top)$ that satisfies the constraint $A\hat{x} = b$.*

A^\dagger is precisely the matrix that “takes b to \hat{x} solution of (10)”.

Proposition 4.5.6. *For a full row rank matrix A , the (unique) solution to (10) is given by the vector $\hat{x} = A^\dagger b$.*

Proof. By using Proposition 4.5.5 we just need to show that $\hat{x} = A^\dagger b$ satisfies $A\hat{x} = b$ and that $\hat{x} = A^\dagger b$ is in $C(A^\top)$. Both these are easy to verify: $A\hat{x} = AA^\dagger b = AA^\top(AA^\top)^{-1}b = b$ and $\hat{x} = A^\dagger b = A^\top((AA^\top)^{-1}b)$ and so $\hat{x} \in C(A^\top)$. \square

Guiding Question 22. We would like to define A^\dagger for all matrices, not just full rank matrices. A natural construction would be to try to define A^\dagger to be the matrix that takes a vector b to the smallest norm solution of the normal equations (3).

To define pseudoinverse of a non full rank matrix A we can do it via de $A = CR$ decomposition (recall from Part I of the course and/or see Appendix A(2)). For $A \in \mathbb{R}^{m \times n}$, with $\text{rank}(A) = r$, the CR decomposition writes $A = CR$ where $C \in \mathbb{R}^{m \times r}$ has the first r linearly independent columns of A and $R \in \mathbb{R}^{r \times n}$ is upper triangular. Note that C is full column rank and R is full row rank.

Definition 4.5.7 (Pseudoinverse for all matrices). *For $A \in \mathbb{R}^{m \times n}$, with $\text{rank}(A) = r$, with CR decomposition $A = CR$ we define the pseudoinverse A^\dagger as*

$$A^\dagger = R^\dagger C^\dagger,$$

which can be rewritten as

$$A^\dagger = R^\top (RR^\top)^{-1} (C^\top C)^{-1} C^\top = R^\top (C^\top CRR^\top)^{-1} C^\top = R^\top (C^\top AR^\top)^{-1} C^\top.$$

The following proposition shows that indeed this definition achieves what was asked in Guiding Question 22.

Proposition 4.5.8. Given $A \in \mathbb{R}^{m \times n}$ and a vector $b \in \mathbb{R}^n$, the (unique) solution to

$$(11) \quad \begin{aligned} \min_{x \in \mathbb{R}^n} \quad & \|x\|^2 \\ \text{s.t.} \quad & A^\top Ax = A^\top b, \end{aligned}$$

is given by $\hat{x} = A^\dagger b$.

Proof. Let r be the rank of A and $A = CR$ with $C \in \mathbb{R}^{m \times r}$ and $R \in \mathbb{R}^{r \times n}$. Then $\hat{x} = A^\dagger b = R^\top (C^\top AR^\top)^{-1} C^\top b$. Thus,

$$A^\top A \hat{x} = A^\top AR^\top (C^\top AR^\top)^{-1} C^\top b = R^\top C^\top AR^\top (C^\top AR^\top)^{-1} C^\top b = R^\top C^\top b = A^\top b.$$

Using Proposition 4.5.5, to show that it is the smallest norm solution we just need to show that $\hat{x} \in C(A^\top A)$, but by Proposition 4.3.2 it is enough to show that $\hat{x} \in C(A^\top)$ and since $C(A^\top) = C(R^\top)$ we have that $\hat{x} = R^\top (C^\top AR^\top)^{-1} C^\top b \in C(A^\top)$. \square

In this proof, the only property of the matrices CR we used is that $A = CR$ and both C and R are full rank. So we have actually shown that we can compute the pseudoinverse from any full rank factorization, not just specifically the CR decomposition. We write it here as a proposition.

Proposition 4.5.9. For $A \in \mathbb{R}^{m \times n}$, with $\text{rank}(A) = r$, and let $S \in \mathbb{R}^{m \times r}$ and $T \in \mathbb{R}^{r \times n}$ such that $A = ST$. Then,

$$A^\dagger = T^\dagger S^\dagger.$$

Remark 4.5.1009. Note that If $A = ST$ and $\text{rank}(A) = r$ then $\text{rank}(S) \geq r$ and $\text{rank}(T) \geq r$ and so the matrices ST in Proposition 4.5.9 are indeed full rank (either full column rank or full row rank).

Proposition 4.5.10. Given $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$, we have

- (1) $(AB)^\dagger = B^\dagger A^\dagger$, as long as $\text{rank}(A) = \text{rank}(B) = n$.
- (2) $(A^\top)^\dagger = (A^\dagger)^\top$,
- (3) AA^\dagger is symmetric, and is the projection matrix for projection on $C(A)$,
- (4) $A^\dagger A$ is symmetric, and is the projection matrix for projection on $C(A^\top)$.

Challenge 23. Prove Proposition 4.5.10. (Hint: use Proposition 4.5.9).

Challenge 1023. Given $A \in \mathbb{R}^{m \times n}$, show that

$$AA^\dagger A = A \quad \text{and} \quad A^\dagger AA^\dagger = A^\dagger.$$

Proposition 4.5.11. Let $A \in \mathbb{R}^{m \times n}$ be a matrix and recall that $C(A)$ and $C(A^\top)$ denote respectively its column and row spaces. When $A : x \rightarrow Ax$ is viewed as a function from $C(A^\top)$ to $C(A)$ it is a bijection. In other words, for all $b \in C(A)$ there is one and only one $x \in C(A^\top)$ such that $Ax = b$.

Challenge 24. Prove Proposition 4.5.11

Further Remark 1024. A different way to define the Pseudo-Inverse of a matrix A is to ask for a matrix A^\dagger that satisfies the conditions in the Challenge 1023 and that both AA^\dagger and $A^\dagger A$ are symmetric. It is nontrivial, but it turns out these conditions are enough to define A^\dagger .

5. LINEAR TRANSFORMATIONS AND DETERMINANTS

Further Remark 25. In this part of the course we slightly deviate from [Str23] and will introduce Linear Transformations, before Determinants. To keep the numbering compatible [Str23] we number the section on Linear Transformations as 5.0.2 (this material is in Chapter 8 of [Str23]).

As we pointed out in Remark 4.4.4, we can view an $m \times n$ matrix A as a function from \mathbb{R}^n to \mathbb{R}^m , that takes $x \in \mathbb{R}^n$ to $Ax \in \mathbb{R}^m$

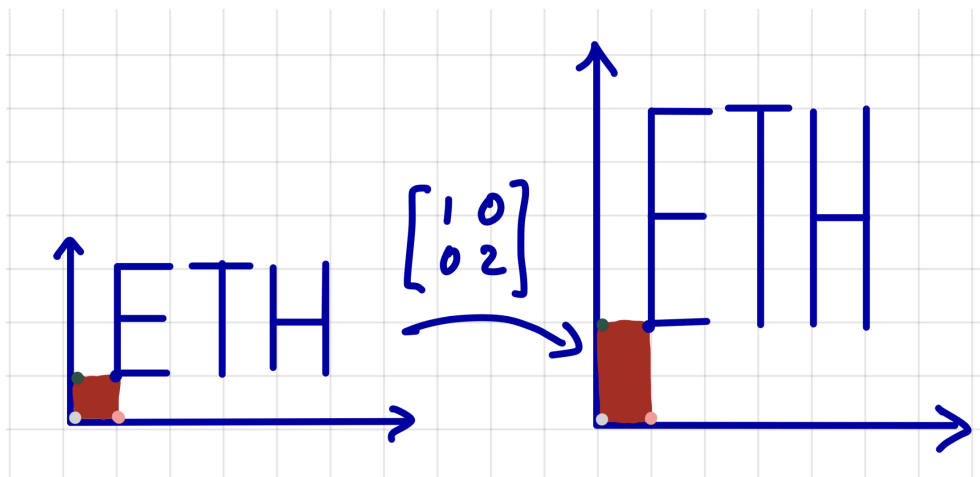
$$\begin{aligned} A : \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ x &\rightarrow Ax. \end{aligned}$$

These functions have a very important property, they are linear: for all $x_1, x_2 \in \mathbb{R}^n$ and any $\alpha \in \mathbb{R}$ we have $A(x_1 + x_2) = Ax_1 + Ax_2$ and $A(\alpha x_1) = \alpha Ax_1$. We will call functions satisfying these properties *Linear Transformations*.

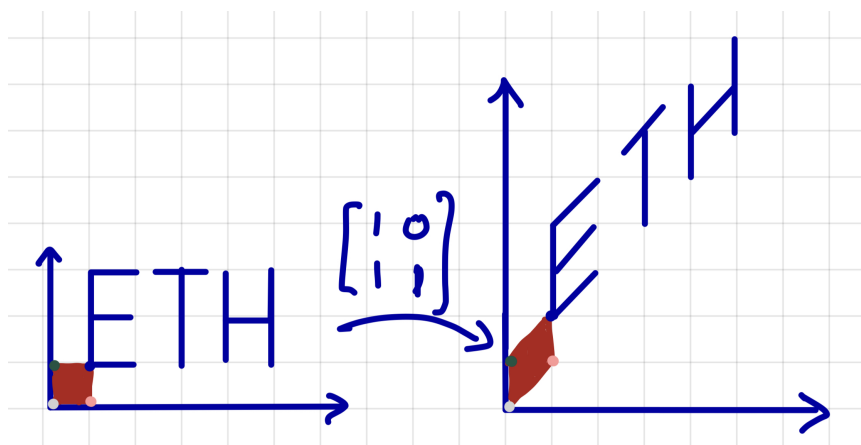
5.0.1. Linear Transformations as transformations of \mathbb{R}^n .

We will now focus on linear transformations from \mathbb{R}^n to itself, corresponding to square matrices $A \in \mathbb{R}^{n \times n}$. Instead of thinking simply of how $x \rightarrow Ax$ maps a vector x to Ax , it is useful to think of this transformation being applied to all of \mathbb{R}^n and to view it as a transformation of the entire \mathbb{R}^n . Let us focus on \mathbb{R}^2 for better visualization and look at a few examples.

Example 5.0.1 (Stretch). The matrix $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ corresponds to stretching by a factor of 2 in the vertical axis. Notice: the first column of A is the image of e_1 and the second the image of e_2 .

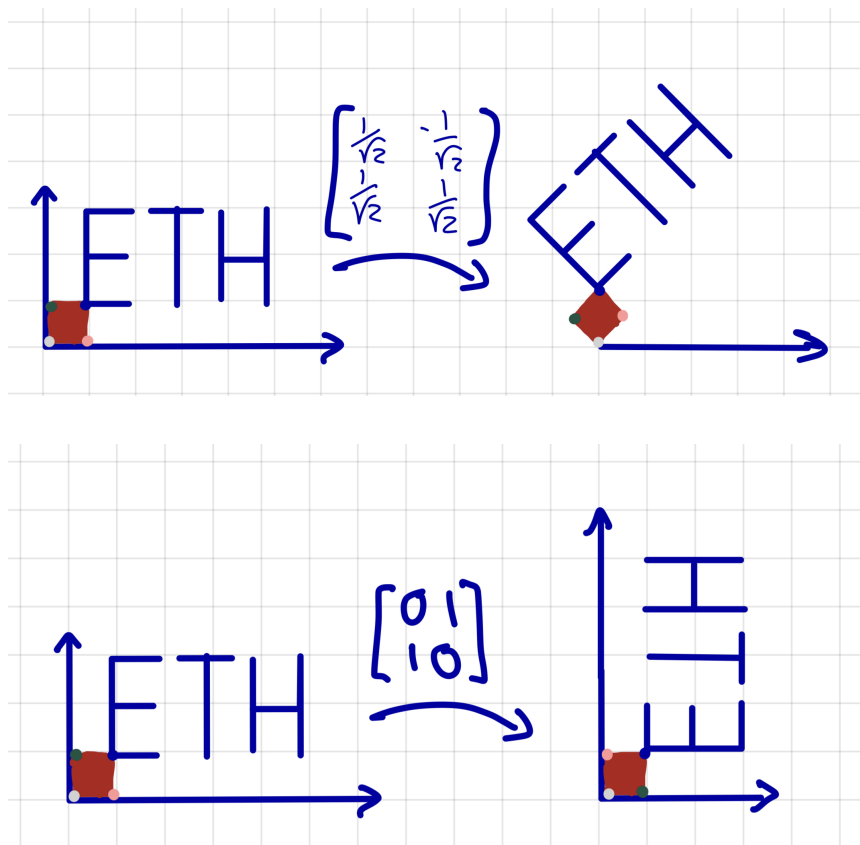


Example 5.0.2 (Shear). The matrix $A = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ corresponds to a shearing transformation given by $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_1 + x_2 \end{bmatrix}$.



Example 5.0.3 (Rotation). The matrix $A = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$ corresponds to a counter-clockwise rotation by $\frac{\pi}{4}$ (or 45°)

Example 5.0.4 (Reflection). The matrix $A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ corresponds to a reflection by the diagonal line $x_2 = x_1$.



Challenge 26. Draw a few more linear transformations. Draw also one corresponding to a non-invertible matrix A . Try to draw one in \mathbb{R}^3 as well.

Challenge 27. Show that a linear transformation takes a triangle to either a triangle, a line segment connecting two points, or a point. What can you say about the rank or invertibility of the corresponding matrix, depending on which of the three objects is the image of a triangle?

5.0.2. Definition of Linear Transformations.

We now treat linear transformations more formally and write the definition of Linear Transformation for any vector space U and V even though in this section we will only treat $U = \mathbb{R}^n$ and $V = \mathbb{R}^m$ (so you can, for now, fully restrict your attention to this setting). Later on, this more general definition will allow us, for example, to discuss linear transformations between subspaces of \mathbb{R}^n and \mathbb{R}^m .

Definition 5.0.5 (Linear Transformation). *Given two vector spaces U and V , a Linear Transformation is a function $T : U \rightarrow V$ such that, for all $u_1, u_2 \in U$ and $\alpha \in \mathbb{R}$ we have*

$$T(u_1 + u_2) = T(u_1) + T(u_2)$$

and

$$T(\alpha u_1) = \alpha T(u_1).$$

Before proceeding let us “collect” a few facts about Linear Transformations.

Proposition 5.0.6. *Let $T : U \rightarrow V$ be a linear transformation and k a positive integer. For all $u_1, \dots, u_k \in U$ and $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ we have*

$$T(\alpha_1 u_1 + \dots + \alpha_k u_k) = \alpha_1 T(u_1) + \dots + \alpha_k T(u_k).$$

Challenge 28. Prove Proposition 5.0.6, iteratively (by induction) using the properties of Linear Transformations.

The most central implication of Proposition 5.0.6 is the fact that the value of T in a basis of U fully determines T .

Proposition 5.0.7. *Let $T : U \rightarrow V$ and $L : U \rightarrow V$ be two linear transformations that take the same value in a basis u_1, \dots, u_n of U . Then $T = L$.*

Proof. Since u_1, \dots, u_n is a basis of U , any $u \in U$ can be written as $u = \alpha_1 u_1 + \dots + \alpha_n u_n$. Using Proposition 5.0.6 we have

$$\begin{aligned} T(u) = T(\alpha_1 u_1 + \dots + \alpha_n u_n) &= \alpha_1 T(u_1) + \dots + \alpha_n T(u_n) = \\ &= \alpha_1 L(u_1) + \dots + \alpha_n L(u_n) = L(\alpha_1 u_1 + \dots + \alpha_n u_n) = L(u) \end{aligned}$$

□

Linear Algebra — A. Bandeira (ETHZ) — Week 10 - 2023.11.24 & 2023.11.29

Please find most up to date notes at: https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_II.pdf

Proposition 5.0.8. *Given a basis u_1, \dots, u_n of U , and any $v_1, \dots, v_n \in V$ there is a Linear Transformation $T : U \rightarrow V$ such that, for all $1 \leq i \leq n$, $T(u_i) = v_i$.*

Challenge 29. Prove Proposition 5.0.8.

Example 5.0.9. *A few examples of linear transformations:*

- (1) *The identity map $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by $T(x) = x$,*
- (2) *For any matrix A , the map $x \rightarrow Ax$,*
- (3) *For a vector $v \in \mathbb{R}^n$ the map $T : \mathbb{R}^n \rightarrow \mathbb{R}$ given by $T(x) = v^\top x$,*
- (4) *The map $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by $T(x) = 0$.*

A few examples of functions that are not linear transformations:

- (1) For a vector $v \in \mathbb{R}^n$ (such that $v \neq 0$) the map $T : \mathbb{R}^n \rightarrow \mathbb{R}$ given by $T(x) = v + x$,
- (2) The map $T : \mathbb{R}^n \rightarrow \mathbb{R}$ given by $T(x) = \|x\|$,
- (3) The map $T : \mathbb{R}^n \rightarrow \mathbb{R}$ given by $T(x) = \|x\|^2$,
- (4) The map $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ given by $T(x) = \frac{1}{\|x\|}x$.

Challenge 30. Show that the first batch of examples above indeed correspond to linear transformations, and that the second does not.

It is easy to see that given an $m \times n$ matrix A , the function $x \rightarrow Ax$ is a Linear Transformation from \mathbb{R}^n to \mathbb{R}^m , what we will show now that the converse is also true.

Proposition 5.0.10. For any Linear Transformation $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$, there exists an $m \times n$ matrix A such that $T(x) = Ax$ for all $x \in \mathbb{R}^n$.

Proof. We will prove this proposition constructively. Let e_1, \dots, e_n be the canonical basis of \mathbb{R}^n (the i -th basis element has a 1 in the i -th entry and zeros elsewhere, or in other words $(e_i)_j = \delta_{ij}$). We write $x = x_1e_1 + \dots + x_n e_n$, then by linearity of T ,

$$T(x) = x_1T(e_1) + \dots + x_nT(e_n) = \begin{bmatrix} | & & | \\ T(e_1) & \cdots & T(e_n) \\ | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = Ax,$$

for $A = \begin{bmatrix} | & & | \\ T(e_1) & \cdots & T(e_n) \\ | & & | \end{bmatrix}$.

□

Challenge 31. (*) Can you describe the linear transformation corresponding to A^\dagger ? (in terms of the linear transformation corresponding to a matrix A)

Exploratory Challenge 32. Can you describe the linear transformation corresponding to A^\top ? (in terms of the linear transformation corresponding to a matrix A)

Further Remark 33. Since we are restricting ourselves to $U = \mathbb{R}^n$ and $V = \mathbb{R}^m$ we are identifying an element $x \in \mathbb{R}^n$ with its coordinates in the canonical basis $x = x_1e_1 + \dots + x_n e_n$. In general, if we use a different basis for U and V we will have a matrix representation for each linear transformation, but it will potentially correspond to a different matrix, it will be the matrix that describes the map from the coordinates in the basis of U to the ones in the basis of V , later in the course we will see some examples, and we will briefly discuss how to go from a matrix representation in one basis to that on another basis, a so-called *change of basis*.

Proposition 5.0.11. *Given two linear transformations $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $L : \mathbb{R}^m \rightarrow \mathbb{R}^p$, with corresponding matrices (as given by Proposition 5.0.10) $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{p \times m}$ the linear transformation $L \circ T$ (given by $L \circ T(x) = L(T(x))$) corresponds to multiplying by the matrix BA . In other words $L \circ T(x) = BAx$.*

Challenge 34. Prove Proposition 5.0.11.

5.1. The Determinant. We will now introduce the notion of determinant $\det(A)$ of a square matrix A . While this has a somewhat involved definition for $n \times n$ matrices, it is useful to first discuss what the determinant geometrically corresponds to, and to focus on small matrices.

In a nutshell, **the determinant of a matrix is a number that corresponds to how much the associated linear transformation inflates space, it corresponds precisely to the volume (or area, in \mathbb{R}^2) of the image of the unit cube (the red square in the pictures above in \mathbb{R}^2); with a negative sign when the orientation changes (in the pictures above in \mathbb{R}^2 , when the order of the colored dots, on the red square, changed).** If we think about the determinant this way, then many of the properties we will list below can be intuitively understood (while it is hard to do so from the formula for the $n \times n$ determinant). For this reason, this section will be somewhat less proof-based, and rather focus on the most relevant properties of the determinant.

Remark 5.1.1. *Grant Sanderson has a website <https://www.3blue1brown.com/> and Youtube channel <https://www.youtube.com/3blue1brown> with excellent animation-heavy explanations of topics in Mathematics, including Linear Algebra. I particularly recommend the video on Determinants, it has also 3 dimensional visualizations that are harder to do on a static medium. You can find it here <https://youtu.be/Ip3X9LOh2dk> or here <https://www.3blue1brown.com/lessons/determinant>. See also Figure 4.*

A calculation of the area of the image of the unit square by left-multiplication by a 2×2 matrix shows (see Figure 4) that

$$\begin{vmatrix} a & b \\ c & d \end{vmatrix} := \det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc.$$

Before we actually formally define determinant for $n \times n$ matrices we will state some of the most important properties of the determinant, you can find the actual definition in Definition 5.1.6.

5.1.1. Determinant and invertibility. Since a square matrix is invertible if the image is full-dimensional, which corresponds to the image of the unit square/cube having non-zero area/volume, then $\det(A) \neq 0$ if and only if A is invertible. This is the following proposition.

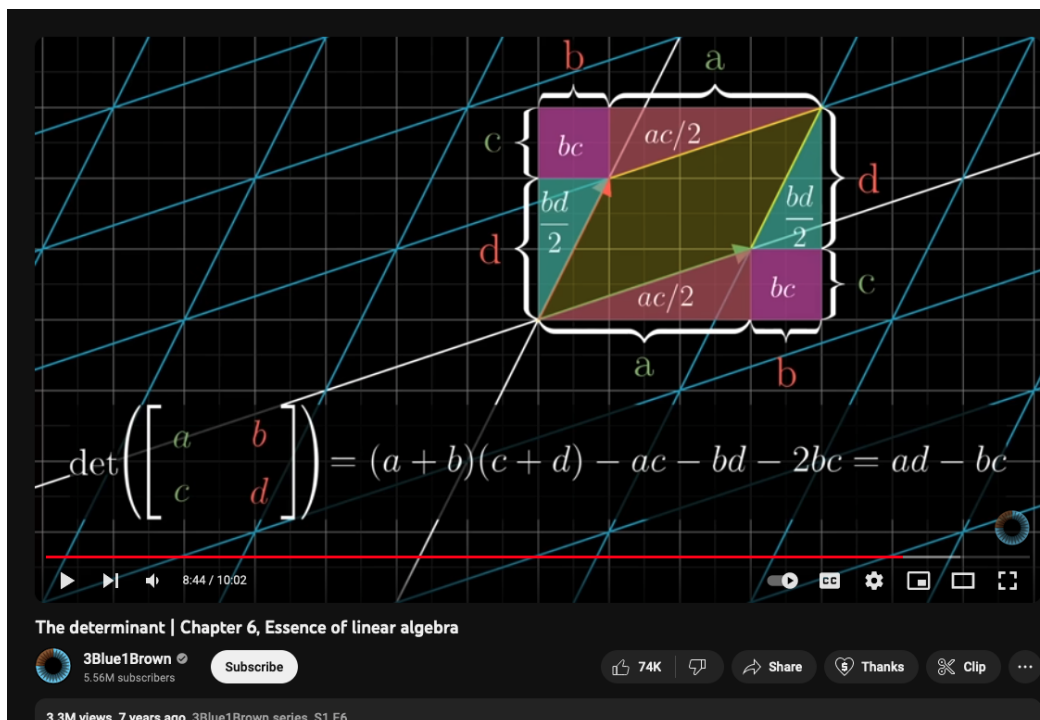


FIGURE 4. Calculation in 3Blue1Brown’s video (see Remark 5.1.1) computing the determinant of a 2×2 matrix as the area of the image of the unit square after a linear transformation (that does not change orientation).

Proposition 5.1.2. *A matrix $A \in \mathbb{R}^{n \times n}$ is invertible if and only if*

$$\det(A) \neq 0.$$

In fact, let us try to invert the matrix $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$, just by naive calculations, not using elimination.¹³

If $a = b = 0$ or $c = d = 0$ then the matrix is not invertible (it has a 0 row).¹⁴ Let’s assume either a or b is non-zero and that either c or d is nonzero. We are looking for a matrix $\begin{bmatrix} w & x \\ y & z \end{bmatrix}$ such that

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} w & x \\ y & z \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Since either a or b is non-zero and either c or d is nonzero, neither of the rows of the matrix are zero. Also, the second/first column of the inverse needs to be orthogonal to the first/second row of the original matrix, and vice-versa.

¹³Elimination is a much better way to do it in general, but bear with me as I am trying to illustrate something, not invert the matrix as efficiently as possible.

¹⁴Note that this is not a necessary condition for non-invertibility, as the all-ones matrix is not invertible while having no zero rows.

We then must have $\begin{bmatrix} x \\ z \end{bmatrix} = \alpha_1 \begin{bmatrix} -b \\ a \end{bmatrix}$ and $\begin{bmatrix} w \\ y \end{bmatrix} = \alpha_2 \begin{bmatrix} d \\ -c \end{bmatrix}$ for some $\alpha_1, \alpha_2 \in \mathbb{R}$.

Since $\begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} w \\ y \end{bmatrix} = 1$ we have: $\alpha_1 \begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} d \\ -c \end{bmatrix} = 1$, which gives $\alpha_1 = \frac{1}{ad-bc}$, note that the denominator is exactly $\det(A)$ which is non-zero when A is invertible. A similar calculation gives $\alpha_2 = \frac{1}{ad-bc} = \frac{1}{\det(A)}$. This gives a formula for the inverse of 2×2 matrices.

Proposition 5.1.3. *Given a 2×2 matrix A with $\det(A) \neq 0$, the inverse is given by*

$$A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}.$$

5.1.2. Determinant and volumes.

Proposition 5.1.4. *Given matrices $A, B \in \mathbb{R}^{n \times n}$ we have*

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

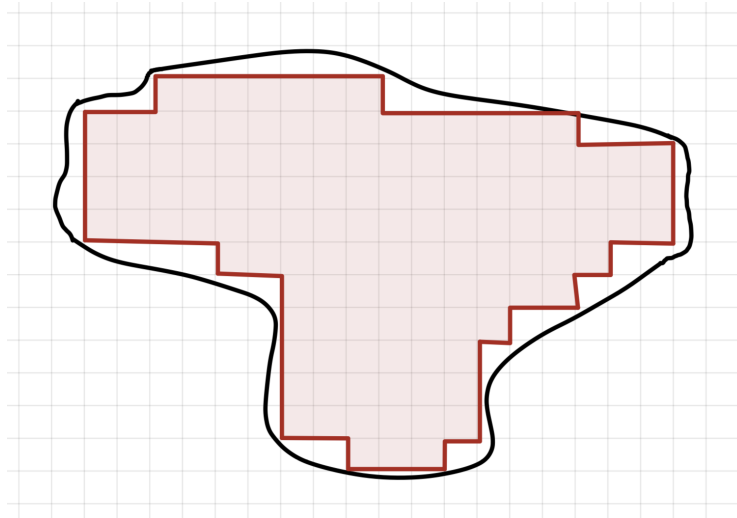


FIGURE 5. Approximation of a region by unit squares

While proving this from the definition of determinant is nontrivial, it is relatively easy to intuitively see why it is true if we recall that determinant measures areas/volumes: Since the area of the image of a square/cube does not change depending on the location of the initial square/cube, and any (nice enough) region of \mathbb{R}^n can be approximated by the union of small squares/cubes (see Figure 5), then the determinant is also the area/volume of the image any (nice enough) unit area/volume 1 region, and so $\det(AB) = \det(A) \det(B)$ (since the image by AB of a unit square/cube is the image by B of the image by A of the same unit square/cube).

5.1.3. *The Definition.* We now give the definition of determinant for $n \times n$ matrices. Before we define determinant we need first to discuss permutations.

Definition 5.1.5 (Sign of Permutation). *Given a permutation $\sigma : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ of n elements, its sign $\text{sgn}(\sigma)$ can be 1 or -1 . The sign counts the parity of the number of pairs of elements that are out of order (sometimes called inversions) after applying the permutation. In other words,*

$$\text{sgn}(\sigma) = \begin{cases} 1 & \text{if } |(i, j) \in \{1, \dots, n\} \times \{1, \dots, n\} \text{ such that } i < j \text{ and } \sigma(i) > \sigma(j)| \text{ is even,} \\ -1 & \text{if } |(i, j) \in \{1, \dots, n\} \times \{1, \dots, n\} \text{ such that } i < j \text{ and } \sigma(i) > \sigma(j)| \text{ is odd.} \end{cases}$$

Exploratory Challenge 35. The sign of a permutation has many nice properties. Try to prove a couple of them:

- (1) The sign of a permutation is multiplicative, i.e.: for two permutations σ, γ we have that $\text{sgn}(\sigma \circ \gamma) = \text{sgn}(\sigma)\text{sgn}(\gamma)$.
- (2) For all $n \geq 2$, exactly half of the permutations have sign 1 and exactly half have sign -1 .

the identity is 1, the sign of a transposition (a permutation that only swaps two elements) is -1 and for two permutations σ, γ we have that $\text{sgn}(\sigma \circ \gamma) = \text{sgn}(\sigma)\text{sgn}(\gamma)$.

Definition 5.1.6. *Given a square matrix $A \in \mathbb{R}^{n \times n}$ the determinant $\det(A)$ is defined as*

$$\det(A) = \sum_{\sigma \in \Pi_n} \text{sgn}(\sigma) \prod_{i=1}^n A_{i, \sigma(i)},$$

where Π_n is the set of all permutations of n elements.

From this Definition one can verify the following propositions.

Proposition 5.1.7. *Given a permutation matrix $P \in \mathbb{R}^{n \times n}$ corresponding to a permutation σ , then $\det(P) = \text{sgn}(\sigma)$. We sometimes also write $\text{sgn}(P)$.*

Proposition 5.1.8. *Given a triangular (either upper- or lower-) matrix $T \in \mathbb{R}^{n \times n}$ we have*

$$\det(T) = \prod_{k=1}^n T_{kk},$$

in particular, $\det(I) = 1$.

Proposition 5.1.9. *Given a matrix $A \in \mathbb{R}^{n \times n}$ we have*

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

The following is a consequence of the propositions above (and the only proof we'll do in this section)

Proposition 5.1.10. *If $Q \in \mathbb{R}^{n \times n}$ is an orthogonal matrix then*

$$\det(Q) = 1 \quad \text{or} \quad \det(Q) = -1.$$

Proof. By Propositions 5.1.8 and 5.1.4 we have $1 = \det(I) = \det(Q^\top Q) = \det(Q^\top) \det(Q)$. by Proposition 5.1.9 we have $1 = \det(Q)^2$ and so $\det(Q)$ is 1 or -1. \square

Following the same line of argument we also have

Proposition 5.1.11. *Given a matrix $A \in \mathbb{R}^{n \times n}$ such that $\det(A) \neq 0$, then A is invertible and*

$$\det(A^{-1}) = \frac{1}{\det(A)}.$$

5.1.4. 3×3 matrices.

If A is a 1×1 matrix, since there is only one permutation of 1 element (the permutation $\sigma(1) = 1$, which has sign 1), we have $\det(A) = A_{11} = A$.

For 2×2 matrices: There are two permutations σ_1 the identity permutation (that doesn't move any element, which has sign 1) and σ_2 the permutation that swaps the two elements (which has sign -1). So, for A a 2×2 matrix, we have

$$\det(A) = \sum_{\sigma \in \Pi_2} \text{sgn}(\sigma) \prod_{i=1}^2 A_{i, \sigma(i)} = (+1) \prod_{i=1}^2 A_{i, \sigma_1(i)} + (-1) \prod_{i=1}^2 A_{i, \sigma_2(i)} = A_{11}A_{22} - A_{12}A_{21}.$$

This corresponds precisely to,

$$\begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc.$$

For 3×3 matrices there are $3! = 6$ permutations, so there will be 6 terms. For A a 3×3 matrix, we can write its determinant as (where an empty entry corresponds to a zero entry)

$$\begin{aligned} \det(A) &= \begin{vmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{vmatrix} \\ &= \begin{vmatrix} A_{11} & & \\ & A_{22} & \\ & & A_{33} \end{vmatrix} + \begin{vmatrix} & A_{12} & \\ A_{21} & & \\ & & A_{33} \end{vmatrix} + \begin{vmatrix} & & A_{12} \\ & & A_{23} \\ A_{31} & & \end{vmatrix} \\ &\quad + \begin{vmatrix} & & A_{13} \\ & A_{22} & \\ A_{31} & & \end{vmatrix} + \begin{vmatrix} & & A_{13} \\ A_{21} & & \\ & A_{32} & \end{vmatrix} + \begin{vmatrix} A_{11} & & \\ & & A_{23} \\ & & A_{32} \end{vmatrix} \\ &= A_{11}A_{22}A_{33} - A_{12}A_{21}A_{33} + A_{12}A_{23}A_{31} - A_{13}A_{22}A_{31} + A_{13}A_{21}A_{32} - A_{11}A_{23}A_{32}. \end{aligned}$$

There is another convenient way of writing this determinant

$$(12) \quad \begin{vmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{vmatrix} = A_{11} \begin{vmatrix} A_{22} & A_{23} \\ A_{32} & A_{33} \end{vmatrix} - A_{12} \begin{vmatrix} A_{21} & A_{23} \\ A_{31} & A_{33} \end{vmatrix} + A_{13} \begin{vmatrix} A_{21} & A_{22} \\ A_{31} & A_{32} \end{vmatrix}.$$

In general, these terms are called the co-factors of A .

Definition 5.1.12. Given $A \in \mathbb{R}^{n \times n}$, for each $1 \leq i, j \leq n$ let \mathcal{A}_{ij} denote the $(n-1) \times (n-1)$ matrix obtained by removing row i and column j from A . Then we define the co-factors of A as

$$C_{ij} = (-1)^{i+j} \det(\mathcal{A}_{ij}).$$

Just as in (12), the determinant can be written in terms of the co-factors.

Proposition 5.1.13. Let $A \in \mathbb{R}^{n \times n}$, for any $1 \leq i \leq n$,

$$\det(A) = \sum_{j=1}^n A_{ij} C_{ij}.$$

The formula we derived above for the inverse of 2×2 matrices (Proposition 5.1.3), also has an analogue in n dimensions.

Proposition 5.1.14. Given $A \in \mathbb{R}^{n \times n}$ with $\det(A) \neq 0$ we have

$$A^{-1} = \frac{1}{\det(A)} C^{\top},$$

where C is the $n \times n$ matrix with the co-factors of A as entries.

One good way to think of this proposition is as the identity $AC^T = \det(A)I$.

Remark 5.1.15. *Computationally speaking, this is not a good way to compute the inverse, as it involves computing many determinants.*

Challenge 36. Verify that Proposition 5.1.14 indeed corresponds to Proposition 5.1.3 when $n = 2$.

Exploratory Challenge 37. Try to prove Proposition 5.1.14 by showing that $AC^T = \det(A)I$. Perhaps start with $n = 3$. You can also use Cramer's Rule (below) to prove this.

5.1.5. *Cramer's Rule.* The determinant also allows us to write a formula for the solution of the linear system of the type $Ax = b$ when $A \in \mathbb{R}^{n \times n}$ and $\det(A) \neq 0$. The idea is simple, we will illustrate it here for $n = 3$.

$$\text{If } \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}, \text{ then we have}$$

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} \begin{bmatrix} x_1 & 0 & 0 \\ x_2 & 1 & 0 \\ x_3 & 0 & 1 \end{bmatrix} = \begin{bmatrix} b_1 & A_{12} & A_{13} \\ b_2 & A_{22} & A_{23} \\ b_3 & A_{32} & A_{33} \end{bmatrix}.$$

Since the determinant is multiplicative, and the determinant of the second matrix in the expression is x_1 , we have

$$\det(A)x_1 = \det(\mathcal{B}_1),$$

where \mathcal{B}_1 is the matrix obtained by A by replacing the first column of A with the vector b .

Since we can do this for any of the columns, we have $x_j = \det(\mathcal{B}_j) / \det(A)$. In general

Proposition 5.1.16 (Cramer's Rule). *Let $A \in \mathbb{R}^{n \times n}$ such that $\det(A) \neq 0$ and $b \in \mathbb{R}^n$ then the solution $x \in \mathbb{R}^n$ of $Ax = b$ is given by*

$$x_j = \frac{\det(\mathcal{B}_j)}{\det(A)},$$

where \mathcal{B}_j is the matrix obtained by A by replacing the j -th column of A with the vector b .

Remark 5.1.17. *As with the formula for the inverse: computationally speaking, this is not a good way to solve linear systems, as it involves computing many determinants.*

5.1.6. *Elimination and the Determinant.* The definition we used for Determinant involves a formula with $n!$ terms, it is computational infeasible for even moderate levels of n (it is faster than exponential! For example, $100!$ has almost 160 digits!), in practice the determinant of a matrix A

is computed by Gaussian Elimination and the matrix decomposition $PA = LU$ (P permutation and so $\det(P) = \text{sgn}(P)$, U is upper triangular and L is lower triangular with only 1s in the diagonal, and so $\det(L) = 1$) and so we would have

$$(13) \quad \det(A) = \frac{1}{\det(P)} \det(L) \det(U) = \text{sgn}(P) \det(U),$$

and since U is a triangular matrix its determinants can be easily computed by Proposition 5.1.8.

Alternatively, one can also think of Gaussian Elimination as directly computing the determinant via the following two propositions

Proposition 5.1.18. *If A is an $n \times n$ matrix and P is a permutation that swaps two elements, meaning that PA corresponds to swapping two rows of A then $\det(PA) = -\det(A)$.*

Proposition 5.1.19. *The determinant is linear in each row (or each column). In other words, for any $a_0, a_1, a_2, \dots, a_n \in \mathbb{R}^n$ and $\alpha_0, \alpha_1 \in \mathbb{R}$ we have*

$$\begin{vmatrix} - & \alpha_0 a_0^\top + \alpha_1 a_1^\top & - \\ - & a_2^\top & - \\ & \vdots & \\ - & a_n^\top & - \end{vmatrix} = \alpha_0 \begin{vmatrix} - & a_0^\top & - \\ - & a_2^\top & - \\ & \vdots & \\ - & a_n^\top & - \end{vmatrix} + \alpha_1 \begin{vmatrix} - & a_1^\top & - \\ - & a_2^\top & - \\ & \vdots & \\ - & a_n^\top & - \end{vmatrix},$$

and

$$\begin{vmatrix} | & | & | & | \\ \alpha_0 a_0 + \alpha_1 a_1 & a_2 & \cdots & a_n \\ | & | & | & | \end{vmatrix} = \alpha_0 \begin{vmatrix} | & | & | & | \\ a_0 & a_2 & \cdots & a_n \\ | & | & | & | \end{vmatrix} + \alpha_1 \begin{vmatrix} | & | & | & | \\ a_1 & a_2 & \cdots & a_n \\ | & | & | & | \end{vmatrix}.$$

Exploratory Challenge 38. The more mathematical way of presenting this material is to define Determinant as a function that goes from $n \times n$ matrices to \mathbb{R} that (i) is linear in each column, (ii) $\det(I) = 1$ and (iii) $\det(A) = 0$ whenever A has two identical columns. It is then possible to prove that the only function satisfying these three properties is the determinant as we defined it. Try to prove it!

Linear Algebra — A. Bandeira (ETHZ) — Week 11 - 2023.12.01 & 2023.12.06

Please find most up to date notes at: https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_II.pdf

6. EIGENVALUES AND EIGENVECTORS

We are (almost) ready for one of the most important concepts (if not the most important one) in Linear Algebra, **eigenvalues and eigenvectors**. In a sense, *it has all been building up to this!*

Guiding Strategy 39. Given a square matrix A , as we will see below, an eigenvalue λ and eigenvector v will be, respectively, a scalar and a non-zero vector satisfying $Av = \lambda v$. This means that $(A - \lambda I)v = 0$ and so $(A - \lambda I)$ is not invertible, or equivalently $\det(A - \lambda I) = 0$. We can look for eigenvalues as solutions of $\det(A - \lambda I) = 0$ which is a polynomial¹⁵ in λ but unfortunately, not all polynomials have real zeros.¹⁶ For example if $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, $\det(A - \lambda I) = 0$ corresponds to $\lambda^2 + 1 = 0$ which only has solutions in \mathbb{C} , the Complex Numbers. For this reason we will start this Chapter with a brief introduction to Complex Numbers. It all starts with asking for a number λ such that $\lambda^2 + 1 = 0$.

Further Reading 40. Complex Analysis is a beautiful topic in Mathematics, what we will cover here is just a tiny peak at it, there is a all bookshelf of excellent books in this topic in our library. I have personally taught a course at ETH on Complex Analysis, and since it was during the COVID pandemic I made videos available online, which are still available at https://www.youtube.com/playlist?list=PLiud-28tsatLRRGqO_Eg_x0S4LVyxuV5p (In particular, the first lecture covers roughly the content here).

6.0. Complex Numbers. If we start with the natural numbers \mathbb{N} and want to solve equations like $x + 10 = 1$, we need negative numbers. This motivates considering the integers \mathbb{Z} . Similarly, rational numbers \mathbb{Q} are needed to solve equations like $10x = 1$ and real numbers \mathbb{R} are needed to solve $x^2 = 2$.¹⁷ Similarly, the Complex Numbers are needed to solve equations such as $x^2 + 1 = 0$. It starts with the introduction of an imaginary number $i \in \mathbb{C}$ such that $i^2 = -1$. You can think of i as $\sqrt{-1}$.

The complex numbers are numbers of the form $z = a + ib$ for $a \in \mathbb{R}$ and $b \in \mathbb{R}$. $\mathbb{C} = \{a + ib : a, b \in \mathbb{R}\}$. Keeping in mind that $i^2 = -1$ we can do operations with complex numbers:

- $(a + ib) + (x + iy) = (a + x) + i(b + y)$,
- $(a + ib)(x + iy) = ax + i(ay + bx) + i^2by = ax + i(ay + bx) - by = (ax - by) + i(ay + bx)$,
- $(a + ib)(a - ib) = a^2 + b^2$,
- $\frac{a+ib}{x+iy} = \frac{(x-iy)(a+ib)}{(x-iy)(x+iy)} = \frac{(ax+by)+i(bx-ay)}{x^2+y^2} = \left(\frac{ax+by}{x^2+y^2}\right) + i\left(\frac{bx-ay}{x^2+y^2}\right)$.

¹⁵This is one of the main reasons we had to cover determinants.

¹⁶A zero of a polynomial P is a point x such that $P(x) = 0$, this is also called a root of the polynomial. In German, it's a "Nullstelle". In fact, a (rather deep) multidimensional version of Theorem 6.0.3, and one of the most important facts in Algebraic Geometry, is called "Hilbert's Nullstellensatz".

¹⁷If you have never seen the proof that there exists no $x \in \mathbb{Q}$ such that $x^2 = 2$ I highly recommend trying to do it: set $x = a/b$ for $a, b \in \mathbb{Z}$ and try to count how many times 2 divides both a and b and find a contradiction.

Given $z \in \mathbb{C}$ with $z = a + ib$ we have the following notation

$$(14) \quad \Re(a + ib) := a \quad \text{called the real part of } z = a + ib,$$

$$(15) \quad \Im(a + ib) := b \quad \text{called the imaginary part of } z = a + ib,$$

$$(16) \quad |z| := \sqrt{a^2 + b^2} \quad \text{called the modulus of } z = a + ib,$$

$$(17) \quad \overline{a + ib} := a - ib \quad \text{called the complex conjugate of } z = a + ib.$$

Note that for $z_1, z_2 \in \mathbb{C}$, we have $|z|^2 = z\bar{z}$, $z_1 z_2 = z_2 z_1$, $\overline{z_1 + z_2} = \bar{z}_1 + \bar{z}_2$, and $\frac{1}{z} = \frac{\bar{z}}{|z|^2}$.

Fact 6.0.1 (Euler's Formula). *Given $\theta \in \mathbb{R}$, we have*

$$(18) \quad e^{i\theta} = \cos \theta + i \sin \theta.$$

This means, in particular, that $e^{i\pi} = -1$. This is usually written as $e^{i\pi} + 1 = 0$.

Further Reading 41. In order to prove Euler's Formula, we need to first define what we mean by $e^{i\theta}$, this can be done, for example, by the Taylor series of the exponential, but this is outside the scope of this course (see Further Reading 40).

Fact 6.0.2 (Polar Coordinates). *A complex number $z \in \mathbb{C}$ can be written as*

$$(19) \quad z = re^{i\theta},$$

where $r \geq 0$ is the modulus of z and $\theta \in \mathbb{R}$ (we can restrict to $\theta \in [0, 2\pi[)$ is an angle, also called the argument of z .

The most important property of Complex Numbers, and what makes them a very natural mathematical object, is that any univariate polynomial equation with complex number coefficients has a (complex) solution, in a certain sense we don't need to extend numbers further, \mathbb{C} is **Algebraically closed**.

Theorem 6.0.3 (Fundamental Theorem of Algebra). *Any degree n non-constant ($n \geq 1$) polynomial $P(z) = \alpha_n z^n + \alpha_{n-1} z^{n-1} + \dots + \alpha_1 z + \alpha_0$ (with $\alpha_n \neq 0$) has a zero: $\lambda \in \mathbb{C}$ such that $P(\lambda) = 0$.*

Further Reading 42. As the name suggests, Theorem 6.0.3 is a central result in Complex Analysis. Proving it is outside the scope of this course.¹⁸ Complex analysis (which leads to the proof of this theorem) is a beautiful example of interaction between analysis, algebra, and geometry. In a nutshell the idea for the classical proof is that differentiable functions in the complex plane

¹⁸But you can see Appendix D for a relatively elementary proof.

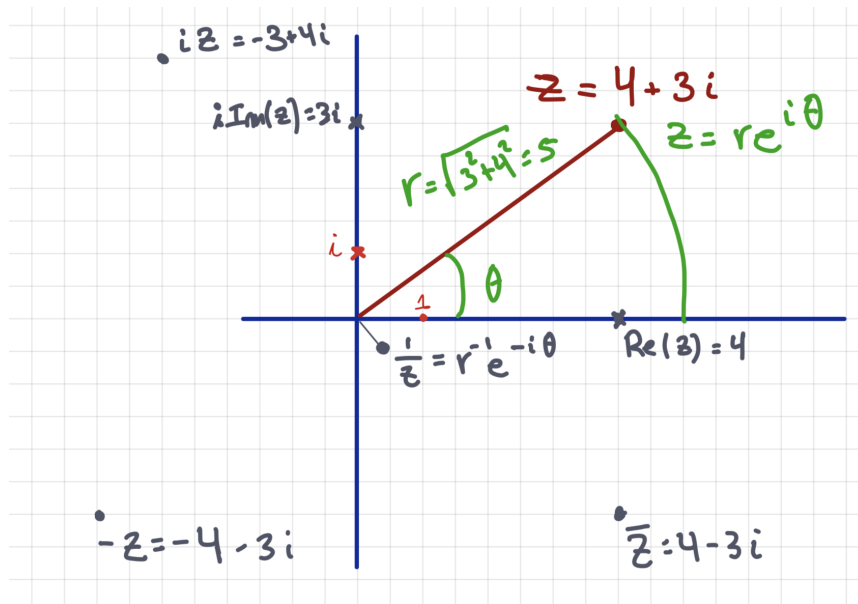


FIGURE 6. A complex number $z = 4 + 3i$ in the Complex plane.

$f : \mathbb{C} \rightarrow \mathbb{C}$ are very special and, in a sense, need to behave like polynomials (this is a deep statement that needs a significant amount of background to properly state and prove). If a polynomial $P(z)$ doesn't have a zero then $1/P(z)$ is a differentiable function that cannot behave like a non-constant polynomial because it does not grow sufficiently far away from zero, and so it must be a constant function which means that $P(z)$ had to be constant, so any non-constant polynomial has a zero. For more on Complex Analysis see Further Reading 40.

Further Remark 43. Once we have λ a zero of $P(z)$, we can divide $P(z)$ by $(z - \lambda)$ to get $P(z) = (z - \lambda)P_1(z)$, then use a zero of P_1 to reiterate, and so on. This argument (carried out carefully) gives the following corollary.

Corollary 6.0.4. Any degree n non-constant ($n \geq 1$) polynomial $P(z) = \alpha_n z^n + \alpha_{n-1} z^{n-1} + \dots + \alpha_1 z + \alpha_0$ (with $\alpha_n \neq 0$) has n zeros: $\lambda_1, \dots, \lambda_n \in \mathbb{C}$, perhaps with repetitions, such that

$$(20) \quad P(z) = \alpha_n (z - \lambda_1)(z - \lambda_2) \cdots (z - \lambda_n).$$

The number of times $\lambda \in \mathbb{C}$ appears in this expansion is called the algebraic multiplicity of the zero.

6.0.1. *Complex-valued Matrices and Vectors.* Analogously to \mathbb{R}^n we also define \mathbb{C}^n as the set of n -dimensional complex valued vectors. We can have complex valued vectors $v \in \mathbb{C}^n$ and matrices $A \in \mathbb{C}^{m \times n}$. The natural operation of “transposing” for complex vectors and matrices is that of

“conjugate transpose” or “hermitian transpose” denoted by A^* , or sometimes A^H ,

$$(21) \quad A^* = \overline{A}^T.$$

Given $v \in \mathbb{C}^n$ we have

$$\|v\|^2 = v^* v = \overline{v}^T v = \sum_{i=1}^n \overline{v_i} v_i = \sum_{i=1}^n |v_i|^2.$$

The inner-product (or dot-product) in \mathbb{C}^n is given by $\langle v, w \rangle = w^* v$.

Similarly to the situation in \mathbb{R}^n , we say $v_1, \dots, v_k \in \mathbb{C}^n$ are linearly independent if there is no (complex valued) non-zero linear combination giving zero, meaning that if $\alpha_1 v_1 + \dots + \alpha_k v_k = 0$ for $\alpha_1, \dots, \alpha_k \in \mathbb{C}$ we must have $\alpha_1 = \dots = \alpha_k = 0$. Also, the span of $v_1, \dots, v_k \in \mathbb{C}^n$ is the set of possible linear combinations $\alpha_1 v_1 + \dots + \alpha_k v_k$ for $\alpha_1, \dots, \alpha_k \in \mathbb{C}$. If v_1, \dots, v_k is a spanning set of a subspace and linearly independent we say it is a basis of that subspace. As with \mathbb{R}^n if we have $v_1, \dots, v_n \in \mathbb{C}^n$ that are either a spanning set of \mathbb{C}^n or linearly independent then they must actually be both (and so are a basis).

Further Reading 44. With these definitions you can already understand the Discrete Fourier Transform (which is the linear transformation corresponding to the DFT matrix, one of the most important complex valued matrices). This is the key object behind signal processing, you can read more about it on the lecture notes of another course I usually teach [BM23]. You can also see a discussion of Fourier Transform, circulant matrices, and signal convolutions in [Str23] (end of Section 6.4).

6.1. Introduction to Eigenvalues and Eigenvectors. Even though the theory can be analogously developed for complex valued matrices, we will focus on real valued matrices.

Guiding Example 45. We will use a guiding example to illustrate both some of the power, and some of the properties, of eigenvalues and eigenvectors. In Guiding Example numbers 45 through 51 we will derive a formula for the n -th Ficonacci Number. The Fibonacci numbers are defined by the recurrence:

$$(22) \quad F_0 = 0, F_1 = 1, \text{ and, for } n \geq 2, F_n = F_{n-1} + F_{n-2}.$$

The recurrence can be rewritten in linear algebraic notation as, for $n \geq 2$,

$$(23) \quad \begin{bmatrix} F_{n+1} \\ F_n \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} F_n \\ F_{n-1} \end{bmatrix}.$$

Defining

$$(24) \quad M = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \text{ and } g_n = \begin{bmatrix} F_{n+1} \\ F_n \end{bmatrix},$$

the recurrence can be rewritten as

$$g_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ and } g_n = M g_{n-1},$$

meaning that

$$(25) \quad g_n = M^n g_0.$$

Definition 6.1.1. Given $A \in \mathbb{R}^{n \times n}$, we say $\lambda \in \mathbb{C}$ is an eigenvalue of A and $v \in \mathbb{C}^n \setminus \{0\}$ is an eigenvector of A , associated with the eigenvalue λ , when the following holds:

$$Av = \lambda v.$$

We call them an eigenvalue-eigenvector pair. If $\lambda \in \mathbb{R}$ then we will call λ a real eigenvalue, and the associated eigenvalue-eigenvector pair a real eigenvalue-eigenvector pair.

Guiding Example 46. Let us try to find eigenvalues (and later the eigenvectors) of $M = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$.

We are looking for $v \in \mathbb{R}^2 \setminus \{0\}$ and $\lambda \in \mathbb{R}$ such that $Mv = \lambda v$, but this can be rewritten as $(M - \lambda I)v = 0$ and since $v \neq 0$ it means that $M - \lambda I$ is non-invertible (also called singular).¹⁹ This is equivalent to $\det(M - \lambda I) = 0$ and so we can find the eigenvalues λ with this equation:

$$(26) \quad 0 = \det(M - \lambda I) = \begin{vmatrix} 1 - \lambda & 1 \\ 1 & 0 - \lambda \end{vmatrix} = (1 - \lambda)(0 - \lambda) - 1 = \lambda^2 - \lambda - 1.$$

By the quadratic formula,²⁰ the solutions to (26) are given by

$$(27) \quad \lambda_1 = \frac{1 + \sqrt{5}}{2} \text{ and } \lambda_2 = \frac{1 - \sqrt{5}}{2}.$$

Further Reading 47 (Golden Ratio). The number $\varphi = \frac{1 + \sqrt{5}}{2}$ is the celebrated Golden Ratio; believed, since the ancient Greeks, to be the ideal aspect ratio for a rectangle.

“Some of the greatest mathematical minds of all ages, from Pythagoras and Euclid in ancient Greece, through the medieval Italian mathematician Leonardo of Pisa and the Renaissance astronomer Johannes Kepler, to present-day scientific figures such as Oxford physicist Roger Penrose, have spent endless hours over this simple ratio and its properties. [...] Biologists, artists, musicians, historians, architects, psychologists, and even mystics have pondered and debated the basis of its ubiquity and appeal. In fact, it is

¹⁹Normally, we would have to look for $\lambda \in \mathbb{C}$ and $v \in \mathbb{C}^n$ but in this case the eigenvalues, as we will see, are real.

²⁰Recall that the quadratic formula says that the zeros of $ax^2 + b + c$ are given by $x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$.

probably fair to say that the Golden Ratio has inspired thinkers of all disciplines like no other number in the history of mathematics.”

— The Golden Ratio: The Story of Phi, the World’s Most Astonishing Number

The following is the original definition which dates back to Euclid around 2300 years ago (they called the number “extreme and mean ratio” back then)

“A straight line is said to have been cut in extreme and mean ratio when, as the whole line is to the greater segment, so is the greater to the lesser”

Guiding Example 48. Now we can try to find the eigenvectors v_1 and v_2 such that $Av_1 = \lambda_1 v_1$ and $Av_2 = \lambda_2 v_2$.

Let us start with v_1 . We are looking for a non-zero element of $N\left(A - \frac{1+\sqrt{5}}{2}I\right)$. In other words

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 - \frac{1+\sqrt{5}}{2} & 1 \\ 1 & -\frac{1+\sqrt{5}}{2} \end{bmatrix} \begin{bmatrix} (v_1)_1 \\ (v_1)_2 \end{bmatrix}.$$

This is an under-determined system and we are looking for a non-zero solution, so let us start by setting $(v_1)_2 = 1$. The second equation gives us $(v_1)_1 = \frac{1+\sqrt{5}}{2}$. Indeed $v_1 = \begin{bmatrix} \frac{1+\sqrt{5}}{2} \\ 1 \end{bmatrix}$ is an eigenvector of M associated to the eigenvalue $\lambda_1 = \frac{1+\sqrt{5}}{2}$.

A similar calculation for $\lambda_2 = \frac{1-\sqrt{5}}{2}$ gives that $v_2 = \begin{bmatrix} \frac{1-\sqrt{5}}{2} \\ 1 \end{bmatrix}$. Indeed

$$(28) \quad \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1+\sqrt{5}}{2} \\ 1 \end{bmatrix} = \frac{1+\sqrt{5}}{2} \begin{bmatrix} \frac{1+\sqrt{5}}{2} \\ 1 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{1-\sqrt{5}}{2} \\ 1 \end{bmatrix} = \frac{1-\sqrt{5}}{2} \begin{bmatrix} \frac{1-\sqrt{5}}{2} \\ 1 \end{bmatrix}$$

Challenge 49. Carry out the calculations in Guiding Example 48 and confirm that we have indeed found two eigenvectors (check the two equalities in (28)).

Further Remark 50. The v_1 and v_2 we constructed in 48 are not the only possible choices, for example any non-zero scalar multiples of these would have also been a possible choice. Normally one picks a unit-norm representative, but in this case we picked vectors that make the calculations the cleanest.

What we carried out in the example above is very general and we now develop the theory for general matrices.

Let λ and v be an eigenvalue-eigenvector pair of a matrix A . Since $v \neq 0$ and $(A - \lambda I)v = Av - \lambda v = 0$ we have that $\det(A - \lambda I) = 0$. Conversely, if $\det(A - \lambda I) = 0$ for some λ , then there

exists $v \in N(A - \lambda I) \setminus \{0\}$ and so λ is an eigenvalue. This gives a procedure to find eigenvalues and eigenvectors: (i) eigenvalues are the solution of $\det(A - \lambda I) = 0$, which is a polynomial equation, and (ii) an associated eigenvector is a non-zero element of $N(A - \lambda I)$.

Let us first formulate this for real eigenvalues and eigenvectors.

Proposition 6.1.2. *Let $A \in \mathbb{R}^{n \times n}$. $\lambda \in \mathbb{R}$ is a (real) eigenvalue of A if and only if $\det(A - \lambda I) = 0$. A vector v is an eigenvector associated with the eigenvalue λ if (and only if) it is a non-zero element of $N(A - \lambda I)$.*

A direct inspection of the formula for the determinant (Definition 5.1.6) gives the following.

Proposition 6.1.3. *$\det(A - \lambda I)$ is a polynomial, in λ , of degree n . The coefficient of the λ^n term is $(-1)^n$.*

The Fundamental Theorem of Algebra (Theorem 6.0.3) immediately implies

Theorem 6.1.4. *Every matrix $A \in \mathbb{R}^{n \times n}$ has an eigenvalue (perhaps complex-valued).*

Remark 6.1.5. *For now we will focus on real eigenvalues, and address complex valued ones later on. Essentially all the properties we will describe below also hold for complex valued eigenvalues (just by replacing \mathbb{R} by \mathbb{C} and doing the appropriate adjustments). For example, Proposition 6.1.2 also holds for complex-valued eigenvalues, one just needs to think of $N(A - \lambda I)$ as a subspace of \mathbb{C}^n , meaning the vectors $v \in \mathbb{C}^n$ such that $(A - \lambda I)v = 0$.*

Guiding Example 51. Let us return to our guiding example. Notice that v_1 and v_2 are linearly independent, and so they are a basis for \mathbb{R}^2 . We can write $g_0 = \alpha_1 v_1 + \alpha_2 v_2$.

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} = g_0 = \alpha_1 v_1 + \alpha_2 v_2 = \begin{bmatrix} \alpha_1 \frac{1+\sqrt{5}}{2} + \alpha_2 \frac{1-\sqrt{5}}{2} \\ \alpha_1 + \alpha_2 \end{bmatrix} = \begin{bmatrix} (\alpha_1 + \alpha_2)\frac{1}{2} + (\alpha_1 - \alpha_2)\frac{\sqrt{5}}{2} \\ \alpha_1 + \alpha_2 \end{bmatrix},$$

and so $\alpha_1 = \frac{1}{\sqrt{5}}$ and $\alpha_2 = -\frac{1}{\sqrt{5}}$.

Recall that $g_n = A^n g_0$ and so

$$g_n = A^n \left(\frac{1}{\sqrt{5}} v_1 - \frac{1}{\sqrt{5}} v_2 \right) = \frac{1}{\sqrt{5}} A^n v_1 - \frac{1}{\sqrt{5}} A^n v_2 = \frac{1}{\sqrt{5}} (A^n v_1 - A^n v_2).$$

Since $Av_1 = \lambda_1 v_1$ we have that $A^2 v_1 = A(\lambda_1 v_1) = \lambda_1^2 v_1$ and iterating this procedure²¹ gives $A^n v_1 = \lambda_1^n v_1$. This means that

$$g_n = \frac{A^n v_1 - A^n v_2}{\sqrt{5}} = \frac{\left(\frac{1+\sqrt{5}}{2}\right)^n v_1 - \left(\frac{1-\sqrt{5}}{2}\right)^n v_2}{\sqrt{5}} = \frac{\left(\frac{1+\sqrt{5}}{2}\right)^n}{\sqrt{5}} \begin{bmatrix} \frac{1+\sqrt{5}}{2} \\ 1 \end{bmatrix} - \frac{\left(\frac{1-\sqrt{5}}{2}\right)^n}{\sqrt{5}} \begin{bmatrix} \frac{1-\sqrt{5}}{2} \\ 1 \end{bmatrix}.$$

²¹A formal proof would use induction

Since F_n is the second coordinate of g_n , we derived a closed formula for the n -th terms of the Fibonacci sequence:

$$(29) \quad F_n = \frac{1}{\sqrt{5}} \left(\frac{1+\sqrt{5}}{2} \right)^n - \frac{1}{\sqrt{5}} \left(\frac{1-\sqrt{5}}{2} \right)^n.$$

An important property that allowed us to do the calculation above was that applying a power of a matrix to an eigenvector was a simple operation, this is the next proposition.

Proposition 6.1.6. *If λ and v are an eigenvalue-eigenvector pair of a matrix A , then, for $k \geq 1$, λ^k and v are an eigenvalue-eigenvector pair of the matrix A^k .*

Proof. Proof by Induction: The base case $k = 1$ is trivial. For the induction step, since λ^k and v are an eigenvalue-eigenvector pair then $A^k v = A(A^{k-1}v) = A(\lambda^{k-1}v) = \lambda^k v$. \square

Proposition 6.1.1006. *Let A be an invertible matrix. If λ and v are an eigenvalue-eigenvector pair of a matrix A , then, $\frac{1}{\lambda}$ and v are an eigenvalue-eigenvector pair of the matrix A^{-1} .*

Proof. Since $Av = \lambda v$ we have $A^{-1}(\lambda v) = v$ and so $\lambda A^{-1}v = v$, which (since $\lambda \neq 0$) is equivalent to $A^{-1}v = \frac{1}{\lambda}v$. \square

Another important property, was that we were able to write a vector as a linear combination of eigenvectors, which was possible because the eigenvectors were linearly independent.

Proposition 6.1.7. *Let $A^{n \times n}$ and let $v_1, \dots, v_k \in \mathbb{R}^n$ be eigenvectors corresponding to eigenvalues $\lambda_1, \dots, \lambda_k \in \mathbb{R}$. If $\lambda_1, \dots, \lambda_k$ are all distinct, the eigenvectors v_1, \dots, v_k are linearly independent.*

Proof. We will prove this by contradiction. Assume that v_1, \dots, v_k are linearly dependent. For $i = 1, \dots, k$, let d_i denote the dimension of the span of v_1, \dots, v_i . Since $v_1 \neq 0$ we have $d_1 = 1$. By the hypothesis $d_k < k$. Let j be the smallest positive integer for which $d_j < j$. Note that, by construction, $d_{j-1} = d_j = j - 1$, this means that v_1, \dots, v_{j-1} are linearly independent but that v_j is in the span of v_1, \dots, v_{j-1} . We can then write

$$(30) \quad v_j = \alpha_1 v_1 + \dots + \alpha_{j-1} v_{j-1}.$$

If we multiply by A both sides we get

$$\lambda_j v_j = Av_j = A(\alpha_1 v_1 + \dots + \alpha_{j-1} v_{j-1}) = \alpha_1 \lambda_1 v_1 + \dots + \alpha_{j-1} \lambda_{j-1} v_{j-1}.$$

Replacing the v_j in the left hand side with the right hand side of (30) we get

$$\lambda_j (\alpha_1 v_1 + \dots + \alpha_{j-1} v_{j-1}) = \alpha_1 \lambda_1 v_1 + \dots + \alpha_{j-1} \lambda_{j-1} v_{j-1},$$

which we can rearrange as

$$(31) \quad \alpha_1 (\lambda_j - \lambda_1) v_1 + \alpha_2 (\lambda_j - \lambda_2) v_2 + \cdots + \alpha_{j-1} (\lambda_{j-1} - \lambda_1) v_{j-1} = 0.$$

Since $\lambda_j - \lambda_i \neq 0$ for all $i \leq j-1$ and not all α_i 's are zero, this is a non-zero linear combination of v_1, \dots, v_{j-1} adding to zero, which would be a contradiction with $d_{j-1} = j-1$. \square

A very important consequence of this is that if a matrix has n distinct real eigenvalues then the eigenvectors form a basis for \mathbb{R}^n .

Theorem 6.1.8. *Let $A \in \mathbb{R}^{n \times n}$ with n distinct real eigenvalues (meaning that the n zeros of $\det(A - \lambda I)$, as described in Corollary 6.0.4, are all distinct) then there is a basis of \mathbb{R}^n , v_1, \dots, v_n , made up of eigenvectors of A .*

Guiding Example 52. Guiding Example 45 is yet to stop providing us with insight into properties of eigenvalues and eigenvectors! Here are a couple of observations, which although outside of the core scope of this course, have significant impact in several areas:

- Notice that since $|\lambda_2| < |\lambda_1|$, the contribution of $\lambda_2^n \alpha_2 v_2$ becomes negligible (when compared to $\lambda_1^n \alpha_1 v_1$) as $n \rightarrow \infty$. This observation can be used in a clever way: we can approximate the eigenvector v_1 by $A^n g_0$ and so if we have a fast way to do matrix-vector multiply, we can approximate eigenvalues and eigenvectors. This is often referred to as the *Power Method*. In a CS Lens I plan to show you how Google's celebrated PageRank algorithm is based on the idea of how eigenvectors can be used for ranking (you can also read more about it here [BSS].²²), calculating the eigenvector using a version of the Power Method is a crucial part of the algorithm.²³
- The vector g_n gets larger and larger as $n \rightarrow \infty$ because $|\lambda_1| > 1$. If both eigenvalues satisfied $|\lambda| < 1$ then $g_n \rightarrow 0$ as $n \rightarrow \infty$. This illustrates the importance of the largest absolute values of the eigenvalues of a matrix in understanding the long term behaviour of systems of the form $A^n g_0$ for some A . If it represents a dynamical system it is related to stability or instability/chaos, if it represents e.g. the evolution of an economical system over time (or the finances of a company) it can be the difference between growth or ruin.²⁴

²²Take a look also at "Landau on Chess Tournaments and Google's PageRank" by Rainer Sinn and Günter M. Ziegler (<https://arxiv.org/pdf/2210.17300.pdf>).

²³An important advantage is that if we already have a good approximation of v_1 , e.g. the page ranks from last week, we can compute a better approximation of v_1 (of this week's rankings) with very few matrix multiplies, you can read more about it here [BSS] and in the references therein.

²⁴Try to modify the Fibonacci recurrence rule so that the new numbers go to zero as $n \rightarrow \infty$. Can you pick a recurrence such that they stabilize as $n \rightarrow \infty$ (without going to ∞ or 0)? Maybe linear algebra students in 2823 years will be studying your sequence!

A few properties of the eigenvalues follow from the fact that, by Corollary 6.0.4,

$$(32) \quad (-1)^n \det(A - zI) = \det(zI - A) = (z - \lambda_1)(z - \lambda_2) \cdots (z - \lambda_n).$$

The polynomial (32) is called the Characteristic Polynomial of the matrix A .²⁵

Proposition 6.1.9. *Given $A \in \mathbb{R}^{n \times n}$ the eigenvalues of A are the same as the ones of A^\top .*

Proof. This follows from (32), and the fact that, for $\det(A - zI) = \det((A - zI)^\top) = \det(A^\top - zI)$. \square

Definition 6.1.10. *Given a matrix $A \in \mathbb{R}^{n \times n}$, the trace of A is defined as*

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

Proposition 6.1.11. *Let $A \in \mathbb{R}^{n \times n}$ and $\lambda_1, \dots, \lambda_n$ its n eigenvalues as they show up in (32) (meaning that a value λ may be repeated, the number of times it shows up is the algebraic multiplicity of the eigenvalue) then*

$$(33) \quad \text{Tr}(A) = \sum_{i=1}^n \lambda_i,$$

$$(34) \quad \det(A) = \prod_{i=1}^n \lambda_i.$$

Remark 6.1.12. *When calculating eigenvalues, Proposition 6.1.11 is very useful to check computations.*

Challenge 1052. Given $B \in \mathbb{R}^{n \times m}$ and $C \in \mathbb{R}^{m \times n}$ show that $\text{Tr}(BC) = \text{Tr}(CB)$.

Challenge 53. Verify Proposition 6.1.11 on your favorite matrix, and also on M in Guiding Example 45.

Challenge 1053. Given $A, B, C \in \mathbb{R}^{n \times n}$ show that we always have $\text{Tr}(ABC) = \text{Tr}(BCA) = \text{Tr}(CAB)$, but that the following is not always true: $\text{Tr}(ABC) = \text{Tr}(ACB)$.

Linear Algebra — A. Bandeira (ETHZ) — Week 12 - 2023.12.08 & 2023.12.13

Please find most up to date notes at: https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_II.pdf

²⁵There is a converse to this in the sense that any monic polynomial can be written as a characteristic polynomial of a matrix, there is a particularly elegant way to build the matrix, if you are interested in learning more, look-up “companion matrix”.

Proof. [of Proposition 6.1.11] To prove (34), simply set $z = 0$ in (32) and note that it gives $(-1)^n \det(A) = (-1)^n \prod_{i=1}^n \lambda_i$. For (33) note that the coefficient of z^{n-1} in the characteristic polynomial (32) is given in the right hand side by $(-\sum_{i=1}^n \lambda_i)$. On the other hand, on the left hand side the coefficient of z^{n-1} is given by the summand in the determinant that multiplies the diagonal elements of $zI - A$, and so it is the coefficient of z^{n-1} of $\prod_{i=1}^n (z - A_{ii})$ which is $-\sum_{i=1}^n A_{ii} = -\text{Tr}(A)$. \square

Caution! 54. We write this caution as Remark 6.1.13 below given how important it is (so that it appears in black font).

Remark 6.1.13. A few *important words of caution*:

- (1) *Even though the eigenvalues of A and A^\top are the same, the eigenvectors are not!*
- (2) *The eigenvalues of $A + B$ are not easily computed from the eigenvalues of A and the ones of B , in particular they are not their sum!*
- (3) *The eigenvalues of AB or BA are not easily computed from the eigenvalues of A and the ones of B , in particular they are not their product!²⁶*
- (4) *Gaussian Elimination doesn't preserve eigenvalues and eigenvectors. The eigenvalues are not the diagonal elements of the U matrix in the $PA = LU$ factorization.²⁷*

While we focus on real eigenvalues on this course, let us see an example of a matrix that has no real eigenvalues and only complex valued ones.

Example 6.1.14. *The eigenvalues of the matrix $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, corresponding to a 90° counter-clockwise rotation, are the solutions to $0 = \det(A - \lambda I) = \lambda^2 + 1$, which are $\lambda_1 = i$ and $\lambda_2 = -i$. The eigenvectors are given by $v_1 = \begin{bmatrix} i \\ 1 \end{bmatrix}$ and $v_2 = \begin{bmatrix} -i \\ 1 \end{bmatrix}$.*

Challenge 55. Try to work out an example for another rotation of \mathbb{R}^2 .

This is a particular case of an orthogonal matrix, whose eigenvalues have a special property.

²⁶Interesting, there is a deep connection between A and B commuting (meaning $AB = BA$) and having the same eigenvectors. This is important in a few fields, in particular in Quantum Physics. If you want to learn more look-up “simultaneously diagonalizable”.

²⁷How to actually compute eigenvalues efficiently is outside of the scope of this course, it turns out that one can do it using the QR decomposition that we learned here as a subroutine. If you want to learn more look-up “QR algorithm”.

Proposition 6.1.15. *Let $Q \in \mathbb{R}^{n \times n}$ be an orthogonal matrix.²⁸ If $\lambda \in \mathbb{C}$ is an eigenvalue of Q , then $|\lambda| = 1$.*

Proof. Let $\lambda \in \mathbb{C}$ be an eigenvalue of Q and $v \in \mathbb{C}^n$ an associated eigenvector. Then $Qv = \lambda v$. Since Q is an orthogonal matrix we have $\|v\|^2 = \|Qv\|^2 = \|\lambda v\|^2 = |\lambda|^2 \|v\|^2$. Since $v \neq 0$ we have $|\lambda| = 1$. \square

Further Fact 56. *If $\lambda \in \mathbb{C}$ is an eigenvalue of a matrix A with real entries, then $\bar{\lambda}$ is also an eigenvalue of A . This can be shown by noticing that $\det(A - \bar{\lambda}I) = \overline{\det(A - \lambda I)}$ (due to the polynomial representation of the determinant in (32), the fact that A has real entries).*

6.1.1. *Repeated eigenvalues.* An important part of the success of the strategy we took in Guiding Example 45 was the fact that we were able to build a basis of \mathbb{R}^2 with eigenvectors of the matrix M . In Theorem 6.1.8 we showed that we can always build a basis of \mathbb{R}^n with eigenvectors of an $n \times n$ matrix A if A has n distinct real eigenvalues. One obstacle could be if some of the eigenvalues are not real valued but, even though we have not focused in complex valued eigenvalues, a straightforward adaption of the proof shows that if A has n distinct eigenvalues (not necessarily real) then there is a basis of \mathbb{C}^n made up of eigenvectors of A . However, repeated eigenvalues can (but doesn't have to) pose a real obstacle to building a basis.

Example 6.1.16. *The matrix $A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ does not have two linearly independent eigenvectors.*

Indeed, $\det(A - \lambda I) = \lambda^2$ which means that $\lambda = 0$ is the only eigenvalue and has algebraic multiplicity 2. However, $N(A - 0I) = N(A)$ only has dimension 1, so there is only one eigenvector (and multiples of it). with

Example 6.1.17. *The zero matrix $A = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$ does have two linearly independent eigenvectors.*

Indeed, $\det(A - \lambda I) = \lambda^2$ which means that $\lambda = 0$ is the only eigenvalue and has algebraic multiplicity 2. But, unlike in Example 6.1.16, $N(A - 0I) = N(A)$ has dimension 2, so there is a basis made up of two eigenvectors (in fact any two linearly independent vectors will be such a basis).

Further Remark 57. Notice that in Example 6.1.16, $N(A^2)$ does have dimension 2. When there exist a positive integer k such that $A^k = 0$ we call A Nilpotent. There is a (rather deep) Theorem that essentially says nilpotency is the only obstacle to getting a complete set of eigenvectors. It

²⁸The \mathbb{C}^n analogue of orthogonal matrices are unitary matrices, $U \in \mathbb{C}^{n \times n}$ that satisfy $U^*U = I$.

roughly says that when there are “missing” eigenvectors they can be found in the Nullspace of powers of $A - \lambda I$, and this gives rise to something called “Jordan Normal Form”.

Definition 6.1.18. *If, given a matrix $A \in \mathbb{R}^{n \times n}$, we can build a basis of \mathbb{R}^n with eigenvectors of A we say that A has a complete set of real eigenvectors.²⁹*

Theorem 6.1.8 states that a matrix with n distinct eigenvalues always has a complete set of real eigenvectors.

Proposition 6.1.19 (Eigenvalues and Eigenvectors of a Projection Matrix). *Let P be the projection matrix on the subspace $U \subseteq \mathbb{R}^n$. Then P has two eigenvalues, 0 and 1, and a complete set of real eigenvectors.*

Proof. Let m be the dimension of U . Let u_1, \dots, u_m be an orthonormal basis of U , and w_1, \dots, w_{n-m} an orthonormal basis of U^\perp . It is easy to see that $Pu_k = 1u_k$ for any $1 \leq k \leq m$ and $Pw_k = 0w_k$ for any $1 \leq k \leq n - m$, so all n vectors are eigenvectors of P (with eigenvalues either 1 or 0). By construction of U^\perp , they form an orthonormal basis. \square

In general when there is an eigenvalue λ with algebraic multiplicity larger than 1, it can be that $N(A - \lambda I)$ is of large enough dimension to find enough linearly independent eigenvectors (as it is the case in projection matrices above, but not in the nilpotent example).

Definition 6.1.20. *Given a matrix $A \in \mathbb{R}^{n \times n}$ and an eigenvalue λ of A we call the dimension of $N(A - \lambda I)$ the geometric multiplicity of λ .*

Further Fact 58. A matrix has a complete set of eigenvectors when the geometric multiplicities are the same as the algebraic multiplicities of all eigenvalues.

Exploratory Challenge 59. Prove Further Fact 58

Example 6.1.21. *For $D \in \mathbb{R}^{n \times n}$ a diagonal matrix, the eigenvalues of D are the diagonal entries of D . The canonical basis e_1, \dots, e_n is a set of eigenvectors of D .*

Challenge 60. Prove the statement in Example 6.1.21

²⁹If the matrix A has complex valued eigenvalues and we can instead build a basis of \mathbb{C}^n we say it has a complete set of eigenvectors. Essentially everything we do below can be (straightforwardly) extended to this case but we will focus on real eigenvalues and eigenvectors for ease of exposition.

Further Fact 1060. The eigenvalues of an $n \times n$ triangular matrix are the n values in the diagonal. However, triangular matrices may not have a complete set of eigenvectors.

Challenge 1061. Prove this fact. Hint: For the positive part use (32), for the negative part recall Example 6.1.16.

Challenge 1062. Let's say a matrix $A \in \mathbb{R}^{n \times n}$ has an LU decomposition (without the need for P in $PA = LU$). Remark 6.1.13 says the eigenvalues of $A \in \mathbb{R}^{n \times n}$ are not the ones of U in the LU decomposition $PA = LU$. The eigenvalues of U are indeed their diagonal entries, and the eigenvalues of L are all 1 (by Further Fact 6.1.1). Why is it the case that the eigenvalues of A are not the diagonal entries of U ?

6.2. Diagonalizing a Matrix and Change of Basis of a Linear Transformation. Let us continue dissecting Guiding Example 45. Essentially, what we did was to write g_0 in the basis v_1, v_2 of eigenvectors of M and then exploit the fact that linear transformation given M had a very simple behaviour when written in the basis v_1, v_2 (the coefficients simply were multiplied by the eigenvalues of M). This motivates us to take a detour in briefly studying linear transformations written in different basis, and to discuss “change of bases”.

6.2.1. Change of basis. For this detour we will briefly consider $m \times n$ matrices, before returning to square matrices when discussing eigenvalues and eigenvectors.

Let $A^{m \times n}$ be a matrix representing a linear transformation $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$ given by $x \in \mathbb{R}^n \rightarrow Ax \in \mathbb{R}^m$, with both input and output written in the canonical bases as $x = \sum_{j=1}^n x_j e_j$ and $Ax = \sum_{i=1}^m (Ax)_i e_i$. Recall (Example 4.4.2) that $(e_i)_j = \delta_{ij}$, and that $(Ax)_i$ is the i -th entry of the vector Ax .

Now, let's say we have a basis for \mathbb{R}^n given by u_1, \dots, u_n and one for \mathbb{R}^m given by v_1, \dots, v_m (neither being necessarily the canonical basis) and we want to understand the the linear transformation L written in this basis. Then L takes a vector $x = \sum_{j=1}^n \alpha_j u_j$ and outputs $L(x) = \sum_{i=1}^m \beta_i v_i$.

We want to compute the matrix B that takes $\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$ to $\begin{bmatrix} \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}$. In other words, such

that $B\alpha = \beta$. Let $U \in \mathbb{R}^{n \times n}$ be the matrix whose columns are the basis elements u_1, \dots, u_n and $V \in \mathbb{R}^{m \times m}$ the matrix whose columns are the basis elements v_1, \dots, v_m . Then, $x = U\alpha$ and $L(x) = V\beta$ and so $\beta = V^{-1}AU\alpha$, the matrix B , corresponding to the linear transformation L written in the new bases is $B = V^{-1}AU$. Note that we can do change of basis between any pair of basis, it needs not be from the canonical basis to another basis, in that case the role of U and V

would be played by the change of basis matrix (the matrix that maps the coefficients of a vector written in the old basis, to its coefficients when written in the new basis).

$$(35) \quad \begin{array}{l} L: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad \text{linear transformation} \\ L\left(\sum_{j=1}^n x_j e_j\right) = \sum_{i=1}^m (Ax)_i e_i \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \\ L\left(\sum_{j=1}^n \alpha_j u_j\right) = \sum_{i=1}^m (B\alpha)_i v_i \quad \alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} \end{array}$$

where $B = V^{-1}AV \in \mathbb{R}^{m \times n}$, $U = \begin{bmatrix} u_1 & \cdots & u_n \end{bmatrix} \in \mathbb{R}^{n \times n}$, $V = \begin{bmatrix} v_1 & \cdots & v_m \end{bmatrix} \in \mathbb{R}^{m \times m}$.

6.2.2. Diagonalizing a Matrix. Let us focus back on square matrices $A \in \mathbb{R}^{n \times n}$. In particular, let A be a matrix with a complete set of real eigenvectors (in the sense of Definition 6.1.18) and let $v_1, \dots, v_n \in \mathbb{R}^{n \times n}$ be a basis formed with eigenvectors of A . The crucial fact we used in Guiding Example 45 also holds in this general situation: if we write a vector $x \in \mathbb{R}^n$ as $x = \sum_{i=1}^n \alpha_i v_i$ then $Ax = \sum_{i=1}^n \lambda_i \alpha_i v_i$ (and also $A^k x = \sum_{i=1}^n \lambda_i^k \alpha_i v_i$, for $k \geq 1$, where λ_i is the eigenvalue associated with the eigenvector v_i). One way to think about this is that the linear transformation corresponding to the matrix A , when written in the basis V is simply a diagonal matrix/transformation. This is the key idea behind *Matrix Diagonalization*. This is one of the most important facts in Linear Algebra.

Theorem 6.2.1. *Let $A \in \mathbb{R}^{n \times n}$ be a matrix with a complete set of real eigenvectors (in the sense of Definition 6.1.18) and let $v_1, \dots, v_n \in \mathbb{R}^{n \times n}$ be a basis formed with eigenvectors of A and let $\lambda_1, \dots, \lambda_n$ be the associated eigenvalues (λ_i associated to v_i). Let V be the matrix whose columns are the eigenvectors v_i , $V = \begin{bmatrix} v_1 & \cdots & v_n \end{bmatrix} \in \mathbb{R}^{n \times n}$. Then,*

$$(36) \quad A = V\Lambda V^{-1},$$

where Λ is a diagonal matrix with $\Lambda_{ii} = \lambda_i$ (and $\Lambda_{ij} = 0$ for all $i \neq j$).

Proof. Since v_1, \dots, v_n is a basis, V is an invertible matrix, so it suffices to prove that

$$(37) \quad V^{-1}AV = \Lambda.$$

This can be done by direct calculation: For any $1 \leq j \leq n$, the j -th column of the matrix $V^{-1}AV$ is given by

$$(V^{-1}AV)_{.j} := (V^{-1}AV) e_j = V^{-1}A v_j = V^{-1} \lambda_j v_j = \lambda_j V^{-1} v_j = \lambda_j e_j,$$

since $V^{-1}v_j = V^{-1}Ve_j = e_j$. Recall that e_j is the vector in \mathbb{R}^n with a 1 in j -th entry and zero elsewhere. Since for any $1 \leq j \leq n$, $\lambda_j e_j$ is also the j -th column of Λ , we have that $V^{-1}AV = \Lambda$. \square

Definition 6.2.2 (Diagonalizable Matrix). *A matrix $A \in \mathbb{R}^{n \times n}$ is called a diagonalizable matrix if there exists an invertible matrix V such that $V^{-1}AV = \Lambda$, where Λ is a diagonal matrix.*

Challenge 63. Most properties of the eigenvalues are very easy to prove by using Theorem 6.2.1 (for the matrices that have a complete set of eigenvectors). Try it!

The eigenvalues of Λ are also $\lambda_1, \dots, \lambda_n$ (recall Example 6.1.21). More generally, for an invertible matrix S we always have that A and $S^{-1}AS$ have the same eigenvalues.

Definition 6.2.3 (Similar Matrices).³⁰ *We say that $A \in \mathbb{R}^{n \times n}$ and $B \in \mathbb{R}^{n \times n}$ are similar matrices if there exists an invertible matrix S such that $B = S^{-1}AS$.*

Proposition 6.2.4. *Similar matrices have the same eigenvalues.*

Challenge 64. Try to show Proposition 6.2.4 via (32).

Challenge 65. Try to show that if λ is an eigenvalue of $S^{-1}AS$ (with associated eigenvector v) then it is also an eigenvalue of A and compute the associated eigenvector in terms of v and S . (without using Proposition 6.2.4 or (32)).

Remark 6.2.5 (Diagonalizing a matrix and finding a good basis for a linear transformation). *If we have a matrix $A \in \mathbb{R}^{n \times n}$ with a complete set of real eigenvectors then Theorem 6.2.1 tells us that the corresponding linear transformation, when viewed in the bases v_1, \dots, v_n is simply a diagonal matrix (recall that in this case $B = \Lambda$, see (35)). This is a remarkable fact: since most matrices have a full set of eigenvectors (in particular all for which the eigenvalues are all distinct do) this says that all the corresponding linear combinations, regardless of how complicated they might seem, are actually just a diagonal operation when viewed in the basis v_1, \dots, v_n .*

6.3. Symmetric Matrices and the Spectral Theorem. This section is devoted to real symmetric matrices,³¹ meaning matrices $A \in \mathbb{R}^{n \times n}$ for which $A^\top = A$ (see Further Remark 69 for a brief discussion of how symmetric matrices appear naturally in several settings). The main goal of this section is to prove the Spectral Theorem.

³⁰The operation $A \rightarrow S^{-1}AS$ is sometimes called conjugation but this is not to be confused with complex conjugation $z \rightarrow \bar{z}$, one term comes from “conjugation” in group theory, the other from “conjugation” in complex analysis.

³¹The same theory can be developed (by a straightforward adaption) to complex matrices but the property of being symmetric is replaced by being Hermitian, which means that a matrix $A = A^* = \bar{A}^\top$. In both situations we say A is self-adjoint.

Theorem 6.3.1 (Spectral Theorem). *Any symmetric³² matrix $A \in \mathbb{R}^{n \times n}$ has n real eigenvalues and an orthonormal basis made of eigenvectors of A .*

Together with Theorem 6.2.1 this implies the following corollary.

Corollary 6.3.2. *For any symmetric matrix $A \in \mathbb{R}^{n \times n}$ there exists an orthogonal matrix $V \in \mathbb{R}^{n \times n}$ (whose columns are eigenvectors of A) such that*

$$A = V\Lambda V^\top,$$

where $\Lambda \in \mathbb{R}^{n \times n}$ is a diagonal matrix with the eigenvalues of A in its diagonal (and $V^\top V = I$).

Remark 6.3.3 (Eigendecomposition). *The decompositions in Corollary 6.3.2 and Theorem 6.2.1 are called Eigendecompositions.*

The following follows easily from the Spectral Theorem.

Corollary 6.3.4. *The rank of a real symmetric matrix A is the number of non-zero eigenvalues (counting repetitions).*

Remark 6.3.5. *For general $n \times n$ (non-symmetric) matrices, the rank is n minus the dimension of the nullspace, so it is n minus the geometric multiplicity of $\lambda = 0$. Since symmetric matrices always have a complete set of eigenvalues and eigenvectors, the geometric multiplicities are always the same as the algebraic multiplicities.*

Proposition 6.3.6. *Let A be a real $n \times n$ symmetric matrix and let v_1, \dots, v_n be an orthonormal basis of eigenvectors of A (the columns of the matrix V in Corollary 6.3.2) and $\lambda_1, \dots, \lambda_n$ the associated eigenvalues. Then*

$$A = \sum_{k=1}^n \lambda_k v_k v_k^\top$$

Proof. Follows directly by Corollary 6.3.2 and the fact that $V^\top V = I$. □

We “build up” to the proof of Theorem 6.3.1 with a few propositions.

Proposition 6.3.7. *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix and $\lambda \in \mathbb{C}$ an eigenvalue of A , then $\lambda \in \mathbb{R}$.*

³²Essentially the same proof shows that any $A \in \mathbb{C}^{n \times n}$ that is Hermitian, also called self-adjoint $A^* = A$, has n real eigenvalues and a complete set of orthonormal eigenvectors in \mathbb{C}^n (the matrix whose columns are these vectors is a Unitary matrix U satisfying $U^*U = I$).

Proof. Let $v \in \mathbb{C}^n$ be an eigenvector associated with the eigenvalue λ . We have $Av = \lambda v$. Recall that, for a matrix (or vector) M , its Hermitian conjugate is given by $M^* = \overline{M}^\top$. Since A is real symmetric we have $A^* = A$. Thus, we have

$$\overline{\lambda} \|v\|^2 = \overline{\lambda} v^* v = (\lambda v)^* v = (Av)^* v = v^* A^* v = v^* Av = v^* \lambda v = \lambda \|v\|^2.$$

Since $v \neq 0$, then $\|v\| \neq 0$ and so $\lambda = \overline{\lambda}$. This implies that $\lambda \in \mathbb{R}$. \square

This, together with Theorem 6.1.4, immediately implies the following.

Corollary 6.3.8. *Every symmetric matrix $A \in \mathbb{R}^{n \times n}$ (satisfying $A = A^\top$) has a real eigenvalue λ .*

Remark 6.3.9. *The fact that two eigenvectors of a real symmetric matrix are orthogonal follows from Theorem 6.3.1 but it is useful to see a simple argument of that (the main difficulty of proving Theorem 6.3.1 is proving that the matrix indeed has a complete set of eigenvectors). Let's say we have $\lambda_1 \neq \lambda_2$ eigenvalues of a real symmetric matrix A and $v_1, v_2 \in \mathbb{R}^n \setminus \{0\}$ corresponding eigenvectors. Then*

$$\lambda_1 v_1^\top v_2 = (Av_1)^\top v_2 = v_1^\top A^\top v_2 = v_1^\top Av_2 = v_1^\top (Av_2) = \lambda_2 v_1^\top v_2,$$

since $\lambda_1 \neq \lambda_2$ we must have that $v_1^\top v_2 = 0$

Further Remark 66. Corollary 6.3.8 is a great example of the usefulness of complex numbers. Even though it is a statement just about real matrices and real eigenvalues we proved it by going through the complex numbers and using results in Complex Analysis. There is an alternative proof without going through the complex numbers but it would need more background, and I find this one more transparent.³³

Linear Algebra — A. Bandeira (ETHZ) — Week 13 - 2023.12.15 & 2023.12.20

Please find most up to date notes at: https://ti.inf.ethz.ch/ew/courses/LA23/notes_part_II.pdf

Proof. [of Theorem 6.3.1]

Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix. We will prove the following by induction, which for $k = n$ implies the theorem we want to show:

- For any $1 \leq k \leq n$ there are k orthogonal eigenvectors of A .

The base case $k = 1$ follows directly by Corollary 6.3.8 as we can always normalize the eigenvector to have norm 1.

³³If you would like to see a nice example of how improving knowledge can lead to much simpler and more transparent methods/models, search “Ptolemaic Epicycle Machine”.

We now assume that the statement is true for k and show it for $k + 1$. We will show that if a real symmetric matrix A has k (with $1 < k < n$) orthonormal eigenvectors then we can build an extra one, orthogonal to the others (to achieve norm 1 we simply need to normalize it).³⁴

Let v_1, \dots, v_k denote k orthonormal eigenvectors of A and $\lambda_1, \dots, \lambda_k$ the respective eigenvalues. Let u_{k+1}, \dots, u_n be an orthonormal basis of the orthogonal complement of the span of v_1, \dots, v_k . Let V_k be the $n \times n$ matrix whose i -th column is v_i if $i \leq k$ and u_i if $i > k$. V_k is an orthogonal matrix. Moreover, let us define $B \in \mathbb{R}^{n \times n}$ as $B = V^\top A V$, then:

$$\begin{aligned} B = V^\top A V &= \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_k^\top & - \\ & u_{k+1}^\top & \\ & \vdots & \\ - & u_n^\top & - \end{bmatrix} \begin{bmatrix} | & & | & | & & | \\ Av_1 & \cdots & Av_k & Au_{k+1} & \cdots & Au_n \\ | & & | & | & & | \end{bmatrix} \\ &= \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_k^\top & - \\ & u_{k+1}^\top & \\ & \vdots & \\ - & u_n^\top & - \end{bmatrix} \begin{bmatrix} | & & | & | & & | \\ \lambda v_1 & \cdots & \lambda v_k & Au_{k+1} & \cdots & Au_n \\ | & & | & | & & | \end{bmatrix} \\ &= \begin{bmatrix} \Lambda_k & 0_{k \times (n-k)} \\ 0_{(n-k) \times k} & C \end{bmatrix}, \end{aligned}$$

where Λ_k is a diagonal matrix with $\lambda_1, \dots, \lambda_k$ in the diagonal, $0_{(n-k) \times k}$ and $0_{k \times (n-k)}$ are zero matrices of size respectively $(n-k) \times k$ and $k \times (n-k)$. C is a $(n-k) \times (n-k)$ symmetric matrix.

Since C is a $(n-k) \times (n-k)$ symmetric matrix, Theorem 6.3.8 implies it has a real eigenvalue λ_{k+1} and a real eigenvector $y \in \mathbb{R}^{n-k}$. Let $w \in \mathbb{R}^n$ be the vector with 0 in the first k coordinates and y in the remaining $n-k$, in other words

$$w_i = \begin{cases} 0 & \text{if } i \leq k \\ y_{i-k} & \text{if } i > k. \end{cases}$$

³⁴In a first reading of the proof I recommend taking $k = 1$ in the induction step, as it is simpler while already containing all the relevant ideas.

We have

$$Bw = \begin{bmatrix} \Lambda_k & 0_{k \times (n-k)} \\ 0_{(n-k) \times k} & C \end{bmatrix} \begin{bmatrix} 0_{k \times 1} \\ y \end{bmatrix} = \begin{bmatrix} 0_{k \times 1} \\ Cy \end{bmatrix} = \begin{bmatrix} 0_{k \times 1} \\ \lambda_{k+1}y \end{bmatrix} = \lambda_{k+1}w.$$

Let $v_{k+1} := Vw$. Since V is orthogonal we have that $A = VB V^\top$. Thus,

$$Av_{k+1} = VB V^\top v_{k+1} = VBw = V\lambda_{k+1}w = \lambda_{k+1}v_{k+1},$$

so v_{k+1} is an eigenvector of A . To see that it is orthogonal to v_1, \dots, v_k note that the inner products $v_i^\top v_{k+1}$ for $i \leq k$ appear in the first k entries of $V^\top v_{k+1} = w$ and that w has its first k coordinates 0 by construction. By normalizing the vector we can have it have unit norm. □

Proposition 6.3.10 (Rayleigh Quotient). *Given a symmetric matrix $A \in \mathbb{R}^{n \times n}$ the Rayleigh Quotient, defined for $x \in \mathbb{R}^n \setminus \{0\}$, as*

$$R(x) = \frac{x^\top Ax}{x^\top x}$$

attains its maximum at $R(v_{\max}) = \lambda_{\max}$ and its minimum at $R(v_{\min}) = \lambda_{\min}$ where λ_{\max} and λ_{\min} are respectively the largest and smallest eigenvalues of A and v_{\max}, v_{\min} their associated eigenvectors.

Proof. It is easy to see that $R(v_{\max}) = \lambda_{\max}$ and $R(v_{\min}) = \lambda_{\min}$ so it suffices to show that, for all $x \in \mathbb{R}^n \setminus \{0\}$ we have $\lambda_{\min} \leq R(x) \leq \lambda_{\max}$. Using Proposition 6.3.6 we can write, for $x \in \mathbb{R}^n \setminus \{0\}$,

$$R(x) = \frac{x^\top (\sum_{k=1}^n \lambda_i v_i v_i^\top) x}{\|x\|^2} = \frac{\sum_{k=1}^n \lambda_i (x^\top v_i)^2}{\|x\|^2},$$

where v_1, \dots, v_n form an orthonormal basis of eigenvectors of A and $\lambda_1, \dots, \lambda_n$ are the associated eigenvalues. Since $(x^\top v_i)^2 \geq 0$ for all $1 \leq i \leq n$ we have that, for all $1 \leq i \leq n$,

$$\lambda_{\min} (x^\top v_i)^2 \leq \lambda_i (x^\top v_i)^2 \leq \lambda_{\max} (x^\top v_i)^2.$$

Collecting all these inequalities we get

$$\lambda_{\min} \frac{\sum_{k=1}^n (x^\top v_i)^2}{\|x\|^2} \leq \frac{\sum_{k=1}^n \lambda_i (x^\top v_i)^2}{\|x\|^2} \leq \lambda_{\max} \frac{\sum_{k=1}^n (x^\top v_i)^2}{\|x\|^2}.$$

To conclude the proof note that, since the v_i 's are orthonormal, the matrix V with the v_i 's as columns is orthogonal and $\sum_{k=1}^n (x^\top v_i)^2 = \|Vx\|^2 = \|x\|^2$ and so $\frac{\sum_{k=1}^n (x^\top v_i)^2}{\|x\|^2} = 1$. □

Definition 6.3.11 (Positive Definite and Positive Semidefinite matrix). A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is said to be Positive Semidefinite (PSD) if all its eigenvalues are non-negative. If all the eigenvalues of A are strictly positive then we say A is Positive Definite (PD).

Exploratory Challenge 67. Even though the eigenvalues of $A + B$ are not easily described by the eigenvalues of A and the ones of B it turns out that if both are PSD (or PD) then so is the sum. Can you show that?

The following follows directly from Proposition 6.3.10.

Proposition 6.3.12. A symmetric matrix $A \in \mathbb{R}^{n \times n}$ is Positive Semidefinite if and only if $x^\top Ax \geq 0$ for all $x \in \mathbb{R}^n$. Analogously, a symmetric matrix $A \in \mathbb{R}^{n \times n}$ is Positive Definite if and only if $x^\top Ax > 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$.

Fact 6.3.13. If two $n \times n$ matrices A and B are PSD (or PD) then their sum is PSD (or PD).

Challenge 68. Exploratory Challenge 67 looked pretty difficult, but now with Proposition 6.3.12 it is much easier, try to prove Fact 6.3.13.

Definition 6.3.14 (Gram Matrix). Given n vectors, v_1, \dots, v_n in \mathbb{R}^m we call their Gram Matrix the $n \times n$ matrix of inner products

$$G_{ij} = v_i^\top v_j.$$

Note that if $V \in \mathbb{R}^{m \times n}$ is the matrix whose columns are the n vectors, then $G = V^\top V$ is the Gram matrix of V .

Remark 6.3.15. Given a matrix $A \in \mathbb{R}^{m \times n}$, as an abuse of notation, we sometimes also call AA^\top also a Gram matrix of A . Notice that, if $a_1, \dots, a_n \in \mathbb{R}^m$ are the columns of A then AA^\top is $m \times m$ and

$$(38) \quad AA^\top = \sum_{i=1}^n a_i a_i^\top.$$

Proposition 6.3.16. Given a real matrix $A \in \mathbb{R}^{m \times n}$, the non-zero eigenvalues of $A^\top A \in \mathbb{R}^{n \times n}$ are the same as the ones of $AA^\top \in \mathbb{R}^{m \times m}$. Both matrices are symmetric and positive semidefinite.

Proof. Let r be the rank of A . We know $\text{rank}(A) = \text{rank}(A^\top) = \text{rank}(A^\top A) = \text{rank}(AA^\top)$ (recall Proposition 4.3.2 and Challenge 9). It is straightforward that both $A^\top A$ and AA^\top are symmetric. Let us prove they are positive semidefinite. We have $x^\top A^\top A x = \|Ax\|^2 \geq 0$ for all x which implies $A^\top A$ is PSD, and the same argument can be used for AA^\top .

Now, both AA^\top and $A^\top A$ have a complete set of real eigenvalues and orthogonal eigenvectors. Let $\lambda_1, \dots, \lambda_r$ be the r non-zero eigenvalues of $A^\top A$ and v_1, \dots, v_r be the corresponding eigenvalues.

We have, for $1 \leq k \leq r$, $A^\top Av_k = \lambda_k v_k$, multiplying by A both sides we get $AA^\top Av_k = \lambda_k Av_k$ and so λ_k is an eigenvalue of AA^\top with eigenvector Av_k . Furthermore, For $j \neq k$ we have $(Av_j)^\top (Av_k) = v_j^\top A^\top Av_k = v_j^\top v_k = 0$ and so the r eigenvectors of AA^\top built this way are orthogonal, and so $\lambda_1, \dots, \lambda_r$ are the nonzero eigenvalues of AA^\top . \square

Proposition 6.3.17. [Cholesky decomposition] Every symmetric positive semidefinite matrix M is a gram matrix of an upper triangular matrix C . $M = C^\top C$ is known as the Cholesky Decomposition.³⁵

Proof. Let M be a symmetric positive semidefinite matrix. Corollary 6.3.2 gives us a decomposition $M = V\Lambda V^\top$ with Λ a diagonal matrix with the eigenvalues of M in the diagonal. Since M is PSD, the diagonal entries of Λ are non-negative and so we can build $\Lambda^{1/2}$ by taking the square root of each diagonal entry of Λ . Then $M = (V\Lambda^{1/2})(V\Lambda^{1/2})^\top$. To make the matrices in the decomposition be upper triangular, simply take the QR decomposition (recall Definition 4.4.12) $(V\Lambda^{1/2})^\top = QR$ with Q such that $Q^\top Q = I$ and R upper triangular. We have $M = (V\Lambda^{1/2})(V\Lambda^{1/2})^\top = (QR)^\top (QR) = R^\top Q^\top QR = R^\top R$. Taking $C = R$ establishes the Proposition.³⁶ \square

Further Remark 69. At first glance, Symmetric matrices look very special (since we must have $A^\top = A$) but they they actually appear very often in both applications and pure mathematics. For example, in my own work, I rarely encounter non-symmetric matrices. There are (at least) two reasons for this: (i) For any matrix B we can form a symmetric matrix $B^\top B$ from which we can study B , this is going to be the key idea being the Singular Value Decomposition. (ii) In many instances, matrices represent relationship between objects — for example, A_{ij} can represent a friendship connection (or a similarity measure) between person (or data point) i and j and in many cases such relationships are symmetric.

7. SINGULAR VALUE DECOMPOSITION; AND SOME OPEN QUESTIONS IN LINEAR ALGEBRA

7.1. The Singular Value Decomposition. We are now reaching what I view as “the ultimate theorem of our class”, the Singular Value Decomposition (SVD). In fact, a mentor of mine once

³⁵To compute an upper triangular matrix C such that $M = C^\top C$ one can use the LU decomposition and needs not to compute an eigendecomposition of the matrix M .

³⁶This is not the classical construction of the Cholesky Decomposition. The classical construction is with Gaussian Elimination, but at this stage of the course I think this is more transparent. Note also that when using Gaussian Elimination C will be a square matrix, while here R can be rectangular if M is not full rank (which makes it a more economical decomposition).

said: “If I could take only one algorithm with me to a desert island, it would be the SVD”. The SVD is a way to generalize the eigendecomposition to non-symmetric, and even non-square, matrices. Instead of eigenvalues we will have singular values and instead of eigenvectors we will have (right and left) singular vectors.

Definition 7.1.1 (SVD — Singular Value Decomposition). *Let $A \in \mathbb{R}^{m \times n}$. There exist orthogonal matrices $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ such that*

$$(39) \quad A = U \Sigma V^\top,$$

where $\Sigma \in \mathbb{R}^{m \times n}$ is a diagonal matrix, in the sense that $\Sigma_{ij} = 0$ when $i \neq j$, and the diagonal elements are non-negative and ordered in descending order. $U^\top U = I$ and $V^\top V = I$.

The columns u_1, \dots, u_m of U are called the left singular vectors of A and are orthonormal. The columns v_1, \dots, v_n of V are called the right singular vectors of A and are orthonormal. The diagonal elements of Σ , $\sigma_i = \Sigma_{ii}$ are called the singular values of A and are ordered as

$$\sigma_1 \geq \dots \geq \sigma_{\min\{m,n\}}.$$

Remark 7.1.2. *If A has rank r we can write the SVD in a more compact form:*

$$(40) \quad A = U_r \Sigma_r V_r^\top,$$

where $U_r \in \mathbb{R}^{m \times r}$ contains the first r left singular vectors, $V_r \in \mathbb{R}^{n \times r}$ contains the first r right singular vectors and $\Sigma_r \in \mathbb{R}^{r \times r}$ is a diagonal matrix with the first r singular values. Notice that storing such a decomposition in the computer requires storing $r \times (m + n + 1)$ real numbers rather than $m \times n$ real numbers which would be required to store A naively. When a matrix has small rank these are crucial savings.³⁷

Oftentimes the subscript is omitted and the compact SVD is simply written as $U \Sigma V^\top$ while specifying the dimensions of the matrices involved to specify which form of the SVD is being considered.

Remark 7.1.3. *Let $A \in \mathbb{R}^{m \times n}$ and $A = U \Sigma V^\top$ be its SVD (as in (39)) then*

$$AA^\top = U \left(\Sigma \Sigma^\top \right) U^\top,$$

and so the left singular vectors of A , the columns of U are the eigenvectors of AA^\top and the singular values of A are the square-root of the eigenvalues of AA^\top (note that $\Sigma \Sigma^\top$ is $m \times m$ diagonal). If $m > n$, A has n singular values and AA^\top has m eigenvalues (which is larger than n), but the “missing” ones are 0.

³⁷Taking this one step forward, when a matrix is well approximated by a low rank matrix, oftentimes one stores only a small rank approximation of a matrix A , this is a crucial idea in tasks ranging from Image Compressions, Numerical Analysis, and Machine Learning, see Section 7.3.

Analogously,

$$A^T A = V \left(\Sigma^T \Sigma \right) V^T,$$

and so the right singular vectors of A , the columns of V are the eigenvectors of $A^T A$ and the singular values of A are the square-root of the eigenvalues of $A^T A$ (note that $\Sigma^T \Sigma$ is $n \times n$ diagonal). If $n > m$, A has m singular values and $A^T A$ has n eigenvalues (which is larger than m), but the “missing” ones are 0.

This observation makes it easier to write the singular values and singular vectors of A in terms of eigenvalues and eigenvectors of AA^T and $A^T A$, which are symmetric matrices (and directly implies, e.g., uniqueness of singular values; and the fact that the rank of a matrix is the number of nonzero singular values). In fact, the proof of the existence SVD will heavily rely on the Spectral Theorem.

An important direct consequence of the SVD, and in particular of (40) is that we can write any rank- r matrix $A \in \mathbb{R}^{m \times n}$ as a sum of r rank-1 matrices:

Proposition 7.1.4. *Let $A \in \mathbb{R}^{m \times n}$ be a matrix with rank r . Let $\sigma_1, \dots, \sigma_r$ be the non-zero singular values of A , u_1, \dots, u_r the corresponding left singular vectors and v_1, \dots, v_r the corresponding right singular vectors. Then*

$$(41) \quad A = \sum_{k=1}^r \sigma_k u_k v_k^T.$$

Challenge 70. The SVD is a powerful tool. Many of the things we did in this course become significantly simpler with the SVD. Now that you have the SVD, try to reread these notes and try to re-interpret the results we derived in terms of the SVD. For example, the Moore-Penrose Pseudoinverse has a very simple description of the SVD, it corresponds to swapping U and V and replacing the non-zero singular values by their inverses, while keeping the zero ones zero. Try to derive this!

Theorem 7.1.5. *[The SVD – the Ultimate Theorem of ETHZ 401-0131-00L]*

Every matrix $A \in \mathbb{R}^{m \times n}$ has an SVD decomposition of the form (39).

In other words:

Every linear transformation is diagonal when viewed in the bases of the singular vectors.

Proof. Let $A^{m \times n}$. Let r be the rank of $A^{m \times n}$. We will build a compact SVD as in (40). It is easy to see that we can get an SVD in the sense of (39) from a compact one by adding singular values that are zero and extending the singular vectors in both U_r and V_r to orthonormal bases.

By Theorem 6.3.1 and Corollary 6.3.2 the matrix AA^\top has a complete set of orthonormal eigenvectors and can be written as

$$(42) \quad AA^\top = U\Lambda U^\top,$$

where $U \in \mathbb{R}^{m \times m}$ is orthogonal and Λ is diagonal. Let us write (42) ordering the diagonal entries of Λ in decreasing order. Furthermore, let us write (42) also in a compact form, by keeping only the r non-zero eigenvalues (and corresponding eigenvectors), so

$$AA^\top = U_r \Lambda_r U_r^\top$$

for $U_r \in \mathbb{R}^{m \times r}$ such that $U_r^\top U_r = I$ and Λ_r is $r \times r$ diagonal with the non-zero eigenvalues of AA^\top . By Proposition 6.3.16 the eigenvalues of AA^\top are non-negative and so the diagonal entries of Λ_r are positive. Let $\Sigma_r \in \mathbb{R}^{r \times r}$ be the diagonal matrix with diagonal entries $\sigma_i := (\Sigma_r)_{ii} = \sqrt{\Lambda_{ii}}$. Our goal is to show that there is a $r \times n$ matrix V_r , with orthonormal columns, such that $A = U_r \Sigma_r V_r^\top$. We would have $\Sigma_r^{-1} U_r^\top A = \Sigma_r^{-1} U_r^\top U_r \Sigma_r V_r^\top = V_r^\top$, or equivalently $V_r = A^\top U_r \Sigma_r^{-1}$. Motivated by this, let's set

$$V_r := A^\top U_r \Sigma_r^{-1},$$

this corresponds to a matrix with columns v_1, \dots, v_r given by $v_k = \frac{1}{\sigma_k} A^\top u_k$. To conclude we need to show that this construction indeed gives a compact SVD, for this we still need to show two things:

(1) $V_r^\top V_r = I$. This can be verified by direct computation, while recalling that $AA^\top = U_r \Lambda_r U_r^\top$:

$$V_r^\top V_r = \left(A^\top U_r \Sigma_r^{-1} \right)^\top A^\top U_r \Sigma_r^{-1} = \Sigma_r^{-1} U_r^\top A A^\top U_r \Sigma_r^{-1} = \Sigma_r^{-1} = \Sigma_r^{-1} U_r^\top U_r \Lambda_r U_r^\top U_r \Sigma_r^{-1} = I$$

(2) $A = U_r \Sigma_r V_r^\top$. Note that

$$U_r \Sigma_r V_r^\top = U_r \Sigma_r \left(A^\top U_r \Sigma_r^{-1} \right)^\top = U_r U_r^\top A,$$

but, by construction, $U_r U_r^\top$ is the projection on $C(AA^\top)$ and so it is also the projection on $C(A)$ (since $C(A) = C(AA^\top)$, use Proposition 4.3.2 for A^\top). Thus $U_r U_r^\top A = A$.

□

7.2. Vector and Matrix Norms. A short section on vector and matrix norms. So far, the norm of a vector $x \in \mathbb{R}^n$ was simply given by $\|x\| = \sqrt{x^\top x}$ but there are instances where it makes sense to measure the “length” of vectors in other ways. One popular way is called the “Manhattan distance” since when traveling in Manhattan one cannot take advantage of Pythagoras Theorem because that would involve cutting through buildings, that norm is given by $\|x\|_1 = \sum_{i=1}^n |x_i|$. In

general, for $1 \leq p \leq \infty$ the ℓ_p norm is given by

$$(43) \quad \|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p},$$

for $p < \infty$, and $\|x\|_\infty = \max_i |x_i|$. Notice that $\|\cdot\|_2$ corresponds to the Euclidean norm that we have used in this course.³⁸

Challenge 71 (★). Prove that, for all $x \in \mathbb{R}^n$, we have $\|x\|_2 \leq \|x\|_1 \leq \sqrt{n}\|x\|_2$.

In several situations one also needs to “measure” the size of matrices (for example, when talking about a matrix being close to another one, we need a notion of distance, or norm of the difference).

Definition 7.2.1 (Two matrix norms). Given a matrix $A \in \mathbb{R}^{m \times n}$ we define two matrix norms:

- $\|A\|_F$, known as the Frobenius norm, is defined as

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2},$$

- $\|A\|_{op}$, known as operator or spectral norm, is defined as

$$\|A\|_{op} = \max_{\substack{x \in \mathbb{R}^n \\ s.t. \|x\|=1}} \|Ax\|.$$

Further Proposition 1071. Given $A \in \mathbb{R}^{m \times n}$ with singular values $\sigma_1 \geq \dots \geq \sigma_{\min\{m,n\}}$. We have

- (1) $\|A\|_F^2 = \text{Tr}(A^T A)$
- (2) $\|A\|_F^2 = \sum_{i=1}^{\min\{m,n\}} \sigma_i^2$
- (3) $\|A\|_{op} = \sigma_1$
- (4) $\|A\|_{op} \leq \|A\|_F \leq \sqrt{\min\{m,n\}} \|A\|_{op}$.

Challenge 72 (★). Prove Further Proposition 1071.³⁹

7.3. Low-Rank Modelling, Images, Data, and Principal Component Analysis.

³⁸The ℓ_1 norm is notable for promoting sparsity when one attempts to minimize it to solve underdetermined linear systems. This is the key idea behind “Compressed Sensing”, and plays a crucial role in many imaging/sensing technologies. You can read more about it in Section 12 of [BM23] or Chapter 10 of [BSS] and references therein.

³⁹Hint: The order I chose for Further Proposition 1071 was so that it is easiest to prove these properties in this order. This is a good exercise to help consolidate your understanding of the SVD, and other concepts in this course.

Further Remark 73. This section serves as a CS Lens. One of the most powerful ideas in Linear Algebra, from an application perspective, is that if one takes a matrix A , writes the SVD of A and then sets of zero all singular values except the largest r , the resulting matrix A_r is good rank- r approximation to the original matrix. In fact, in several senses it is the best possible rank- r approximation. The celebrated Eckart-Young-Minsky Theorem says that it is the best approximation in the sense that $\|A - A_r\|$ is the smallest possible in a large class of matrix norms $\|\cdot\|$. This idea is key for recommendation systems (and played a central role in the “Netflix prize”), it is the basis for Principal Component Analysis (PCA), and can even do image compression. I will describe some of these things as a CS Lens but you can also read an exposition of it in lecture notes from another course. With your background you can now read Sections 3–5 of [BM23], or you can wait to learn all of this in other classes in your degree.⁴⁰

7.10. Some Mathematical Open Problems. Now that we have covered the notion of eigenvalues and eigenvectors, there are a few fascinating open questions we can state. These are questions (or conjectures), that we currently do not know the answer to (or that we are not sure they are true). I have a list of 42 open problems in some lecture notes [Ban16] I wrote almost a decade ago, some have been solved in the meantime, but many remain open. I write below a few for which you have all the necessary background to understand: Let me know if you solve any of them; regardless of whether its days or decades from now, I will be very happy to hear about your solution!

Conjecture 74 (Hadamard conjecture). For any n multiple of 4 there exists an Hadamard matrix H that is $n \times n$.

An Hadamard matrix $H \in \mathbb{R}^{n \times n}$ is a matrix with only entries 1 or -1 that is a multiple of an orthogonal matrix. In other words $H_{i,j} = \pm 1$ for all i, j and $H^T H = nI$. Yet in other words: the columns of H are an orthogonal basis for \mathbb{R}^n formed with only vectors with entries ± 1 .

Open Problem 75 (Mutually Unbiased Bases). See Open Problem 6.2 in [Ban16].

Conjecture 76 (Zauner’s Conjecture). See Open Problem 6.3 in [Ban16].

Conjecture 77 (Komlos Conjecture). See Open Problem 0.1 in [Ban16].

⁴⁰A particularly striking example of PCA is the “Genes mirror geography within Europe” experiment by J. Novembre et al, available at <https://www.nature.com/articles/nature07331> and <https://faculty.eeb.ucla.edu/Novembre/Novembreetal2008Nature.pdf>. Also discussed on National Geographic <https://www.nationalgeographic.com/science/article/european-genes-mirror-european-geography>.

Conjecture 78 (Matrix Spencer Conjecture). See Open Problem 4.3 in [Ban16].

Open Problem 79 (Rank of the Matrix Multiplication Tensor). What is the rank of the Matrix Multiplication Tensor corresponding to multiplication of 3×3 matrices.

A $d_1 \times d_2 \times d_3$ tensor T is what we can think of as a cubic matrix. It has $d_1 d_2 d_3$ entries given by T_{ijl} . We say T has rank r if r is the smallest integer such that we can write

$$T = \sum_{k=1}^r a_k \otimes b_k \otimes c_k,$$

for $a_k \in \mathbb{R}^{d_1}, b_k \in \mathbb{R}^{d_2}, c_k \in \mathbb{R}^{d_3}$, for $k = 1, \dots, r$. In other words

$$T_{ijl} = \sum_{k=1}^r (a_k)_i (b_k)_j (c_k)_l.$$

Recall Proposition 7.1.4 to see why for matrices this corresponds to the notion of rank we have been using. While computing the rank of a matrix is computationally easy, doing so for tensors is notoriously difficult (because they lack a spectral theory of eigenvalues and eigenvectors).

There is a way to think of Strassen's algorithm (that you saw in a CS Lens in Part I of the course) in terms of a decomposition of a certain Tensor in terms of rank-1 tensors. In this description we focus on $n \times n$ matrices, but the same thing can be done for rectangular matrices.

The $n \times n$ matrix multiplication tensor is a $n^2 \times n^2 \times n^2$ tensor, where each dimension is indexed by pairs $(i_1, i_2), (j_1, j_2), (l_1, l_2)$ and T is given by

$$T_{(i_1, i_2), (j_1, j_2), (l_1, l_2)} = \begin{cases} 1 & \text{if } i_1 = j_1, j_2 = l_1, l_2 = i_2 \\ 0 & \text{o.w.} \end{cases}.$$

Strassen's algorithm can be viewed as the fact that the rank of the 2×2 matrix multiplication tensor (a $4 \times 4 \times 4$ tensor) is ≤ 7 . The rank of the 3×3 matrix multiplication tensor (a $9 \times 9 \times 9$ tensor) remains unknown.^{41 42}

⁴¹To the best of my knowledge, the current "world records" for lower and upper bounds are 19 and 23, see for example <https://mathoverflow.net/questions/249256/best-known-bounds-on-border-ranks-of-small-matrix-multiplication-tensors?noredirect=1&lq=1> or <https://mathoverflow.net/questions/151058/best-known-bounds-on-tensor-rank-of-matrix-multiplication-of-3x3-matrices>.

⁴²See <https://www.youtube.com/watch?v=fDAPJ7rvUw> for a very nice description of how AI methods found a better low rank decomposition for the matrix multiplication tensor for some dimensions, and <https://www.nature.com/articles/s41586-022-05172-4> for the paper the video discusses (make sure to take a look at the table in Figure 3 in that article).

Conjecture 80 (The Paley ETF Conjecture). See Open Problem 6.4 in [Ban16] (see also Open Problem 5.1. in the same reference).⁴³

Acknowledgements. Many thanks to all the Linear Algebra teaching team! A particularly shout out to many of you that caught countless typos on these notes and gave countless suggestions that improve them! Thanks (a partial list): Bernd Gärtner, Sebastian Haslebacher, Felix Breuer, Till Schnabel, Mattia Taiana, Tanguiy Magne, Aviv Segall, Matthieu Croci, Sergey Prokudin, Mia Filic, Sofia Giampietro, Seyedmorteza Sadat, Mark Sosman. Felix and Till in particular are responsible for countless improvements in the notes, thanks! Thanks also to Dustin Mixon (OSU) for inspiring discussions on how to teach linear algebra, and for showing us `eigenquiz`!

And, most importantly, a special thanks to the ETH D-INFK students of Linear Algebra in Fall 2023 that made teaching this course an absolute pleasure! I hope you also enjoyed it!

APPENDIX A. SOME IMPORTANT PRELIMINARIES AND REMARKS ON NOTATION

To follow these notes the reader needs to be familiar with basics of vector and matrix operations and manipulations; understand what is dimension of a subspace, and in particular that is well-defined (that every basis of a subspace has the same size); and understand what is the rank of a matrix (and in particular that the dimension of the column space and the row space are the same). Even though Gaussian Elimination is not a core ingredient of this part of the course, we still assume that the reader is familiar with it. The students of 401-0131-00L are familiar with all this via Part I of this course.

Some further important preliminaries and/or remarks:

- (1) The dot product $x \cdot y$ between two real valued vectors is sometimes also called inner product and written as $\langle x, y \rangle$ (it is equal to $x^\top y$). For \mathbb{C}^n the inner product is given by $\langle x, y \rangle = y^* x$.
- (2) Matrix Factorization for A an $m \times n$ matrix with rank r :
 $A = CR$,
 C is $m \times r$ with linearly independent columns (they are the first r linearly independent columns of A). R is $r \times n$, it is upper triangular (i.e. $R_{ij} = 0$ if $i > j$), and it has an $r \times r$ identity as a submatrix, corresponding to the locations of the first r linearly independent columns of A .
- (3) For V a subspace (or a vector space) with dimension n the following holds:
 - Any basis of V has size n .

⁴³This one needs some (light) Number Theory background to understand. I posed this one (with collaborators) and spent many hours trying to make progress on it...

- Any spanning set of V has size $\geq n$.
- Any spanning set of V with size n is also a basis.
- Any set of linearly independent vectors in V has size $\leq n$.
- Any set of linearly independent vectors in V with size n is also a basis.

APPENDIX B. WEEKLY SCHEDULE

Numbers represent Fr-Wed weeks of 4×45 min lectures (except first and last).

CS lenses are generally not in the script (nor do they appear in this schedule).

- (7) 8.11.2023: [Introductions](#); [4.2. Projections](#)
- (8) [4.3. Least Squares & Fitting a line](#), [4.4 Orthonormal bases](#)
- (9) [QR decomposition](#), [4.5. and Intro to Linear Transformations](#)
- (10) [Determinants](#). [Finish 5](#). [Gentle intro to 6.1](#)
- (11) [Complex Numbers](#). [Start of Eigenvalues/Eigenvectors \(6\)](#).
- (12) [Continue with 6](#), [matrix diagonalization](#), [Change of basis for LT](#), [start Spectral Thm](#).
- (13) [Finish Spectral Theorem](#). [The Singular Value Decomposition](#). [Matrix norms \(brief intro\)](#).
- (14) 22.12.2023: [Entire lecture of “CS lenses”](#), a sort of technical “Ask Me Anything” session that won’t cover core material of the course. If you have any particular topic you would like to me cover, suggest it on the forum!

APPENDIX C. CS LENS LECTURES (NOT PART OF CHAPTER 7)

- [Kernel Methods](https://ti.inf.ethz.ch/ew/courses/LA23/slides/CS_Lens_kernels.pdf): https://ti.inf.ethz.ch/ew/courses/LA23/slides/CS_Lens_kernels.pdf
- [Graphs, Networks, and Linear Algebra](https://ti.inf.ethz.ch/ew/courses/LA23/slides/CS_Lens_Graphs1.pdf) https://ti.inf.ethz.ch/ew/courses/LA23/slides/CS_Lens_Graphs1.pdf. I also cover this in some classes I have taught, you can see manuscripts here: [BSS, BM23].
- [Google’s PageRank algorithm](#). I also cover this in a class I have taught, see [BSS].
- See Chapter 7 for more.

APPENDIX D. A “SIMPLE PROOF” OF FUNDAMENTAL THEOREM OF ALGEBRA

In this appendix I present a brief sketch of a (relatively) simple proof of the Fundamental Theorem of Algebra I learned from Alessio Figalli. It uses some analysis that you haven’t seen yet, but I will make a comment when it appears.

Let $P(z)$ be a polynomial of degree n . Without loss of generality we can assume it is monic $P(z) = z^n + \alpha_{n-1}z^{n-1} + \dots + \alpha_0$. Suppose $P(z)$ has no zeros/roots. There is a $r \in \mathbb{R}$ large enough such that the infimum of $|P(z)|$ inside the close disc D_r of radius r centered at zero is smaller than that outside the disc D_r (because far from the origin the term z^n dominates and forces $|P(z)|$ to be large outside of D_r). Since D_r is a compact set and $|P(z)|$ is continuous it needs to attain its minimum⁴⁴ at a point $z_0 \in D_r$. Note that $P(z_0) \neq 0$. Write $Q(z) = P(z - z_0)$, it is also a polynomial of degree n , $Q(z) = \beta_0 + \beta_1z + \beta_2z^2 + \dots + z^n$. Notice that $\beta_0 = P(z_0)$. let k be the first coefficient of $Q(z)$ (not including β_0) that is nonzero (meaning that $\beta_k \neq 0$ but $\beta_i = 0$ for all $0 < i < k$). Then $Q(z) = \beta_0 + \beta_kz^k + \beta_{k+1}z^{k+1} + \dots$. Take $\varepsilon > 0$ arbitrarily small and consider $Q\left(\varepsilon\left(-\frac{\beta_0}{\beta_k}\right)^{\frac{1}{k}}\right)$. It is not difficult to see that for ε small enough the higher order terms are negligible and the term $\beta_0 + \beta_kz^k$ has smaller modulus and so one can pick ε such that $\left|Q\left(\varepsilon\left(-\frac{\beta_0}{\beta_k}\right)^{\frac{1}{k}}\right)\right| < |Q(0)|$ which is a contradiction with the fact that $|P(z_0)|$ was minimum.

References⁴⁵

- [Ban16] Afonso S. Bandeira. Ten lectures and forty-two open problems in the mathematics of data science. Available online at: <https://people.math.ethz.ch/~abandeira/TenLecturesFortyTwoProblems.pdf>. See also <https://ocw.mit.edu/courses/18-s096-topics-in-mathematics-of-data-science-fall-2015/>, 2016.
- [BM23] Afonso S. Bandeira and Antoine Maillard. Mathematics of signals, networks, and learning. Available online at: https://anmaillard.github.io/teaching/msnl_spring_2023.pdf. Videos from an earlier version of the course available at <https://youtube.com/playlist?list=PLiud-28tsatL0MbfJFQQS7MYkrFrujCYp>, 2023.
- [BSS] A. S. Bandeira, A. Singer, and T. Strohmer. Mathematics of data science. Book draft available at <https://people.math.ethz.ch/~abandeira/BandeiraSingerStrohmer-MDS-draft.pdf>. Videos available at: <https://www.youtube.com/playlist?list=PLiud-28tsatIKUitdoH3EEUZL-9i516IL>.
- [Str23] Gilbert Strang. *Introduction to Linear Algebra*. (Table of contents available at <https://math.mit.edu/~gs/linearalgebra/ila6/indexila6.html>). Wellesley - Cambridge Press., sixth edition, 2023.

⁴⁴This is the part that needs more analysis/topology background: the fact that continuous functions on a compact (think closed and bounded) set needs to attain a minimum. You will learn about this in Analysis. Perhaps not surprisingly, the “other proof” of Corollary 6.3.8 that does not involve complex numbers, also needs this fact.

⁴⁵In some PDF viewers the \sim in the urls above does not show as the correct character, if the link appears broken delete the \sim and write a new \sim on the url.