

Linear Algebra Reformed for 21st C Application

A. J. Roberts
University of Adelaide
South Australia *

July 26, 2016

Please inform me of typographical glitches, writing infelicities, and errors.

* <http://www.maths.adelaide.edu.au/anthony.roberts>

Contents

1 Vectors	10
1.1 Vectors have magnitude and direction	12
1.2 Adding and stretching vectors	21
1.3 The dot product determines angles and lengths	34
1.4 The cross product	57
1.5 Use Matlab/Octave for vector computation	72
2 Systems of linear equations	83
2.1 Introduction to systems of linear equations	86
2.2 Directly solve linear systems	95
2.3 Linear combinations span sets	122
3 Matrices encode system interactions	132
3.1 Matrix operations and algebra	134
3.2 The inverse of a matrix	164
3.3 Factorise to the singular value decomposition	205
3.4 Subspaces, basis and dimension	240
3.5 Project to solve inconsistent equations	276
3.6 Introducing linear transformations	336
4 Eigenvalues and eigenvectors of symmetric matrices	375
4.1 Introduction to eigenvalues and eigenvectors	377
4.2 Beautiful properties for symmetric matrices	403
5 Approximate matrices	437
5.1 Measure changes to matrices	438
5.2 Regularise linear equations	482

6 Determinants distinguish matrices	506
6.1 Geometry underlies determinants	507
6.2 Laplace expansion theorem for determinants	523
7 Eigenvalues and eigenvectors in general	551
7.1 Find eigenvalues and eigenvectors of matrices	559
7.2 Linear independent vectors may form a basis	605
7.3 Diagonalisation identifies the transformation	637

V0-1P

Preface

Traditional courses in linear algebra make considerable use of the reduced row echelon form (RREF), but the RREF is an unreliable tool for computation in the face of inexact data and arithmetic. The [Singular Value Decomposition] SVD can be regarded as a modern, computationally powerful replacement for the RREF.¹

Cleve Moler, MathWorks (2006)

The Singular Value Decomposition (SVD) is sometimes called the *jewel in the crown* of linear algebra. Traditionally the SVD is introduced and explored at the end of several linear algebra courses. Question: Why were students required to wait until the end of the course, if at all, to be introduced to beauty and power of this jewel? Answer: limitations of hand calculation.

This book establishes a new route through linear algebra, one that reaches the SVD jewel in linear algebra's crown very early, in Section 3.3. Thereafter its beautiful power both explores many modern applications and also develops traditional linear algebra concepts, theory, and methods. No rigour is lost in this new route: indeed, this book demonstrates that most theory is better proved with an SVD rather than with the traditional RREF. This new route through linear algebra becomes available by the ready availability of ubiquitous computing in the 21st century.

As so many other disciplines use the SVD, it is not only important that mathematicians understand what it is, but also teach it thoroughly in linear algebra and matrix analysis courses. *(Turner et al. 2015, p.30)*

Aims for students

How should mathematical sciences departments reshape their curricula to suit the needs of a well-educated workforce in the twenty-first century?

... The mathematical sciences themselves are changing as the needs of big data and the challenges of modeling

¹ <http://au.mathworks.com/company/newsletters/articles/professor-svd.html> [9 Jan 2015]

complex systems reveal the limits of traditional curricula.
Bressoud et al. (2014)

Linear algebra is packed with compelling results for application in science, engineering and computing, and with answers for the twenty-first century needs of big data and complex systems. This book provides the conceptual understanding of the essential linear algebra of vectors and matrices for modern engineering and science. The traditional linear algebra course has been reshaped herein to meet modern demands.

Crucial is to inculcate the terms and corresponding relationships that you will most often encounter later in professional life, often when using professional software. For example, the manual for the engineering software package Fluent most often invokes the linear algebra terms of diagonal, dot product, eigenvalue, least square, orthogonal, projection, principal axes, symmetric, unit vector. Engineers need to know these terms. What such useful terms mean, their relationships, and use in applications are central to the mathematical development in this book: you will see them introduced early and used often.

For those who proceed on to do higher mathematics, the development also provides a great solid foundation of key concepts, relationships and transformations necessary for higher mathematics.

Important for all is to develop facility in manipulating, interpreting and transforming between visualisations, algebraic forms, and vector-matrix representations—of problems, working, solutions and interpretation. In particular, one overarching aim of the book is to encourage your formulation, thinking and operation at the crucial system-wide level of matrix/vector operations.

In view of ubiquitous computing, throughout this book explicitly integrates computer support for developing concepts and their relations. The central computational tools to understand are the operation `A\` for solving straightforward linear equations; the function `svd()` for difficult linear equation; the function `svd()` for approximation; and function `eig()` for probing structures. This provides a framework to understand key computational tools to effectively utilise the so-called third arm of science: namely, computation.

Throughout the book examples, many graphical, introduce and illustrate the concepts and relationships between the concepts. Working through these will help form the mathematical relationships essential for application. Also included are various applications, described to varying levels of details, to empower you to see how the mathematics will empower you to answer many practical challenges

in engineering and science.

The main contribution of mathematics to the natural sciences is not in formal computations . . . , but in the investigation of those non-formal questions where the exact setting of the question (what are we searching for and what specific models must be used) usually constitute half the matter.

. . . Examples teach no less than rules, and errors, more than correct but abstruse proofs. Looking at the pictures in this book, the reader will understand more than learning by rote dozens of axioms

(*Arnold 2014, p.xiii*)

Background for teachers

Depending upon the background of your students, your class should pick up the story somewhere in the first two chapters or so. Some students will have previously learnt some vector material in just 2D and 3D, in which case refresh the concepts in n D as presented in the first chapter.

As a teacher you can use this book in several ways.

- One way is as a reasonably rigorous mathematical development of concepts and interconnections by invoking its definitions, theorems, and proofs, all interleaved with examples.
- Another way is as the development of practical techniques and insight for application orientated science and engineering students via the motivating examples to appropriate definitions, theorems and applications.
- Or any mix of these two.

The concept of linear independence does not appear until relatively late, namely in Chapter 7. This is good for several reasons. First, orthogonality is much more commonly invoked in science and engineering than is linear independence. Second, it is well documented that students struggle with linear independence:

there is ample talk in the math ed literature of classes hitting a ‘brick wall’, when linear (in)dependence is studied in the middle of such a course *Uhlig (2002)*

Consequently, here we learn the more specific orthogonality before the more abstract linear independence. Many modern applications are made available by the relatively early introduction of orthogonality.

Indeed one of the aims of this book is to organise the development so that if a student only studies part of the material, then s/he still obtains a powerful and useful body of knowledge for science or engineering.

The book typically introduces concepts in low dimensional cases, and subsequently develops the general theory of the concept. This is to help focus the learning, empowered by visualisation, while also making a preformal connection to be strengthened subsequently. People are not one dimensional; knowledge is not linear. Cross-references make many connections, explicitly recalling earlier learning (although sometimes forward to material not yet ‘covered’).

information that is recalled grows stronger with each retrieval . . . spaced practice is preferable to massed practice. *Halpern & Hakel (2003) [p.38]*

On visualisation

All Linear Algebra courses should stress visualization and geometric interpretation of theoretical ideas in 2- and 3-dimensional spaces. Doing so highlights “algebraic and geometric” as “contrasting but complementary points of view,” *(Schumacher et al. 2015, p.38)*

Throughout, this book also integrates visualisation. This visualisation reflects the fundamentally geometric nature of linear algebra. It also empowers learners to utilise different parts of their brain and integrate the knowledge together from the different perspectives. Visualisation also facilitates greater skills at interpretation and modelling so essential in applications. Lastly, a visual exercise question develops learner’s understanding without them being able to defer the challenge to online tools, as yet.

Visual representations are effective because they tap into the capabilities of the powerful and highly parallel human visual system. We like receiving information in visual form and can process it very efficiently: around a quarter of our brains are devoted to vision, more than all our other senses combined . . . researchers (especially those from mathematic backgrounds) see visual notations as being informal, and that serious analysis can only take place at the level of their semantics. However, this is a misconception: visual languages are no less formal than textual ones *Moody (2009)*

On integrated computation

Cowen argued that because “no serious application of linear algebra happens without a computer,” computation should be part of every beginning Linear Algebra course. . . . While the increasing applicability of linear algebra does not require that we stop teaching theory, Cowen argues that “it should encourage us to see the role of the theory in the subject as it is applied.”

(Schumacher et al. 2015, p.38)

We need to empower students to use computers to improve their understanding, learning and application of mathematics; not only integrated in their study but also in their later professional career.

One often expects it should be easy to sprinkle a few computational tips and tools throughout a mathematics course. This is not so—extra computing is difficult. There are two reasons for the difficulty: first, the number of computer language details that have to be learned is surprisingly large; second, for students it is a genuine intellectual overhead to learn and relate both the mathematics and the computations.

Consequently, this book chooses a computing language where it is as simple as reasonably possible to perform linear algebra operations: Matlab/Octave appears to answer this criteria.² Further, we are as ruthless as possible in invoking herein the smallest feasible set of commands and functions from Matlab/Octave so that students have the minimum to learn. Most teachers will find many of their favourite commands are missing—this omission is all to the good in focussing upon useful mathematical development aided by only essential integrated computation.

This book does not aim to teach computer programming: there is no flow control, no looping, no recursion, nor function definitions. The aim herein is to use short sequences of declarative assignment statements, coupled with the power of vector and matrix data structures, to learn core mathematical concepts, applications and their relationships in linear algebra.

The internet is now ubiquitous and pervasive. So too is computing power: students can execute Matlab/Octave not only on laptops, but also on tablets and smart phones, perhaps using university or public servers, octave-online.net, Matlab-Online or Matlab-Mobile. We no longer need separate computer laboratories. Instead, expect students to access computational support simply by reaching

² To compare popular packages, just look at the length of expressions students have to type in order to achieve core computations: Matlab/Octave is almost always the shortest ([Nakos & Joyner 1998](#), e.g.). (Of course be wary of this metric: e.g., APL would surely be too concise!)

into their pockets or bags.

long after Riemann had passed away, historians discovered that he had developed advanced techniques for calculating the Riemann zeta function and that his formulation of the Riemann hypothesis—often depicted as a triumph of pure thought—was actually based on painstaking numerical work.

Donoho & Stodden (2015)

Acknowledgements

I acknowledge with thanks the work of many others who inspired much design and details here, including the stimulating innovations of calculus reform ([Hughes-Hallett et al. 2013](#), e.g.), the comprehensive efforts behind recent reviews of undergraduate mathematics teaching ([Alpers et al. 2013](#), [Bressoud et al. 2014](#), [Turner et al. 2015](#), [Schumacher et al. 2015](#), [Bliss et al. 2016](#), e.g.), and the books of [Anton & Rorres \(1991\)](#), [Davis & Uhl \(1999\)](#), [Holt \(2013\)](#), [Larson \(2013\)](#), [Lay \(2012\)](#), [Nakos & Joyner \(1998\)](#), [Poole \(2015\)](#), [Will \(2004\)](#). I also thank the entire \LaTeX team, especially Knuth, Lamport, Feuersänger, and the AMS.

1 Vectors

Chapter Contents

1.1	Vectors have magnitude and direction	12
1.1.1	Exercises	18
1.2	Adding and stretching vectors	21
1.2.1	Basic operations	21
1.2.2	Parametric equation of a line	25
1.2.3	Manipulation requires algebraic properties . .	28
1.2.4	Exercises	32
1.3	The dot product determines angles and lengths . .	34
1.3.1	Work done involves the dot product	40
1.3.2	Algebraic properties of the dot product	41
1.3.3	Orthogonal vectors are at right-angles	46
1.3.4	Normal vectors and equations of a plane	48
1.3.5	Exercises	53
1.4	The cross product	57
1.4.1	Exercises	67
1.5	Use Matlab/Octave for vector computation	72
1.5.1	Exercises	79

This chapter is a relatively concise introduction to vectors, their properties, and a little computation with Matlab/Octave. Skim or study as needed.

Mathematics started with counting. The natural numbers $1, 2, 3, \dots$ quantify how many objects have been counted. Eventually, via many existential arguments over centuries about whether meaningful, negative numbers and the zero were included to form the **integers** $\dots, -2, -1, 0, 1, 2, \dots$. In the mean time people needed to quantify fractions such as two and half a bags, or a third of a cup which led to the rational numbers such as $2\frac{1}{2} = \frac{5}{2}$ or $\frac{1}{3}$, and now defined as all numbers writeable in the form $\frac{p}{q}$ for integers p and q (q non-zero). Roughly two thousand years ago, Pythagoras was forced to recognise that for many triangles the length of a side could not be rational, and hence there must be more numbers in the world about us than rationals could provide. To cope with non-rational numbers such as $\sqrt{2} = 1.41421\dots$ and $\pi = 3.14159\dots$, mathematicians

define the **real numbers** to be all numbers which in principle can be written as a decimal expansion such as

$$\frac{9}{7} = 1.285714285714\cdots \quad \text{or} \quad e = 2.718281828459\cdots.$$

Such decimal expansions may terminate or repeat or may need to continue on indefinitely (as denoted by the three dots, called an ellipsis). The frequently invoked symbol \mathbb{R} denotes the set of all possible real numbers.

In the sixteenth century Gerolamo Cardano¹ developed a procedure to solve cubic polynomial equations. But the procedure involved manipulating $\sqrt{-1}$ which seemed a crazy figment of imagination. Nonetheless the procedure worked. Subsequently, many practical uses were found for $\sqrt{-1}$, now denoted by i . Consequently, many areas of modern science and engineering use **complex numbers** which are those of the form $a + bi$ for real numbers a and b . The symbol \mathbb{C} denotes the set of all possible complex numbers. This book mostly uses integers and real numbers, but eventually we need the marvellous complex numbers.

In some places this book uses the term **scalar** to denote a number that could be integer, real or complex. The term ‘scalar’ arises because such numbers are often used to scale the length of a ‘vector’.

¹ Considered one of the great mathematicians of the Renaissance, Cardano was one of the key figures in the foundation of probability and the earliest introducer of the binomial coefficients and the binomial theorem in the western world. He wrote more than 200 works on medicine, mathematics, physics, chemistry, biology, astronomy, philosophy, religion, and music. . . . he is well-known for his achievements in algebra. He made the first systematic use of negative numbers, published with attribution the solutions of other mathematicians for the cubic and quartic equations, and acknowledged the existence of imaginary numbers. (Wikipedia, 2015)

1.1 Vectors have magnitude and direction

Section Contents

1.1.1 Exercises	18
---------------------------	----

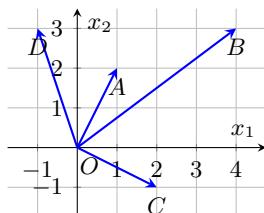
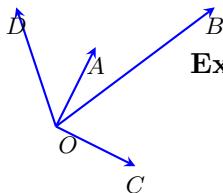
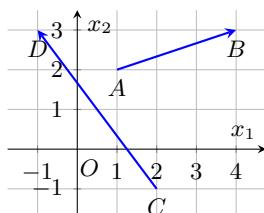
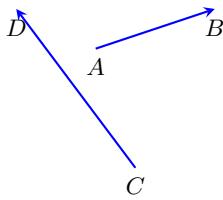
There are more things in heaven and earth, Horatio,
than are dreamt of in your philosophy.

(*Hamlet I.5:159–167*)

In the eighteenth century, astronomers needed to describe both the position and velocity of the planets. Such a description required quantities which have both a magnitude and a direction. Step outside, a wind blowing at 8 m/s from the south-west also has both a magnitude and direction. Quantities that have both the properties of a magnitude and a direction are called **vectors** (from the Latin for *carrier*).

Example 1.1.1 (displacement vector). An important class of vectors are the so-called **displacement vectors**. Given two points in space, say A and B , the displacement vector \overrightarrow{AB} is the directed line segment from the point A to the point B —as illustrated by the two displacement vectors \overrightarrow{AB} and \overrightarrow{CD} in the margin. For example, if your home is at position A and your school at position B , then travelling from home to school is moving by the amount of the displacement vector \overrightarrow{AB} .

To be able to manipulate vectors we describe them with numbers from a coordinate system. So choose an origin for the coordinate system, usually denoted O , and draw coordinate axes in the plane (or space), as illustrated for the above two displacement vectors. Here the displacement vector \overrightarrow{AB} goes three units to the right and one unit up, so we denote it by the ordered pair of numbers $\overrightarrow{AB} = (3, 1)$. Whereas the displacement vector \overrightarrow{CD} goes three units to the left and four units up, so we denote it by the ordered pair of numbers $\overrightarrow{CD} = (-3, 4)$. ■



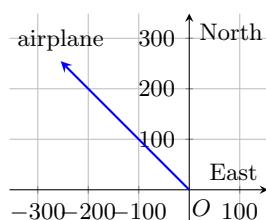
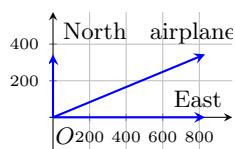
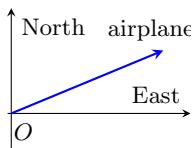
Example 1.1.2 (position vector). The next important class of vectors are the **position vectors**. Given some chosen fixed origin in space, then \overrightarrow{OA} is the position vector of the point A . The marginal picture illustrates the position vectors of four points in the plane, given a chosen origin O .

Again, to be able to manipulate such vectors we describe them with numbers from a coordinate system. So draw coordinate axes in the plane (or space), as illustrated for the above four position vectors. Here the position vector \overrightarrow{OA} goes one unit to the right and two units up so we denote it by $\overrightarrow{OA} = (1, 2)$. Similarly, the

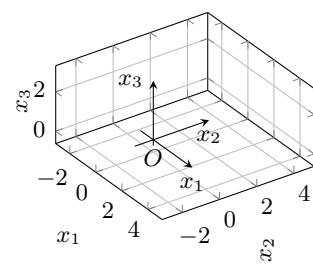
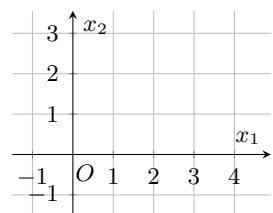
position vectors $\overrightarrow{OB} = (4, 3)$, $\overrightarrow{OC} = (2, -1)$, and $\overrightarrow{OB} = (-1, 3)$. Of course you recognise that the ordered pairs of numbers in the position vectors are exactly the coordinates of each of the specified end-points. ■

Example 1.1.3 (velocity vector).

Consider an airplane in level flight at 900 km/hr to the east-north-east. Choosing coordinate axes oriented to the East and the North, the direction of the airplane is at an angle 22.5° from the East, as illustrated in the margin. Trigonometry then tells us that the Eastward part of the speed of the airplane is $900 \cos(22.5^\circ) = 831.5$ km/hr, whereas the Northward part of the speed is $900 \sin(22.5^\circ) = 344.4$ km/hr (as indicated in the margin). Further, the airplane is in level flight, not going up or down, so in the third direction of space (vertically) its speed component is zero. Putting these together forms the velocity vector $(831.5, 344.4, 0)$ in km/hr in space.



Another airplane takes off from an airport at 360 km/hr to the northwest and climbs at 2 m/s. The direction northwest is 45° to the East-West lines and 45° to the North-South lines. Trigonometry then tells us that the Westward speed of the airplane is $360 \cos(45^\circ) = 360 \cos(\frac{\pi}{4}) = 254.6$ km/hr, whereas the Northward speed is $360 \sin(45^\circ) = 360 \sin(\frac{\pi}{4}) = 254.6$ km/hr as illustrated in the margin. But West is the opposite direction to East, so if we choose to write East as positive, then West must be negative. Consequently, together with the climb in the vertical, the velocity vector is $(-254.6 \text{ km/hr}, 254.6 \text{ km/hr}, 2 \text{ m/s})$. But it is best to avoid mixing units within a vector, so here convert all speeds to m/s: here 360 km/hr upon dividing by 3600 secs/hr and multiplying by 1000 m/km gives $360 \text{ km/hr} = 100 \text{ m/s}$. Then the North and West speeds are $100 \cos(\frac{\pi}{4}) = 70.71 \text{ m/s}$. Consequently, the velocity vector of the climbing airplane should be described as $(-70.71, 70.71, 2)$ in m/s. ■



In applications, as these examples illustrate, the ‘physical’ vector exists before the coordinate system. It is only when we choose a specific coordinate system that a ‘physical’ vector gets expressed by numbers. Throughout, unless otherwise specified, this book assumes that vectors are expressed in what is called a **standard coordinate system**.

- In the two dimensions of the plane the standard coordinate system has two coordinate axes, one horizontal and one vertical at right-angles to each other, often labelled x_1 and x_2 respectively (as illustrated in the margin), although labels x and y are also common.
- In the three dimensions of space the standard coordinate system has three coordinate axes, two horizontal and one

vertical all at right-angles to each other, often labelled x_1 , x_2 and x_3 respectively (as illustrated in the margin), although labels x , y and z are also common.

- Correspondingly, in so-called ‘ n dimensions’ the standard coordinate system has n coordinate axes, all at right-angles to each other, and often labelled x_1, x_2, \dots, x_n , respectively.

Definition 1.1.4. *Given a standard coordinate system with n coordinate axes, all at right-angles to each other, a **vector** is an ordered n -tuple of real numbers x_1, x_2, \dots, x_n equivalently written either as a row in parentheses or as a column in brackets,*

$$(x_1, x_2, \dots, x_n) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

(they mean the same, it is just more convenient to usually use a row in parentheses in text, and a column in brackets in displayed mathematics). The real numbers x_1, x_2, \dots, x_n are called the **components** of the vector, and the number of components is termed its **size** (here n). The components are determined such that letting X be the point with coordinates (x_1, x_2, \dots, x_n) then the position vector \overrightarrow{OX} has the same magnitude and direction as the vector denoted (x_1, x_2, \dots, x_n) . Two vectors of the same size are **equal**, $=$, if all their corresponding components are equal (vectors with different sizes are never equal).

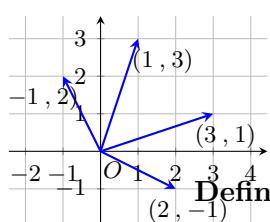
Robert Recorde invented the equal sign circa 1557 “bicause noe 2 thynges can be moare equalle”.

Examples 1.1.1 and 1.1.2 introduced some vectors and wrote them as a row in parentheses, such as $\overrightarrow{AB} = (3, 1)$. In this book exactly the same thing is meant by the columns in brackets: for example,

$$\overrightarrow{AB} = (3, 1) = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \quad \overrightarrow{CD} = (-3, 4) = \begin{bmatrix} -3 \\ 4 \end{bmatrix},$$

$$\overrightarrow{OC} = (2, -1) = \begin{bmatrix} 2 \\ -1 \end{bmatrix}, \quad (-70.71, 70.71, 2) = \begin{bmatrix} -70.71 \\ 70.71 \\ 2 \end{bmatrix}.$$

However, as defined subsequently, a row of numbers within brackets is to be quite different: $(3, 1) \neq [3 \ 1]$, and $(831, 344, 0) \neq [831 \ 344 \ 0]$.



The *ordering* of the components is very important. For example, as illustrated in the margin, the vector $(3, 1)$ is very different from the vector $(1, 3)$; similarly, the vector $(2, -1)$ is very different from the vector $(-1, 2)$.

Definition 1.1.5. *The set of all vectors with n components is denoted \mathbb{R}^n . The vector with all components zero, $(0, 0, \dots, 0)$, is called the **zero vector** and denoted by **0**.*

Example 1.1.6.

- All the vectors we can draw and imagine in the two dimensional plane form \mathbb{R}^2 . Sometimes we write that \mathbb{R}^2 is the plane because of this very close connection.
- All the vectors we can draw and imagine in three dimensional space form \mathbb{R}^3 . Again, sometimes we write that \mathbb{R}^3 is three dimensional space because of the close connection.
- The set \mathbb{R}^1 is the set of all vectors with one component, and that one component is measured along one axis. Hence \mathbb{R}^1 is effectively the same as the set of real numbers labelling that axis.

■

As just introduced for the zero vector $\mathbf{0}$, this book generally denotes vectors by a bold letter (except for displacement vectors). The other common notation you may see elsewhere is to denote vectors by a small over-arrow such as in the “zero vector $\vec{0}$ ”. Less commonly, some books and articles use an over- or under-tilde (\sim) to denote vectors. Be aware of this different notation in reading other books.

Question: why do we need vectors with n components, in \mathbb{R}^n , when the world around us is only three dimensional? Answer: because vectors can encode much more than spatial structure as in the next example.

Example 1.1.7 (linguistic vectors). Consider the following four sentences.

- (a) The dog sat on the mat.
- (b) The cat scratched the dog.
- (c) The cat and dog sat on the mat.
- (d) The dog scratched.

These four sentences involve up to three objects, cat, dog and mat, and two actions, sat and scratched. Some characteristic of the sentences is captured simply by counting the number of times each of these three objects and two actions appear in each sentence, and then forming a vector from the counts. Let's use vectors $\mathbf{w} = (N_{\text{cat}}, N_{\text{dog}}, N_{\text{mat}}, N_{\text{sat}}, N_{\text{scratched}})$ where the various N are the counts of each word (\mathbf{w} for words). The previous statement implicitly specifies that we use five coordinate axes, perhaps labelled “cat”, “dog”, “mat”, “sat” and “scratched”, and that distance along each axis represents the number of times the corresponding word is used. These word vectors are in \mathbb{R}^5 . Then

- (a) “The dog sat on the mat” is summarised by the vector $\mathbf{w} = (0, 1, 1, 1, 0)$.
- (b) “The cat scratched the dog” is summarised by the vector $\mathbf{w} = (1, 1, 0, 0, 1)$.

- (c) “The cat and dog sat on the mat” is summarised by the vector $\mathbf{w} = (1, 1, 1, 1, 0)$.
- (d) “The dog scratched” is summarised by the vector $\mathbf{w} = (0, 1, 0, 0, 1)$.
- (e) An empty sentence is the zero vector $\mathbf{w} = (0, 0, 0, 0, 0)$.
- (f) Together, the two sentences “The dog sat on the mat. The cat scratched the dog.” are summarised by the vector $\mathbf{w} = (1, 2, 1, 1, 1)$.

Using such crude summary representations of some text, even of entire documents, empowers us to use powerful mathematical techniques to relate documents together, compare and contrast, express similarities, look for type clusters, and so on. In application we would not just count words for objects (nouns) and actions (verbs), but also qualifications (adjectives and adverbs).²

People generally know and use thousands of words. Consequently, in practice, such word vectors typically have thousands of components corresponding to coordinate axes of thousands of distinct words. To cope with such vectors of many components, modern linear algebra has been developed to powerfully handle problems involving vectors with thousands, millions or even an ‘infinite number’ of components.

■

Definition 1.1.8 (Pythagoras). *Given any vector $\mathbf{v} = (v_1, v_2, \dots, v_n)$ in \mathbb{R}^n , define the **length** (or **magnitude**) of vector \mathbf{v} to be the real number (≥ 0)*

$$|\mathbf{v}| = \sqrt{v_1^2 + v_2^2 + \cdots + v_n^2}.$$

*A vector of length one is called a **unit vector**. (Be aware that many books, especially advanced books, denote the length of a vector with a pair of double lines, as in $\|\mathbf{v}\|$.)*

Example 1.1.9. Find the lengths of the following vectors.

(a) $\mathbf{a} = (-3, 4)$

(b) $\mathbf{b} = (3, 3)$

Solution: $|\mathbf{a}| =$

$$\sqrt{(-3)^2 + 4^2} = \sqrt{25} = 5.$$

Solution:

$$|\mathbf{b}| = \sqrt{3^2 + 3^2} = \sqrt{18} = 3\sqrt{2}.$$

(c) $\mathbf{c} = (1, -2, 3)$

(d) $\mathbf{d} = (1, -1, -1, 1)$

Solution: $|\mathbf{c}| =$

$$\sqrt{1^2 + (-2)^2 + 3^2} = \sqrt{14}.$$

Solution: $|\mathbf{d}| =$

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

■

² Look up Latent Semantic Indexing, such as at https://en.wikipedia.org/wiki/Latent_semantic_indexing [April 2015]

King – man + woman = queen

Computational linguistics has dramatically changed the way researchers study and understand language. The ability to number-crunch huge amounts of words for the first time has led to entirely new ways of thinking about words and their relationship to one another.

This number-crunching shows exactly how often a word appears close to other words, an important factor in how they are used. So the word Olympics might appear close to words like running, jumping, and throwing but less often next to words like electron or stegosaurus. This set of relationships can be thought of as a multidimensional vector that describes how the word Olympics is used within a language, which itself can be thought of as a vector space.

And therein lies this massive change. This new approach allows languages to be treated like vector spaces with precise mathematical properties. Now the study of language is becoming a problem of vector space mathematics. ^a *Technology Review, 2015*

^a <http://www.technologyreview.com/view/541356> [Oct 2015]

Example 1.1.10. Write down three different vectors, all three with the same number of components, that are (a) of length 5, (b) of length 3, and (c) of length -2 .

Solution: (a) Humans knew of the $3 : 4 : 5$ right-angled triangle thousands of years ago, so perhaps one answer could be $(3, 4)$, $(-4, 3)$ and $(5, 0)$.

- (b) One answer might be $(3, 0, 0)$, $(0, 3, 0)$ and $(0, 0, 3)$. A more interesting answer might arise from knowing $1^2 + 2^2 + 2^2 = 3^2$ leading to an answer of $(1, 2, 2)$, $(2, -1, 2)$ and $(-2, 2, 1)$.
- (c) Since the length of a vector is $\sqrt{\dots}$ which is always positive or zero, the length cannot be negative, so there is no possible answer to this last request.

■

Theorem 1.1.11. *The zero vector is the only vector of length zero: $|\mathbf{v}| = 0$ if and only if $\mathbf{v} = \mathbf{0}$.*

Proof. First establish the zero vector has length zero. From Definition 1.1.8, in \mathbb{R}^n ,

$$|\mathbf{0}| = \sqrt{0^2 + 0^2 + \cdots + 0^2} = \sqrt{0} = 0.$$

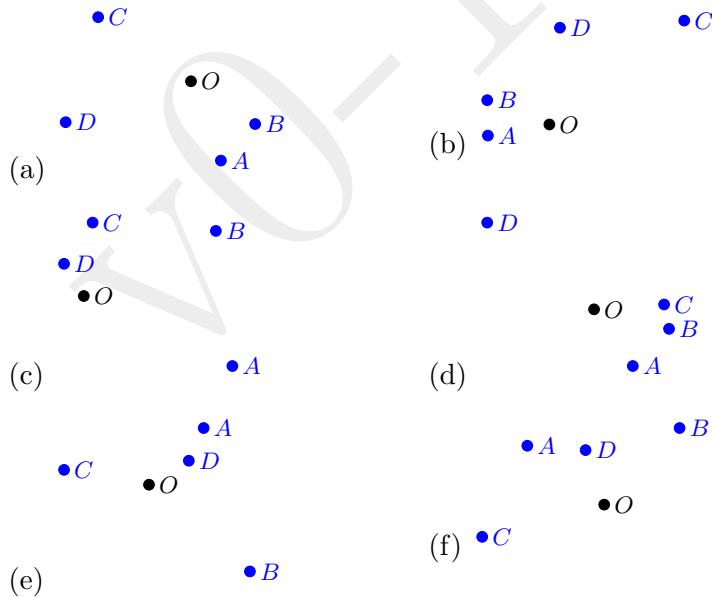
Second, if a vector has length zero then it must be the zero vector. Let vector $\mathbf{v} = (v_1, v_2, \dots, v_n)$ in \mathbb{R}^n have zero length. By squaring both sides of the Definition 1.1.8 for length we then know that

$$\underbrace{v_1^2}_{\geq 0} + \underbrace{v_2^2}_{\geq 0} + \cdots + \underbrace{v_n^2}_{\geq 0} = 0.$$

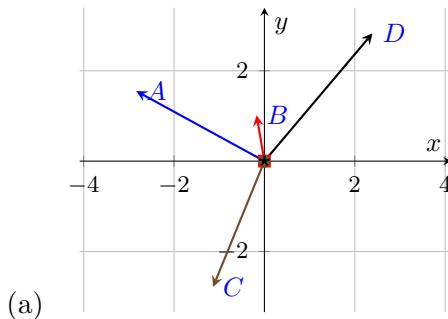
Being squares, all terms on the left are non-negative, so the only way they can all add to zero is if they are all zero. That is, $v_1 = v_2 = \cdots = v_n = 0$. Hence, the vector \mathbf{v} must be the zero vector $\mathbf{0}$. \square

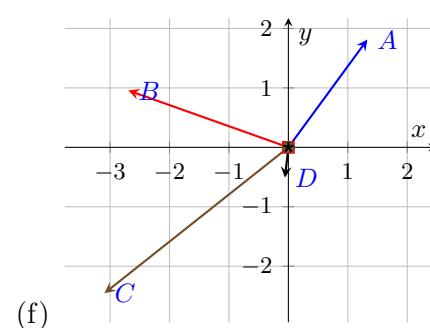
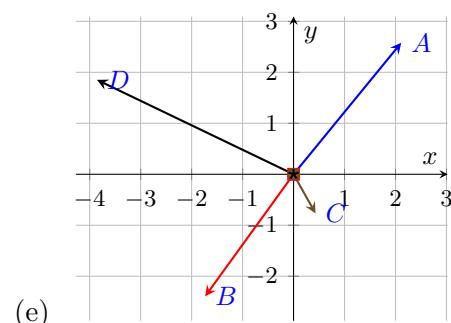
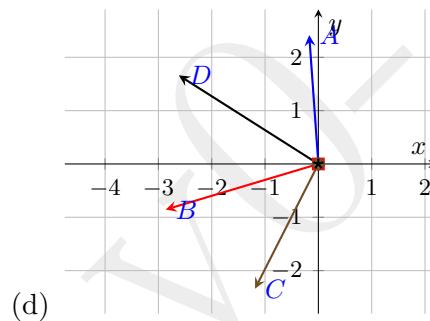
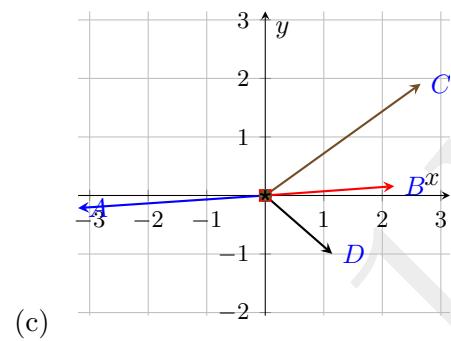
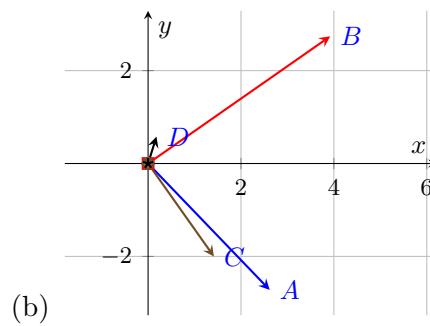
1.1.1 Exercises

Exercise 1.1.1. For each case: on the plot draw the displacement vectors \vec{AB} and \vec{CD} , and the position vectors of the points A and D .



Exercise 1.1.2. For each case: roughly estimate (to say ± 0.2) each of the two components of the four position vectors of the points A , B , C and D .





Exercise 1.1.3. For each case plotted in Exercise 1.1.2: from your estimated

components of each of the four position vectors, calculate the length (or magnitude) of the four vectors. Also use a ruler (or otherwise) to directly measure an estimate of the length of each vector. Confirm your calculated lengths reasonably approximate your measured lengths.

Exercise 1.1.4. Below are the titles of eight books that The Society of Industrial and Applied Mathematics (SIAM) reviewed recently.

- (a) Introduction to Finite and Spectral Element Methods using MATLAB
- (b) Derivative Securities and Difference Methods
- (c) Iterative Methods for Linear Systems: Theory and Applications
- (d) Singular Perturbations: Introduction to System Order Reduction Methods with Applications
- (e) Risk and Portfolio Analysis: Principles and Methods
- (f) Differential Equations: Theory, Technique, and Practice
- (g) Contract Theory in Continuous-Time Models
- (h) Stochastic Chemical Kinetics: Theory and Mostly Systems Biology Applications

Make a list of the five significant words that appear more than once in this list (not including the common nontechnical words such as “and” and “for”, and not distinguishing between words with a common root). Being consistent about the order of words, represent each of the eight titles by a word vector in \mathbb{R}^7 .

1.2 Adding and stretching vectors

Section Contents

1.2.1	Basic operations	21
1.2.2	Parametric equation of a line	25
1.2.3	Manipulation requires algebraic properties . .	28
1.2.4	Exercises	32

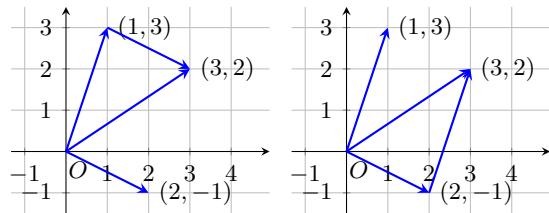
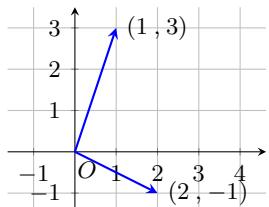
Useful operations on vectors are those which are physically meaningful. Then algebraic manipulations derives powerful results in applications. The two basic operations are addition and scalar multiplication.

These operations then make sense of statements such as “king – man + women = queen”.

1.2.1 Basic operations

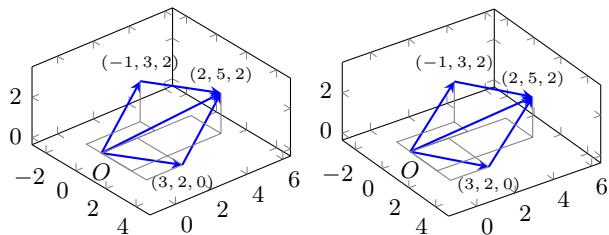
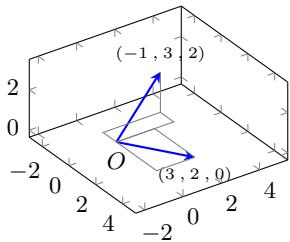
Example 1.2.1. Vectors of the same size are added component-wise. Equivalently, obtain the same result by geometrically joining the two vectors ‘head-to-tail’ and drawing the vector from the start to the finish.

- (a) $(1, 3) + (2, -1) = (1 + 2, 3 + (-1)) = (3, 2)$ as illustrated below where (given the two vectors plotted in the margin) the vector $(2, -1)$ is drawn from the end of $(1, 3)$, and the end point of the result determines the vector addition $(3, 2)$, as shown below-left.



This result $(3, 2)$ is the same if the vector $(1, 3)$ is drawn from the end of $(2, -1)$ as shown above-right. That is, $(2, -1) + (1, 3) = (1, 3) + (2, -1)$. That the order of addition is immaterial is the commutative law of vector addition that is established in general by Theorem 1.2.13a.

- (b) $(3, 2, 0) + (-1, 3, 2) = (3 + (-1), 2 + 3, 0 + 2) = (2, 5, 2)$ as illustrated below where (given the two vectors as plotted in the margin) the vector $(-1, 3, 2)$ is drawn from the end of $(3, 2, 0)$, and the end point of the result determines the vector addition $(2, 5, 2)$. As below, find the same result by drawing the vector $(3, 2, 0)$ from the end of $(-1, 3, 2)$.

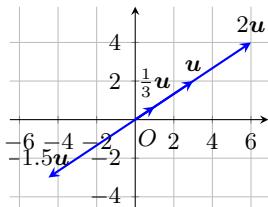


As drawn above, many of the three-D plots in this book are **stereo pairs**: drawing the plot from two slightly different viewpoints: relax your eye muscles to ‘look through’ the page, and then focus on the pair of plots to see the three-D effect. With practice viewing such three-D stereo pairs becomes less difficult.

- (c) The addition $(1, 3) + (3, 2, 0)$ is not defined and cannot be done as the two vectors have a different number of components, different sizes. ■

Example 1.2.2. To multiply a vector by a scalar, multiply each component by the scalar. Equivalently, obtain the result through stretching the vector by a factor of the scalar.

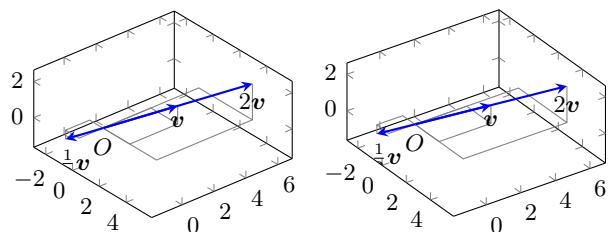
- (a) Let the vector $\mathbf{u} = (3, 2)$ then, as illustrated in the margin,



- (b) Let the vector $\mathbf{v} = (2, 3, 1)$ then, as illustrated below in stereo,

$$2\mathbf{v} = 2 \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \cdot 2 \\ 2 \cdot 3 \\ 2 \cdot 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 6 \\ 2 \end{bmatrix},$$

$$\left(-\frac{1}{2}\right)\mathbf{v} = -\frac{1}{2} \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} = \begin{bmatrix} -\frac{1}{2} \cdot 2 \\ -\frac{1}{2} \cdot 3 \\ -\frac{1}{2} \cdot 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -\frac{3}{2} \\ -\frac{1}{2} \end{bmatrix}.$$



Definition 1.2.3. Let two vectors in \mathbb{R}^n be $\mathbf{u} = (u_1, u_2, \dots, u_n)$ and $\mathbf{v} = (v_1, v_2, \dots, v_n)$, and let c be a scalar. Then the **sum** or **addition** of \mathbf{u} and \mathbf{v} , denoted $\mathbf{u} + \mathbf{v}$, is the vector obtained by joining \mathbf{v} to \mathbf{u} 'head-to-tail', and is computed as

$$\mathbf{u} + \mathbf{v} = (u_1 + v_1, u_2 + v_2, \dots, u_n + v_n).$$

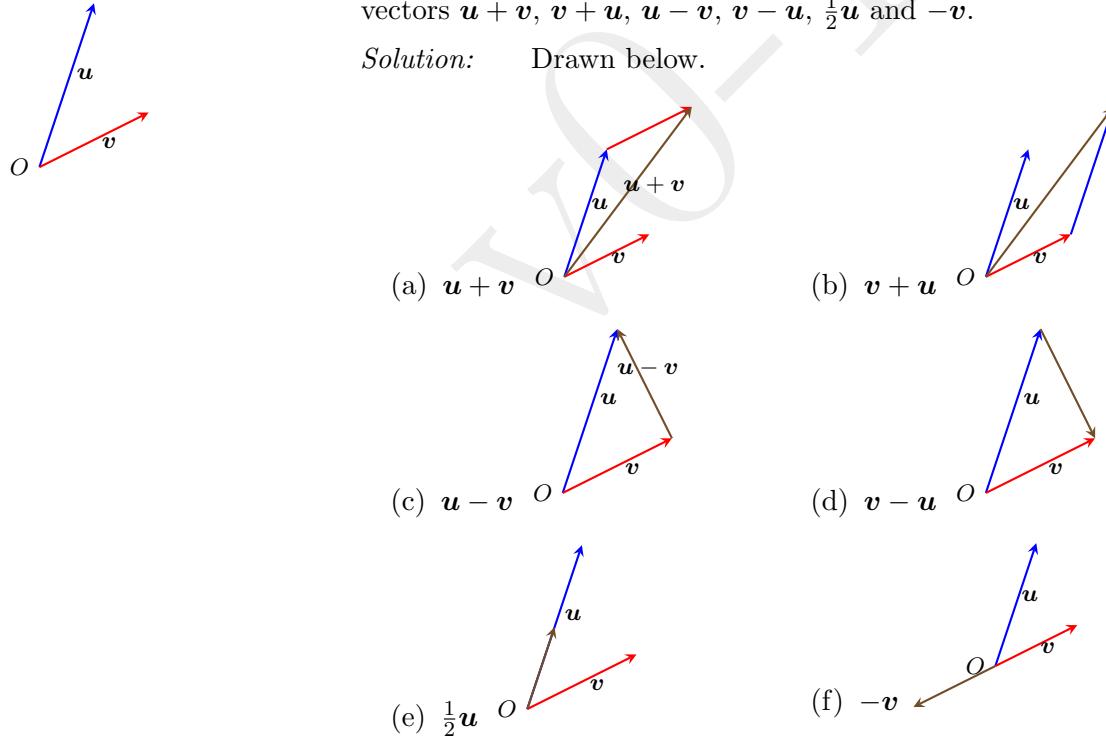
The **scalar multiplication** of \mathbf{u} by c , denoted $c\mathbf{u}$, is the vector of length $|c|\|\mathbf{u}\|$ in the direction of \mathbf{u} when $c > 0$ but in the opposite direction when $c < 0$, and is computed as

$$c\mathbf{u} = (cu_1, cu_2, \dots, cu_n).$$

The **negative** of \mathbf{u} denoted $-\mathbf{u}$, is defined as the scalar multiple $-\mathbf{u} = (-1)\mathbf{u}$, and is a vector of the same length as \mathbf{u} but in exactly the opposite direction. The **difference** $\mathbf{u} - \mathbf{v}$ is defined as $\mathbf{u} + (-\mathbf{v})$ and is equivalently the vector drawn from the end of \mathbf{v} to the end of \mathbf{u} .

Example 1.2.4. For the vectors \mathbf{u} and \mathbf{v} shown in the margin, draw the vectors $\mathbf{u} + \mathbf{v}$, $\mathbf{v} + \mathbf{u}$, $\mathbf{u} - \mathbf{v}$, $\mathbf{v} - \mathbf{u}$, $\frac{1}{2}\mathbf{u}$ and $-\mathbf{v}$.

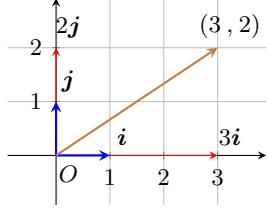
Solution: Drawn below.



Using vector addition and scalar multiplication we often write vectors in terms of so-called standard unit vectors. In the plane and drawn in the margin are the two unit vectors \mathbf{i} and \mathbf{j} (length one) in the direction of the two coordinate axes. Then, for example,

$$(3, 2) = (3, 0) + (0, 2) \quad (\text{by addition})$$

$$\begin{aligned}
 &= 3(1, 0) + 2(0, 1) \quad (\text{by scalar mult}) \\
 &= 3\mathbf{i} + 2\mathbf{j}.
 \end{aligned}$$

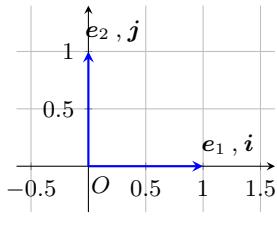


Similarly, in three dimensional space we often write vectors in terms of the three unit vectors \mathbf{i} , \mathbf{j} and \mathbf{k} aligned along the three coordinate axes. For example,

$$\begin{aligned}
 (2, 3, -1) &= (2, 0, 0) + (0, 3, 0) + (0, 0, -1) \quad (\text{by addition}) \\
 &= 2(1, 0, 0) + 3(0, 1, 0) - (0, 0, 1) \quad (\text{by scalar mult}) \\
 &= 2\mathbf{i} + 3\mathbf{j} - \mathbf{k}.
 \end{aligned}$$

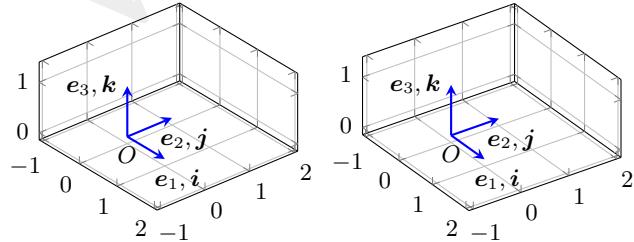
The next definition generalises these standard unit vectors to vectors in \mathbb{R}^n .

Definition 1.2.5. *Given a standard coordinate system with n coordinate axes, all at right-angles to each other, the **standard unit vectors** $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ are the vectors of length one in the direction of the corresponding coordinate axis (as illustrated in the margin for \mathbb{R}^2 and below for \mathbb{R}^3). That is,*



$$\mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots \quad \mathbf{e}_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}.$$

In \mathbb{R}^2 and \mathbb{R}^3 the symbols \mathbf{i} , \mathbf{j} and \mathbf{k} are often used as synonyms for \mathbf{e}_1 , \mathbf{e}_2 and \mathbf{e}_3 , respectively (as also illustrated).

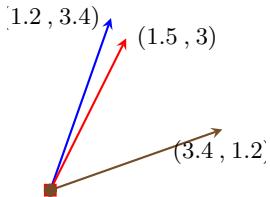


That is, for three examples, the following are equivalent ways of writing the same vector:

$$(3, 2) = \begin{bmatrix} 3 \\ 2 \\ 0 \end{bmatrix} = 3\mathbf{i} + 2\mathbf{j} = 3\mathbf{e}_1 + 2\mathbf{e}_2;$$

$$(2, 3, -1) = \begin{bmatrix} 2 \\ 3 \\ -1 \end{bmatrix} = 2\mathbf{i} + 3\mathbf{j} - \mathbf{k} = 2\mathbf{e}_1 + 3\mathbf{e}_2 - \mathbf{e}_3;$$

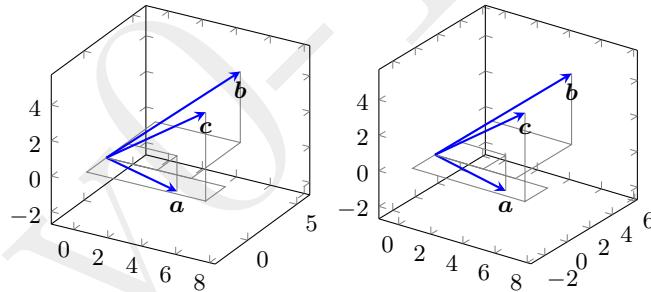
$$(0, -3.7, 0, 0.1, -3.9) = \begin{bmatrix} 0 \\ -3.7 \\ 0 \\ 0.1 \\ -3.9 \end{bmatrix} = -3.7\mathbf{e}_2 + 0.1\mathbf{e}_4 - 3.9\mathbf{e}_5.$$



Distance Defining a ‘distance’ between vectors empowers us to compare vectors concisely. For example we would like to say that $(1.2, 3.4) \approx (1.5, 3)$ to an error 0.5 (as illustrated in the margin). Why 0.5? Because the difference between the vectors $(1.5, 3) - (1.2, 3) = (0.3, -0.4)$ has length $\sqrt{0.3^2 + (-0.4)^2} = 0.5$. Conversely, we would like to recognise that vectors $(1.2, 3.4)$ and $(3.4, 1.2)$ are very different (as also illustrated)—there is a large ‘distance’ between them. Why is there a large ‘distance’? Because the difference between the vectors $(1.2, 3.4) - (3.4, 1.2) = (-2.2, 2.2)$ has length $\sqrt{(-2.2)^2 + 2.2^2} = 2.2\sqrt{2} = 3.1113$ which is relatively large. This concept of distance between two vectors \mathbf{u} and \mathbf{v} , directly analogous to the distance between two points, is the length $|\mathbf{u} - \mathbf{v}|$.

Definition 1.2.6. *The **distance** between vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n is the length of their difference, $|\mathbf{u} - \mathbf{v}|$.*

Example 1.2.7. Given three vectors $\mathbf{a} = 3\mathbf{i} + 2\mathbf{j} - 2\mathbf{k}$, $\mathbf{b} = 5\mathbf{i} + 5\mathbf{j} + 4\mathbf{k}$ and $\mathbf{c} = 7\mathbf{i} - 2\mathbf{j} + 5\mathbf{k}$ (shown below in stereo): which pair are the closest to each other? and which pair are furthest from each other?



Solution: Compute the distances between each pair.

- $|\mathbf{b} - \mathbf{a}| = |2\mathbf{i} + 3\mathbf{j} + 6\mathbf{k}| = \sqrt{2^2 + 3^2 + 6^2} = \sqrt{49} = 7$.
- $|\mathbf{c} - \mathbf{a}| = |4\mathbf{i} - 4\mathbf{j} + 7\mathbf{k}| = \sqrt{4^2 + (-4)^2 + 7^2} = \sqrt{81} = 9$.
- $|\mathbf{c} - \mathbf{b}| = |2\mathbf{i} - 7\mathbf{j} - \mathbf{k}| = \sqrt{2^2 + (-7)^2 + (-1)^2} = \sqrt{54} = 7.3485$.

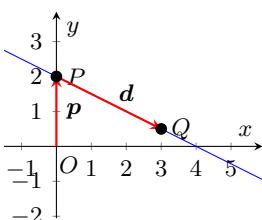
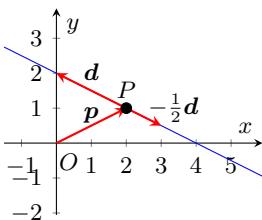
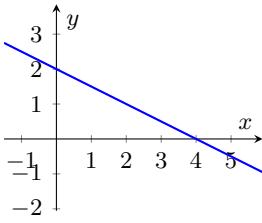
The smallest distance of 7 is between \mathbf{a} and \mathbf{b} so these two are the closest pair of vectors. The largest distance of 9 is between \mathbf{a} and \mathbf{c} so these two are the furthest pair of vectors.

■

1.2.2 Parametric equation of a line

We are familiar with lines in the plane, and equations that describe them. Let’s now consider such equations from a vector view. The insights empower us to generalise the descriptions to lines in space, and then in any number of dimensions.

Example 1.2.8. Consider the line drawn in the margin in some chosen coordinate system. Recall one way to find an equation of the line is to find the intercepts with the axes, here at $x = 4$ and $y = 2$, then write down $\frac{x}{4} + \frac{y}{2} = 1$ as an equation of the line. Then algebraic rearrangement gives various other forms, such as $x + 2y = 4$ or $y = 2 - x/2$.



The alternative is to describe the line with vectors. Choose any point P on the line, such as $(2, 1)$ as drawn in the margin. Then view any other point on the line as having position vector that is the vector sum of \overrightarrow{OP} and a vector aligned along the line. Denote \overrightarrow{OP} by \mathbf{p} as drawn. Then, for example, the point $(0, 2)$ on the line has position vector $\mathbf{p} + \mathbf{d}$ for vector $\mathbf{d} = (-2, 1)$ because $\mathbf{p} + \mathbf{d} = (2, 1) + (-2, 1) = (0, 2)$. Other points on the line are also given using the same vectors, \mathbf{p} and \mathbf{d} : for example, the point $(3, \frac{1}{2})$ has position vector $\mathbf{p} - \frac{1}{2}\mathbf{d}$ (as drawn) because $\mathbf{p} - \frac{1}{2}\mathbf{d} = (2, 1) - \frac{1}{2}(-2, 1) = (3, \frac{1}{2})$; and the point $(-2, 3)$ has position vector $\mathbf{p} + 2\mathbf{d} = (2, 1) + 2(-2, 1)$. In general, every point on the line may be expressed as $\mathbf{p} + t\mathbf{d}$ for some scalar t .

For any given line, there are many possible choices of \mathbf{p} and \mathbf{d} in such a vector representation. A different looking, but equally valid form is obtained from any pair of points on the line. For example, one could choose point P to be $(0, 2)$ and point Q to be $(3, \frac{1}{2})$, as drawn in the margin. Let position vector $\mathbf{p} = \overrightarrow{OP} = (0, 2)$ and the vector $\mathbf{d} = \overrightarrow{PQ} = (3, -\frac{3}{2})$, then every point on the line has position vector $\mathbf{p} + t\mathbf{d}$ for some scalar t :

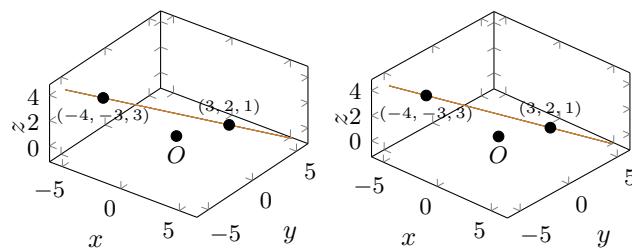
$$(2, 1) = (0, 2) + (2, -1) = (0, 2) + \frac{2}{3}(3, -\frac{3}{2}) = \mathbf{p} + \frac{2}{3}\mathbf{d}; \\ (6, -1) = (0, 2) + (6, -3) = (0, 2) + 2(3, -\frac{3}{2}) = \mathbf{p} + 2\mathbf{d}; \\ (-1, \frac{5}{2}) = (0, 2) + (-1, \frac{1}{2}) = (0, 2) - \frac{1}{3}(3, -\frac{3}{2}) = \mathbf{p} - \frac{1}{3}\mathbf{d}.$$

Other choices of points P and Q give other valid vector equations for a given line. ■

Definition 1.2.9. A **parametric equation** of a line is $\mathbf{x} = \mathbf{p} + t\mathbf{d}$ where \mathbf{p} is the position vector of some point on the line, the so-called **direction vector** \mathbf{d} is parallel to the line ($\mathbf{d} \neq \mathbf{0}$), and the scalar **parameter** t varies over all real values to give all position vectors on the line.

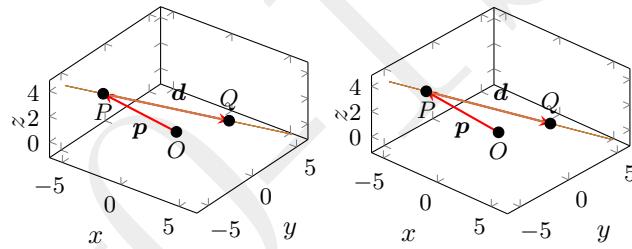
Beautifully, this definition applies for lines in any number of dimensions by using vectors with the corresponding number of components.

Example 1.2.10. Given that the line drawn below in space goes through points $(-4, -3, 3)$ and $(3, 2, 1)$, find a parametric equation of the line.



Solution: Let's call the points $(-4, -3, 3)$ and $(3, 2, 1)$ as P and Q , respectively, and as shown below. First, choose a point on the line, say P , and set its position vector $\mathbf{p} = \overrightarrow{OP} = (-4, -3, 3) = -4\mathbf{i} - 3\mathbf{j} + 3\mathbf{k}$, as drawn. Second, choose a direction vector to be, say, $\mathbf{d} = \overrightarrow{PQ} = (3, 2, 1) - (-4, -3, 3) = 7\mathbf{i} + 5\mathbf{j} - 2\mathbf{k}$, also drawn. A parametric equation of the line is then $\mathbf{x} = \mathbf{p} + t\mathbf{d}$, specifically

$$\begin{aligned}\mathbf{x} &= (-4\mathbf{i} - 3\mathbf{j} + 3\mathbf{k}) + t(7\mathbf{i} + 5\mathbf{j} - 2\mathbf{k}) \\ &= (-4 + 7t)\mathbf{i} + (-3 + 5t)\mathbf{j} + (3 - 2t)\mathbf{k}.\end{aligned}$$



■

Example 1.2.11. Given the parametric equation of a line in space is $\mathbf{x} = (-4 + 2t, 3 - t, -1 - 4t)$, find the value of the parameter t that gives each of the following points on the line: $(-1.6, 1.8, -5.8)$, $(-3, 2.5, -3)$, and $(-6, 4, 4)$.

Solution: • For the point $(-1.6, 1.8, -5.8)$ we need to find the parameter value t such that $-4 + 2t = -1.6$, $3 - t = 1.8$ and $-1 - 4t = -5.8$. The first of these requires $t = (-1.6 + 4)/2 = 1.2$, the second requires $t = 3 - 1.8 = 1.2$, and the third requires $t = (-1 + 5.8)/4 = 1.2$. All three agree that choosing parameter $t = 1.2$ gives the required point.

- For the point $(-3, 2.5, -3)$ we need to find the parameter value t such that $-4 + 2t = -3$, $3 - t = 2.5$ and $-1 - 4t = -3$. The first of these requires $t = (-3 + 4)/2 = 0.5$, the second requires $t = 3 - 2.5 = 0.5$, and the third requires $t = (-1 + 3)/4 = 0.5$. All three agree that choosing parameter $t = 0.5$ gives the required point.
- For the point $(-6, 4, 4)$ we need to find the parameter value t such that $-4 + 2t = -6$, $3 - t = 4$ and $-1 - 4t = 4$. The first of these requires $t = (-6 + 4)/2 = -1$, the second requires $t = 3 - 4 = -1$, and the third requires $t = (-1 - 4)/4 = -1.25$.

Since these three require different values of t , namely -1 and -1.25 , it means that there is no single value of the parameter t that gives the required point. Hence the task is impossible because the point $(-6, 4, 4)$ cannot be on the line.
³

■

1.2.3 Manipulation requires algebraic properties

It seems to be nothing other than that art which they call by the barbarous name of ‘algebra’, if only it could be disentangled from the multiple numbers and inexplicable figures that overwhelm it . . . *Descartes*

To unleash the power of algebra on vectors, we need to know the properties of vector operations. Many of the following properties are familiar as they directly correspond to familiar properties of arithmetic operations on scalars. Moreover, the proofs show the vector properties follow directly from the familiar properties of arithmetic operations on scalars.

Example 1.2.12. Let vectors $\mathbf{u} = (1, 2)$, $\mathbf{v} = (3, 1)$, and $\mathbf{w} = (-2, 3)$, and let scalars $a = -\frac{1}{2}$ and $b = \frac{5}{2}$. Verify the following properties hold:

(a) $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$ (commutative law);

Solution: $\mathbf{u} + \mathbf{v} = (1, 2) + (3, 1) = (1+3, 2+1) = (4, 3)$, whereas $\mathbf{v} + \mathbf{u} = (3, 1) + (1, 2) = (3+1, 1+2) = (4, 3)$ is the same.

(b) $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$ (associative law);

Solution: $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = (4, 3) + (-2, 3) = (2, 6)$, whereas $\mathbf{u} + (\mathbf{v} + \mathbf{w}) = \mathbf{u} + ((3, 1) + (-2, 3)) = (1, 2) + (1, 4) = (2, 6)$ is the same.

(c) $\mathbf{u} + \mathbf{0} = \mathbf{u}$;

Solution: $\mathbf{u} + \mathbf{0} = (1, 2) + (0, 0) = (1+0, 2+0) = (1, 2) = \mathbf{u}$.

(d) $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$;

Solution: Recall $-\mathbf{u} = (-1)\mathbf{u} = (-1)(1, 2) = (-1, -2)$, and so $\mathbf{u} + (-\mathbf{u}) = (1, 2) + (-1, -2) = (1-1, 2-2) = (0, 0) = \mathbf{0}$.

(e) $a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$ (a distributive law);

Solution: $a(\mathbf{u} + \mathbf{v}) = -\frac{1}{2}(4, 3) = (-\frac{1}{2} \cdot 4, -\frac{1}{2} \cdot 3) = (-2, -\frac{3}{2})$, whereas $a\mathbf{u} + a\mathbf{v} = -\frac{1}{2}(1, 2) + (-\frac{1}{2})(3, 1) = (-\frac{1}{2}, -1) + (-\frac{3}{2}, -\frac{1}{2}) = (-\frac{1}{2} - \frac{3}{2}, -1 - \frac{1}{2}) = (-2, -\frac{3}{2})$ which is the same.

³ Section 3.5 develops how to treat such inconsistent information in order to ‘best’ solve such impossible tasks.

(f) $(a + b)\mathbf{u} = a\mathbf{u} + b\mathbf{u}$ (a distributive law);

Solution: $(a + b)\mathbf{u} = (-\frac{1}{2} + \frac{5}{2})(1, 2) = 2(1, 2) = (2 \cdot 1, 2 \cdot 2) = (2, 4)$, whereas $a\mathbf{u} + b\mathbf{u} = (-\frac{1}{2})(1, 2) + \frac{5}{2}(1, 2) = (-\frac{1}{2}, -1) + (\frac{5}{2}, 5) = (-\frac{1}{2} + \frac{5}{2}, -1 + 5) = (-2, 4)$ which is the same.

(g) $(ab)\mathbf{u} = a(b\mathbf{u})$;

Solution: $(ab)\mathbf{u} = (-\frac{1}{2} \cdot 5)(1, 2) = (-\frac{5}{4})(1, 2) = (-\frac{5}{4}, -\frac{5}{2})$, whereas $a(b\mathbf{u}) = a(\frac{5}{2}(1, 2)) = (-\frac{1}{2})(\frac{5}{2}, 5) = (-\frac{5}{4}, -\frac{5}{2})$ which is the same.

(h) $1\mathbf{u} = \mathbf{u}$;

Solution: $1\mathbf{u} = 1(1, 2) = (1 \cdot 1, 1 \cdot 2) = (1, 2) = \mathbf{u}$.

(i) $0\mathbf{u} = \mathbf{0}$;

Solution: $0\mathbf{u} = 0(1, 2) = (0 \cdot 1, 0 \cdot 2) = (0, 0) = \mathbf{0}$.

(j) $|a\mathbf{u}| = |a| \cdot |\mathbf{u}|$.

Solution: Now $|a| = |-\frac{1}{2}| = \frac{1}{2}$, and the length $|\mathbf{u}| = \sqrt{1^2 + 2^2} = \sqrt{5}$ (Definition 1.1.8). Consequently, $|a\mathbf{u}| = |(-\frac{1}{2})(1, 2)| = |(-\frac{1}{2}, -1)| = \sqrt{(-\frac{1}{2})^2 + (-1)^2} = \sqrt{\frac{1}{4} + 1} = \sqrt{\frac{5}{4}} = \frac{1}{2}\sqrt{5} = |a| \cdot |\mathbf{u}|$ as required. ■

Theorem 1.2.13. Let \mathbf{u} , \mathbf{v} and \mathbf{w} be any vectors with n components (in \mathbb{R}^n), and let a and b be any scalars. Then the following properties hold:

(a) $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$ (commutative law);

(b) $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$ (associative law);

(c) $\mathbf{u} + \mathbf{0} = \mathbf{0} + \mathbf{u} = \mathbf{u}$;

(d) $\mathbf{u} + (-\mathbf{u}) = (-\mathbf{u}) + \mathbf{u} = \mathbf{0}$;

(e) $a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$ (a distributive law);

(f) $(a + b)\mathbf{u} = a\mathbf{u} + b\mathbf{u}$ (a distributive law);

(g) $(ab)\mathbf{u} = a(b\mathbf{u})$;

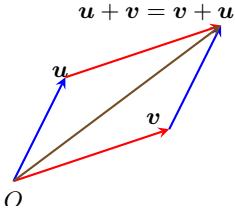
(h) $1\mathbf{u} = \mathbf{u}$;

(i) $0\mathbf{u} = \mathbf{0}$;

(j) $|a\mathbf{u}| = |a| \cdot |\mathbf{u}|$.

Proof. We prove property 1.2.13a, and leave the proof of other properties as exercises. The approach is to establish the properties of vector operations using the known properties of scalar operations.

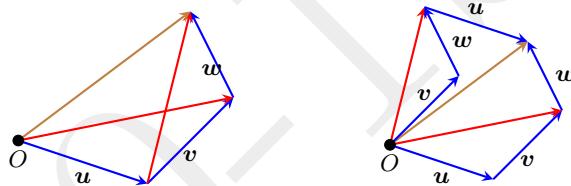
Property 1.2.13a is the commutativity of vector addition. Example 1.2.1a shows graphically how $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$ in one case, and the margin here shows another case. In general, let vectors $\mathbf{u} = (u_1, u_2, \dots, u_n)$ and $\mathbf{v} = (v_1, v_2, \dots, v_n)$ then



$$\begin{aligned}
 & \mathbf{u} + \mathbf{v} \\
 &= (u_1, u_2, \dots, u_n) + (v_1, v_2, \dots, v_n) \\
 &= (u_1 + v_1, u_2 + v_2, \dots, u_n + v_n) \quad (\text{by Defn. 1.2.3}) \\
 &= (v_1 + u_1, v_2 + u_2, \dots, v_n + u_n) \quad (\text{commutative scalar add}) \\
 &= (v_1, v_2, \dots, v_n) + (u_1, u_2, \dots, u_n) \quad (\text{by Defn. 1.2.3}) \\
 &= \mathbf{v} + \mathbf{u}.
 \end{aligned}$$

□

Example 1.2.14. Which of the following two diagrams best illustrates the associative law 1.2.13b? Give reasons.



Solution: The left diagram.

- In the left diagram, the two red vectors represent $\mathbf{u} + \mathbf{v}$ (left) and $\mathbf{v} + \mathbf{w}$ (right). Thus the left-red followed by the blue \mathbf{w} represents $(\mathbf{u} + \mathbf{v}) + \mathbf{w}$, whereas the \mathbf{u} followed by the right-red represents $\mathbf{u} + (\mathbf{v} + \mathbf{w})$. The brown vector shows they are equal: $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$.
- The right-hand diagram invokes the commutative law as well. The top-left part of the diagram shows $(\mathbf{v} + \mathbf{w}) + \mathbf{u}$, whereas the bottom-right part shows $(\mathbf{u} + \mathbf{v}) + \mathbf{w}$. That these are equal, the brown vector, requires both the commutative and associative laws.

■

We frequently use these algebraic properties, such as in rearranging and solving vector equations.

Example 1.2.15. Find the vector \mathbf{x} such that $3\mathbf{x} - 2\mathbf{u} = 6\mathbf{v}$.

Solution: Using Theorem 1.2.13, all the following equations are equivalent:

$$\begin{aligned}
 & 3\mathbf{x} - 2\mathbf{u} = 6\mathbf{v}; \\
 & (3\mathbf{x} - 2\mathbf{u}) + 2\mathbf{u} = 6\mathbf{v} + 2\mathbf{u} \quad (\text{add } 2\mathbf{u} \text{ to both sides}); \\
 & 3\mathbf{x} + (-2\mathbf{u} + 2\mathbf{u}) = 6\mathbf{v} + 2\mathbf{u} \quad (\text{by 1.2.13b, associativity});
 \end{aligned}$$

$$\begin{aligned}
 3\mathbf{x} + \mathbf{0} &= 6\mathbf{v} + 2\mathbf{u} \quad (\text{by 1.2.13d}); \\
 3\mathbf{x} &= 6\mathbf{v} + 2\mathbf{u} \quad (\text{by 1.2.13c}); \\
 \frac{1}{3}(3\mathbf{x}) &= \frac{1}{3}(6\mathbf{v} + 2\mathbf{u}) \quad (\text{multiply both sides by } \frac{1}{3}); \\
 \frac{1}{3}(3\mathbf{x}) &= \frac{1}{3}(6\mathbf{v}) + \frac{1}{3}(2\mathbf{u}) \quad (\text{by 1.2.13e, distributivity}); \\
 (\frac{1}{3} \cdot 3)\mathbf{x} &= (\frac{1}{3} \cdot 6)\mathbf{v} + (\frac{1}{3} \cdot 2)\mathbf{u} \quad (\text{by 1.2.13g}); \\
 1\mathbf{x} &= 2\mathbf{v} + \frac{2}{3}\mathbf{u} \quad (\text{by scalar operations}); \\
 \mathbf{x} &= 2\mathbf{v} + \frac{2}{3}\mathbf{u} \quad (\text{by 1.2.13h}).
 \end{aligned}$$

Generally we do not write down all such details. Generally the following shorter derivation is acceptable. The following are equivalent:

$$\begin{aligned}
 3\mathbf{x} - 2\mathbf{u} &= 6\mathbf{v}; \\
 3\mathbf{x} &= 6\mathbf{v} + 2\mathbf{u} \quad (\text{adding } 2\mathbf{u} \text{ to both sides}); \\
 \mathbf{x} &= 2\mathbf{v} + \frac{2}{3}\mathbf{u} \quad (\text{dividing both sides by 3}).
 \end{aligned}$$

But exercises and examples in this section often explicitly require full details and justification.

■

Example 1.2.16. Rearrange $3\mathbf{x} - \mathbf{a} = 2(\mathbf{a} + \mathbf{x})$ for vector \mathbf{x} in terms of \mathbf{a} : giving excruciating detail of the justification using Theorem 1.2.13.

Solution: Using Theorem 1.2.13, the following statements are equivalent:

$$\begin{aligned}
 3\mathbf{x} - \mathbf{a} &= 2(\mathbf{a} + \mathbf{x}) \\
 3\mathbf{x} - \mathbf{a} &= 2\mathbf{a} + 2\mathbf{x} \quad (\text{by 1.2.13e, distributivity}); \\
 (3\mathbf{x} - \mathbf{a}) + \mathbf{a} &= (2\mathbf{a} + 2\mathbf{x}) + \mathbf{a} \quad (\text{adding } \mathbf{a} \text{ to both sides}); \\
 3\mathbf{x} + (-\mathbf{a} + \mathbf{a}) &= 2\mathbf{a} + (2\mathbf{x} + \mathbf{a}) \quad (\text{by 1.2.13b, associativity}); \\
 3\mathbf{x} + \mathbf{0} &= 2\mathbf{a} + (\mathbf{a} + 2\mathbf{x}) \quad (\text{by 1.2.13d and 1.2.13a}); \\
 3\mathbf{x} &= (2\mathbf{a} + \mathbf{a}) + 2\mathbf{x} \quad (\text{by 1.2.13c and 1.2.13b}); \\
 3\mathbf{x} &= (2\mathbf{a} + 1\mathbf{a}) + 2\mathbf{x} \quad (\text{by 1.2.13h}); \\
 3\mathbf{x} &= (2 + 1)\mathbf{a} + 2\mathbf{x} \quad (\text{by 1.2.13f, distributivity}); \\
 3\mathbf{x} + (-2)\mathbf{x} &= 3\mathbf{a} + 2\mathbf{x} + (-2)\mathbf{x} \quad (\text{sub. } 2\mathbf{x} \text{ from both sides}); \\
 (3 + (-2))\mathbf{x} &= 3\mathbf{a} + (2 + (-2))\mathbf{x} \quad (\text{by 1.2.13f, distributivity}); \\
 1\mathbf{x} &= 3\mathbf{a} + 0\mathbf{x} \quad (\text{by scalar arithmetic}); \\
 \mathbf{x} &= 3\mathbf{a} + \mathbf{0} \quad (\text{by 1.2.13h and 1.2.13i}); \\
 \mathbf{x} &= 3\mathbf{a} \quad (\text{by 1.2.13c}).
 \end{aligned}$$

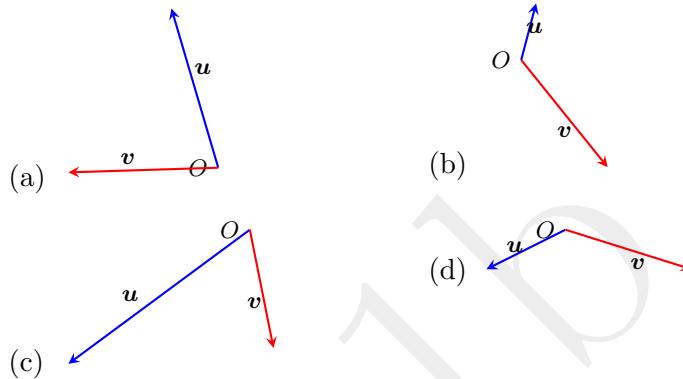
If the question had not requested full details, then the following would be enough. The following statements are equivalent:

$$\begin{aligned}
 3\mathbf{x} - \mathbf{a} &= 2(\mathbf{a} + \mathbf{x}) \quad (\text{distribute the multiplication}) \\
 3\mathbf{x} &= 2\mathbf{a} + 2\mathbf{x} + \mathbf{a} \quad (\text{adding } \mathbf{a} \text{ to both sides}); \\
 \mathbf{x} &= 3\mathbf{a} \quad (\text{subtracting } 2\mathbf{x} \text{ from both sides}).
 \end{aligned}$$

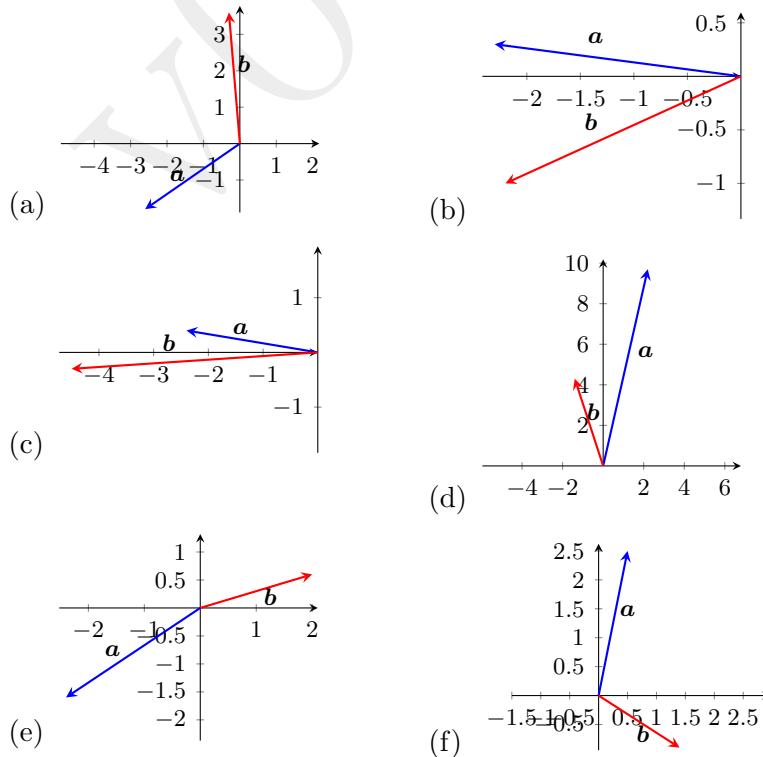


1.2.4 Exercises

Exercise 1.2.1. For each of the pairs of vectors \mathbf{u} and \mathbf{v} shown below, draw the vectors $\mathbf{u} + \mathbf{v}$, $\mathbf{v} + \mathbf{u}$, $\mathbf{u} - \mathbf{v}$, $\mathbf{v} - \mathbf{u}$, $\frac{1}{2}\mathbf{u}$ and $-\mathbf{v}$.



Exercise 1.2.2. For each of the following pairs of vectors shown below, use a ruler (or other measuring stick) to directly measure the distance between the pair of vectors.



Exercise 1.2.3. For each of the following groups of vectors, use the distance between vectors to find which pair in the group are closest to each other, and which pair in the group are furthest from each other.

(a) $\mathbf{u} = (-5, 0, 3)$, $\mathbf{v} = (1, -6, 10)$, $\mathbf{w} = (-4, 4, 11)$

- (b) $\mathbf{u} = (2, 2, -1)$, $\mathbf{v} = (3, 6, -9)$, $\mathbf{w} = (1, -2, -9)$
 (c) $\mathbf{u} = (1, 1, -3)$, $\mathbf{v} = (7, 7, -10)$, $\mathbf{w} = (-1, 4, -9)$
 (d) $\mathbf{u} = 3\mathbf{i}$, $\mathbf{v} = 4\mathbf{i} - 2\mathbf{j} + 2\mathbf{k}$, $\mathbf{w} = 4\mathbf{i} + 2\mathbf{j} + 2\mathbf{k}$
 (e) $\mathbf{u} = (-5, 3, 5, 6)$, $\mathbf{v} = (-6, 1, 3, 10)$, $\mathbf{w} = (-4, 6, 2, 15)$
 (f) $\mathbf{u} = (-4, -1, -1, 2)$, $\mathbf{v} = (-5, -2, -2, 1)$, $\mathbf{w} = (-3, -2, -2, 1)$
- (g) $\mathbf{u} = 5\mathbf{e}_1 + \mathbf{e}_3 + 5\mathbf{e}_4$, $\mathbf{v} = 6\mathbf{e}_1 - 2\mathbf{e}_2 + 3\mathbf{e}_3 + \mathbf{e}_4$, $\mathbf{w} = 7\mathbf{e}_1 - 2\mathbf{e}_2 - 3\mathbf{e}_3$

(h) $\mathbf{u} = 2\mathbf{e}_1 + 4\mathbf{e}_1 - \mathbf{e}_3 + 5\mathbf{e}_4$, $\mathbf{v} = -2\mathbf{e}_1 + 8\mathbf{e}_2 - 6\mathbf{e}_3 - 3\mathbf{e}_4$,
 $\mathbf{w} = -6\mathbf{e}_3 + 11\mathbf{e}_4$

Exercise 1.2.4. Find a parametric equation of the line through the given two points.

- (a) $(-11, 0, 3)$, $(-3, -2, 2)$ (b) $(-4, 1, -2)$, $(3, -5, 5)$
 (c) $(2.4, 5.5, -3.9)$,
 $(1.5, -5.4, -0.5)$ (d) $(0.2, -7.2, -4.6, -2.8)$,
 $(3.3, -1.1, -0.4, -0.3)$
 (e) $(2.2, 5.8, 4, 3, 2)$,
 $(-1.1, 2.2, -2.4, -3.2, 0.9)$ (f) $(1.8, -3.1, -1, -1.3, -3.3)$,
 $(-1.4, 0.8, -2.6, 3.1, -0.8)$

Exercise 1.2.5. Verify the algebraic properties of Theorem 1.2.13 for each of the following sets of vectors and scalars.

- (a) $\mathbf{u} = 2.4\mathbf{i} - 0.3\mathbf{j}$, $\mathbf{v} = -1.9\mathbf{i} + 0.5\mathbf{j}$, $\mathbf{w} = -3.5\mathbf{i} - 1.8\mathbf{j}$, $a = 0.4$ and $b = 1.4$.
 (b) $\mathbf{u} = (1/3, 14/3)$, $\mathbf{v} = (4, 4)$, $\mathbf{w} = (2/3, -10/3)$, $a = -2/3$ and $b = -1$.
 (c) $\mathbf{u} = -\frac{1}{2}\mathbf{j} + \frac{3}{2}\mathbf{k}$, $\mathbf{v} = 2\mathbf{i} - \mathbf{j}$, $\mathbf{w} = 2\mathbf{i} - \mathbf{k}$, $a = -3$ and $b = \frac{1}{2}$.
 (d) $\mathbf{u} = (2, 1, 4, -2)$, $\mathbf{v} = (-3, -2, 0, -1)$, $\mathbf{w} = (-6, 5, 4, 2)$, $a = -4$ and $b = 3$.

Exercise 1.2.6. Prove in detail some algebraic properties chosen from Theorem 1.2.13b–1.2.13j on vector addition and scalar multiplication.

Exercise 1.2.7. For each of the following vectors equations, rearrange the equations to get vector \mathbf{x} in terms of the other vectors. Give excruciating detail of the justification using Theorem 1.2.13.

- (a) $\mathbf{x} + \mathbf{a} = \mathbf{0}$.
 (b) $2\mathbf{x} - \mathbf{b} = 3\mathbf{b}$.
 (c) $3(\mathbf{x} + \mathbf{a}) = \mathbf{x} + (\mathbf{a} - 2\mathbf{x})$.
 (d) $-4\mathbf{b} = \mathbf{x} + 3(\mathbf{a} - \mathbf{x})$.

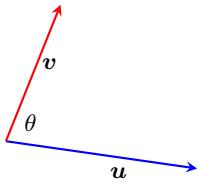
1.3 The dot product determines angles and lengths

Section Contents

1.3.1	Work done involves the dot product	40
1.3.2	Algebraic properties of the dot product	41
1.3.3	Orthogonal vectors are at right-angles	46
1.3.4	Normal vectors and equations of a plane	48
1.3.5	Exercises	53

The previous Section 1.2 discussed how to add, subtract and stretch vectors. Question: can we multiply two vectors? The answer is yes, but different. ‘Vector multiplication’ has major differences to multiplication of scalars. This section introduces the so-called dot product that, among other attributes, gives us a way of determining the angle between two vectors.

Example 1.3.1. Consider the two vectors $\mathbf{u} = (7, -1)$ and $\mathbf{v} = (2, 5)$ as plotted in the margin. What is the angle θ between the two vectors?

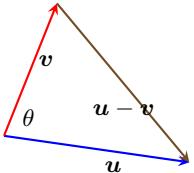


Solution: Form a triangle with the vector $\mathbf{u} - \mathbf{v} = (5, -6)$ going from the tip of \mathbf{v} to the tip of \mathbf{u} , as shown in the margin. The sides of the triangles are of length $|\mathbf{u}| = \sqrt{7^2 + (-1)^2} = \sqrt{50} = 5\sqrt{2}$, $|\mathbf{v}| = \sqrt{2^2 + 5^2} = \sqrt{29}$, and $|\mathbf{u} - \mathbf{v}| = \sqrt{5^2 + (-6)^2} = \sqrt{61}$. By the cosine rule for triangles

$$|\mathbf{u} - \mathbf{v}|^2 = |\mathbf{u}|^2 + |\mathbf{v}|^2 - 2|\mathbf{u}||\mathbf{v}|\cos\theta$$

Here this rule rearranges to

$$\begin{aligned} |\mathbf{u}||\mathbf{v}|\cos\theta &= \frac{1}{2}(|\mathbf{u}|^2 + |\mathbf{v}|^2 - |\mathbf{u} - \mathbf{v}|^2) \\ &= \frac{1}{2}(50 + 29 - 61) \\ &= 9. \end{aligned}$$



Dividing by the product of the lengths then gives $\cos\theta = 9/(5\sqrt{58}) = 0.2364$ so the angle $\theta = \arccos(0.2364) = 1.3322 = 76.33^\circ$ as is reasonable from the plots. ■

The interest in this Example 1.3.1 is the number nine on the right-hand side of $|\mathbf{u}||\mathbf{v}|\cos\theta = 9$. The reason is that 9 just happens to be $14 - 5$, which in turn just happens to be $7 \cdot 2 + (-1) \cdot 5$, and it is no coincidence that this expression is the same as $u_1v_1 + u_2v_2$ in terms of vector components $\mathbf{u} = (u_1, u_2) = (7, -1)$ and $\mathbf{v} = (v_1, v_2) = (2, 5)$. Repeat this example for any pair of vectors \mathbf{u} and \mathbf{v} to find that always $|\mathbf{u}||\mathbf{v}|\cos\theta = u_1v_1 + u_2v_2$ (Exercise 1.3.1). This equality suggests that the sum of products of corresponding components of \mathbf{u} and \mathbf{v} is closely connected to the angle between the vectors.

Definition 1.3.2. For any two vectors in \mathbb{R}^n , $\mathbf{u} = (u_1, u_2, \dots, u_n)$ and $\mathbf{v} = (v_1, v_2, \dots, v_n)$, define the **dot product** (or **inner product**), denoted by a dot between the two vectors, as the scalar

$$\mathbf{u} \cdot \mathbf{v} = u_1 v_1 + u_2 v_2 + \cdots + u_n v_n.$$

The dot product of two vectors gives a scalar result, not a vector.

When writing the vector dot product, the dot between the two vectors is essential. We sometimes also denote the scalar product by such a dot (to clarify a product) and sometimes omit the dot between the scalars, for example $a \cdot b = ab$ for scalars. But for the vector dot product the dot must not be omitted: ' $\mathbf{u}\mathbf{v}$ ' is meaningless.

Example 1.3.3. Compute the dot product between the following pairs of vectors.

(a) $\mathbf{u} = (-2, 5, -2)$, $\mathbf{v} = (3, 3, -2)$

Solution: $\mathbf{u} \cdot \mathbf{v} = (-2)3 + 5 \cdot 3 + (-2)(-2) = 13$. Alternatively, $\mathbf{v} \cdot \mathbf{u} = 3(-2) + 3 \cdot 5 + (-2)(-2) = 13$. That these give the same result is a commutative law, Theorem 1.3.10a, and so in the following we compute the dot product only one way around.

(b) $\mathbf{u} = (1, -3, 0)$, $\mathbf{v} = (1, 2)$

Solution: There is no answer: the dot product cannot be computed as the two vectors are of different sizes.

(c) $\mathbf{a} = (-7, 3, 0, 2, 2)$, $\mathbf{b} = (-3, 4, -4, 2, 0)$

Solution: $\mathbf{a} \cdot \mathbf{b} = (-7)(-3) + 3 \cdot 4 + 0(-4) + 2 \cdot 2 + 2 \cdot 0 = 37$.

(d) $\mathbf{p} = (-0.1, -2.5, -3.3, 0.2)$, $\mathbf{q} = (-1.6, 1.1, -3.4, 2.2)$

Solution: $\mathbf{p} \cdot \mathbf{q} = (-0.1)(-1.6) + (-2.5)1.1 + (-3.3)(-3.4) + 0.2 \cdot 2.2 = 9.07$.

■

Theorem 1.3.4. For any two non-zero vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n , the **angle** θ between the vectors is determined by

$$\cos \theta = \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}| |\mathbf{v}|}, \quad 0 \leq \theta \leq \pi \quad (0^\circ \leq \theta \leq 180^\circ).$$

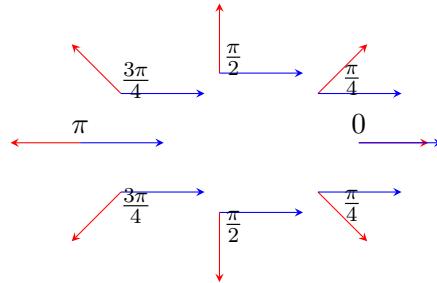


Table 1.1: when a cosine is one of these special values, then we know the corresponding angle exactly. In other cases we usually use a calculator (\arccos or \cos^{-1}) or computer (`acos()`) to compute the angle numerically.

θ	θ	$\cos \theta$	$\cos \theta$
0	0°	1	
$\pi/6$	30°	$\sqrt{3}/2$	0.8660
$\pi/4$	45°	$1/\sqrt{2}$	0.7071
$\pi/3$	60°	$1/2$	0.5
$\pi/2$	90°	0	
$2\pi/3$	120°	$-1/2$	-0.5
$3\pi/4$	135°	$-1/\sqrt{2}$	-0.7071
$5\pi/6$	150°	$-\sqrt{3}/2$	-0.8660
π	180°	-1	

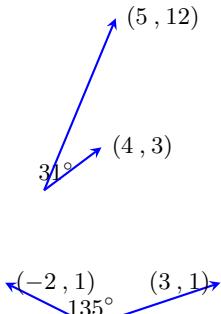
Example 1.3.5. Determine the angle between the following pairs of vectors.

- (a) $(4, 3)$ and $(5, 12)$

Solution: These vectors (shown in the margin) have length $\sqrt{4^2 + 3^2} = \sqrt{25} = 5$ and $\sqrt{5^2 + 12^2} = \sqrt{169} = 13$, respectively. Their dot product $(4, 3) \cdot (5, 12) = 20 + 36 = 56$. Hence $\cos \theta = 56/(5 \cdot 13) = 0.8615$ and so angle $\theta = \arccos(0.8615) = 0.5325 = 30.51^\circ$.

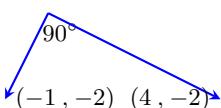
- (b) $(3, 1)$ and $(-2, 1)$

Solution: These vectors (shown in the margin) have length $\sqrt{3^2 + 1^2} = \sqrt{10}$ and $\sqrt{(-2)^2 + 1^2} = \sqrt{5}$, respectively. Their dot product $(3, 1) \cdot (-2, 1) = -6 + 1 = -5$. Hence $\cos \theta = -5/(\sqrt{10} \cdot \sqrt{5}) = -1/\sqrt{2} = -0.7071$ and so angle $\theta = \arccos(-1/\sqrt{2}) = 2.3562 = \frac{3}{4}\pi = 135^\circ$ (Table 1.1).



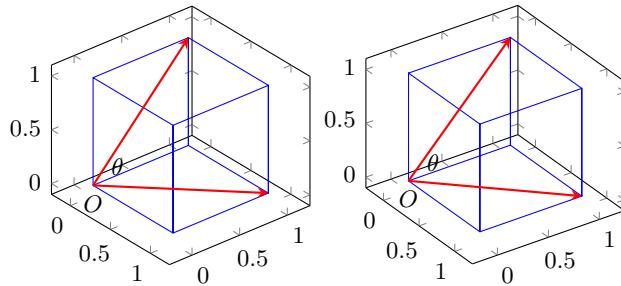
- (c) $(4, -2)$ and $(-1, -2)$

Solution: These vectors (shown in the margin) have length $\sqrt{4^2 + (-2)^2} = \sqrt{20} = 2\sqrt{5}$ and $\sqrt{(-1)^2 + (-2)^2} = \sqrt{5}$, respectively. Their dot product $(4, -2) \cdot (-1, -2) = -4 + 4 = 0$. Hence $\cos \theta = 0/(2\sqrt{5} \cdot \sqrt{5}) = 0$ and so angle $\theta = \frac{1}{2}\pi = 90^\circ$ (Table 1.1). ■



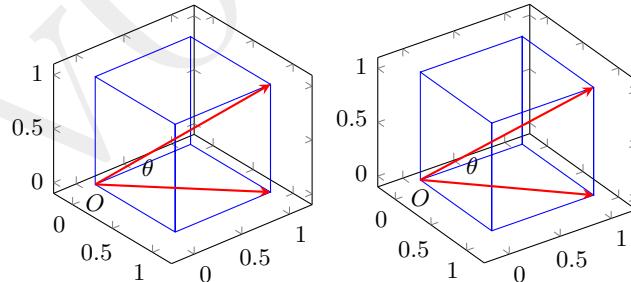
Example 1.3.6. In chemistry one computes the angles between bonds in molecules and crystals. In engineering one needs the angles between beams and struts in complex structures. The dot product determines such angles.

- (a) Consider the cube drawn in stereo below, and compute the angle between the diagonals on two adjacent faces.



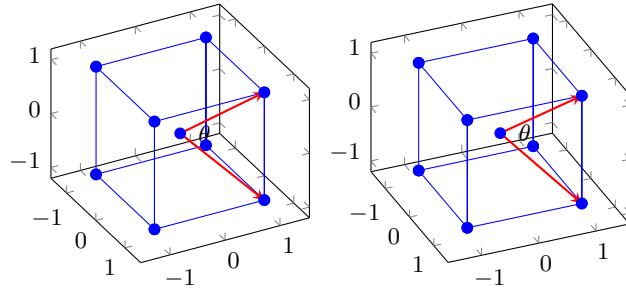
Solution: Draw two vectors along adjacent diagonals: the above pair of vectors are $(1, 1, 0)$ and $(0, 1, 1)$. They both have the same length as $|(1, 1, 0)| = \sqrt{1^2 + 1^2 + 0^2} = \sqrt{2}$ and $|(0, 1, 1)| = \sqrt{0^2 + 1^2 + 1^2} = \sqrt{2}$. The dot product is $(1, 1, 0) \cdot (0, 1, 1) = 0 + 1 + 0 = 1$. Hence the cosine $\cos \theta = 1/(\sqrt{2} \cdot \sqrt{2}) = 1/2$. Table 1.1 gives the angle $\theta = \frac{\pi}{3} = 60^\circ$.

- (b) Consider the cube drawn in stereo below: what is the angle between a diagonal on a face and a diagonal of the cube?



Solution: Draw two vectors along the diagonals: the above pair of vectors are $(1, 1, 0)$ and $(1, 1, 1)$. The face-diagonal has length $|(1, 1, 0)| = \sqrt{1^2 + 1^2 + 0^2} = \sqrt{2}$ whereas the cube diagonal has length $|(1, 1, 1)| = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$. The dot product is $(1, 1, 0) \cdot (1, 1, 1) = 1 + 1 + 0 = 2$. Hence $\cos \theta = 2/(\sqrt{2} \cdot \sqrt{3}) = \sqrt{2}/\sqrt{3} = 0.8165$. Then a calculator (or Matlab/Octave, see Section 1.5) gives the angle $\theta = \arccos(0.8165) = 0.6155 = 35.26^\circ$.

- (c) A body-centered cubic lattice (such as that formed by caesium chloride crystals) has one lattice point in the center of the unit cell as well as the eight corner points. Consider the body-centered cube of atoms drawn in stereo below with the center of the cube at the origin: what is the angle between the center atom and two adjacent corner atoms?



Solution: Draw two corresponding vectors from the center atom: the above pair of vectors are $(1, 1, 1)$ and $(1, 1, -1)$. These have the same length $|(1, 1, 1)| = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$ and $|(1, 1, -1)| = \sqrt{1^2 + 1^2 + (-1)^2} = \sqrt{3}$. The dot product is $(1, 1, 1) \cdot (1, 1, -1) = 1 + 1 - 1 = 1$. Hence $\cos \theta = 1/(\sqrt{3} \cdot \sqrt{3}) = 1/3 = 0.3333$. Then a calculator (or Matlab/Octave, see Section 1.5) gives the angle $\theta = \arccos(1/3) = 1.2310 = 70.53^\circ$.

■

Example 1.3.7 (semantic similarity). Recall that Example 1.1.7 introduced the encoding of sentences and documents as word count vectors. In the example, a word vector has five components, $(N_{\text{cat}}, N_{\text{dog}}, N_{\text{mat}}, N_{\text{sat}}, N_{\text{scratched}})$ where the various N are the counts of each word in any sentence or document. For example,

- (a) “The dog sat on the mat” has word vector $\mathbf{a} = (0, 1, 1, 1, 0)$.
- (b) “The cat scratched the dog” has word vector $\mathbf{b} = (1, 1, 0, 0, 1)$.
- (c) “The cat and dog sat on the mat” has word vector $\mathbf{c} = (1, 1, 1, 1, 0)$.

Use the angle between these three word vectors to characterise the similarity of sentences: a small angle means the sentences are somehow close; a large angle means the sentences are disparate.

Solution: First, these word vectors have lengths $|\mathbf{a}| = |\mathbf{b}| = \sqrt{3}$ and $|\mathbf{c}| = 2$. Second, the ‘angles’ between these sentences are the following.

- The angle θ_{ab} between “The dog sat on the mat” and “The cat scratched the dog” satisfies

$$\cos \theta_{ab} = \frac{\mathbf{a} \cdot \mathbf{b}}{|\mathbf{a}| |\mathbf{b}|} = \frac{0 + 1 + 0 + 0 + 0}{\sqrt{3} \cdot \sqrt{3}} = \frac{1}{3}.$$

A calculator (or Matlab/Octave, see Section 1.5) then gives the angle $\theta_{ab} = \arccos(1/3) = 1.2310 = 70.53^\circ$ so the sentences are quite dissimilar.

- The angle θ_{ac} between “The dog sat on the mat” and “The cat and dog sat on the mat” satisfies

$$\cos \theta_{ac} = \frac{\mathbf{a} \cdot \mathbf{c}}{|\mathbf{a}| |\mathbf{c}|} = \frac{0 + 1 + 1 + 1 + 0}{\sqrt{3} \cdot 2} = \frac{3}{2\sqrt{3}} = \frac{\sqrt{3}}{2}.$$

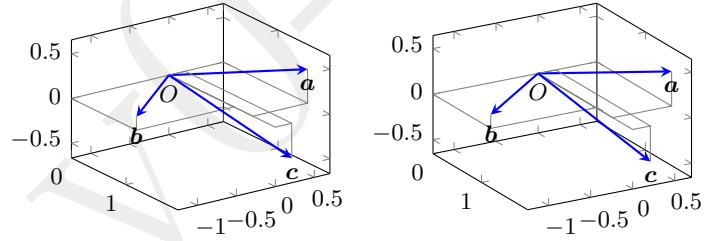
Table 1.1 gives the angle $\theta_{ac} = \frac{\pi}{6} = 30^\circ$ so the sentences are roughly similar.

- The angle θ_{bc} between “The cat scratched the dog” and “The cat and dog sat on the mat” satisfies

$$\cos \theta_{bc} = \frac{\mathbf{b} \cdot \mathbf{c}}{|\mathbf{b}| |\mathbf{c}|} = \frac{1 + 1 + 0 + 0 + 0}{\sqrt{3} \cdot 2} = \frac{2}{2\sqrt{3}} = \frac{1}{\sqrt{3}}.$$

A calculator (or Matlab/Octave, see Section 1.5) then gives the angle $\theta_{bc} = \arccos(1/\sqrt{3}) = 0.9553 = 54.74^\circ$ so the sentences are moderately dissimilar.

The following stereo plot schematically draws these three vectors at the correct angles from each other, and with correct lengths, in some abstract coordinate system (Section 3.4 gives the techniques to do such plots systematically).

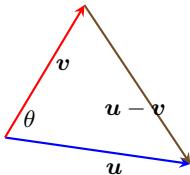


■

Proof. To prove the angle Theorem 1.3.4, form a triangle from vectors \mathbf{u} , \mathbf{v} and $\mathbf{u} - \mathbf{v}$ as illustrated in the margin. Recall and apply the cosine rule for triangles

$$|\mathbf{u} - \mathbf{v}|^2 = |\mathbf{u}|^2 + |\mathbf{v}|^2 - 2|\mathbf{u}||\mathbf{v}| \cos \theta.$$

In \mathbb{R}^n this rule rearranges to



$$\begin{aligned} 2|\mathbf{u}||\mathbf{v}| \cos \theta &= |\mathbf{u}|^2 + |\mathbf{v}|^2 - |\mathbf{u} - \mathbf{v}|^2 \\ &= u_1^2 + u_2^2 + \cdots + u_n^2 + v_1^2 + v_2^2 + \cdots + v_n^2 \\ &\quad - (u_1 - v_1)^2 - (u_2 - v_2)^2 - \cdots - (u_n - v_n)^2 \\ &= u_1^2 + u_2^2 + \cdots + u_n^2 + v_1^2 + v_2^2 + \cdots + v_n^2 \\ &\quad - u_1^2 + 2u_1v_1 - v_1^2 - u_2^2 + 2u_2v_2 - v_2^2 \\ &\quad - \cdots - u_n^2 + 2u_nv_n - v_n^2 \\ &= 2u_1v_1 + 2u_2v_2 + \cdots + 2u_nv_n \\ &= 2(u_1v_1 + u_2v_2 + \cdots + u_nv_n) \end{aligned}$$

$$= 2\mathbf{u} \cdot \mathbf{v}.$$

Dividing both sides by $2|\mathbf{u}||\mathbf{v}|$ gives $\cos \theta = \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}||\mathbf{v}|}$ as required. \square

1.3.1 Work done involves the dot product

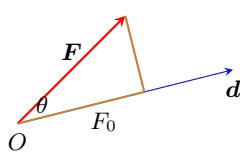
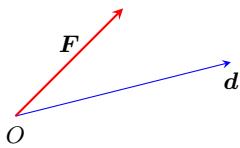
In physics and engineering, “work” has a precise meaning related to energy: when a force of magnitude F acts on a body and that body moves a distance d , then the work done by the force is $W = Fd$. This formula applies only for one dimensional force and displacement, the case when the force and the displacement are all in the same direction. For example if a 5 kg barbell drops downwards 2 m under the force of gravity (9.8 newtons/kg), then the work done by gravity on the barbell during the drop is

$$W = F \times d = (5 \times 9.8) \times 2 = 98 \text{ joules.}$$

This work done goes to the kinetic energy of the falling barbell. The kinetic energy dissipates when the barbell hits the floor.

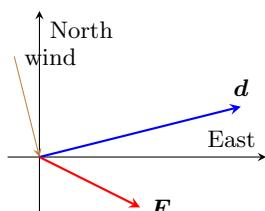
In general, the applied force and the displacement are not in the same direction (as illustrated in the margin). Consider the general case when a vector force \mathbf{F} acts on a body which moves a displacement \mathbf{d} . Then the work done by the force on the body is the length of the displacement times the component of the force in the direction of the displacement—the component of the force at right-angles to the displacement does no work.

As illustrated in the margin, draw a right-angled triangle to decompose the force \mathbf{F} into the component F_0 in the direction of the displacement, and an unnamed component at right-angles. Then by the scalar formula, the work done is $W = F_0|\mathbf{d}|$. As drawn, the force \mathbf{F} makes an angle θ to the displacement \mathbf{d} : the dot product determines this angle via $\cos \theta = (\mathbf{F} \cdot \mathbf{d})/(|\mathbf{F}||\mathbf{d}|)$ (Theorem 1.3.4). By basic trigonometry, the adjacent side of the force triangle has length $F_0 = |\mathbf{F}| \cos \theta = |\mathbf{F}| \frac{\mathbf{F} \cdot \mathbf{d}}{|\mathbf{F}||\mathbf{d}|} = \frac{\mathbf{F} \cdot \mathbf{d}}{|\mathbf{d}|}$. Finally, the work done $W = F_0|\mathbf{d}| = \frac{\mathbf{F} \cdot \mathbf{d}}{|\mathbf{d}|}|\mathbf{d}| = \mathbf{F} \cdot \mathbf{d}$, the dot product of the vector force and vector displacement.



Example 1.3.8. A sailing boat travels a distance of 40 m East and 10 m North, as drawn in the margin. The wind from abeam of strength and direction $(1, -4)$ m/s generates a force $\mathbf{F} = (20, -10)$ (newtons) on the sail, as drawn. What is the work done by the wind.

Solution: The direction of the wind is immaterial except for the force it generates. The displacement vector $\mathbf{d} = (40, 10)$ m. Then the work done is $W = \mathbf{F} \cdot \mathbf{d} = (40, 10) \cdot (20, -10) = 800 - 100 = 700$ joules. \blacksquare



Finding components of vectors in various directions, called projection and surprisingly common in applications, is developed much further by section 3.5.3.

1.3.2 Algebraic properties of the dot product

To manipulate the dot product in algebraic expressions, we need to know its basic algebraic rules. The following rules of Theorem 1.3.10 are analogous to well known rules for scalar multiplication.

Example 1.3.9. Given vectors $\mathbf{u} = (-2, 5, -2)$, $\mathbf{v} = (3, 3, -2)$ and $\mathbf{w} = (2, 0, -5)$, and scalar $a = 2$, verify that (properties 1.3.10c and 1.3.10d)

- $a(\mathbf{u} \cdot \mathbf{v}) = (a\mathbf{u}) \cdot \mathbf{v} = \mathbf{u} \cdot (a\mathbf{v})$ (a form of associativity);
- $(\mathbf{u} + \mathbf{v}) \cdot \mathbf{w} = \mathbf{u} \cdot \mathbf{w} + \mathbf{v} \cdot \mathbf{w}$ (distributivity).

Solution: • For the first:

$$\begin{aligned} a(\mathbf{u} \cdot \mathbf{v}) &= 2((-2, 5, -2) \cdot (3, 3, -2)) \\ &= 2((-2)3 + 5 \cdot 3 + (-2)(-2)) \\ &= 2 \cdot 13 = 26; \\ (a\mathbf{u}) \cdot \mathbf{v} &= (-4, 10, -4) \cdot (3, 3, -2) \\ &= (-4)3 + 10 \cdot 3 + (-4)(-2) \\ &= 26; \\ \mathbf{u} \cdot (a\mathbf{v}) &= (-2, 5, -2) \cdot (6, 6, -4) \\ &= (-2)6 + 5 \cdot 6 + (-2)(-4) \\ &= 26. \end{aligned}$$

These three are equal.

- For the second:

$$\begin{aligned} (\mathbf{u} + \mathbf{v}) \cdot \mathbf{w} &= (1, 8, -4) \cdot (2, 0, -5) \\ &= 1 \cdot 2 + 8 \cdot 0 + (-4)(-5) \\ &= 22; \\ \mathbf{u} \cdot \mathbf{w} + \mathbf{v} \cdot \mathbf{w} &= (-2, 5, -2) \cdot (2, 0, -5) \\ &\quad + (3, 3, -2) \cdot (2, 0, -5) \\ &= [(-2)2 + 5 \cdot 0 + (-2)(-5)] \\ &\quad + [3 \cdot 2 + 3 \cdot 0 + (-2)(-5)] \\ &= 6 + 16 = 22. \end{aligned}$$

These are both equal.

■

Theorem 1.3.10 (dot properties). *Let \mathbf{u} , \mathbf{v} and \mathbf{w} be any vectors in \mathbb{R}^n , and let a be any scalar. The following properties hold:*

- (a) $\mathbf{u} \cdot \mathbf{v} = \mathbf{v} \cdot \mathbf{u}$ (commutative law);
- (b) $\mathbf{u} \cdot \mathbf{0} = \mathbf{0} \cdot \mathbf{u} = 0$;
- (c) $a(\mathbf{u} \cdot \mathbf{v}) = (au) \cdot \mathbf{v} = \mathbf{u} \cdot (av)$;
- (d) $(\mathbf{u} + \mathbf{v}) \cdot \mathbf{w} = \mathbf{u} \cdot \mathbf{w} + \mathbf{v} \cdot \mathbf{w}$ (distributive law);
- (e) $\mathbf{u} \cdot \mathbf{u} \geq 0$, and further $\mathbf{u} \cdot \mathbf{u} = 0$ if and only if $\mathbf{u} = \mathbf{0}$.

Proof. Here prove only the commutative law 1.3.10a and the inequality 1.3.10e. Exercise 1.3.6 asks you to analogously prove the other properties. At the core of each proof is the definition of the dot product which empowers us to deduce a property via the corresponding property for scalars.

- To prove the commutative law 1.3.10a consider

$$\begin{aligned}\mathbf{u} \cdot \mathbf{v} &= u_1 v_1 + u_2 v_2 + \cdots + u_n v_n \quad (\text{by Defn. 1.3.2}) \\ &\quad (\text{using each scalar mult. is commutative}) \\ &= v_1 u_1 + v_2 u_2 + \cdots + v_n u_n \\ &= \mathbf{v} \cdot \mathbf{u} \quad (\text{by Defn. 1.3.2}).\end{aligned}$$

- To prove the inequality 1.3.10e consider

$$\begin{aligned}\mathbf{u} \cdot \mathbf{u} &= u_1 u_1 + u_2 u_2 + \cdots + u_n u_n \quad (\text{by Defn. 1.3.2}) \\ &= u_1^2 + u_2^2 + \cdots + u_n^2 \\ &\geq 0 + 0 + \cdots + 0 \quad (\text{as each scalar term is } \geq 0) \\ &= 0.\end{aligned}$$

To prove the further part, first consider the zero vector. From Definition 1.3.2, in \mathbb{R}^n ,

$$\mathbf{0} \cdot \mathbf{0} = 0^2 + 0^2 + \cdots + 0^2 = 0.$$

Second, let vector $\mathbf{u} = (u_1, u_2, \dots, u_n)$ in \mathbb{R}^n satisfy $\mathbf{u} \cdot \mathbf{u} = 0$. Then we know that

$$\underbrace{u_1^2}_{\geq 0} + \underbrace{u_2^2}_{\geq 0} + \cdots + \underbrace{u_n^2}_{\geq 0} = 0.$$

Being squares, all terms on the left are non-negative, so the only way they can all add to zero is if they are all zero. That is, $u_1 = u_2 = \cdots = u_n = 0$. Hence, the vector \mathbf{u} must be the zero vector $\mathbf{0}$.

□

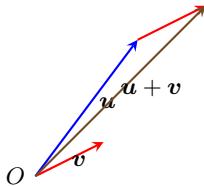
The last part of this proof, that $\mathbf{u} \cdot \mathbf{u} = 0$ if and only if $\mathbf{u} = \mathbf{0}$, may look uncannily familiar. The reason is that this last part is essentially the same as the proof of Theorem 1.1.11 that the zero vector is the only vector of length zero. The next Theorem 1.3.13 establishes that this connection between dot products and lengths is no coincidence.

Example 1.3.11. For the two vectors $\mathbf{u} = (3, 4)$ and $\mathbf{v} = (2, 1)$ verify the following three properties:

- (a) $\sqrt{\mathbf{u} \cdot \mathbf{u}} = |\mathbf{u}|$, the length of \mathbf{u} ;
- (b) $|\mathbf{u} \cdot \mathbf{v}| \leq |\mathbf{u}||\mathbf{v}|$ (Cauchy–Schwarz inequality);
- (c) $|\mathbf{u} + \mathbf{v}| \leq |\mathbf{u}| + |\mathbf{v}|$ (triangle inequality).

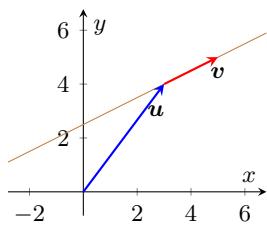
Solution: (a) Here $\sqrt{\mathbf{u} \cdot \mathbf{u}} = \sqrt{3 \cdot 3 + 4 \cdot 4} = \sqrt{25} = 5$, whereas the length $|\mathbf{u}| = \sqrt{3^2 + 4^2} = \sqrt{25} = 5$ (Definition 1.1.8). These expressions are equal.

- (b) Here $|\mathbf{u} \cdot \mathbf{v}| = |3 \cdot 2 + 4 \cdot 1| = 10$, whereas $|\mathbf{u}||\mathbf{v}| = 5\sqrt{2^2 + 1^2} = 5\sqrt{5} = 11.180$. Hence $|\mathbf{u} \cdot \mathbf{v}| = 10 \leq 11.180 = |\mathbf{u}||\mathbf{v}|$.
- (c) Here $|\mathbf{u} + \mathbf{v}| = |(5, 5)| = \sqrt{5^2 + 5^2} = \sqrt{50} = 7.071$, whereas $|\mathbf{u}| + |\mathbf{v}| = 5 + \sqrt{5} = 7.236$. Hence $|\mathbf{u} + \mathbf{v}| = 7.071 \leq 7.236 = |\mathbf{u}| + |\mathbf{v}|$. This is called the triangle inequality because the vectors \mathbf{u} , \mathbf{v} and $\mathbf{u} + \mathbf{v}$ may be viewed as forming a triangle, as illustrated in the margin, and this inequality follows because the length of any side of a triangle must be less than the sum of the other two sides.



The Cauchy–Schwarz inequality is one point of distinction between this ‘vector multiplication’ and scalar multiplication: for scalars $|ab| = |a||b|$, but the dot product of vectors is typically less, $|\mathbf{u} \cdot \mathbf{v}| \leq |\mathbf{u}||\mathbf{v}|$.

Example 1.3.12. The general proof of the Cauchy–Schwarz inequality involves a trick, so let’s introduce the trick using the vectors of Example 1.3.11. Let vectors $\mathbf{u} = (3, 4)$ and $\mathbf{v} = (2, 1)$ and consider the line given parametrically (Definition 1.2.9) as the position vectors $\mathbf{x} = \mathbf{u} + t\mathbf{v} = (3+2t, 4+t)$ for scalar parameter t —illustrated in the margin. The position vector \mathbf{x} of any point on the line has length ℓ (Definition 1.1.8) where



$$\begin{aligned}\ell^2 &= (3+2t)^2 + (4+t)^2 \\ &= 9 + 12t + 4t^2 + 16 + 8t + t^2 \\ &= \underbrace{25}_c + \underbrace{20}_b t + \underbrace{5}_a t^2,\end{aligned}$$

a quadratic polynomial in t . We know that the length $\ell > 0$ (the line does not pass through the origin so no \mathbf{x} is zero). Hence the quadratic in t cannot have any zeros. By the known properties of quadratic equations it follows that the discriminant $b^2 - 4ac < 0$. Indeed it is: here $b^2 - 4ac = 20^2 - 4 \cdot 5 \cdot 25 = 400 - 500 = -100 < 0$. Usefully, here $a = 5 = |\mathbf{v}|^2$, $c = 25 = |\mathbf{u}|^2$ and $b = 20 = 2 \cdot 10 = 2(\mathbf{u} \cdot \mathbf{v})$. So $b^2 - 4ac < 0$, written as $\frac{1}{4}b^2 < ac$, becomes the statement that $\frac{1}{4}[2(\mathbf{u} \cdot \mathbf{v})]^2 = (\mathbf{u} \cdot \mathbf{v})^2 < |\mathbf{v}|^2|\mathbf{u}|^2$. Taking the square-root of both

sides verifies the Cauchy–Schwarz inequality. The proof of the next theorem establishes it in general. \blacksquare

Theorem 1.3.13. *For all vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n the following properties hold:*

- (a) $\sqrt{\mathbf{u} \cdot \mathbf{u}} = |\mathbf{u}|$, the length of \mathbf{u} ;
- (b) $|\mathbf{u} \cdot \mathbf{v}| \leq |\mathbf{u}||\mathbf{v}|$ (Cauchy–Schwarz inequality);
- (c) $|\mathbf{u} \pm \mathbf{v}| \leq |\mathbf{u}| + |\mathbf{v}|$ (triangle inequality).

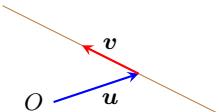
Proof. Each property depends upon the previous.

1.3.13a

$$\begin{aligned}\sqrt{\mathbf{u} \cdot \mathbf{u}} &= \sqrt{u_1 u_1 + u_2 u_2 + \cdots + u_n u_n} \quad (\text{by Defn. 1.3.2}) \\ &= \sqrt{u_1^2 + u_2^2 + \cdots + u_n^2} \\ &= |\mathbf{u}| \quad (\text{by Defn. 1.1.8}).\end{aligned}$$

1.3.13b To prove the Cauchy–Schwarz inequality between vectors \mathbf{u} and \mathbf{v} first consider the trivial case when $\mathbf{v} = \mathbf{0}$: then the left-hand side $|\mathbf{u} \cdot \mathbf{v}| = |\mathbf{u} \cdot \mathbf{0}| = |\mathbf{0}| = 0$; whereas the right-hand side $|\mathbf{u}||\mathbf{v}| = |\mathbf{u}||\mathbf{0}| = |\mathbf{u}|0 = 0$; and so the inequality $|\mathbf{u} \cdot \mathbf{v}| \leq |\mathbf{u}||\mathbf{v}|$ is satisfied in this case.

Second, for the case when $\mathbf{v} \neq \mathbf{0}$, consider the line given parametrically by $\mathbf{x} = \mathbf{u} + t\mathbf{v}$ for (real) scalar parameter t , as illustrated in the margin. The distance ℓ of a point on the line from the origin is the length of its position vector, and by property 1.3.13a

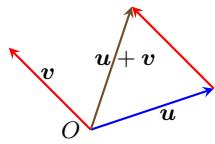


$$\begin{aligned}\ell^2 &= \mathbf{x} \cdot \mathbf{x} \\ &= (\mathbf{u} + t\mathbf{v}) \cdot (\mathbf{u} + t\mathbf{v}) \\ &\quad (\text{then using distributivity 1.3.10d}) \\ &= \mathbf{u} \cdot (\mathbf{u} + t\mathbf{v}) + (t\mathbf{v}) \cdot (\mathbf{u} + t\mathbf{v}) \\ &\quad (\text{again using distributivity 1.3.10d}) \\ &= \mathbf{u} \cdot \mathbf{u} + \mathbf{u} \cdot (t\mathbf{v}) + (t\mathbf{v}) \cdot \mathbf{u} + (t\mathbf{v}) \cdot (t\mathbf{v}) \\ &\quad (\text{using scalar mult. property 1.3.10c}) \\ &= \mathbf{u} \cdot \mathbf{u} + t(\mathbf{u} \cdot \mathbf{v}) + t(\mathbf{v} \cdot \mathbf{u}) + t^2(\mathbf{v} \cdot \mathbf{v}) \\ &\quad (\text{using 1.3.13a and commutativity 1.3.10a}) \\ &= |\mathbf{u}|^2 + 2(\mathbf{u} \cdot \mathbf{v})t + |\mathbf{v}|^2t^2 \\ &= at^2 + bt + c,\end{aligned}$$

a quadratic in t , with coefficients $a = |\mathbf{v}|^2 > 0$, $b = 2(\mathbf{u} \cdot \mathbf{v})$, and $c = |\mathbf{u}|^2$. Since $\ell^2 \geq 0$ (it may be zero if the line goes through the origin), then this quadratic in t has either no zeros or just one zero. By the properties of quadratic equations, the

discriminant $b^2 - 4ac \leq 0$, that is, $\frac{1}{4}b^2 \leq ac$. Substituting the particular coefficients here gives $\frac{1}{4}[2(\mathbf{u} \cdot \mathbf{v})]^2 = (\mathbf{u} \cdot \mathbf{v})^2 \leq |\mathbf{v}|^2|\mathbf{u}|^2$. Taking the square-root of both sides then establishes the Cauchy–Schwarz inequality $|\mathbf{u} \cdot \mathbf{v}| \leq |\mathbf{u}||\mathbf{v}|$.

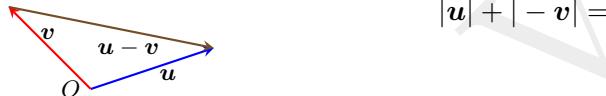
- 1.3.13c** To prove the triangle inequality between vectors \mathbf{u} and \mathbf{v} first observe the Cauchy–Schwarz inequality implies $(\mathbf{u} \cdot \mathbf{v}) \leq |\mathbf{u}||\mathbf{v}|$ since the left-hand side has magnitude \leq the right-hand side. Then consider (analogous to the $t = 1$ case of the above)



$$\begin{aligned}
 |\mathbf{u} + \mathbf{v}|^2 &= (\mathbf{u} + \mathbf{v}) \cdot (\mathbf{u} + \mathbf{v}) \\
 &\quad (\text{then using distributivity 1.3.10d}) \\
 &= \mathbf{u} \cdot (\mathbf{u} + \mathbf{v}) + \mathbf{v} \cdot (\mathbf{u} + \mathbf{v}) \\
 &\quad (\text{again using distributivity 1.3.10d}) \\
 &= \mathbf{u} \cdot \mathbf{u} + \mathbf{u} \cdot \mathbf{v} + \mathbf{v} \cdot \mathbf{u} + \mathbf{v} \cdot \mathbf{v} \\
 &\quad (\text{using 1.3.13a and commutativity 1.3.10a}) \\
 &= |\mathbf{u}|^2 + 2(\mathbf{u} \cdot \mathbf{v}) + |\mathbf{v}|^2 \\
 &\quad (\text{using Cauchy–Schwarz inequality}) \\
 &\leq |\mathbf{u}|^2 + 2|\mathbf{u}||\mathbf{v}| + |\mathbf{v}|^2 \\
 &= (|\mathbf{u}| + |\mathbf{v}|)^2.
 \end{aligned}$$

Take the square-root of both sides to establish the triangle inequality $|\mathbf{u} + \mathbf{v}| \leq |\mathbf{u}| + |\mathbf{v}|$.

The minus case follows because $|\mathbf{u} - \mathbf{v}| = |\mathbf{u} + (-\mathbf{v})| \leq |\mathbf{u}| + |-\mathbf{v}| = |\mathbf{u}| + |\mathbf{v}|$.



□

- Example 1.3.14.** Verify the Cauchy–Schwarz inequality (+ case) and the triangle inequality for the vectors $\mathbf{a} = (-1, -2, 1, 3, -2)$ and $\mathbf{b} = (-3, -2, 10, 2, 2)$.

Solution: We need the length of the vectors:

$$\begin{aligned}
 |\mathbf{a}| &= \sqrt{(-1)^2 + (-2)^2 + 1^2 + 3^2 + (-2)^2} \\
 &= \sqrt{19} = 4.3589, \\
 |\mathbf{b}| &= \sqrt{(-3)^2 + (-2)^2 + 10^2 + 2^2 + 2^2} \\
 &= \sqrt{121} = 11.
 \end{aligned}$$

The dot product

$$\begin{aligned}
 \mathbf{a} \cdot \mathbf{b} &= (-1)(-3) + (-2)(-2) + 1 \cdot 10 + 3 \cdot 2 + (-2)2 \\
 &= 19.
 \end{aligned}$$

Hence $|\mathbf{a} \cdot \mathbf{b}| = 19 < 47.948 = |\mathbf{a}||\mathbf{b}|$, which verifies the Cauchy–Schwarz inequality.

Now, the length of the sum

$$|\mathbf{a} + \mathbf{b}| = |(-4, -4, 11, 5, 0)|$$

$$\begin{aligned}
 &= \sqrt{(-4)^2 + (-4)^2 + 11^2 + 5^2 + 0^2} \\
 &= \sqrt{178} = 13.342.
 \end{aligned}$$

Hence $|\mathbf{a} + \mathbf{b}| = 13.342 < 15.359 = 11 + \sqrt{19} = |\mathbf{a}| + |\mathbf{b}|$, which verifies the triangle inequality. ■

1.3.3 Orthogonal vectors are at right-angles

Of all the angles that vectors can make with each other, the two most important angles are, firstly, when the vectors are aligned with each other, and secondly, when the vectors are at right-angles to each other. Recall Theorem 1.3.4 gives the angle θ between two vectors via $\cos \theta = \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}||\mathbf{v}|}$. For vectors at right-angles $\theta = 90^\circ$ and so $\cos \theta = 0$ and hence non-zero vectors are at right-angles only when the dot product $\mathbf{u} \cdot \mathbf{v} = 0$. We give a special name to vectors at right-angles.

Definition 1.3.15. Two vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n are termed **orthogonal** (or **perpendicular**) if and only if their dot product $\mathbf{u} \cdot \mathbf{v} = 0$. The term ‘orthogonal’ derives from the Greek for ‘right-angled’.

By convention the zero vector $\mathbf{0}$ is orthogonal to all other vectors. However, in practice, we almost always use the notion of orthogonality only in connection with *non-zero* vectors. Often the requirement that the orthogonal vectors are non-zero is explicitly made, but beware that sometimes the requirement may be implicit in the problem.

Example 1.3.16. The standard unit vectors (Definition 1.2.5) are orthogonal to each other. For example, consider the standard unit vectors \mathbf{i} , \mathbf{j} and \mathbf{k} in \mathbb{R}^3 :

- $\mathbf{i} \cdot \mathbf{j} = (1, 0, 0) \cdot (0, 1, 0) = 0 + 0 + 0 = 0$;
- $\mathbf{j} \cdot \mathbf{k} = (0, 1, 0) \cdot (0, 0, 1) = 0 + 0 + 0 = 0$;
- $\mathbf{k} \cdot \mathbf{i} = (0, 0, 1) \cdot (1, 0, 0) = 0 + 0 + 0 = 0$.

By Definition 1.3.15 these are orthogonal to each other. ■

Example 1.3.17. Which pairs of the following vectors, if any, are perpendicular to each other? $\mathbf{u} = (-1, 1, -3, 0)$, $\mathbf{v} = (2, 4, 2, -6)$ and $\mathbf{w} = (-1, 6, -2, 3)$.

Solution: Use the dot product.

- $\mathbf{u} \cdot \mathbf{v} = (-1, 1, -3, 0) \cdot (2, 4, 2, -6) = -2 + 4 - 6 + 0 = -4 \neq 0$ so this pair are not perpendicular.
- $\mathbf{u} \cdot \mathbf{w} = (-1, 1, -3, 0) \cdot (-1, 6, -2, 3) = 1 + 6 + 6 + 0 = 13 \neq 0$ so this pair are not perpendicular.

- $\mathbf{v} \cdot \mathbf{w} = (2, 4, 2, -6) \cdot (-1, 6, -2, 3) = -2 + 24 - 4 - 18 = 0$ so this pair are the only two vectors perpendicular to each other.

■

Example 1.3.18. Find the number b such that vectors $\mathbf{a} = \mathbf{i} + 4\mathbf{j} + 2\mathbf{k}$ and $\mathbf{b} = \mathbf{i} + b\mathbf{j} - 3\mathbf{k}$ are at right-angles.

Solution: For vectors to be at right-angles, their dot product must be zero. Hence find b such that

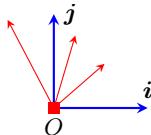
$$0 = \mathbf{a} \cdot \mathbf{b} = (\mathbf{i} + 4\mathbf{j} + 2\mathbf{k}) \cdot (\mathbf{i} + b\mathbf{j} - 3\mathbf{k}) = 1 + 4b - 6 = 4b - 5.$$

Solving gives $b = 5/4$. That is, $\mathbf{i} + \frac{5}{4}\mathbf{j} - 3\mathbf{k}$ is at right-angles to $\mathbf{i} + 4\mathbf{j} + 2\mathbf{k}$.

■

Key properties The next couple of innocuous looking theorems are vital keys to important results in subsequent chapters.

To introduce the first theorem, consider the 2D plane and try to draw a non-zero vector at right-angles to both the two standard unit vectors \mathbf{i} and \mathbf{j} . The red vectors in the margin illustrate three failed attempts to draw a vector at right-angles to both \mathbf{i} and \mathbf{j} . It cannot be done. No vector in the plane can be at right angles to both the standard unit vectors in the plane.

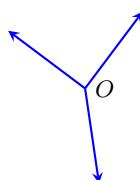


Theorem 1.3.19. *There is no non-zero vector orthogonal to all n standard unit vectors in \mathbb{R}^n .*

Proof. Let $\mathbf{u} = (u_1, u_2, \dots, u_n)$ be a vector in \mathbb{R}^n that is orthogonal to all n standard unit vectors. Then by Definition 1.3.15 of orthogonality:

- $0 = \mathbf{u} \cdot \mathbf{e}_1 = (u_1, u_2, \dots, u_n) \cdot (1, 0, \dots, 0) = u_1 + 0 + \dots + 0 = u_1$, and so the first component must be zero;
- $0 = \mathbf{u} \cdot \mathbf{e}_2 = (u_1, u_2, \dots, u_n) \cdot (0, 1, \dots, 0) = 0 + u_2 + 0 + \dots + 0 = u_2$, and so the second component must be zero; and so on to
- $0 = \mathbf{u} \cdot \mathbf{e}_n = (u_1, u_2, \dots, u_n) \cdot (0, 0, \dots, 1) = 0 + 0 + \dots + u_n = u_n$, and so the last component must be zero.

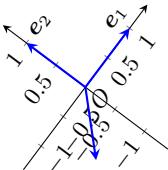
Since $u_1 = u_2 = \dots = u_n = 0$ the only vector that is orthogonal to all the standard unit vectors is $\mathbf{u} = \mathbf{0}$, the zero vector. □



To introduce the second theorem, imagine trying to draw three unit vectors in any orientation in the 2D plane such that all three are at right-angles to each other. The margin illustrates one attempt. It cannot be done. There are at most two vectors in 2D that are all at right-angles to each other.

Theorem 1.3.20 (orthogonal completeness). *In a set of orthogonal unit vectors in \mathbb{R}^n , there can be no more than n vectors in the set.*⁴

Proof. Use contradiction. Suppose there are more than n orthogonal unit vectors in the set. Define a coordinate system for \mathbb{R}^n using the first n unit vectors as the n standard unit vectors (as illustrated in the margin). Theorem 1.3.19 then says there cannot be any more non-zero vectors orthogonal than these n standard unit vectors. This contradicts there being more than n orthogonal unit vectors. To avoid this contradiction the supposition must be wrong; that is, there cannot be more than n orthogonal unit vectors in \mathbb{R}^n . \square



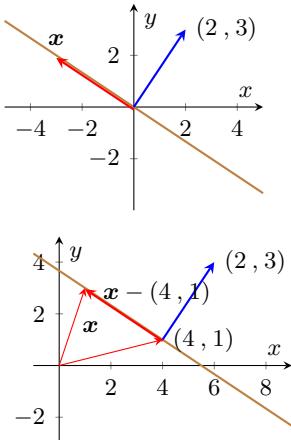
1.3.4 Normal vectors and equations of a plane

This section uses the dot product to find equations of a plane in 3D. The key is to write points in the plane as those at right-angles to the direction perpendicular to the plane, called a normal. Let's start with an example of the idea in 2D.

Example 1.3.21. First find the equation of the line that is perpendicular to the vector $(2, 3)$ and that passes through the origin. Second, find is the equation of the line that passes through the point $(4, 1)$, instead of the origin?

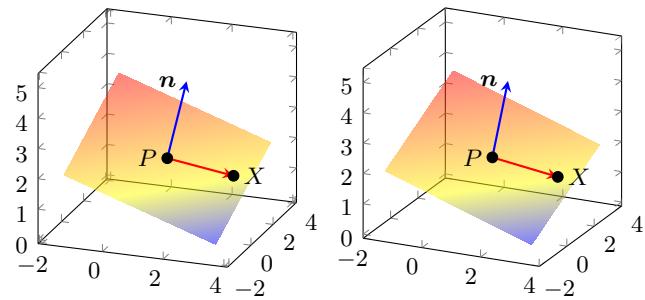
Solution: Recall that vectors at right-angles have a zero dot product (section 1.3.3). Thus the position vector \mathbf{x} of any point in the line satisfies the dot product $\mathbf{x} \cdot (2, 3) = 0$. For $\mathbf{x} = (x, y)$, as illustrated in the margin, $\mathbf{x} \cdot (2, 3) = 2x + 3y$ so the equation of the line is $2x + 3y = 0$.

When the line goes through $(4, 1)$, then it is the displacement $\mathbf{x} - (4, 1)$ that must be orthogonal to $(2, 3)$, as illustrated. That is, the equation of the line is $(\mathbf{x} - (4, 1)) \cdot (2, 3) = 0$. Evaluating the dot product gives $2(x-4) + 3(y-1) = 0$; that is, $2x + 3y = 2 \cdot 4 + 3 \cdot 1 = 11$ is an equation of the line. \blacksquare



Now use the same approach to finding an equation of a plane in 3D. The problem is to find the equation of the plane that goes through a given point P and is perpendicular to a given vector \mathbf{n} , called a **normal vector**. As illustrated below, that means to find all points X such that \overrightarrow{PX} is orthogonal to \mathbf{n} .

⁴ For the pure at heart, this property forms part of the definition of what we mean by \mathbb{R}^n . The representation of a vector in \mathbb{R}^n by n components (here Definition 1.1.4) then follows as a consequence, instead of vice-versa.

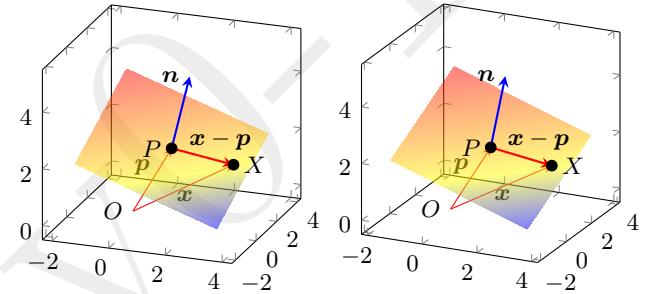


Denote the position vector of P by $\mathbf{p} = (x_0, y_0, z_0)$, the position vector of X by $\mathbf{x} = (x, y, z)$, and let the normal vector be $\mathbf{n} = (a, b, c)$. Then, as drawn below, the displacement vector $\overrightarrow{PX} = \mathbf{x} - \mathbf{p} = (x - x_0, y - y_0, z - z_0)$ and so for \overrightarrow{PX} to be orthogonal to \mathbf{n} requires $\mathbf{n} \cdot (\mathbf{x} - \mathbf{p}) = 0$; that is, an **equation of the plane** is

$$a(x - x_0) + b(y - y_0) + c(z - z_0) = 0,$$

equivalently

$$ax + by + cz = d \quad \text{for constant } d = ax_0 + by_0 + cz_0.$$



Example 1.3.22. Find an equation of the plane through point $P = (1, 1, 2)$ that has normal vector $\mathbf{n} = (1, -1, 3)$. (This is the case drawn in the above illustrations.) Hence write down three distinct points on the plane.

Solution: Letting $\mathbf{x} = (x, y, z)$ be the coordinates of a point in the plane, the above argument asserts an equation of the plane is $\mathbf{n} \cdot (\mathbf{x} - \overrightarrow{OP}) = 0$ which becomes $1(x - 1) - 1(y - 1) + 3(z - 2) = 0$; that is, $x - 1 - y + 1 + 3z - 6 = 0$, which rearranged is $x - y + 3z = 6$.

To find some points in the plane, rearrange this equation to $z = 2 - x/3 + y/3$ and then substitute any values for x and y : $x = y = 0$ gives $z = 2$ so $(0, 0, 2)$ is on the plane; $x = 3$ and $y = 0$ gives $z = 1$ so $(3, 0, 1)$ is on the plane; $x = 2$ and $y = -2$ gives $z = 2/3$ so $(2, -2, 2/3)$ is on the plane; and so on. ■

Example 1.3.23. Write down a normal vector to each of the following planes:

- (a) $3x - 6y + z = 4$; (b) $z = 0.2x - 3.3y - 1.9$.

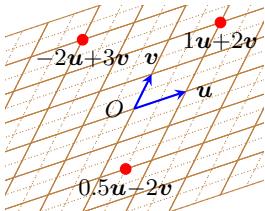
Solution:

- (a) In this standard form $3x - 6y + 2z = 4$ a normal vector is the coefficients of the variables, $\mathbf{n} = (3, -6, 2)$ (or any scalar multiple).
- (b) Rearrange $z = 0.2x - 3.3y - 1.9$ to standard form $-0.2x + 3.3y + z = -1.9$ then a normal is $\mathbf{n} = (-0.2, 3.3, 1)$ (or any multiple).

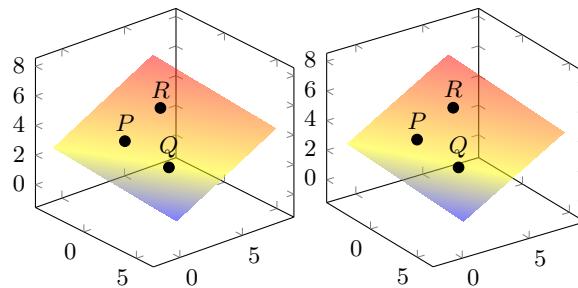
■

Parametric equation of a plane An alternative way of describing a plane is via a parametric equation analogous to the parametric equation of a line (section 1.2.2). Such a parametric representation generalises to any dimensions (Section 2.3).

The basic idea, as illustrated in the margin, is that given any plane (through the origin for the moment), then choosing almost any two vectors in the plane allows us to write all points in the plane as a sum of multiples of the two vectors. With the given vectors \mathbf{u} and \mathbf{v} shown in the margin, illustrated are the points $\mathbf{u} + 2\mathbf{v}$, $\frac{1}{2}\mathbf{u} - 2\mathbf{v}$ and $-2\mathbf{u} + 3\mathbf{v}$. Similarly, all points in the plane have a position vector in the form $s\mathbf{u} + t\mathbf{v}$ for some scalar parameters s and t . The grid shown in the margin illustrates the sum of integral and half-integral multiples. The formula $\mathbf{x} = s\mathbf{u} + t\mathbf{v}$ for parameters s and t is called a parametric equation of the plane.



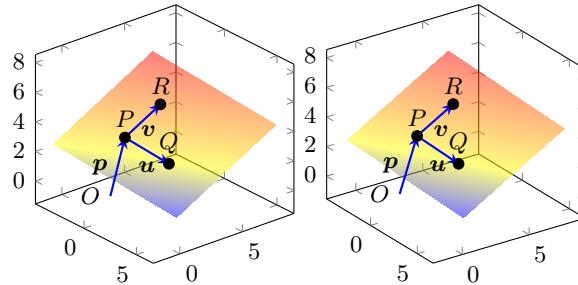
Example 1.3.24. Find a parametric equation of the plane that passes through the three points $P = (-1, 2, 3)$, $Q = (2, 3, 2)$ and $R = (0, 4, 5)$, drawn below in stereo.



Solution: This plane does not pass through the origin, so we first choose a point and make the description relative to that point: say choose point P with position vector $\mathbf{p} = \overrightarrow{OP} = -\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}$. Then, as illustrated below, two vectors parallel to the required plane are

$$\begin{aligned}\mathbf{u} &= \overrightarrow{PQ} = \overrightarrow{OQ} - \overrightarrow{OP} \\ &= (2\mathbf{i} + 3\mathbf{j} + 2\mathbf{k}) - (-\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}) \\ &= 3\mathbf{i} + \mathbf{j} - \mathbf{k}, \\ \mathbf{v} &= \overrightarrow{PR} = \overrightarrow{OR} - \overrightarrow{OP}\end{aligned}$$

$$\begin{aligned}
 &= (4\mathbf{j} + 5\mathbf{k}) - (-\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}) \\
 &= \mathbf{i} + 2\mathbf{j} + 2\mathbf{k}.
 \end{aligned}$$



Lastly, every point in the plane is the sum of the displacement vector \mathbf{p} and arbitrary multiples of the parallel vectors \mathbf{u} and \mathbf{v} . That is, a parametric equation of the plane is $\mathbf{x} = \mathbf{p} + s\mathbf{u} + t\mathbf{v}$ which here is

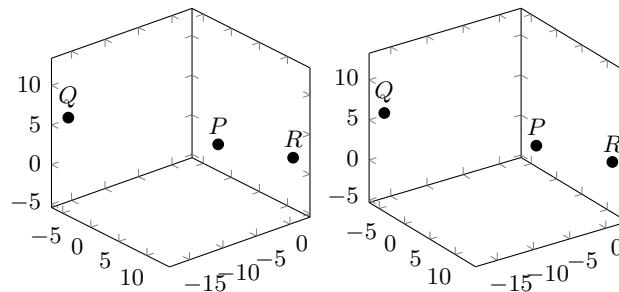
$$\begin{aligned}
 \mathbf{x} &= (-\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}) + s(3\mathbf{i} + \mathbf{j} - \mathbf{k}) + t(\mathbf{i} + 2\mathbf{j} + 2\mathbf{k}) \\
 &= (-1 + 3s + t)\mathbf{i} + (2 + s + 2t)\mathbf{j} + (3 - s + 2t)\mathbf{k}.
 \end{aligned}$$

■

Definition 1.3.25. A *parametric equation* of a plane is $\mathbf{x} = \mathbf{p} + s\mathbf{u} + t\mathbf{v}$ where \mathbf{p} is the position vector of some point in the plane, the two vectors \mathbf{u} and \mathbf{v} are parallel to the plane ($\mathbf{u}, \mathbf{v} \neq \mathbf{0}$ and are at an angle to each other), and the scalar **parameters** s and t vary over all real values to give position vectors of all points in the plane.

The beauty of this definition is that it applies for planes in any number of dimensions by using vectors with the corresponding number of components.

Example 1.3.26. Find a parametric equation of the plane that passes through the three points $P = (6, -4, 3)$, $Q = (-4, -18, 7)$ and $R = (11, 3, 1)$, drawn below in stereo.



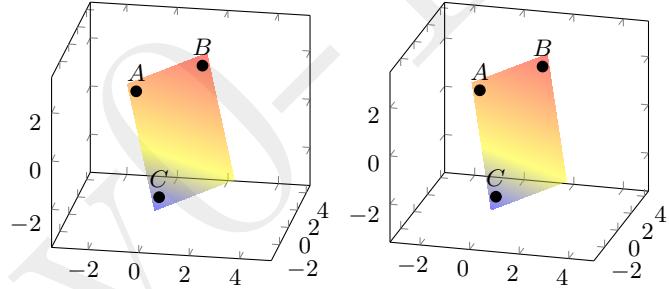
Solution: First choose a point and make the description relative to that point: say choose point P with position vector $\mathbf{p} = \overrightarrow{OP} = 6\mathbf{i} - 4\mathbf{j} + 3\mathbf{k}$. Then, as illustrated below, two vectors parallel to the required plane are

$$\mathbf{u} = \overrightarrow{PQ} = \overrightarrow{OQ} - \overrightarrow{OP}$$

$$\begin{aligned}
&= (-4\mathbf{i} - 18\mathbf{j} + 7\mathbf{k}) - (6\mathbf{i} - 4\mathbf{j} + 3\mathbf{k}) \\
&= -10\mathbf{i} - 14\mathbf{j} + 4\mathbf{k}, \\
\mathbf{v} &= \overrightarrow{PR} = \overrightarrow{OR} - \overrightarrow{OP} \\
&= (11\mathbf{i} + 3\mathbf{j} + \mathbf{k}) - (6\mathbf{i} - 4\mathbf{j} + 3\mathbf{k}) \\
&= 5\mathbf{i} + 7\mathbf{j} - 2\mathbf{k}.
\end{aligned}$$

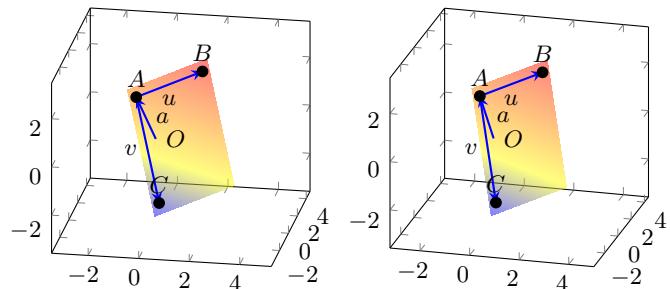
Oops: notice that $\mathbf{u} = -2\mathbf{v}$ so the vectors \mathbf{u} and \mathbf{v} are not at a nontrivial angle; instead they are aligned along a line because the three points P , Q and R are collinear. There are an infinite number of planes passing through such collinear points. Hence we cannot answer the question which requires “the plane”. ■

Example 1.3.27. Find a parametric equation of the plane that passes through the three points $A = (-1.2, 2.4, 0.8)$, $B = (1.6, 1.4, 2.4)$ and $C = (0.2, -0.4, -2.5)$, drawn below in stereo.



Solution: First choose a point and make the description relative to that point: say choose point A with position vector $\mathbf{a} = \overrightarrow{OA} = -1.2\mathbf{i} + 2.4\mathbf{j} + 0.8\mathbf{k}$. Then, as illustrated below, two vectors parallel to the required plane are

$$\begin{aligned}
\mathbf{u} &= \overrightarrow{AB} = \overrightarrow{OB} - \overrightarrow{OA} \\
&= (1.6\mathbf{i} + 1.4\mathbf{j} + 2.4\mathbf{k}) - (-1.2\mathbf{i} + 2.4\mathbf{j} + 0.8\mathbf{k}) \\
&= 2.8\mathbf{i} - \mathbf{j} + 1.6\mathbf{k}, \\
\mathbf{v} &= \overrightarrow{AC} = \overrightarrow{OC} - \overrightarrow{OA} \\
&= (0.2\mathbf{i} - 0.4\mathbf{j} - 2.5\mathbf{k}) - (-1.2\mathbf{i} + 2.4\mathbf{j} + 0.8\mathbf{k}) \\
&= 1.4\mathbf{i} - 2.8\mathbf{j} - 3.3\mathbf{k}.
\end{aligned}$$



Lastly, every point in the plane is the sum of the displacement vector \mathbf{a} , and arbitrary multiples of the parallel vectors \mathbf{u} and \mathbf{v} . That is, a parametric equation of the plane is $\mathbf{x} = \mathbf{a} + s\mathbf{u} + t\mathbf{v}$ which here is

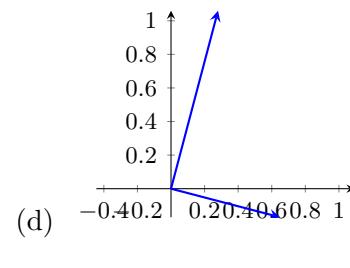
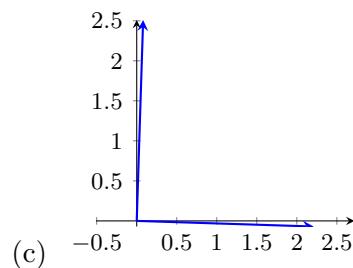
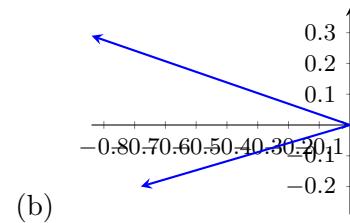
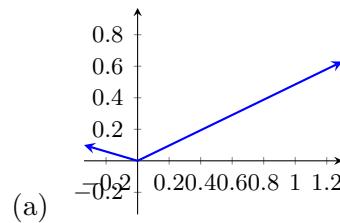
$$\begin{aligned}\mathbf{x} &= \begin{bmatrix} -1.2 \\ 2.4 \\ 0.8 \end{bmatrix} + s \begin{bmatrix} 2.8 \\ -1 \\ 1.6 \end{bmatrix} + t \begin{bmatrix} 1.4 \\ -2.8 \\ -3.3 \end{bmatrix} \\ &= \begin{bmatrix} -1.2 + 2.8s + 1.4t \\ 2.4 - s - 2.8t \\ 0.8 + 1.6s - 3.3t \end{bmatrix}.\end{aligned}$$

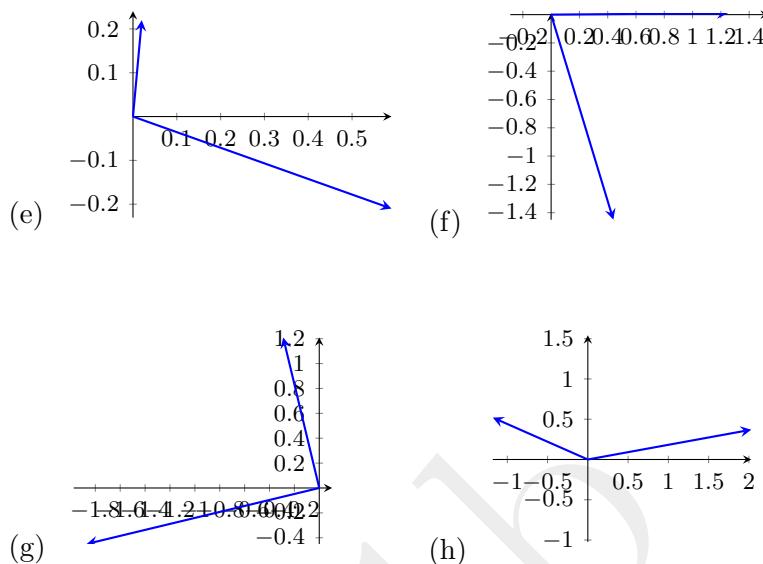
1.3.5 Exercises

Exercise 1.3.1. Following Example 1.3.1, use the cosine rule for triangles to find the angle between the following pairs of vectors. Confirm that $|\mathbf{u}||\mathbf{v}| \cos \theta = \mathbf{u} \cdot \mathbf{v}$ in each case.

- | | |
|------------------------------|--|
| (a) (6, 5) and (-3, 1) | (b) (6, 2, 2) and (-1, -2, 5) |
| (c) (2, 2.9) and (-1.4, 0.8) | (d) (-3.6, 0, -0.7) and
(1.2, -0.9, -0.6) |

Exercise 1.3.2. Which of the following pairs of vectors appear orthogonal?





Exercise 1.3.3. Recall that Example 1.1.7 represented the following sentences by word vectors $\mathbf{w} = (N_{\text{cat}}, N_{\text{dog}}, N_{\text{mat}}, N_{\text{sat}}, N_{\text{scratched}})$.

- “The cat and dog sat on the mat” is summarised by the vector $\mathbf{a} = (1, 1, 1, 1, 0)$.
 - “The dog scratched” is summarised by the vector $\mathbf{b} = (0, 1, 0, 0, 1)$.
 - “The dog sat on the mat; the cat scratched the dog.” is summarised by the vector $\mathbf{c} = (1, 2, 1, 1, 1)$.

Find the similarity between pairs of these sentences by calculating the angle between each pair of word vectors. What is the most similar pair of sentences?

Exercise 1.3.4. Recall Exercise 1.1.4 found word vectors in \mathbb{R}^7 for the titles of eight books that The Society of Industrial and Applied Mathematics (SIAM) reviewed recently. The following four titles have more than one word counted in the word vectors.

- (a) Introduction to Finite and Spectral Element Methods using MATLAB
 - (b) Iterative Methods for Linear Systems: Theory and Applications
 - (c) Singular Perturbations: Introduction to System Order Reduction Methods with Applications
 - (d) Stochastic Chemical Kinetics: Theory and Mostly Systems Biology Applications

Find the similarity between pairs of these titles by calculating the angle between each pair of corresponding word vectors in \mathbb{R}^7 . What is the most similar pair of titles? What is the most dissimilar titles?

Exercise 1.3.5. Suppose two non-zero word vectors are orthogonal. Explain what such orthogonality means in terms of the words of the original sentences.

Exercise 1.3.6. For the properties of the dot product, Theorem 1.3.10, prove some properties chosen from 1.3.10b–1.3.10d.

Exercise 1.3.7. Verify the Cauchy–Schwarz inequality (+ case) and also the triangle inequality for the following pairs of vectors.

$$(a) (2, -4, 4) \text{ and } (6, 7, 6) \quad (b) (1, -2, 2) \text{ and } (-3, 6, -6)$$

$$(c) (-2, -3, 6) \text{ and } (3, 1, 2) \quad (d) (3, -5, -1, -1) \text{ and } (1, -1, -1, -1)$$

$$(e) \begin{bmatrix} -0.2 \\ 0.8 \\ -3.8 \\ -0.3 \end{bmatrix} \text{ and } \begin{bmatrix} 2.4 \\ -5.2 \\ 5.0 \\ 1.9 \end{bmatrix} \quad (f) \begin{bmatrix} 0.8 \\ 0.8 \\ 6.6 \\ -1.5 \end{bmatrix} \text{ and } \begin{bmatrix} 4.4 \\ -0.6 \\ 2.1 \\ 2.2 \end{bmatrix}$$

Exercise 1.3.8. Find an equation of the plane with the given normal vector \mathbf{n} and through the given point P .

$$(a) P = (1, 2, -3), \quad (b) P = (5, -4, -13), \\ \mathbf{n} = (2, -5, -2). \quad \mathbf{n} = (-1, 0, -1).$$

$$(c) P = (10, -4, -1), \quad (d) P = (2, -5, -1), \\ \mathbf{n} = (-2, 4, 5). \quad \mathbf{n} = (4, 9, -4).$$

$$(e) P = (1.7, -4.2, 2.2), \quad (f) P = (3, 5, -2), \\ \mathbf{n} = (1, 0, 4). \quad \mathbf{n} = (-2.5, -0.5, 0.4).$$

$$(g) P = (-7.3, -1.6, 5.8), \quad (h) P = (0, -1.2, 2.2), \\ \mathbf{n} = (-2.8, -0.8, 4.4). \quad \mathbf{n} = (-1.4, -8.1, -1.5).$$

Exercise 1.3.9. Write down a normal vector to the plane described by each of the following equations.

$$(a) 2x + 3y + 2z = 6 \quad (b) -7x - 2y + 4 = -5z$$

$$(c) -12x_1 + 2x_2 + 2x_3 - 8 = 0 \quad (d) 2x_3 = 8x_1 + 5x_2 + 1$$

$$(e) 0.1x = 1.5y + 1.1z + 0.7 \quad (f) \\ -5.5x_1 + 1.6x_2 = 6.7x_3 - 1.3$$

Exercise 1.3.10. For each case, find a parametric equation of the plane through the three given points.

$$(a) (0, 5, -4), (-3, -2, 2), \quad (b) (0, -1, -1), (-4, 1, -5), \\ (5, 1, -3). \quad (0, -3, -2).$$

$$(c) (2, 2, 3), (2, 3, 3), \quad (d) (-1, 2, 2), (0, 1, -1), \\ (3, 1, 0). \quad (1, 0, -4).$$

$$\begin{array}{l} \text{(e)} \quad (0.4, -2.2, 8.7), \\ (-2.2, 1.3, -4.9), \\ (-1.4, 3.2, -0.4). \end{array}$$

$$\begin{array}{l} \text{(f)} \quad (2.2, -6.7, 2), \\ (-2.6, -1.6, -0.5), \\ (2.9, 5.4, -0.6). \end{array}$$

$$\begin{array}{l} \text{(g)} \quad (-5.6, -2.2, -6.8), \\ (-1.8, 4.3, -3.9), \\ (2.5, -3.5, -1.7), \end{array}$$

$$\begin{array}{l} \text{(h)} \quad (1.8, -0.2, -0.7), \\ (-1.6, 2, -3.7), \\ (1.4, -0.5, 0.5), \end{array}$$

Exercise 1.3.11. For each case of Exercise 1.3.10 that you have done, find two other parametric equations of the plane.

V0-1D

1.4 The cross product

Section Contents

Area of a parallelogram	57
Normal vector to a plane	58
Definition of a cross product	59
Geometry of a cross product	60
Algebraic properties of a cross product	63
Volume of a parallelepiped	65
1.4.1 Exercises	67

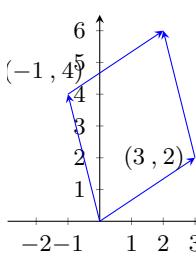
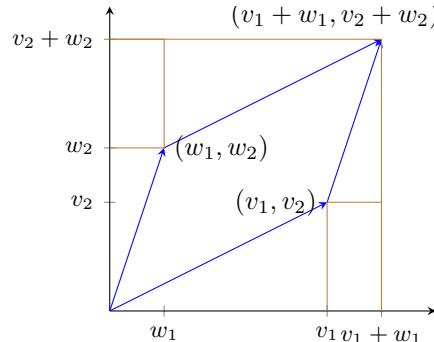
This section is optional, but is vital in many topics of science and engineering.

In the three dimensions of the world we live in, as well as the dot product there is another way to multiply vectors, called the cross product. For more than three dimensions, qualitatively different techniques are developed in subsequent chapters.

Area of a parallelogram

Consider the parallelogram drawn in blue. It has sides given by vectors $\mathbf{v} = (v_1, v_2)$ and $\mathbf{w} = (w_1, w_2)$ as shown. Let's determine the area of the parallelogram by that of the containing rectangle less the two small rectangles and the four small triangles. The two small rectangles have the same area, namely $w_1 v_2$. The two small triangles on the left and the right also have the same area, namely $\frac{1}{2} w_1 w_2$. The two small triangles on the top and the bottom similarly have the same area, namely $\frac{1}{2} v_1 v_2$. Thus, the parallelogram has

$$\begin{aligned} \text{area} &= (v_1 + w_1)(v_2 + w_2) - 2w_1 v_2 - 2 \cdot \frac{1}{2} w_1 w_2 - 2 \cdot \frac{1}{2} v_1 v_2 \\ &= v_1 v_2 + v_1 w_2 + w_1 v_2 + w_1 w_2 - 2w_1 v_2 - w_1 w_2 - v_1 v_2 \\ &= v_1 w_2 - v_2 w_1. \end{aligned}$$



In application, sometimes this right-hand side expression is negative because vectors \mathbf{v} and \mathbf{w} are the ‘wrong way’ around. Thus in general the parallelogram area = $|v_1 w_2 - v_2 w_1|$.

Example 1.4.1. What is the area of the parallelogram (illustrated in the margin) whose edges are formed by the vectors $(3, 2)$ and $(-1, 4)$?

Solution: The parallelogram area = $|3 \cdot 4 - 2 \cdot (-1)| = |12 + 2| = 14$. The illustration indicates that this area must be about right as with imagination one could cut the area and move it about to form a rectangle roughly 3 by 5, and hence the area should be roughly 15.

■

Interestingly, we meet this expression for area, $v_1w_2 - v_2w_1$, in another context: that of equations for a plane and its normal vector.

Normal vector to a plane

Recall section 1.3.4 introduced that we describe planes either via an equation such as $x - y + 3z = 6$ or via a parametric description such as $\mathbf{x} = (1, 1, 2) + (1, 1, 0)s + (0, 3, 1)t$. These determine the same plane, just different algebraic descriptions. One converts between these two descriptions using the cross product.

Example 1.4.2. Derive that the plane described parametrically by $\mathbf{x} = (1, 1, 2) + (1, 1, 0)s + (0, 3, 1)t$ has normal equation $x - y + 3z = 6$.

Solution: The key to deriving the normal equation is to find that a normal vector to the plane is $(1, -1, 3)$. This normal vector comes from the two vectors that multiply the parameters in the parametric form, $(1, 1, 0)$ and $(0, 3, 1)$. (Those who have seen 3×3 determinants will recognise the following has the same pattern—see Chapter 6.) Write the vectors as two consecutive columns, following a first column of the *symbols* of the standard unit vectors \mathbf{i} , \mathbf{j} and \mathbf{k} , in

$$\mathbf{n} = \begin{vmatrix} \mathbf{i} & 1 & 0 \\ \mathbf{j} & 1 & 3 \\ \mathbf{k} & 0 & 1 \end{vmatrix}$$

(cross out 1st column and each row, multiplying each by common entry, with alternating sign)

$$\begin{aligned} &= \mathbf{i} \begin{vmatrix} 1 & 0 \\ 1 & 3 \end{vmatrix} - \mathbf{j} \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} + \mathbf{k} \begin{vmatrix} 1 & 0 \\ 1 & 3 \end{vmatrix} \\ &= \mathbf{i} \begin{vmatrix} 1 & 3 \\ 0 & 1 \end{vmatrix} - \mathbf{j} \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} + \mathbf{k} \begin{vmatrix} 1 & 0 \\ 1 & 3 \end{vmatrix} \end{aligned}$$

(draw diagonals, then subtract product of red diagonal from product of the blue)

$$\begin{aligned} &= \mathbf{i} \begin{vmatrix} 1 & 3 \\ 0 & 1 \end{vmatrix} - \mathbf{j} \begin{vmatrix} 1 & 0 \\ 0 & 1 \end{vmatrix} + \mathbf{k} \begin{vmatrix} 1 & 0 \\ 1 & 3 \end{vmatrix} \\ &= \mathbf{i}(1 \cdot 1 - 0 \cdot 3) - \mathbf{j}(1 \cdot 1 - 0 \cdot 0) + \mathbf{k}(1 \cdot 3 - 1 \cdot 0) \\ &= \mathbf{i} - \mathbf{j} + 3\mathbf{k}. \end{aligned}$$

Using this normal vector, the equation of the plane must be of the form $x - y + 3z = \text{constant}$. Since the plane goes through

point $(1, 1, 2)$, the constant $= 1 - 1 + 3 \cdot 2 = 6$; that is, the plane is $x - y + 3z = 6$ (as given). ■

Definition of a cross product

General formula The procedure used in Example 1.4.2 to derive a normal vector leads to an algebraic formula. Let's apply the same procedure to two general vectors $\mathbf{v} = (v_1, v_2, v_3)$ and $\mathbf{w} = (w_1, w_2, w_3)$. The procedure computes

$$\mathbf{n} = \begin{vmatrix} \mathbf{i} & v_1 & w_1 \\ \mathbf{j} & v_2 & w_2 \\ \mathbf{k} & v_3 & w_3 \end{vmatrix}$$

(cross out 1st column and each row, multiplying each by common entry, with alternating sign)

$$= \mathbf{i} \begin{vmatrix} \mathbf{i} & v_1 & w_1 \\ \mathbf{j} & v_2 & w_2 \\ \mathbf{k} & v_3 & w_3 \end{vmatrix} - \mathbf{j} \begin{vmatrix} \mathbf{i} & v_1 & w_1 \\ \mathbf{j} & v_2 & w_2 \\ \mathbf{k} & v_3 & w_3 \end{vmatrix} + \mathbf{k} \begin{vmatrix} \mathbf{i} & v_1 & w_1 \\ \mathbf{j} & v_2 & w_2 \\ \mathbf{k} & v_3 & w_3 \end{vmatrix}$$

$$= \mathbf{i} \begin{vmatrix} v_2 & w_2 \\ v_3 & w_3 \end{vmatrix} - \mathbf{j} \begin{vmatrix} v_1 & w_1 \\ v_3 & w_3 \end{vmatrix} + \mathbf{k} \begin{vmatrix} v_1 & w_1 \\ v_2 & w_2 \end{vmatrix}$$

(draw diagonals, then subtract product of red diagonal from product of the blue)

$$= \mathbf{i} \begin{vmatrix} v_2 & w_2 \\ v_3 & w_3 \end{vmatrix} - \mathbf{j} \begin{vmatrix} v_1 & w_1 \\ v_3 & w_3 \end{vmatrix} + \mathbf{k} \begin{vmatrix} v_1 & w_1 \\ v_2 & w_2 \end{vmatrix}$$

$$= \mathbf{i}(v_2w_3 - v_3w_2) - \mathbf{j}(v_1w_3 - v_3w_1) + \mathbf{k}(v_1w_2 - v_2w_1).$$

We use this formula to define the cross product algebraically, and then see what it means geometrically.

Definition 1.4.3. Let $\mathbf{v} = (v_1, v_2, v_3)$ and $\mathbf{w} = (w_1, w_2, w_3)$ be any two vectors in \mathbb{R}^3 . The **cross product** (or **vector product**) $\mathbf{v} \times \mathbf{w}$ is defined algebraically as

$$\mathbf{v} \times \mathbf{w} := \mathbf{i}(v_2w_3 - v_3w_2) + \mathbf{j}(v_1w_3 - v_3w_1) + \mathbf{k}(v_1w_2 - v_2w_1).$$

Example 1.4.4. Among the standard unit vectors, derive that

- | | |
|---|--|
| (a) $\mathbf{i} \times \mathbf{j} = \mathbf{k}$, | (b) $\mathbf{j} \times \mathbf{i} = -\mathbf{k}$, |
| (c) $\mathbf{j} \times \mathbf{k} = \mathbf{i}$, | (d) $\mathbf{k} \times \mathbf{j} = -\mathbf{i}$, |
| (e) $\mathbf{k} \times \mathbf{i} = \mathbf{j}$, | (f) $\mathbf{i} \times \mathbf{k} = -\mathbf{j}$, |
| (g) $\mathbf{i} \times \mathbf{i} = \mathbf{j} \times \mathbf{j} = \mathbf{k} \times \mathbf{k} = \mathbf{0}$. | |

Solution: Using Definition 1.4.3:

$$\begin{aligned} \mathbf{i} \times \mathbf{j} &= (1, 0, 0) \times (0, 1, 0) \\ &= \mathbf{i}(0 \cdot 0 - 0 \cdot 1) + \mathbf{j}(0 \cdot 0 - 1 \cdot 0) + \mathbf{k}(1 \cdot 1 - 0 \cdot 0) \end{aligned}$$

$$\begin{aligned}
&= \mathbf{k}; \\
\mathbf{j} \times \mathbf{i} &= (0, 1, 0) \times (1, 0, 0) \\
&= \mathbf{i}(1 \cdot 0 - 0 \cdot 0) + \mathbf{j}(0 \cdot 1 - 0 \cdot 0) + \mathbf{k}(1 \cdot 1 - 1 \cdot 0) \\
&= -\mathbf{k}; \\
\mathbf{i} \times \mathbf{i} &= (1, 0, 0) \times (1, 0, 0) \\
&= \mathbf{i}(0 \cdot 0 - 0 \cdot 0) + \mathbf{j}(0 \cdot 1 - 1 \cdot 0) + \mathbf{k}(1 \cdot 0 - 0 \cdot 1) \\
&= \mathbf{0}.
\end{aligned}$$

Exercise 1.4.1 asks you to correspondingly establish the other six identities.

These example cross products most clearly demonstrate the orthogonality of a cross product to its two argument vectors (Theorem 1.4.6a) and that the direction is in the right-hand sense (Theorem 1.4.6b). ■

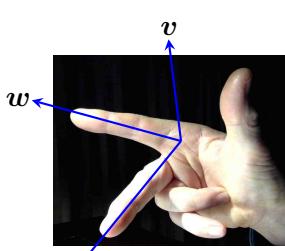
Geometry of a cross product

Example 1.4.5 (parallelogram area). Let's revisit the introduction to this section. Consider the parallelogram in the x_1x_2 -plane with edges formed by the \mathbb{R}^3 vectors $\mathbf{v} = (v_1, v_2, 0)$ and $\mathbf{w} = (w_1, w_2, 0)$. At the start of this Section 1.4 we derived that the parallelogram formed by these vectors has area $= |v_1w_2 - v_2w_1|$. Compare this area with the cross product

$$\begin{aligned}
\mathbf{v} \times \mathbf{w} &= \mathbf{i}(v_2 \cdot 0 - 0 \cdot w_2) + \mathbf{j}(0 \cdot w_1 - v_1 \cdot 0) + \mathbf{k}(v_1w_2 - v_2w_1) \\
&= \mathbf{i}0 + \mathbf{j}0 + \mathbf{k}(v_1w_2 - v_2w_1) \\
&= \mathbf{k}(v_1w_2 - v_2w_1).
\end{aligned}$$

Consequently, the length of this cross product equals the area of the parallelogram formed by \mathbf{v} and \mathbf{w} (Theorem 1.4.6d). (Also the direction of the cross product, $\pm\mathbf{k}$, is orthogonal to the x_1x_2 -plane containing the two vectors—Theorem 1.4.6a). ■

Theorem 1.4.6 (cross product geometry). *Let \mathbf{v} and \mathbf{w} be any two vectors in \mathbb{R}^3 .*



- (a) *the vector $\mathbf{v} \times \mathbf{w}$ is orthogonal to both \mathbf{v} and \mathbf{w} ;*
- (b) *the direction of $\mathbf{v} \times \mathbf{w}$ is in the right-hand sense in that if \mathbf{v} is in the direction of your thumb, and \mathbf{w} is in the direction of your straight index finger, then $\mathbf{v} \times \mathbf{w}$ is in the direction of your bent second/longest finger—all on your right-hand as illustrated in the margin;*
- (c) *$|\mathbf{v} \times \mathbf{w}| = |\mathbf{v}| |\mathbf{w}| \sin \theta$ where θ is the angle between vectors \mathbf{v} and \mathbf{w} ($0 \leq \theta \leq \pi$, equivalently $0^\circ \leq \theta \leq 180^\circ$); and*

- (d) the length $|\mathbf{v} \times \mathbf{w}|$ is the area of the parallelogram with edges \mathbf{v} and \mathbf{w} .

Proof. Let $\mathbf{v} = (v_1, v_2, v_3)$ and $\mathbf{w} = (w_1, w_2, w_3)$.

- 1.4.6a Recall that two vectors are orthogonal if their dot product is zero (Definition 1.3.15). To determine orthogonality between \mathbf{v} and the cross product $\mathbf{v} \times \mathbf{w}$, consider

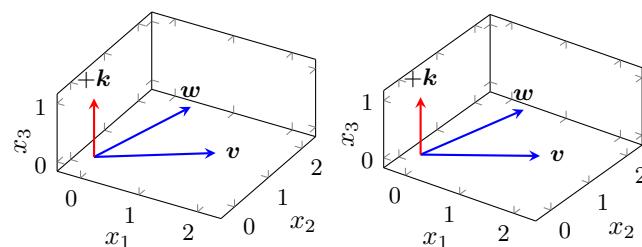
$$\begin{aligned}\mathbf{v} \cdot (\mathbf{v} \times \mathbf{w}) &= (v_1\mathbf{i} + v_2\mathbf{j} + v_3\mathbf{k}) \cdot [\mathbf{i}(v_2w_3 - v_3w_2) \\ &\quad + \mathbf{j}(v_3w_1 - v_1w_3) + \mathbf{k}(v_1w_2 - v_2w_1)] \\ &= v_1(v_2w_3 - v_3w_2) + v_2(v_3w_1 - v_1w_3) \\ &\quad + v_3(v_1w_2 - v_2w_1) \\ &= v_1v_2w_3 - v_1v_3w_2 + v_2v_3w_1 \\ &\quad - v_1v_2w_3 + v_1v_3w_2 - v_2v_3w_1 \quad = 0\end{aligned}$$

as each term in the penultimate line cancels with the term underneath in the last line. Since the dot product is zero, the cross product $\mathbf{v} \times \mathbf{w}$ is orthogonal to vector \mathbf{v} .

Similarly, $\mathbf{v} \times \mathbf{w}$ is orthogonal to \mathbf{w} (Exercise 1.4.5).

- 1.4.6b This right-handed property follows from the convention that the standard unit vectors \mathbf{i} , \mathbf{j} and \mathbf{k} are right-handed: that if \mathbf{i} is in the direction of your thumb, and \mathbf{j} is in the direction of your straight index finger, then \mathbf{k} is in the direction of your bent second/longest finger—all on your right-hand.

We prove only for the case of vectors in the x_1x_2 -plane, in which case $\mathbf{v} = (v_1, v_2, 0)$ and $\mathbf{w} = (w_1, w_2, 0)$, and when both $v_1, w_1 > 0$. One example is in stereo below.



Example 1.4.5 derived the cross product $\mathbf{v} \times \mathbf{w} = \mathbf{k}(v_1w_2 - v_2w_1)$. Consequently, this cross product is in the $+\mathbf{k}$ direction only when $v_1w_2 - v_2w_1 > 0$ (it is in the $-\mathbf{k}$ direction in the complementary case when $v_1w_2 - v_2w_1 < 0$). This inequality for $+\mathbf{k}$ rearranges to $v_1w_2 > v_2w_1$. Dividing by the positive v_1w_1 requires $\frac{w_2}{w_1} > \frac{v_2}{v_1}$. That is, in the x_1x_2 -plane the ‘slope’ of vector \mathbf{w} must greater than the ‘slope’ of vector \mathbf{v} . In this case, if \mathbf{v} is in the direction of your thumb on your right-hand, and \mathbf{w} is in the direction of your straight index finger, then your bent second/longest finger is in the direction $+\mathbf{k}$ as required by the cross-product $\mathbf{v} \times \mathbf{w}$.

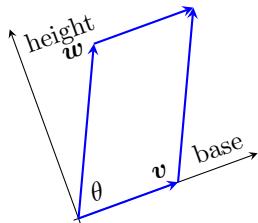
1.4.6c Exercise 1.4.6 establishes the identity $|\mathbf{v} \times \mathbf{w}|^2 = |\mathbf{v}|^2|\mathbf{w}|^2 - (\mathbf{v} \cdot \mathbf{w})^2$. From Theorem 1.3.4 substitute $\mathbf{v} \cdot \mathbf{w} = |\mathbf{v}||\mathbf{w}| \cos \theta$ into this identity:

$$\begin{aligned} |\mathbf{v} \times \mathbf{w}|^2 &= |\mathbf{v}|^2|\mathbf{w}|^2 - (\mathbf{v} \cdot \mathbf{w})^2 \\ &= |\mathbf{v}|^2|\mathbf{w}|^2 - (|\mathbf{v}||\mathbf{w}| \cos \theta)^2 \\ &= |\mathbf{v}|^2|\mathbf{w}|^2 - |\mathbf{v}|^2|\mathbf{w}|^2 \cos^2 \theta \\ &= |\mathbf{v}|^2|\mathbf{w}|^2(1 - \cos^2 \theta) \\ &= |\mathbf{v}|^2|\mathbf{w}|^2 \sin^2 \theta. \end{aligned}$$

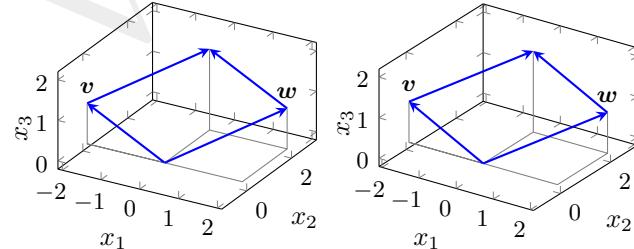
Take the square-root of both sides to determine $|\mathbf{v} \times \mathbf{w}| = \pm |\mathbf{v}||\mathbf{w}| \sin \theta$. But $\sin \theta \geq 0$ since the angle $0 \leq \theta \leq \pi$, and all the lengths are also ≥ 0 , so only the plus case applies. That is, the length $|\mathbf{v} \times \mathbf{w}| = |\mathbf{v}||\mathbf{w}| \sin \theta$ as required.

1.4.6d Consider the plane containing the vectors \mathbf{v} and \mathbf{w} , and hence containing the parallelogram formed by these vectors—as illustrated in the margin. Using vector \mathbf{v} as the base of the parallelogram, with length $|\mathbf{v}|$, by basic trigonometry the height of the parallelogram is then $|\mathbf{w}| \sin \theta$. Hence the area of the parallelogram is the product base \times height $= |\mathbf{v}||\mathbf{w}| \sin \theta = |\mathbf{v} \times \mathbf{w}|$ by the previous part 1.4.6c.

□



Example 1.4.7. Find the area of the parallelogram with edges formed by vectors $\mathbf{v} = (-2, 0, 1)$ and $\mathbf{w} = (2, 2, 1)$ —as in stereo below.



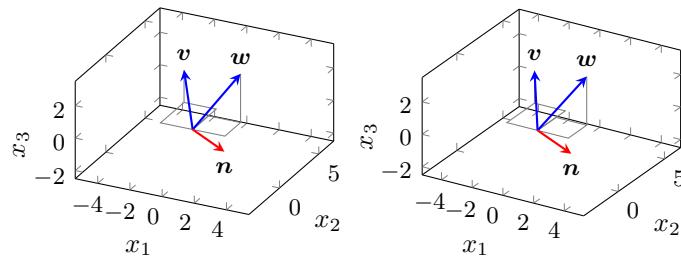
Solution: The area is the length of the cross product

$$\begin{aligned} \mathbf{v} \times \mathbf{w} &= \mathbf{i}(0 \cdot 1 - 1 \cdot 2) + \mathbf{j}(1 \cdot 2 - (-2) \cdot 1) + \mathbf{k}((-2) \cdot 2 - 0 \cdot 2) \\ &= -2\mathbf{i} + 4\mathbf{j} - 4\mathbf{k}. \end{aligned}$$

Then the parallelogram area $|\mathbf{v} \times \mathbf{w}| = \sqrt{(-2)^2 + 4^2 + (-4)^2} = \sqrt{4 + 16 + 16} = \sqrt{36} = 6$.

■

Example 1.4.8. Find a normal vector to the plane containing the two vectors $\mathbf{v} = -2\mathbf{i} + 3\mathbf{j} + 2\mathbf{k}$ and $\mathbf{w} = 2\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}$ —illustrated below. Hence find an equation of the plane given parametrically as $\mathbf{x} = -2\mathbf{i} - \mathbf{j} + 3\mathbf{k} + (-2\mathbf{i} + 3\mathbf{j} + 2\mathbf{k})s + (2\mathbf{i} + 2\mathbf{j} + 3\mathbf{k})t$.



Solution: Use Definition 1.4.3 of the cross-product to find a normal vector:

$$\begin{aligned}\mathbf{v} \times \mathbf{w} &= \mathbf{i}(3 \cdot 3 - 2 \cdot 2) + \mathbf{j}(2 \cdot 2 - (-2) \cdot 3) + \mathbf{k}((-2) \cdot 2 - 3 \cdot 2) \\ &= 5\mathbf{i} + 10\mathbf{j} - 10\mathbf{k}.\end{aligned}$$

A normal vector is any vector proportional to this, so we could divide by five and choose normal vector $\mathbf{n} = \mathbf{i} + 2\mathbf{j} - 2\mathbf{k}$ (as illustrated above).

An equation of the plane through $-2\mathbf{i} - \mathbf{j} + 3\mathbf{k}$ is then given by the dot product

$$\begin{aligned}(\mathbf{i} + 2\mathbf{j} - 2\mathbf{k}) \cdot [(x+2)\mathbf{i} + (y+1)\mathbf{j} + (z-3)\mathbf{k}] &= 0, \\ \text{that is, } x+2+2y+2-2z+6 &= 0, \\ \text{that is, } x+2y-2z+10 &= 0\end{aligned}$$

is the required normal equation of the plane.

■

Algebraic properties of a cross product

Exercises 1.4.9–1.4.11 establish three of the following four useful algebraic properties of the cross product.

Theorem 1.4.9 (cross product properties). *Let \mathbf{u} , \mathbf{v} and \mathbf{w} be any vectors in \mathbb{R}^3 , and c be any scalar:*

- (a) $\mathbf{v} \times \mathbf{v} = \mathbf{0}$;
- (b) $\mathbf{w} \times \mathbf{v} = -(\mathbf{v} \times \mathbf{w})$ (not commutative);
- (c) $(c\mathbf{v}) \times \mathbf{w} = c(\mathbf{v} \times \mathbf{w}) = \mathbf{v} \times (c\mathbf{w})$;
- (d) $\mathbf{u} \times (\mathbf{v} + \mathbf{w}) = \mathbf{u} \times \mathbf{v} + \mathbf{u} \times \mathbf{w}$ (distributive law).

Proof. Let's prove property 1.4.9a two ways—algebraically and geometrically. Exercises 1.4.9–1.4.11 ask you to prove the other properties.

- Algebraically: with vector $\mathbf{v} = (v_1, v_2, v_3)$, Definition 1.4.3 gives

$$\begin{aligned}\mathbf{v} \times \mathbf{v} &= \mathbf{i}(v_2 v_3 - v_3 v_2) + \mathbf{j}(v_3 v_1 - v_1 v_3) + \mathbf{k}(v_1 v_2 - v_2 v_1) \\ &= 0\mathbf{i} + 0\mathbf{j} + 0\mathbf{k} = \mathbf{0}.\end{aligned}$$

- Geometrically: from Theorem 1.4.6d, $|\mathbf{v} \times \mathbf{v}|$ is the area of the parallelogram with edges \mathbf{v} and \mathbf{v} . But such a parallelogram has zero area, so $|\mathbf{v} \times \mathbf{v}| = 0$. Since the only vector of length zero is the zero vector (Theorem 1.1.11), $\mathbf{v} \times \mathbf{v} = \mathbf{0}$.

□

Example 1.4.10. As an example of Theorem 1.4.9b, Example 1.4.4 showed that $\mathbf{i} \times \mathbf{j} = \mathbf{k}$, whereas reversing the order of the cross product gives the negative $\mathbf{j} \times \mathbf{i} = -\mathbf{k}$. Given Example 1.4.8 derived $\mathbf{v} \times \mathbf{w} = 5\mathbf{i} + 10\mathbf{j} - 10\mathbf{k}$ in the case when $\mathbf{v} = -2\mathbf{i} + 3\mathbf{j} + 2\mathbf{k}$ and $\mathbf{w} = 2\mathbf{i} + 2\mathbf{j} + 3\mathbf{k}$, what is $\mathbf{w} \times \mathbf{v}$?

Solution: By Theorem 1.4.9b, $\mathbf{w} \times \mathbf{v} = -(\mathbf{v} \times \mathbf{w}) = -5\mathbf{i} - 10\mathbf{j} + 10\mathbf{k}$.

■

Example 1.4.11. Given $(\mathbf{i} + \mathbf{j} + \mathbf{k}) \times (-2\mathbf{i} - \mathbf{j}) = \mathbf{i} - 2\mathbf{j} + \mathbf{k}$, what is $(3\mathbf{i} + 3\mathbf{j} + 3\mathbf{k}) \times (-2\mathbf{i} - \mathbf{j})$?

Solution: The first vector is $3(\mathbf{i} + \mathbf{j} + \mathbf{k})$ so by Theorem 1.4.9c,

$$\begin{aligned} & (3\mathbf{i} + 3\mathbf{j} + 3\mathbf{k}) \times (-2\mathbf{i} - \mathbf{j}) \\ &= [3(\mathbf{i} + \mathbf{j} + \mathbf{k})] \times (-2\mathbf{i} - \mathbf{j}) \\ &= 3[(\mathbf{i} + \mathbf{j} + \mathbf{k}) \times (-2\mathbf{i} - \mathbf{j})] \\ &= 3[\mathbf{i} - 2\mathbf{j} + \mathbf{k}] = 3\mathbf{i} - 6\mathbf{j} + 3\mathbf{k}. \end{aligned}$$

■

Example 1.4.12. The properties of Theorem 1.4.9 empower algebraic manipulation. Use such algebraic manipulation, and the identities among standard unit vectors of Example 1.4.4, compute the cross product $(\mathbf{i} - \mathbf{j}) \times (4\mathbf{i} + 2\mathbf{k})$.

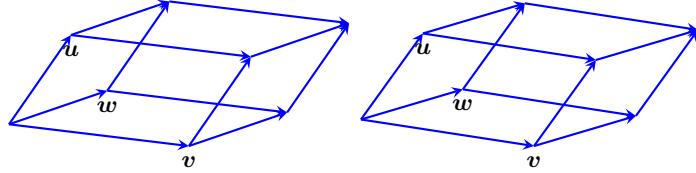
Solution: In full detail:

$$\begin{aligned} & (\mathbf{i} - \mathbf{j}) \times (4\mathbf{i} + 2\mathbf{k}) \\ &= (\mathbf{i} - \mathbf{j}) \times (4\mathbf{i}) + (\mathbf{i} - \mathbf{j}) \times (2\mathbf{k}) \quad (\text{by Thm 1.4.9d}) \\ &= 4(\mathbf{i} - \mathbf{j}) \times \mathbf{i} + 2(\mathbf{i} - \mathbf{j}) \times \mathbf{k} \quad (\text{by Thm 1.4.9c}) \\ &= -4\mathbf{i} \times (\mathbf{i} - \mathbf{j}) - 2\mathbf{k} \times (\mathbf{i} - \mathbf{j}) \quad (\text{by Thm 1.4.9b}) \\ &= -4[\mathbf{i} \times \mathbf{i} + \mathbf{i} \times (-\mathbf{j})] - 2[\mathbf{k} \times \mathbf{i} + \mathbf{k} \times (-\mathbf{j})] \quad (\text{by Thm 1.4.9d}) \\ &= -4[\mathbf{i} \times \mathbf{i} - \mathbf{i} \times \mathbf{j}] - 2[\mathbf{k} \times \mathbf{i} - \mathbf{k} \times \mathbf{j}] \quad (\text{by Thm 1.4.9c}) \\ &= -4[\mathbf{0} - \mathbf{k}] - 2[\mathbf{j} - (-\mathbf{i})] \quad (\text{by Example 1.4.4}) \\ &= -2\mathbf{i} - 2\mathbf{j} + 4\mathbf{k}. \end{aligned}$$

■

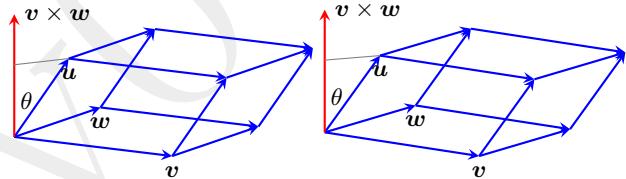
Volume of a parallelepiped

Consider the parallelepiped with edges formed by three vectors \mathbf{u} , \mathbf{v} and \mathbf{w} in \mathbb{R}^3 , as illustrated in stereo below. Our challenge is to derive that the volume of the parallelepiped is $|\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})|$.



The volume of the parallelepiped is the area of its base times its height.

- The base of the parallelepiped is the parallelogram formed with edges \mathbf{v} and \mathbf{w} . Hence the base has area $|\mathbf{v} \times \mathbf{w}|$ (Theorem 1.4.6d).
- The height of the parallelepiped is then that part of \mathbf{u} in the direction of a normal vector to \mathbf{v} and \mathbf{w} . We know that $\mathbf{v} \times \mathbf{w}$ is orthogonal to both \mathbf{v} and \mathbf{w} (Theorem 1.4.6a), so by trigonometry the height must be $|\mathbf{u}| \cos \theta$ for angle θ between \mathbf{u} and $\mathbf{v} \times \mathbf{w}$, as illustrated below.



To cater for cases where $\mathbf{v} \times \mathbf{w}$ points in the opposite direction to that shown, the height is $|\mathbf{u}| |\cos \theta|$. The dot product determines this cosine (Theorem 1.3.4):

$$\cos \theta = \frac{\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})}{|\mathbf{u}| |\mathbf{v} \times \mathbf{w}|}.$$

The height of the parallelepiped is then

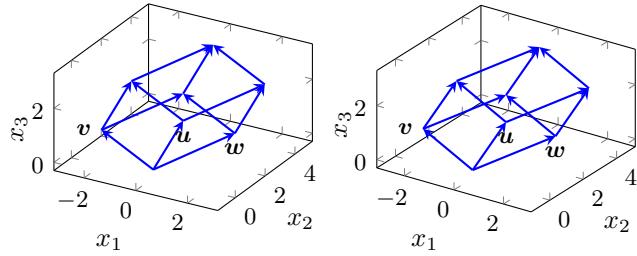
$$|\mathbf{u}| |\cos \theta| = |\mathbf{u}| \frac{|\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})|}{|\mathbf{u}| |\mathbf{v} \times \mathbf{w}|} = \frac{|\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})|}{|\mathbf{v} \times \mathbf{w}|}.$$

Consequently, the volume of the parallelepiped equals

$$\text{base} \cdot \text{height} = |\mathbf{v} \times \mathbf{w}| \frac{|\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})|}{|\mathbf{v} \times \mathbf{w}|} = |\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})|.$$

Definition 1.4.13. *For any three vectors \mathbf{u} , \mathbf{v} and \mathbf{w} in \mathbb{R}^3 , the **scalar triple product** is $\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})$.*

Example 1.4.14. Use the scalar triple product to find the area of the parallelepiped formed by vectors $\mathbf{u} = (0, 2, 1)$, $\mathbf{v} = (-2, 0, 1)$ and $\mathbf{w} = (2, 2, 1)$ —as illustrated in stereo below.



Solution: Example 1.4.7 found the cross product $\mathbf{v} \times \mathbf{w} = -2\mathbf{i} + 4\mathbf{j} - 4\mathbf{k}$. So the scalar triple product $\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w}) = (2\mathbf{j} + \mathbf{k}) \cdot (-2\mathbf{i} + 4\mathbf{j} - 4\mathbf{k}) = 8 - 4 = 4$. Hence the volume of the parallelepiped is 4 (cubic units).

The order of the vectors in a scalar triple product only affects the sign of the result. For example, we also find the volume of this parallelepiped via $\mathbf{v} \cdot (\mathbf{u} \times \mathbf{w})$. Returning to the procedure of Example 1.4.2 to find the cross product gives

$$\begin{aligned}\mathbf{u} \times \mathbf{w} &= \begin{vmatrix} \mathbf{i} & 0 & 2 \\ \mathbf{j} & 2 & 2 \\ \mathbf{k} & 1 & 1 \end{vmatrix} \\ &= \mathbf{i} \begin{vmatrix} \mathbf{i} & 0 & 2 \\ \mathbf{j} & 2 & 2 \\ \mathbf{k} & 1 & 1 \end{vmatrix} - \mathbf{j} \begin{vmatrix} \mathbf{i} & 0 & 2 \\ \mathbf{j} & 2 & 2 \\ \mathbf{k} & 1 & 1 \end{vmatrix} + \mathbf{k} \begin{vmatrix} \mathbf{i} & 0 & 2 \\ \mathbf{j} & 2 & 2 \\ \mathbf{k} & 1 & 1 \end{vmatrix} \\ &= \mathbf{i} \begin{vmatrix} 2 & 2 \\ 1 & 1 \end{vmatrix} - \mathbf{j} \begin{vmatrix} 0 & 2 \\ 1 & 1 \end{vmatrix} + \mathbf{k} \begin{vmatrix} 0 & 2 \\ 2 & 2 \end{vmatrix} \\ &= \mathbf{i} \cancel{\begin{vmatrix} 2 & 2 \\ 1 & 1 \end{vmatrix}} - \mathbf{j} \cancel{\begin{vmatrix} 0 & 2 \\ 1 & 1 \end{vmatrix}} + \mathbf{k} \cancel{\begin{vmatrix} 0 & 2 \\ 2 & 2 \end{vmatrix}} \\ &= \mathbf{i}(2 \cdot 1 - 1 \cdot 2) - \mathbf{j}(0 \cdot 1 - 1 \cdot 2) + \mathbf{k}(0 \cdot 2 - 2 \cdot 2) \\ &= 2\mathbf{j} - 4\mathbf{k}.\end{aligned}$$

Then the triple product $\mathbf{v} \cdot (\mathbf{u} \times \mathbf{w}) = (-2\mathbf{i} + \mathbf{k}) \cdot (2\mathbf{j} - 4\mathbf{k}) = 0 + 0 - 4 = -4$. Hence the volume of the parallelepiped is $|-4| = 4$ as before.

■

Using the procedure of Example 1.4.2 to find a scalar triple product establishes a strong connection to the determinants of Chapter 6. In the second solution to the previous Example 1.4.14, in finding $\mathbf{u} \times \mathbf{w}$, the unit vectors \mathbf{i} , \mathbf{j} and \mathbf{k} just acted as place holding symbols to eventually ensure a multiplication by the correct component of \mathbf{v} in the dot product. We could seamlessly combine the two products by replacing the symbols \mathbf{i} , \mathbf{j} and \mathbf{k} directly with the corresponding component of \mathbf{v} :

$$\mathbf{v} \cdot (\mathbf{u} \times \mathbf{w}) = \begin{vmatrix} -2 & 0 & 2 \\ 0 & 2 & 2 \\ 1 & 1 & 1 \end{vmatrix}$$

$$\begin{aligned}
&= -2 \begin{vmatrix} -2 & 0 & 2 \\ 0 & 2 & 2 \\ 1 & 1 & 1 \end{vmatrix} - 0 \begin{vmatrix} -2 & 0 & 2 \\ 0 & 2 & 2 \\ 1 & 1 & 1 \end{vmatrix} + 1 \begin{vmatrix} -2 & 0 & 2 \\ 0 & 2 & 2 \\ 1 & 1 & 1 \end{vmatrix} \\
&= -2 \begin{vmatrix} 2 & 2 \\ 1 & 1 \end{vmatrix} - 0 \begin{vmatrix} 0 & 2 \\ 1 & 1 \end{vmatrix} + 1 \begin{vmatrix} 0 & 2 \\ 2 & 2 \end{vmatrix} \\
&= -2 \begin{vmatrix} 2 & 2 \\ 1 & 1 \end{vmatrix} - 0 \begin{vmatrix} 0 & 2 \\ 1 & 1 \end{vmatrix} + 1 \begin{vmatrix} 0 & 2 \\ 2 & 2 \end{vmatrix} \\
&= -2(2 \cdot 1 - 1 \cdot 2) - 0(0 \cdot 1 - 1 \cdot 2) + 1(0 \cdot 2 - 2 \cdot 2) \\
&= -2 \cdot 0 - 0(-2) + 1(-4) = -4.
\end{aligned}$$

Hence the parallelepiped formed by \mathbf{u} , \mathbf{v} and \mathbf{w} has volume $| -4 |$, as before. Here the volume follows from the above manipulations of the matrix of numbers formed with columns of the matrix being the vectors \mathbf{u} , \mathbf{v} and \mathbf{w} . Chapter 6 shows that this computation of volume generalises to determining, via analogous matrices of vectors, the ‘volume’ of objects formed by vectors with any number of components.

1.4.1 Exercises

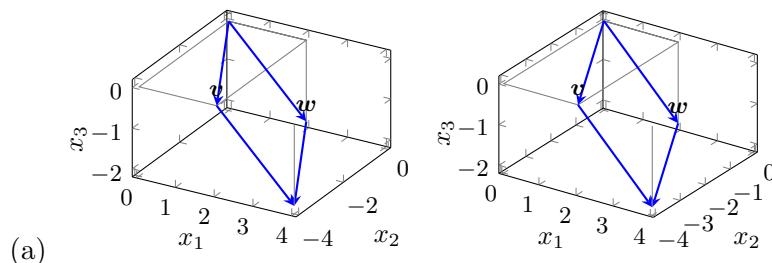
Exercise 1.4.1. Use Definition 1.4.3 to establish some of the standard unit vector identities in Example 1.4.4:

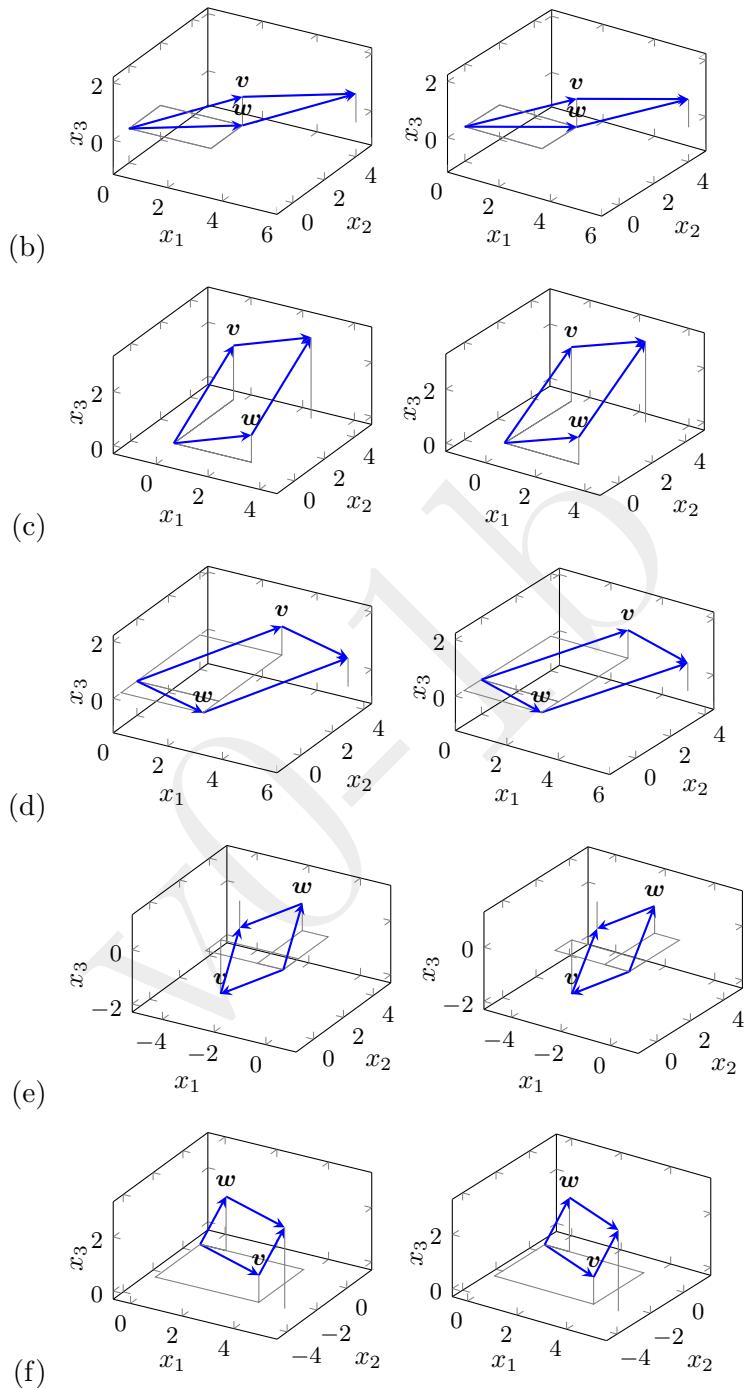
- (a) $\mathbf{j} \times \mathbf{k} = \mathbf{i}$, $\mathbf{k} \times \mathbf{j} = -\mathbf{i}$, $\mathbf{j} \times \mathbf{j} = \mathbf{0}$;
- (b) $\mathbf{k} \times \mathbf{i} = \mathbf{j}$, $\mathbf{i} \times \mathbf{k} = -\mathbf{j}$, $\mathbf{k} \times \mathbf{k} = \mathbf{0}$.

Exercise 1.4.2. Use Definition 1.4.3, perhaps via the procedure used in Example 1.4.2, to determine the following cross products. Confirm each cross product is orthogonal to the two vectors in the given product. Show your details.

- | | |
|---|--|
| (a) $(3\mathbf{i} + \mathbf{j}) \times (3\mathbf{i} - 3\mathbf{j} - 2\mathbf{k})$ | (b) $(3\mathbf{i} + \mathbf{k}) \times (5\mathbf{i} + 6\mathbf{k})$ |
| (c) $(2\mathbf{i} - \mathbf{j} - 3\mathbf{k}) \times (3\mathbf{i} + 2\mathbf{k})$ | (d) $(\mathbf{i} - \mathbf{j} + 2\mathbf{k}) \times (3\mathbf{i} + 3\mathbf{k})$ |
| (e) $(-1, 3, 2) \times (3, -5, 1)$ | (f) $(3, 0, 4) \times (5, 1, 2)$ |
| (g) $(4, 1, 3) \times (3, 2, -1)$ | (h) $(3, -7, 3) \times (2, 1, 0)$ |

Exercise 1.4.3. For each of the stereo pictures below, estimate the area of the pictured parallelogram by estimating the edge vectors \mathbf{v} and \mathbf{w} (all components are integers), then computing their cross product.





Exercise 1.4.4. Each of the following equations describes a plane in 3D.
Find a normal vector to each of the planes.

- $\mathbf{x} = (-1, 0, 1) + (-5, 2, -1)s + (2, -4, 0)t$
- $2x + 2y + 4z = 20$
- $x_1 - x_2 + x_3 + 2 = 0$
- $\mathbf{x} = 6\mathbf{i} - 3\mathbf{j} + (3\mathbf{i} - 3\mathbf{j} - 2\mathbf{k})s - (\mathbf{i} + \mathbf{j} + \mathbf{k})t$
- $\mathbf{x} = \mathbf{j} + 2\mathbf{k} + (\mathbf{i} - \mathbf{k})s + (-5\mathbf{i} + \mathbf{j} - 3\mathbf{k})t$

- (f) $3y = x + 2z + 4$
 (g) $3p + 8q - 9 = 4r$
 (h) $\mathbf{x} = (-2, 2, -3) + (-3, 2, 0)s + (-1, 3, 2)t$

Exercise 1.4.5. Use Definition 1.4.3 to prove that for all vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^3$, the cross product $\mathbf{v} \times \mathbf{w}$ is orthogonal to \mathbf{w} .

Exercise 1.4.6. For all vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^3$ prove the identity $|\mathbf{v} \times \mathbf{w}|^2 = |\mathbf{v}|^2|\mathbf{w}|^2 - (\mathbf{v} \cdot \mathbf{w})^2$ (an identity invoked in the proof of Theorem 1.4.6c). Use the algebraic Definitions 1.3.2 and 1.4.3 of the dot and cross products to expand both sides of the identity and show both sides expand to the same complicated expression.

Exercise 1.4.7. Using Theorem 1.4.9, and the identities among standard unit vectors of Example 1.4.4, compute the following cross products. Record and justify each step in detail.

- | | |
|--|---|
| (a) $\mathbf{i} \times (3\mathbf{j})$ | (b) $(4\mathbf{j} + 3\mathbf{k}) \times \mathbf{k}$ |
| (c) $(4\mathbf{k}) \times (\mathbf{i} + 6\mathbf{j})$ | (d) $\mathbf{j} \times (3\mathbf{i} + 2\mathbf{k})$ |
| (e) $(2\mathbf{i} + 2\mathbf{k}) \times (\mathbf{i} + \mathbf{j})$ | (f) $(\mathbf{i} - 5\mathbf{j}) \times (-\mathbf{j} + 3\mathbf{k})$ |

Exercise 1.4.8. You are given that three specific vectors \mathbf{u} , \mathbf{v} and \mathbf{w} in \mathbb{R}^3 have the following cross products:

$$\begin{aligned}\mathbf{u} \times \mathbf{v} &= -\mathbf{j} + \mathbf{k}, \\ \mathbf{u} \times \mathbf{w} &= \mathbf{i} - \mathbf{k}, \\ \mathbf{v} \times \mathbf{w} &= -\mathbf{i} + 2\mathbf{j}.\end{aligned}$$

Use Theorem 1.4.9 to compute the following cross products. Record and justify each step in detail.

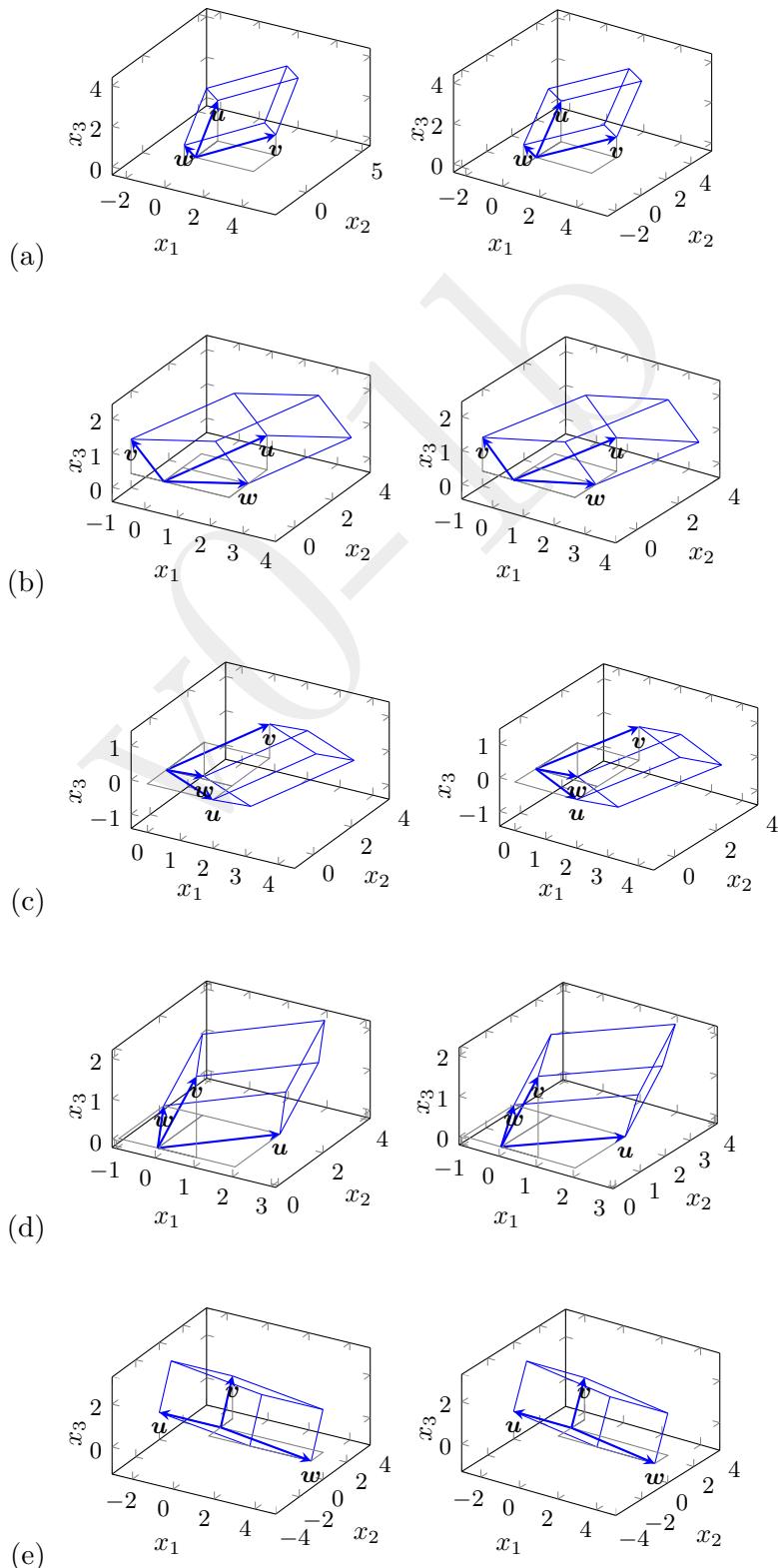
- | | |
|---|--|
| (a) $(\mathbf{u} + \mathbf{v}) \times \mathbf{w}$ | (b) $(3\mathbf{u} + \mathbf{w}) \times (2\mathbf{u})$ |
| (c) $(3\mathbf{v}) \times (\mathbf{u} + \mathbf{v})$ | (d) $(2\mathbf{v} + \mathbf{w}) \times (\mathbf{u} + 3\mathbf{v})$ |
| (e) $(2\mathbf{v} + 3\mathbf{w}) \times (\mathbf{u} + 2\mathbf{w})$ | (f) $(\mathbf{u} + 4\mathbf{v} + 2\mathbf{w}) \times \mathbf{w}$ |

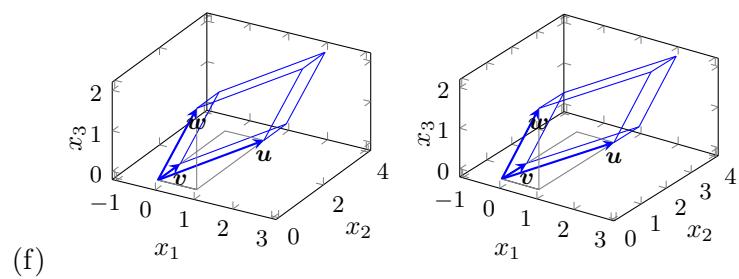
Exercise 1.4.9. Use Definition 1.4.3 to algebraically prove Theorem 1.4.9b—the property that $\mathbf{w} \times \mathbf{v} = -(\mathbf{v} \times \mathbf{w})$. Explain how this property also follows from the basic geometry of the cross product (Theorem 1.4.6).

Exercise 1.4.10. Use Definition 1.4.3 to algebraically prove Theorem 1.4.9c—the property that $(c\mathbf{v}) \times \mathbf{w} = c(\mathbf{v} \times \mathbf{w}) = \mathbf{v} \times (c\mathbf{w})$. Explain how this property also follows from the basic geometry of the cross product (Theorem 1.4.6)—consider $c > 0$, $c = 0$ and $c < 0$ separately.

Exercise 1.4.11. Use Definition 1.4.3 to algebraically prove Theorem 1.4.9d—the distributive property that $\mathbf{u} \times (\mathbf{v} + \mathbf{w}) = \mathbf{u} \times \mathbf{v} + \mathbf{u} \times \mathbf{w}$.

Exercise 1.4.12. For each of the following illustrated parallelepipeds: estimate the edge vectors \mathbf{u} , \mathbf{v} and \mathbf{w} (all components are integers); then use the scalar triple product to estimate the volume of the parallelepiped.





1.5 Use Matlab/Octave for vector computation

Section Contents

1.5.1 Exercises	79
---------------------------	----

It is the science of *calculation*,—which becomes continually more necessary at each step of our progress, and which must ultimately govern the whole of the applications of science to the arts of life.

Charles Babbage, 1832

Subsequent chapters invoke the computer packages Matlab/Octave to perform calculations that would be tedious and error prone when done by hand. This section introduces Matlab/Octave so that you can start to become familiar with it on small problems. You should directly compare the computed answer with your calculation by hand. The aim is to develop some basic confidence with Matlab/Octave before later using it to save considerable time in longer tasks.

- Matlab is commercial software available from Mathworks.⁵ It is also available over the internet as Matlab-Online or Matlab-Mobile.
- Octave is free software, that for our limited purposes is almost the same as Matlab, and downloadable over the internet.⁶
- Alternatively, your home institution may provide Matlab/Octave via a web service that is useable from smart phones, tablets and computers.

Example 1.5.1. Use the Matlab/Octave command `norm()` to compute the length/magnitude of the following vectors (Definition 1.1.8).

- (a) $(2, -1)$

Solution: Start Matlab/Octave. After a prompt, “>>” in Matlab or “octave:>” in Octave, type a command, followed by the Return/Enter key to get it executed. As indicated by Table 1.2 the numbers with brackets separated by semi-colons forms a vector, and the = character assigns the result to variable for subsequent use.

⁵ <http://mathworks.com>

⁶ <https://www.gnu.org/software/octave/>

Table 1.2: Use Matlab/Octave to help compute vector results with the following basics. This and subsequent tables throughout the book summarise the main Matlab/Octave information.

- Real numbers are limited to being zero or of magnitude from 10^{-323} to 10^{+308} , both positive and negative (called the **floating point** numbers). Real numbers are computed and stored to a maximum precision of nearly sixteen significant digits.^a
- Matlab/Octave potentially uses complex numbers (\mathbb{C}), but mostly we stay with real numbers (\mathbb{R}).
- Each Matlab/Octave command is usually typed on one line by itself.
- `[. ; . ; .]` where each dot denotes a number, forms vectors in \mathbb{R}^3 (or use newlines instead of the semi-colons). Use more numbers separated by semi-colons for vectors in other \mathbb{R}^n .
- `=` assigns the value of the expression on the right to the variable name on the left.
If the computation of an expression is not assigned to anything explicitly, then by default it is assigned to the variable `ans` (as denoted by “`ans =`” in Matlab/Octave).
- Variable names are alphanumeric starting with a letter.
- `size(v)` returns the number of components of the vector (Definition 1.1.4): if the vector v is in \mathbb{R}^m , then `size(v)` returns $[m \ 1]$.
- `norm(v)` computes the length/magnitude of the vector v (Definition 1.1.8).
- `+,-,*` is vector/scalar addition, subtraction, and multiplication, but only provided the sizes of the two vectors are the same. Parentheses () control the order of operations.
- `/x` divides by a scalar x . However, be warned that `/v` for a vector v typically gives strange results as Matlab/Octave interprets it to mean you want to (approximately) solve some linear equation.
- `dot(u,v)` computes the dot product of vectors u and v (Definition 1.3.2)—if they have the same size.
- `acos(q)` computes the arc-cos, the inverse cosine, of the scalar q in radians. To find the angle in degrees use `acos(q)*180/pi`.
- `quit` terminates the Matlab/Octave session.

^aIf desired, ‘computer algebra’ software provides us with an arbitrary level of precision, even exact. Current computer algebra software includes the free Sage, Maxima and Reduce, and the commercial Maple, Mathematica and (via Matlab) MuPad.

Assign the vector to a variable \mathbf{a} by the command $\mathbf{a}=[2;-1]$. Then executing $\mathbf{norm}(\mathbf{a})$ reports $\mathbf{ans} = 2.2361$ as shown in the dialogue to the right.

```
>> a=[2;-1]
a =
    2
   -1
>> norm(a)
ans = 2.2361
```

This computes the answer $|(2, -1)| = \sqrt{2^2 + (-1)^2} = \sqrt{5} = 2.2361$ (to five significant digits which we take to be practically exact).

The qr-code appearing in the margin here encodes these Matlab/Octave commands. You may scan such qr-codes with your favourite app⁷, and then copy and paste the code direct into a Matlab/Octave client. Alternatively, if reading an electronic version of this book, then you may copy and paste the commands; however, be warned that the quote character ' usually needs correcting. Although here the saving in typing is negligible, later you can save considerable typing

(b) $(-1, 1, -5, 4)$

Solution: In Matlab/Octave:

Assign the vector to a variable with $\mathbf{b}=[-1;1;-5;4]$ as shown to the right. Then execute $\mathbf{norm}(\mathbf{b})$ and find that Matlab/Octave reports $\mathbf{ans} = 6.5574$

```
>> b=[-1;1;-5;4]
b =
   -1
    1
   -5
    4
>> norm(b)
ans = 6.5574
```

Hence 6.5574 is the length of $(-1, 1, -5, 4)$ (to five significant digits which we take to be practically exact).

(c) $(-0.3, 4.3, -2.5, -2.8, 7, -1.9)$

Solution: In Matlab/Octave:

i. assign the vector with the command
 $\mathbf{c}=[-0.3;4.3;-2.5;-2.8;7;-1.9]$

ii. execute $\mathbf{norm}(\mathbf{c})$ and find that Matlab/Octave reports
 $\mathbf{ans} = 9.2347$

Hence the length of vector $(-0.3, 4.3, -2.5, -2.8, 7, -1.9)$ is 9.2347 (to five significant digits).



⁷ At the time of writing, qr-code scanning applications for smart-phones include *WaspScan* and *QRReader*—but I have no expertise to assess their quality.



Example 1.5.2. Use Matlab/Octave operators $+$, $-$, $*$ to compute the value of the expressions $\mathbf{u} + \mathbf{v}$, $\mathbf{u} - \mathbf{v}$, $3\mathbf{u}$ for vectors $\mathbf{u} = (-4.1, 1.7, 4.1)$ and $\mathbf{v} = (2.9, 0.9, -2.4)$ (Definition 1.2.3).

Solution: In Matlab/Octave type the commands, each followed by Return/Enter key.

Assign the named vectors with the commands
 $\mathbf{u} = [-4.1; 1.7; 4.1]$ and
 $\mathbf{v} = [2.9; 0.9; -2.4]$ to see the two steps in the dialogue to the right.

```
>> u=[-4.1;1.7;4.1]
u =
-4.1000
1.7000
4.1000
>> v=[2.9;0.9;-2.4]
v =
2.9000
0.9000
-2.4000
```

Execute $\mathbf{u} + \mathbf{v}$ to find from the dialogue on the right that the sum $\mathbf{u} + \mathbf{v} = (-1.2, 2.6, 1.7)$.

```
>> u+v
ans =
-1.2000
2.6000
1.7000
```

Execute $\mathbf{u} - \mathbf{v}$ to find from the dialogue on the right that the difference $\mathbf{u} - \mathbf{v} = (-7, 0.8, 6.5)$.

```
>> u-v
ans =
-7.0000
0.8000
6.5000
```

Execute $3\mathbf{u}$ to find from the dialogue on the right that the scalar multiple $3\mathbf{u} = (-12.3, 5.1, 12.3)$ (the asterisk is essential to compute multiplication).

```
>> 3*u
ans =
-12.3000
5.1000
12.3000
```

Example 1.5.3. Use Matlab/Octave to confirm that $2(2\mathbf{p} - 3\mathbf{q}) + 6(\mathbf{q} - \mathbf{p}) = -2\mathbf{p}$ for vectors $\mathbf{p} = (1, 0, 2, -6)$ and $\mathbf{q} = (2, 4, 3, 5)$.

Solution: In Matlab/Octave

Assign the first vector with
 $p=[1;0;2;-6]$ as shown to the right.



```
>> p=[1;0;2;-6]
p =
    1
    0
    2
   -6
```

Assign the other vector with
 $q=[2;4;3;5]$.

```
>> q=[2;4;3;5]
q =
    2
    4
    3
    5
```

Compute $2(2p - 3q) + 6(q - p)$ with the command
 $2*(2*p-3*q)+6*(q-p)$ as shown to the right, and see the result is evidently $-2p$.

```
>> 2*(2*p-3*q)+6*(q-p)
ans =
    -2
     0
    -4
    12
```

Confirm it is $-2p$ by adding $2p$ to the above result with the command `ans+2*p` as shown to the right, and see the zero vector result.

```
>> ans+2*p
ans =
    0
    0
    0
    0
```

Example 1.5.4. Use Matlab/Octave to confirm the commutative law (Theorem 1.2.13a) $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$ for vectors $\mathbf{u} = (8, -6, -4, -2)$ and $\mathbf{v} = (4, 3, -1)$.

Solution: In Matlab/Octave



Assign $\mathbf{u} = (8, -6, -4, -2)$
with command
 $\mathbf{u} = [8; -6; -4; -2]$ as shown to
the right.

```
>> u=[8;-6;-4;-2]
u =
    8
   -6
   -4
   -2
```

Assign $\mathbf{v} = (4, 3, -1)$ with the
command $\mathbf{v} = [4; 3; -1]$.

```
>> v=[4;3;-1]
v =
    4
    3
   -1
```

Compute $\mathbf{u} + \mathbf{v}$ with the
command $\mathbf{u} + \mathbf{v}$ as shown to the
right. Matlab prints an error
message because the vectors \mathbf{u}
and \mathbf{v} are of different sizes and
so cannot be added together.

```
>> u+v
Error using +
Matrix dimensions must agree.
```

Check the sizes of the vectors
in the sum using `size(u)` and
`size(v)` to confirm \mathbf{u} is in \mathbb{R}^4
whereas \mathbf{v} is in \mathbb{R}^3 . Hence the
two vectors cannot be added
(Definition 1.2.3).

```
>> size(u)
ans =
    4    1
>> size(v)
ans =
    3    1
```

Alternatively, Octave gives the following error message in which
“nonconformant arguments” means of wrong sizes.

```
error: operator +: nonconformant arguments
(op1 is 4x1, op2 is 3x1)
```



Example 1.5.5. Use Matlab/Octave to compute the angles between the pair of vectors $(4, 3)$ and $(5, 12)$ (Theorem 1.3.4).

Solution: In Matlab/Octave



Because each vector is used twice in the formula

$$\cos \theta = (\mathbf{u} \cdot \mathbf{v}) / (|\mathbf{u}| |\mathbf{v}|),$$

give each a name as shown to the right.

```
>> u=[4;3]
u =
    4
    3
>> v=[5;12]
v =
    5
    12
```

Then invoke `dot()` to compute the dot product in the formula for the cosine of the angle.

Lastly, invoke `acos()` for the arc-cosine, then convert the radians to degrees to find the angle $\theta = 30.510^\circ$.

```
>> cost=dot(u,v)/norm(u)/norm(v)
cost = 0.8615
```

```
>> theta=acos(cost)*180/pi
theta = 30.510
```

■

Example 1.5.6. Verify the distributive law for the dot product $(\mathbf{u}+\mathbf{v}) \cdot \mathbf{w} = \mathbf{u} \cdot \mathbf{w} + \mathbf{v} \cdot \mathbf{w}$ (Theorem 1.3.10d) for vectors $\mathbf{u} = (-0.1, -3.1, -2.9, -1.3)$, $\mathbf{v} = (-3, 0.5, 6.4, -0.9)$ and $\mathbf{w} = (-1.5, -0.2, 0.4, -3.1)$.

Solution: In Matlab/Octave



Assign vector

$\mathbf{u} = (-0.1, -3.1, -2.9, -1.3)$ with the command
 $\mathbf{u} = [-0.1; -3.1; -2.9; -1.3]$ as shown to the right.

```
>> u=[-0.1;-3.1;-2.9;-1.3]
u =
    -0.1000
    -3.1000
    -2.9000
    -1.3000
```

Assign vector

$\mathbf{v} = (-3, 0.5, 6.4, -0.9)$ with the command
 $\mathbf{v} = [-3; 0.5; 6.4; -0.9]$ as shown to the right.

```
>> v=[-3;0.5;6.4;-0.9]
v =
    -3.0000
    0.5000
    6.4000
    -0.9000
```

Assign vector

$$\mathbf{w} = (-1.5, -0.2, 0.4, -3.1)$$

with the command

$$\mathbf{w} = [-1.5; -0.2; 0.4; -3.1]$$

as shown to the right.

```
>> w=[-1.5;-0.2;0.4;-3.1]
```

w =

-1.5000

-0.2000

0.4000

-3.1000

Compute the dot product

$$(\mathbf{u} + \mathbf{v}) \cdot \mathbf{w}$$

with the command

`dot(u+v,w)`

to find the answer

is 13.390.

```
>> dot(u+v,w)
```

ans = 13.390

Compare this with the dot product expression

$$\mathbf{u} \cdot \mathbf{w} + \mathbf{v} \cdot \mathbf{w}$$

via the command

`dot(u,w)+dot(v,w)`

to find the answer is 13.390.

```
>> dot(u,w)+dot(v,w)
```

ans = 13.390

That the two answers agree verifies the distributive law for dot products.

■

On two occasions I have been asked [by members of Parliament!], “Pray, Mr. Babbage, if you put into the machine wrong figures, will the right answers come out?”

I am not able rightly to apprehend the kind of confusion of ideas that could provoke such a question.

Charles Babbage

1.5.1 Exercises

Exercise 1.5.1. Use Matlab/Octave to compute the length of each of the following vectors (the first five have integer lengths).

$$(a) \begin{bmatrix} 2 \\ 3 \\ 6 \end{bmatrix} \quad (b) \begin{bmatrix} 4 \\ -4 \\ 7 \end{bmatrix} \quad (c) \begin{bmatrix} -2 \\ 6 \\ 9 \end{bmatrix} \quad (d) \begin{bmatrix} 0.5 \\ 0.1 \\ -0.5 \\ 0.7 \end{bmatrix} \quad (e) \begin{bmatrix} 8 \\ -4 \\ 5 \\ -8 \end{bmatrix}$$

$$(f) \begin{bmatrix} 1.1 \\ 1.7 \\ -4.2 \\ -3.8 \\ 0.9 \end{bmatrix} \quad (g) \begin{bmatrix} 2.6 \\ -0.1 \\ 3.2 \\ -0.6 \\ -0.2 \end{bmatrix} \quad (h) \begin{bmatrix} 1.6 \\ -1.1 \\ -1.4 \\ 2.3 \\ -1.6 \end{bmatrix}$$

Exercise 1.5.2. Use Matlab/Octave to determine which are wrong out of the following identities and relations for vectors $\mathbf{p} = (0.8, -0.3, 1.1, 2.6, 0.1)$ and $\mathbf{q} = (1, 2.8, 1.2, 2.3, 2.3)$.

- | | |
|--|--|
| (a) $3(\mathbf{p} - \mathbf{q}) = 3\mathbf{p} - 3\mathbf{q}$ | (b) $2(\mathbf{p} - 3\mathbf{q}) + 3(2\mathbf{q} - \mathbf{p}) = \mathbf{p}$ |
| (c) $\frac{1}{2}(\mathbf{p} - \mathbf{q}) + \frac{1}{2}(\mathbf{p} + \mathbf{q}) = \mathbf{p}$ | (d) $ \mathbf{p} + \mathbf{q} \leq \mathbf{p} + \mathbf{q} $ |
| (e) $ \mathbf{p} - \mathbf{q} \leq \mathbf{p} + \mathbf{q} $ | (f) $ \mathbf{p} \cdot \mathbf{q} \leq \mathbf{p} \mathbf{q} $ |

Exercise 1.5.3. Use Matlab/Octave to find the angles between pairs of vectors in each of the following groups.

(a) $\mathbf{p} = (2, 3, 6)$, $\mathbf{q} = (6, 2, -3)$, $\mathbf{r} = (3, -6, 2)$



(b) $\mathbf{u} = \begin{bmatrix} -1 \\ -7 \\ -1 \\ -7 \end{bmatrix}$, $\mathbf{v} = \begin{bmatrix} -1 \\ 4 \\ 4 \\ 4 \end{bmatrix}$, $\mathbf{w} = \begin{bmatrix} 1 \\ -4 \\ -4 \\ -4 \end{bmatrix}$



(c) $\mathbf{u} = \begin{bmatrix} 5 \\ 1 \\ -1 \\ 3 \end{bmatrix}$, $\mathbf{v} = \begin{bmatrix} -6 \\ 4 \\ 2 \\ -5 \end{bmatrix}$, $\mathbf{w} = \begin{bmatrix} -4 \\ 2 \\ 1 \\ -2 \end{bmatrix}$



(d) $\mathbf{u} = \begin{bmatrix} -9 \\ -8 \\ 4 \\ 8 \end{bmatrix}$, $\mathbf{v} = \begin{bmatrix} 9 \\ 3 \\ 1 \\ -3 \end{bmatrix}$, $\mathbf{w} = \begin{bmatrix} -6 \\ 5 \\ -2 \\ -4 \end{bmatrix}$



(e) $\mathbf{a} = \begin{bmatrix} -4.1 \\ 9.8 \\ 0.3 \\ 1.4 \\ 2.7 \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} -0.6 \\ 2.6 \\ -1.2 \\ -0.2 \\ -0.9 \end{bmatrix}$, $\mathbf{c} = \begin{bmatrix} -2.8 \\ -0.9 \\ -6.2 \\ -2.3 \\ -4.7 \end{bmatrix}$, $\mathbf{d} = \begin{bmatrix} 1.8 \\ -3.4 \\ -8.6 \\ 1.4 \\ 1.8 \end{bmatrix}$



(f) $\mathbf{a} = \begin{bmatrix} -0.5 \\ 2.0 \\ -3.4 \\ 1.8 \\ 0.1 \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} 5.4 \\ 7.4 \\ 0.5 \\ 0.7 \\ 1.3 \end{bmatrix}$, $\mathbf{c} = \begin{bmatrix} -0.2 \\ -1.5 \\ -0.3 \\ 1.1 \\ 2.5 \end{bmatrix}$, $\mathbf{d} = \begin{bmatrix} 1.0 \\ 2.0 \\ -1.3 \\ -4.4 \\ -2.0 \end{bmatrix}$

Answers to selected exercises

1.2.2a : 5.9

1.2.2c : 2.2

1.2.2e : 4.9

1.2.3a : \mathbf{u} and \mathbf{w} are closest; \mathbf{v} and \mathbf{w} are furthest.

1.2.3c : \mathbf{u} and \mathbf{w} are closest; \mathbf{v} and \mathbf{u} are furthest.

1.2.3e : \mathbf{u} and \mathbf{v} are closest; \mathbf{u} and \mathbf{w} are furthest.

1.2.3g : \mathbf{u} and \mathbf{v} are closest; \mathbf{u} and \mathbf{w} are furthest.

1.2.4a : One possibility is $\mathbf{x} = (-11 + 8t, -2t, 3 - t)$

1.2.4c : One possibility is $\mathbf{x} = (2.4 - 0.9t, 5.5 - 10.9t, -3.9 + 3.4t)$

1.2.4e : One possibility is $\mathbf{x} = (2.2 - 3.3t, 5.8 - 3.6t, 4 - 6.4t, 3 - 6.2t, 2 - 1.1t)$

1.3.2a : Not orthogonal.

1.3.2c : Orthogonal.

1.3.2e : Not orthogonal.

1.3.2g : Orthogonal.

1.3.3 : $\theta_{ab} = 69.30^\circ$, $\theta_{ac} = 27.89^\circ$, $\theta_{bc} = 41.41^\circ$. The first and third sentences have smallest angle and so are most similar.

1.3.10a : $(-3, -2, 2) + (3, 7, -6)s + (8, 3, -5)t$

1.3.10c : $(2, 3, 3) + (0, -1, 0)s + (1, -2, -3)t$

1.3.10e : $(-2.2, 1.3, -4.9) + (2.6, -3.5, 13.6)s + (0.8, 1.9, 4.5)t$

1.3.10g : $(-1.8, 4.3, -3.9) + (-3.8, -6.5, -2.9)s + (4.3, -7.8, 2.2)t$

1.4.2a : $-2\mathbf{i} + 6\mathbf{j} - 12\mathbf{k}$

1.4.2c : $-2\mathbf{i} - 13\mathbf{j} + 3\mathbf{k}$

1.4.2e : $(13, 7, -4)$

1.4.2g : $(-7, 13, 5)$

1.4.3a : 12

1.4.3c : 14

1.4.3e : $\sqrt{138} = 11.75$

1.4.4a : $\propto (-2, -1, 8)$

1.4.4c : $\propto \mathbf{i} - \mathbf{j} + \mathbf{k}$

1.4.4e : $\propto \mathbf{i} + 8\mathbf{j} + \mathbf{k}$

1.4.4g : $(p, q, r) \propto (3, 8, -4)$

1.4.7a : $3\mathbf{k}$

1.4.7c : $-24\mathbf{i} + 4\mathbf{j}$

1.4.7e : $-2\mathbf{i} + 2\mathbf{j} + 2\mathbf{k}$

1.4.8a : $2\mathbf{j} - \mathbf{k}$

1.4.8c : $3\mathbf{j} - 3\mathbf{k}$

1.4.8e : $-7\mathbf{i} + 10\mathbf{j} + \mathbf{k}$

1.4.12a : 12

1.4.12c : 10

1.4.12e : 6

1.5.1a : 7

1.5.1c : 11

1.5.1e : 13

1.5.1g : 4.1725

1.5.3a : They are all orthogonal.

1.5.3c : $\theta_{uv} = 142.78^\circ$, $\theta_{uw} = 146.44^\circ$, $\theta_{vw} = 12.10^\circ$

1.5.3e : $\theta_{ab} = 42.82^\circ$, $\theta_{ac} = 99.11^\circ$, $\theta_{ad} = 109.89^\circ$, $\theta_{bc} = 64.32^\circ$,
 $\theta_{bd} = 92.89^\circ$, $\theta_{cd} = 61.70^\circ$

2 Systems of linear equations

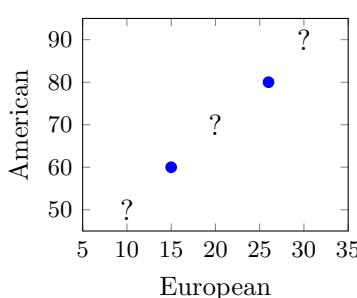
Chapter Contents

2.1	Introduction to systems of linear equations	86
2.1.1	Exercises	92
2.2	Directly solve linear systems	95
2.2.1	Compute a system's solution	95
2.2.2	Algebraic manipulation solves systems	106
2.2.3	Three possible numbers of solutions	113
2.2.4	Exercises	116
2.3	Linear combinations span sets	122
2.3.1	Exercises	127

Linear relationships are commonly identified in science and engineering, and are commonly expressed as linear equations. One of the reasons is that scientists and engineers can do amazingly powerful algebraic transformations with linear equations. Such transformations and their practical implications are the subject of this book.

One vital use in science and engineering is in the scientific task of taking scattered experimental data and inferring a general algebraic relation between the quantities measured. In computing science this task is often called ‘data mining’, ‘knowledge discovery’ or even ‘artificial intelligence’—although the algebraic relation is then typically discussed as a computational procedure. But appearing within these tasks is always linear equations to be solved.

I am sure you can guess where we are going with this example, but let's pretend we do not know.



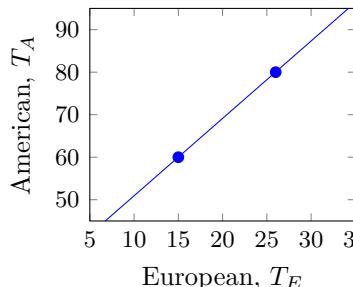
Example 2.0.1 (scientific inference). Two colleagues, an American and a European, discuss the weather; in particular, they discuss the temperature. The American says “yesterday the temperature was 80° but today is much cooler at 60° ”. The European says, “that's not what I heard, I heard the temperature was 26° and today is 15° ”. (The marginal figure plots these two data points.) “Hmmm, we must be using a different temperature scale”, they say. Being scientists they start to use linear algebra to *infer*, from the two days of temperature data, a general relation between their temperature scales—a relationship valid over a wide range of temperatures (denoted by the question marks in the marginal figure). Let's assume that, in terms of the European temperature T_E , the

American temperature $T_A = cT_E + d$ for some constants c and d they and we aim to find. The two days of data then give that

$$80 = c26 + d \quad \text{and} \quad 60 = c15 + d.$$

To find c and d :

- subtract the second equation from the first to deduce $80 - 60 = 26c + d - 15c - d$ which simplifies to $20 = 11c$, that is, $c = 20/11 = 1.82$ to two decimal places;
- use this value of c in either equation, say the second, gives $60 = \frac{20}{11}15 + d$ which rearranges to $d = 360/11 = 32.73$ to two decimal places.



We deduce that the temperature relationship is $T_A = 1.82 T_E + 32.73$ (as plotted in the marginal figure). The two colleagues now *predict* that they will be able to use this formula to translate their temperature into that of the other, and vice versa.

You may quite rightly object that the two colleagues *assumed* a linear relation, they do *not know* it is linear. You may also object that the predicted relation is erroneous as it should be $T_A = \frac{9}{5}T_E + 32$ (the relation between Celsius and Fahrenheit). Absolutely, you should object. Scientifically, the deduced relation $T_A = 1.82 T_E + 32.73$ is only a conjecture that fits the known data. More data and more linear algebra together empower us to both confirm the linearity (or not as the case may be), and also to improve the accuracy of the coefficients. This is fundamental scientific methodology—and central to it is the algebra of linear equations. ■

Linear algebra and equations are also crucial for nonlinear relationships. Figure 2.1 shows four plots of the same nonlinear function, but on successively smaller scales. Zooming in on the point $(0, 1)$ we see the curve looks straighter and straighter until on the microscale it is effectively a straight line. The same is true for anywhere on any smooth nonlinear curve: we discover that the curve looks like a straight line on the microscale. Thus we may view any smooth nonlinear function as roughly being made up of lots of microscale straight line segments. Linear equations and their algebra on this microscale empower our understanding of nonlinear relations—for example, microscale linearity underwrites all of calculus.

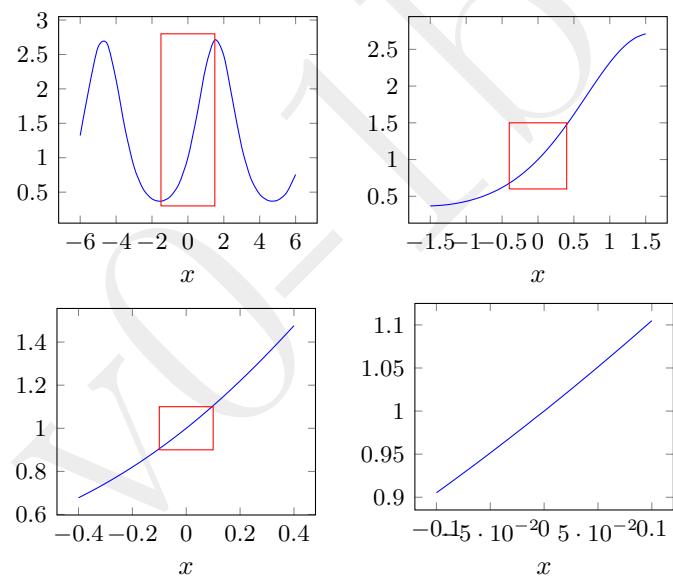


Figure 2.1: zoom in anywhere on any smooth nonlinear curve, such as the plotted $f(x)$, and we discover that the curve looks like a straight line on the microscale. The (red) rectangles show the region plotted in the next graph in the sequence.

Table 2.1: examples of linear equations, and equations that are not linear (called nonlinear equations).

linear	nonlinear
$-3x + 2 = 0$	$x^2 - 3x + 2 = 0$
$2x - 3y = -1$	$2xy = 3$
$-1.2x_1 + 3.4x_2 - x_3 = 5.6$	$x_1^2 + 2x_2^2 = 4$
$r - 5s = 2 - 3s + 2t$	$r/s = 2 + t$
$\sqrt{3}t_1 + \frac{\pi}{2}t_2 - t_3 = 0$	$3\sqrt{t_1} + t_2^3/t_3 = 0$
$(\cos \frac{\pi}{6})x + e^2y = 1.23$	$x + e^{2y} = 1.23$

2.1 Introduction to systems of linear equations

Section Contents

2.1.1 Exercises	92
---------------------------	----

The great aspect of linear equations is that we can straightforwardly manipulate them algebraically to deduce results: some results are not only immensely useful in applications but also in further theory.

Example 2.1.1 (simple algebraic manipulation). Following Example 2.0.1, recall that the temperature in Fahrenheit $T_F = \frac{9}{5}T_C + 32$ in terms of the temperature in Celsius, T_C . Straightforward algebra answers the following questions.

- What is a formula for the Celsius temperature as a function of the temperature in Fahrenheit? Answer by rearranging the equation: subtract 32 from both sides, $T_F - 32 = \frac{9}{5}T_C$; multiply both sides by $\frac{5}{9}$, then $\frac{5}{9}(T_F - 32) = T_C$; that is, $T_C = \frac{5}{9}T_F - \frac{160}{9}$.
- What temperature has the same *numerical value* in the two scales? That is, when is $T_F = T_C$? Answer by algebra: we want $T_C = T_F = \frac{9}{5}T_C + 32$; subtract $\frac{9}{5}T_C$ from both sides to give $-\frac{4}{5}T_C = 32$; multiply both sides by $-\frac{5}{4}$, then $T_C = -\frac{5}{4} \times 32 = -40$.

■

Linear equations are characterised by each unknown never being multiplied or divided by another unknown, or itself, nor inside ‘curvaceous’ functions. Table 2.1 lists examples of both. The power of linear algebra is especially important for large numbers of unknown variables: generally we say there are n variables.

Definition 2.1.2. A *linear equation* in the n variables x_1, x_2, \dots, x_n is an equation that can be written in the form

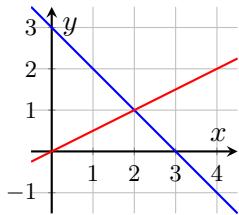
$$a_1x_1 + a_2x_2 + \cdots + a_nx_n = b$$

where the **coefficients** a_1, a_2, \dots, a_n and the **constant term** b are given scalar constants. An equation that cannot be written in this form is called a **nonlinear equation**. A **system** of linear equations is a set of one or more linear equations in one or more variables (usually more than one).

Example 2.1.3 (two equations in two variables). Graphically and algebraically solve each of the following systems.

$$(a) \begin{aligned} x + y &= 3 \\ 2x - 4y &= 0 \end{aligned}$$

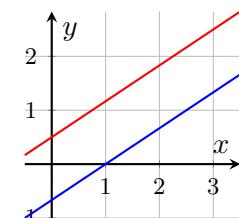
Solution: To draw the graphs seen in the marginal plot, rearrange the linear equations as $y = 3 - x$ and $y = x/2$. They intersect at the point $(2, 1)$ so $x = 2$ and $y = 1$ is the unique solution.



Algebraically, one could add twice the first equation to half of the second equation: $2(x + y) + \frac{1}{2}(2x - 4y) = 2 \cdot 3 + \frac{1}{2} \cdot 0$ which simplifies to $3x = 6$ as the y terms cancel; hence $x = 2$. Then say consider the second equation, $2x - 4y = 0$, which now becomes $2 \cdot 2 - 4y = 0$, that is, $y = 1$. This algebra gives the same solution $(x, y) = (2, 1)$ as graphically.

$$(b) \begin{aligned} 2x - 3y &= 2 \\ -4x + 6y &= 3 \end{aligned}$$

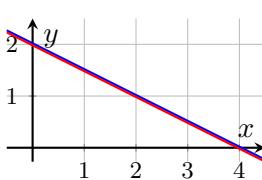
Solution: To draw the graphs seen in the marginal plot, rearrange the linear equations as $y = \frac{2}{3}x - \frac{2}{3}$ and $y = \frac{2}{3}x + \frac{1}{2}$. Evidently these lines never intersect, they are parallel, so there appears to be no solution.



Algebraically, one could add twice the first equation to the second equation: $2(2x - 3y) + (-4x + 6y) = 2 \cdot 2 + 3$ which, as the x and y terms cancel, simplifies to $0 = 7$. This equation is a contradiction as zero is not equal to seven. Thus there are no solutions to the system.

$$(c) \begin{aligned} x + 2y &= 4 \\ 2x + 4y &= 8 \end{aligned}$$

Solution: To draw the graphs seen in the marginal plot, rearrange the linear equations as $y = 2 - x/2$ and $y = 2 - x/2$. They are the same line so every point on this line is a solution of the system. There are an infinite number of possible solutions.

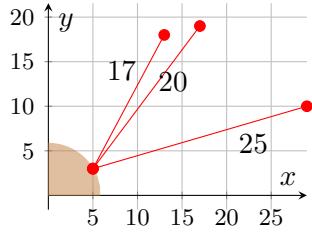


Algebraically, the rearrangement of both equations to exactly the same $y = 2 - x/2$ establishes an infinite number of solutions, here parametrised by x .

Example 2.1.4 (Global Positioning System). The Global Positioning System (GPS) is a network of 24 satellites orbiting the Earth. Each satellite knows very accurately its position at all times, and broadcasts this position by radio. Receivers, such as smart-phones, pick up these signals and from the time taken for the signals to arrive know the distance to those satellites within ‘sight’. The receivers solve a system of equations and inform you of their precise position.

Let’s solve a definite example problem, but in two dimensions for simplicity. Suppose you and your smart-phone are at some unknown location (x, y) in the 2D-plane, on the Earth’s surface where the Earth has radius about 6 Mm (here all distances are measured in units of Megametres, Mm, thousands of km). But your smart-phone picks up the broadcast from three GPS satellites, and then determines their distance from you. From the broadcast and the timing, suppose you then know that a satellite at $(29, 10)$ is 25 away, one at $(17, 19)$ is 20 away, and one at $(13, 18)$ is 17 away (as drawn in the margin). Find your location (x, y) .

Solution: From these three sources of information, Pythagoras and the length of displacement vectors gives the three equations



$$\begin{aligned}(x - 29)^2 + (y - 10)^2 &= 25^2, \\ (x - 17)^2 + (y - 19)^2 &= 20^2, \\ (x - 13)^2 + (y - 18)^2 &= 17^2.\end{aligned}$$

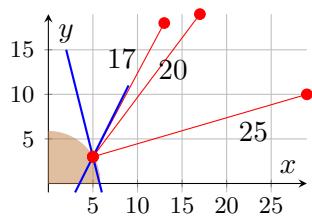
These three equations constrain your as yet unknown location (x, y) . Expanding the squares in these equations gives the equivalent system

$$\begin{aligned}x^2 - 58x + 841 + y^2 - 20y + 100 &= 625, \\ x^2 - 34x + 289 + y^2 - 38y + 361 &= 400, \\ x^2 - 26x + 169 + y^2 - 36y + 324 &= 289.\end{aligned}$$

Involving squares of the unknowns, these are a nonlinear system of equations and so appear to lie outside the remit of this book. However, straightforward algebra transforms these three nonlinear equations into a system of two linear equations which we solve.

Let’s subtract the third equation from each of the other two, then the nonlinear squared terms cancel giving a system of two linear equations in two variables:

$$\begin{aligned}-32x + 672 + 16y - 224 &= 336 \iff -2x + y = -7, \\ -8x + 120 - 2y + 37 &= 111 \iff -4x - y = -23.\end{aligned}$$



Graphically, include these two lines to the picture, namely $y = -7 + 2x$ and $y = 23 - 4x$, and then their intersection gives your location.

Algebraically, one could add the two equations together: $(-2x + y) + (-4x - y) = -7 - 23$ which reduces to $-6x = -30$, that is, $x = 5$.

Then either equation, say the first, determines $y = -7 + 2x = -7 + 2 \cdot 5 = 3$. That is, your location is $(x, y) = (5, 3)$ (in Mm), as drawn.

If the x -axis is a line through the equator, and the y -axis goes through the North pole, then your location would be at latitude $\tan^{-1} \frac{3}{5} = 30.96^\circ\text{N}$. ■

Example 2.1.5 (three equations in three variables). Graph the surfaces and algebraically solve the system

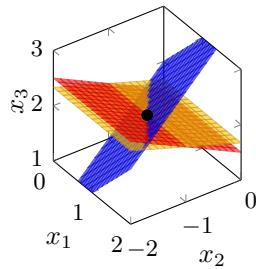
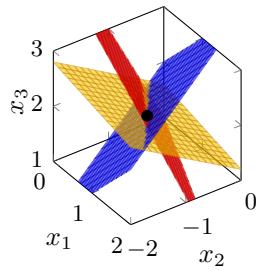
$$\begin{aligned}x_1 + x_2 - x_3 &= -2, \\x_1 + 3x_2 + 5x_3 &= 8, \\x_1 + 2x_2 + x_3 &= 1.\end{aligned}$$

Solution: The marginal plot shows the three planes represented by the given equations (in the order blue, brown, red), and plots the (black) point we seek of intersection of all three planes.

Algebraically we combine and manipulate the equations in a sequence of steps designed to simplify the form of the system. *By doing the same manipulation to the whole of each of the equations, we ensure the validity of the result.*

- (a) Subtract the first equation from each of the other two equations to deduce (as illustrated)

$$\begin{aligned}x_1 + x_2 - x_3 &= -2, \\2x_2 + 6x_3 &= 10, \\x_2 + 2x_3 &= 3.\end{aligned}$$

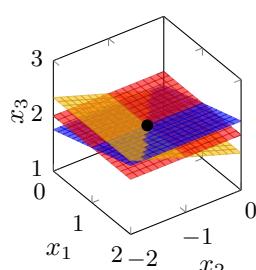


- (b) Divide the second equation by two:

$$\begin{aligned}x_1 + x_2 - x_3 &= -2, \\x_2 + 3x_3 &= 5, \\x_2 + 2x_3 &= 3.\end{aligned}$$

- (c) Subtract the second equation from each of the other two (as illustrated):

$$\begin{aligned}x_1 - 4x_3 &= -7, \\x_2 + 3x_3 &= 5, \\-x_3 &= -2.\end{aligned}$$

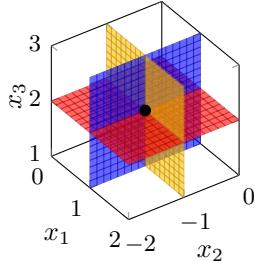


- (d) Multiply the third equation by (-1) :

$$\begin{aligned}x_1 - 4x_3 &= -7, \\x_2 + 3x_3 &= 5, \\x_3 &= 2.\end{aligned}$$

- (e) Add four times the third equation to the first, and subtract three times it from the second (as illustrated):

$$\begin{aligned}x_1 &= 1, \\x_2 &= -1, \\x_3 &= 2.\end{aligned}$$



Thus the only solution to this system of three linear equations in three variables is $(x_1, x_2, x_3) = (1, -1, 2)$. ■

The sequence of marginal graphs in the previous Example 2.1.5 illustrate the equations at each main step in the algebraic manipulations. Apart from keeping the solution point fixed, the sequence of graphs looks rather chaotic. Indeed there is no particular geometric pattern or interpretation of the steps in this algebra. One feature of Section 3.3 is that we discover how the so-called ‘singular value decomposition’ solves linear equations via a sound method with a strong geometric interpretation. This geometric interpretation then empowers further methods useful in applications.

Transform into abstract setting Linear algebra has an important aspect crucial in applications. A crucial skill in applying linear algebra is that it takes an application problem and transforms it into an abstract setting. Example 2.0.1 transformed the problem of inferring a line through two data points into solving two linear equations. The next Example 2.1.6 similarly transforms the problem of inferring a plane through three data points into solving three linear equations. The original application is often not easily recognisable in the abstract version. Nonetheless, it is the abstraction by linear algebra that empowers immense results.

Example 2.1.6 (infer a surface through three points). This example illustrates the previous paragraph. Given a geometric problem of inferring what plane passes through three given points, we transform this problem into the linear algebra task of finding the intersection point of three given planes. This task we do.

Suppose we observe that at some given temperature and humidity we get some rainfall: let’s find a formula that predicts the rainfall from temperature and humidity measurements. In some *completely artificial units*, Table 2.2 lists measured temperature (‘temp’), humidity (‘humid’), and rainfall (‘rain’).

Solution: To infer a relation to hold generally—to fill in the gaps between the known measurements, seek ‘rainfall’ to be predicted by the linear formula

$$(\text{‘rain’}) = x_1 + x_2(\text{‘temp’}) + x_3(\text{‘humid’}),$$

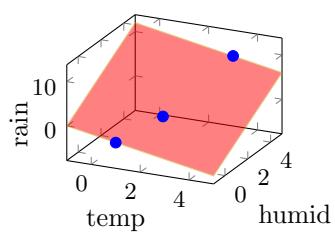


Table 2.2: in some artificial units, this table lists measured temperature, humidity, and rainfall.

‘temp’	‘humid’	‘rain’
1	-1	-2
3	5	8
2	1	1

for some coefficients x_1 , x_2 and x_3 to be determined. The measured data of Table 2.2 constrains and determines these coefficients: substitute each triple of measurements to require

$$\begin{aligned} -2 &= x_1 + x_2(1) + x_3(-1), & x_1 + x_2 - x_3 &= -2, \\ 8 &= x_1 + x_2(3) + x_3(5), & \iff & x_1 + 3x_2 + 5x_3 = 8, \\ 1 &= x_1 + x_2(2) + x_3(1), & x_1 + 2x_2 + x_3 &= 1. \end{aligned}$$

The previous Example 2.1.5 solves this set of three linear equations in three unknowns to determine the solution that the coefficients $(x_1, x_2, x_3) = (1, -1, 2)$. That is, the requisite formula to infer rain from any given temperature and humidity is

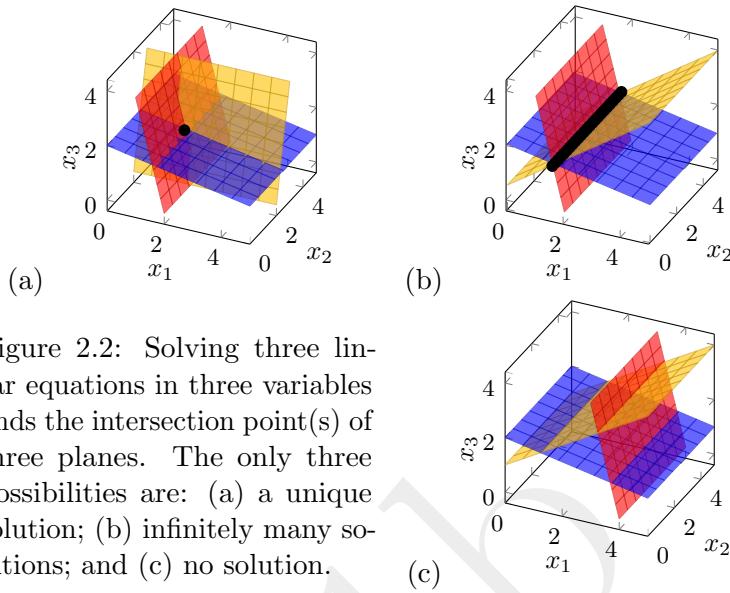
$$(\text{‘rain’}) = 1 - (\text{‘temp’}) + 2(\text{‘humid’}).$$

This example illustrates that the geometry of fitting a plane to three points (as plotted) translates into the abstract geometry of finding the intersection of three planes (plotted in the previous example). The linear algebra procedure for this latter abstract problem then gives the required ‘physical’ solution.

■

The solution of three linear equations in three variables leads to finding the intersection point of three planes. Figure 2.1 illustrates the three general possibilities: a unique solution (as in Example 2.1.5), or infinitely many solutions, or no solution. The solution of two linear equations in two variables also has the same three possibilities—as deduced and illustrated in Example 2.1.3. The next section establishes the general key property of a system of any number of linear equations in any number of variables: the system has either

- a unique solution (a consistent system), or
- infinitely many solutions (a consistent system), or
- no solutions (an inconsistent system).



2.1.1 Exercises

Exercise 2.1.1. Graphically and algebraically solve each of the following systems.

$$(a) \begin{array}{l} x - 2y = -3 \\ -4x = -4 \end{array}$$

$$(b) \quad \begin{aligned} x + 2y &= 5 \\ 6x - 2y &= 2 \end{aligned}$$

$$(c) \quad \begin{aligned} x - y &= 2 \\ -2x + 7y &= -4 \end{aligned}$$

$$(d) \quad \begin{aligned} 3x - 2y &= 2 \\ -3x + 2y &= -2 \end{aligned}$$

$$(e) \quad \begin{aligned} 3x - 2y &= 1 \\ 6x - 4y &= -2 \end{aligned}$$

$$(f) \quad \begin{aligned} 4x - 3y &= -1 \\ -5x + 4y &= 1 \end{aligned}$$

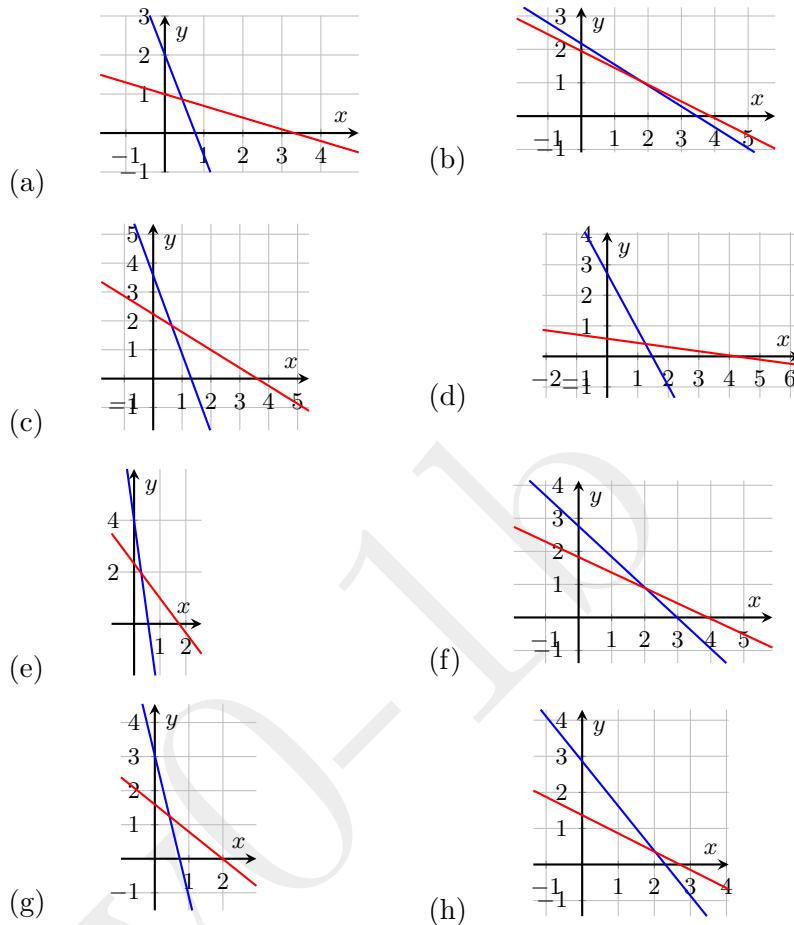
$$\begin{array}{l} \text{(g)} \quad p + q = 3 \\ \qquad -p - q = 2 \end{array}$$

$$(h) \quad \begin{aligned} p - q &= 1 \\ -3p + 5q &= -4 \end{aligned}$$

$$(1) \quad u - v = -1$$

$$(\text{J}) \quad -u - v = 1$$

Exercise 2.1.2. For each of the following graphs: estimate the equations of the pair of lines; solve the pair of equations algebraically; and confirm the algebraic solution is reasonably close to the intersection of the pair of lines.



Exercise 2.1.3. Graphically and algebraically solve each of the following systems of three equations for the two unknowns.

$$(a) \begin{aligned} 4x + y &= 8 \\ 3x - 3y &= -\frac{3}{2} \\ -4x + 2y &= -2 \end{aligned} \quad (b) \begin{aligned} -4x + 3y &= \frac{7}{2} \\ 7x + y &= -3 \\ x - 2y &= \frac{3}{2} \end{aligned}$$

$$(c) \begin{aligned} 2x + 2y &= 2 \\ -3x - 3y &= -3 \\ x + y &= 1 \end{aligned} \quad (d) \begin{aligned} -2x - 4y &= 3 \\ x + 2y &= 3 \\ -4x - 8y &= -6 \end{aligned}$$

$$(e) \begin{aligned} 3x + 2y &= 4 \\ -2x - 4y &= -4 \\ 4x + 2y &= 5 \end{aligned} \quad (f) \begin{aligned} -2x + 3y &= -3 \\ -5x + 2y &= -9 \\ 3x + 3y &= 6 \end{aligned}$$

Exercise 2.1.4 (Global Positioning System in 2D). For each case below, and in two dimensions, suppose you know from three GPS satellites that you and your GPS receiver are given distances away from the given locations of each of the three satellites (locations and distance are in Mm). Following Example 2.1.4, determine your position.

- | | |
|---|---|
| (a) 26 from (10 , 30)
29 from (20 , 27) | (b) 25 from (11 , 29)
26 from (28 , 15)
20 from (16 , 21) |
| (c) 20 from (22 , 12)
26 from (16 , 24)
29 from (26 , 21) | (d) 17 from (12 , 21)
25 from (10 , 29)
26 from (27 , 15) |

In which of these cases: are you at the ‘North Pole’? flying high above the Earth? the measurement data is surely in error?

VO-1b

2.2 Directly solve linear systems

Section Contents

2.2.1	Compute a system's solution	95
2.2.2	Algebraic manipulation solves systems	106
2.2.3	Three possible numbers of solutions	113
2.2.4	Exercises	116

The previous Section 2.1 solved some example systems of linear equations by hand algebraic manipulation. We continue to do so for small systems. However, such by-hand solutions are tedious for systems bigger than say four equations in four unknowns. For bigger systems—which are typical in applications—we use computers to find solutions because computers are ideal for tedious manipulations.

2.2.1 Compute a system's solution

It is unworthy of excellent persons to lose hours like
slaves in the labour of calculation.

Gottfried Wilhelm von Leibniz

Computers primarily deal with numbers, not algebraic equations, so we have to abstract the coefficients of a system into a numerical data structure. We use matrices and vectors.

Example 2.2.1. The first system of Example 2.1.3a

$$\begin{aligned} x + y &= 3 \\ 2x - 4y &= 0 \end{aligned} \quad \text{is written} \quad \underbrace{\begin{bmatrix} 1 & 1 \\ 2 & -4 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x \\ y \end{bmatrix}}_x = \underbrace{\begin{bmatrix} 3 \\ 0 \end{bmatrix}}_b.$$

That is, the system $\begin{cases} x + y = 3 \\ 2x - 4y = 0 \end{cases}$ is equivalent to $A\mathbf{x} = \mathbf{b}$ for the so-called coefficient matrix $A = \begin{bmatrix} 1 & 1 \\ 2 & -4 \end{bmatrix}$, right-hand side vector $\mathbf{b} = (3, 0)$, and vector of variables $\mathbf{x} = (x, y)$. ■

In this chapter, the two character symbol ' $A\mathbf{x}$ ' is just a shorthand for all the left-hand sides of the linear equations in a system. However, Section 3.1 defines a useful multiplicative meaning to this composite symbol.

The beauty of the form $A\mathbf{x} = \mathbf{b}$ is that the numbers involved in the system are abstracted into the matrix A and vector \mathbf{b} : Matlab/Octave handles such numerical matrices and vectors. For some of you, writing a system in this matrix-vector form $A\mathbf{x} = \mathbf{b}$ (Definition 2.2.2 below) will appear to be just some mystic rearrangement of symbols—such an interpretation is sufficient for this chapter. However, those of you who have met matrix multiplication will recognise that $A\mathbf{x} = \mathbf{b}$ is an expression involving natural operations

for matrices and vectors: Section 3.1 defines and explores such useful operations.

Definition 2.2.2 (matrix-vector form). *For a given system of m linear equations in n variables*

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1, \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2, \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &= b_m, \end{aligned}$$

its **matrix-vector form** is $A\mathbf{x} = \mathbf{b}$ for the $m \times n$ matrix of coefficients

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix},$$

and vectors $\mathbf{x} = (x_1, x_2, \dots, x_n)$ and $\mathbf{b} = (b_1, b_2, \dots, b_m)$. If $m = n$ (the number of equations is the same as the number of variables), then A is called a **square matrix** (the number of rows is the same as the number of columns).

Example 2.2.3 (matrix-vector form).

Write the following systems in matrix-vector form.

$$(a) \begin{aligned} x_1 + x_2 - x_3 &= -2, \\ x_1 + 3x_2 + 5x_3 &= 8, \\ x_1 + 2x_2 + x_3 &= 1. \end{aligned} \quad (b) \begin{aligned} -2r + 3s &= 6, \\ s - 4t &= -\pi. \end{aligned}$$

Solution: (a) The first system, that of Example 2.1.5, is of three equations in three variables ($m = n = 3$) and is written in the form $A\mathbf{x} = \mathbf{b}$ as

$$\underbrace{\begin{bmatrix} 1 & 1 & -1 \\ 1 & 3 & 5 \\ 1 & 2 & 1 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_x = \underbrace{\begin{bmatrix} -2 \\ 8 \\ 1 \end{bmatrix}}_b$$

for square matrix A .

- (b) The second system has three variables called r , s and t and two equations. Variables ‘missing’ from an equation are represented as zero times that variable, thus the system

$$\begin{aligned} -2r + 3s + 0t &= 6, \\ 0r + s - 4t &= -\pi, \end{aligned} \quad \text{is} \quad \underbrace{\begin{bmatrix} -2 & 3 & 0 \\ 0 & 1 & -4 \end{bmatrix}}_A \underbrace{\begin{bmatrix} r \\ s \\ t \end{bmatrix}}_x = \underbrace{\begin{bmatrix} 6 \\ -\pi \end{bmatrix}}_b$$

for 2×3 matrix A .



Procedure 2.2.4 (unique solution). *In Matlab/Octave, to solve the system $Ax = b$ for a square matrix A , use commands listed in Table 1.2 and 2.3 to:*

1. *form matrix A and column vector b ;*
2. *check $\text{rcond}(A)$ exists and is not too small, $1 \geq \text{good} > 10^{-2} > \text{poor} > 10^{-4} > \text{bad} > 10^{-8} > \text{terrible}$, (rcond is always between zero and one inclusive);*
3. *if rcond both exists and is acceptable, then execute $x=A\backslash b$ to compute the solution vector x .*

Checking $\text{rcond}()$ avoids gross mistakes such as “Lies My Computer Told Me” (Poole 2015, p.83). Section 3.3.2 discovers what rcond is, and why rcond avoids mistakes.¹ In practice, decisions about acceptability are rarely black and white, and so the qualitative ranges of rcond reflects practical reality.

In theory, there is no difference between theory and practice. But, in practice, there is.

Jan L. A. van de Snepscheut

Example 2.2.5. Use Matlab/Octave to solve the system (from Example 2.1.5)

$$\begin{aligned} x_1 + x_2 - x_3 &= -2, \\ x_1 + 3x_2 + 5x_3 &= 8, \\ x_1 + 2x_2 + x_3 &= 1. \end{aligned}$$

Solution: Begin by writing the system in the abstract matrix-vector form $Ax = b$ as already done by Example 2.2.3. Then the three steps of Procedure 2.2.4 are the following.

- (a) Form matrix A and column vector b with the Matlab/Octave assignments

```
A=[1 1 -1; 1 3 5; 1 2 1]
b=[-2;8;1]
```

Table 2.3 summarises that in Matlab/Octave: each line is one command; the $=$ symbol assigns the value of the right-hand expression to the variable name of the left-hand side²; and the brackets $[]$ construct matrices and vectors.

¹ Interestingly, there are incredibly rare pathological matrices for which $\text{rcond}()$ and $A\backslash$ fails us (Driscoll & Maki 2007). For example, among 32×32 matrices the probability is about 10^{-22} of encountering a matrix for which $\text{rcond}()$ misleads us by more than a factor of a hundred in using $A\backslash$.

² Beware that the symbol “=” in Matlab/Octave is a procedural assignment of a value, which is quite different in nature to the “=” in algebra which denotes equality.

Table 2.3: To realise Procedure 2.2.4, and other procedures, we need these basics of Matlab/Octave as well as that of Table 1.2.

- The floating point numbers are extended by `Inf`, denoting ‘infinity’, and `NaN`, denoting ‘not a number’ such as the indeterminate $0/0$.
- `[... ; ... ; ...]` forms both matrices and vectors, or use newlines instead of the semi-colons.
- `rcond(A)` of a square matrix A estimates the reciprocal of the so-called condition number (defined precisely by Definition 3.3.14).
- `x=A\b` computes an ‘answer’ to $Ax = b$ —but it may not be a solution unless `rcond(A)` exists and is not small;
- Change an element of an array or vector by assigning a new value with assignments `A(i,j)=...` or `b(i)=...` where `i` and `j` denote some indices.
- For a vector (or matrix) `t` and an exponent `p`, the operation `t.^p` computes the p th power of each element in the vector; for example, if `t=[1;2;3;4;5]` then `t.^2` computes `[1;4;9;16;25]`.
- The function `ones(m,1)` gives a column vector of m ones, $(1, 1, \dots, 1)$.
- Lastly, remember that ‘the answer’ by a computer is not necessarily ‘the solution’ of your problem.

- (b) Check `rcond(A)`: here it is 0.018 which is in the good range.
 (c) Since `rcond` is acceptable, then execute `x=A\b` to compute the solution vector $\mathbf{x} = (1, -1, 2)$ (and assign it to the variable `x`, see Table 2.3).

All together that is the four commands

```
A=[1 1 -1; 1 3 5; 1 2 1]
b=[-2;8;1]
rcond(A)
x=A\b
```

Such qr-codes in the margin encodes these commands for you to possibly scan, copy and paste into Matlab/Octave.



Example 2.2.6. Following the previous Example 2.2.5, solve each of the two systems:

$$(a) \begin{aligned} x_1 + x_2 - x_3 &= -2, \\ x_1 + 3x_2 + 5x_3 &= 5, \\ x_1 - 3x_2 + x_3 &= 1; \end{aligned} \quad (b) \begin{aligned} x_1 + x_2 - x_3 &= -2, \\ x_1 + 3x_2 - 2x_3 &= 5, \\ x_1 - 3x_2 + x_3 &= 1. \end{aligned}$$

Solution: Begin by writing, or at least by imaging, each system

in matrix-vector form:

$$\underbrace{\begin{bmatrix} 1 & 1 & -1 \\ 1 & 3 & 5 \\ 1 & -3 & 1 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_x = \underbrace{\begin{bmatrix} -2 \\ 5 \\ 1 \end{bmatrix}}_b ; \quad \underbrace{\begin{bmatrix} 1 & 1 & -1 \\ 1 & 3 & -2 \\ 1 & -3 & 1 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_x = \underbrace{\begin{bmatrix} -2 \\ 5 \\ 1 \end{bmatrix}}_b .$$

As the matrices and vectors are modifications of the previous Example 2.2.5 we reduce typing by modifying the matrix and vector of the previous example (using the ability to change each element in a matrix, see Table 2.3).

- (a) For the first system execute $A(3,2)=-3$ and $b(2)=5$ to see the matrix and vector are now

$$\begin{aligned} A &= \\ &\begin{matrix} 1 & 1 & -1 \\ 1 & 3 & 5 \\ 1 & -3 & 1 \end{matrix} \\ b &= \\ &\begin{matrix} -2 \\ 5 \\ 1 \end{matrix} \end{aligned}$$

Check: $\text{rcond}(A)$ is 0.14 which is good, so obtain the solution from $x=A\backslash b$, namely

$$\begin{aligned} x &= \\ &\begin{matrix} -0.6429 \\ -0.1429 \\ 1.2143 \end{matrix} \end{aligned}$$

That is, the solution $x = (-0.64, -0.14, 1.21)$ to two decimal places (2 d.p.).³

- (b) For the second system now execute $A(2,3)=-2$ to see the new matrix is the required

$$\begin{aligned} A &= \\ &\begin{matrix} 1 & 1 & -1 \\ 1 & 3 & -2 \\ 1 & -3 & 1 \end{matrix} \end{aligned}$$

Check: find that $\text{rcond}(A)$ is zero which is classified as terrible. Consequently we cannot compute a solution of this second system of linear equations (as in Figure 2.1(c)).

If we were to try $x=A\backslash b$ in this second system, then Matlab/Octave would report

³The four or five significant digits printed by Matlab/Octave is effectively exact for most practical purposes. This text often reports two significant digits as two is enough for most human readable purposes. When a numerical result is reported to two decimal places, the text indicates this truncation with “(2 d.p.)”.

Warning: Matrix is singular to working precision.

However, one cannot rely on Matlab/Octave producing such useful messages: use `rcond` to almost always avoid mistakes.

■

Example 2.2.7. Use Matlab/Octave to solve the system

$$\begin{aligned}x_1 - 2x_2 + 3x_3 + x_4 + 2x_5 &= 7, \\-2x_1 - 6x_2 - 3x_3 - 2x_4 + 2x_5 &= -1, \\2x_1 + 3x_2 - 2x_5 &= -9, \\-2x_1 + x_2 &= -3, \\-2x_1 - 2x_2 + x_3 + x_4 - 2x_5 &= 5.\end{aligned}$$

Solution: Following Procedure 2.2.4, form the corresponding matrix and vector as

```
A=[1 -2 3 1 2
   -2 -6 -3 -2 2
   2 3 0 0 -2
   -2 1 0 0 0
   -2 -2 1 1 -2]
b=[7;-1;-9;-3;5]
```

Check: find `rcond(A)` is acceptably 0.020 so compute the solution via `x=A\b` to find

```
x =
  0.8163
 -1.3673
 -6.7551
 17.1837
  3.2653
```



that is, the solution $\mathbf{x} = (0.82, -1.37, -6.76, 17.18, 3.27)$ (2 d.p.).

■

Example 2.2.8. What system of linear equations are represented by the following matrix-vector expression? and what is the result of using Procedure 2.2.4 for this system?

$$\begin{bmatrix} -7 & 3 \\ 7 & -5 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix} = \begin{bmatrix} 3 \\ -2 \\ 1 \end{bmatrix}.$$

Solution: The corresponding system of linear equations is

$$\begin{aligned}-7y + 3z &= 3, \\7y - 5z &= -2, \\y - 2z &= 1.\end{aligned}$$

Invoking Procedure 2.2.4:

- (a) form matrix A and column vector \mathbf{b} with

$$\mathbf{A} = [-7 \ 3; \ 7 \ -5; \ 1 \ -2]$$

$$\mathbf{b} = [3; -2; 1]$$

- (b) check `rcond(A)`: Matlab/Octave gives the message

Error using `rcond`

Input must be a square matrix.

As `rcond` does not exist, the procedure cannot give a solution.

The reason for the procedure not leading to a solution is that a system of three equations in two variables, as here, generally does not have a solution.⁴

■

Example 2.2.9 (partial fraction decomposition). Recall that mathematical analysis sometimes needs to separate a rational function into a sum of simpler ‘partial’ fractions. For example, for some purposes $\frac{3}{(x-1)(x+2)}$ needs to be written as $\frac{1}{x-1} - \frac{1}{x+2}$. Solving linear equations helps: here pose that $\frac{3}{(x-1)(x+2)} = \frac{A}{x-1} + \frac{B}{x+2}$ for some unknown A and B ; then write the right-hand side over the common denominator,

$$\frac{A}{x-1} + \frac{B}{x+2} = \frac{A(x+2) + B(x-1)}{(x-1)(x+2)} = \frac{(A+B)x + (2A-B)}{(x-1)(x+2)}$$

and this equals $\frac{3}{(x-1)(x+2)}$ only if $A+B=0$ and $2A-B=3$; solving these two linear equations gives the required $A=1$ and $B=-1$ to determine the decomposition $\frac{3}{(x-1)(x+2)} = \frac{1}{x-1} - \frac{1}{x+2}$.

Now find the partial fraction decomposition of $\frac{-4x^3+8x^2-5x+2}{x^2(x-1)^2}$.

Solution: Recalling that repeated factors require extra terms in the decomposition, seek a decomposition of the form

$$\begin{aligned} & \frac{A}{x^2} + \frac{B}{x} + \frac{C}{(x-1)^2} + \frac{D}{x-1} \\ &= \frac{A(x-1)^2 + Bx(x-1)^2 + Cx^2 + Dx^2(x-1)}{x^2(x-1)^2} \\ &= \frac{(B+D)x^3 + (A-2B+C-D)x^2 + (-2A+B)x + (A)}{x^2(x-1)^2} \\ &= \frac{-4x^3 + 8x^2 - 5x + 2}{x^2(x-1)^2}. \end{aligned}$$

⁴ If one were to execute `x=A\b`, then you would find Matlab/Octave gives the ‘answer’ $\mathbf{x} = (-0.77, -0.73)$ (2 d.p.). But this answer is not a solution. Instead this answer has another meaning, often sensibly useful, which is explained by Section 3.5. Using `rcond` avoids us confusing such an answer with a solution.

For this last equality to hold for all x the coefficients of various powers of x must be equal: this leads to the linear equation system

$$\begin{aligned} B + D &= -4 \\ A - 2B + C - D &= 8 \\ -2A + B &= -5 \\ A &= 2 \end{aligned} \iff \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & -2 & 1 & -1 \\ -2 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} A \\ B \\ C \\ D \end{bmatrix} = \begin{bmatrix} -4 \\ 8 \\ -5 \\ 2 \end{bmatrix}.$$

Either solve by hand or by computer.

- By hand, the last equation gives $A = 2$ so the third equation then gives $B = -1$. Then the first gives $D = -3$. Lastly, the second then gives $C = 8 - 2 + 2(-1) + (-3) = 1$. That is, the decomposition is

$$\frac{-4x^3 + 8x^2 - 5x + 2}{x^2(x-1)^2} = \frac{2}{x^2} - \frac{1}{x} + \frac{1}{(x-1)^2} - \frac{3}{x-1}$$

- Using Matlab/Octave, form the matrix and right-hand side with

```
a=[0 1 0 1
    1 -2 1 -1
    -2 1 0 0
    1 0 0 0]
b=[-4;8;-5;2]
```



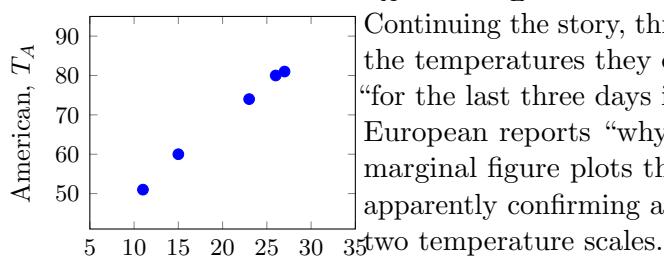
Then solve by checking `rcond(a)`, which at 0.04 is good, and so `ABCD=a\b` finds the answer $(2, -1, 1, -3)$. As before, these coefficients give the decomposition as

$$\frac{-4x^3 + 8x^2 - 5x + 2}{x^2(x-1)^2} = \frac{2}{x^2} - \frac{1}{x} + \frac{1}{(x-1)^2} - \frac{3}{x-1}$$

■

Example 2.2.10 (rcond avoids disaster). In Example 2.0.1 an American and European compared temperatures and using two days temperatures discovered the approximation that the American temperature $T_A = 1.82 T_E + 32.73$ where T_E denotes the European temperature.

Continuing the story, three days later they again meet and compare the temperatures they experienced: the American reporting that “for the last three days it has been 51° , 74° and 81° ”, whereas the European reports “why, I recorded it as 11° , 23° and 27° ”. The marginal figure plots this data with the original two data points, apparently confirming a reasonable linear relationship between the



Let’s fit a polynomial to this temperature data.

Solution: There are five data points which will each give an equation to be satisfied. This suggests we use linear algebra to determine five coefficients in a formula. Let's fit the data with the quartic polynomial

$$T_A = c_1 + c_2 T_E + c_3 T_E^2 + c_4 T_E^3 + c_5 T_E^4, \quad (2.1)$$

and use the data to determine the coefficients c_1, c_2, \dots, c_5 . Substituting each of the five pairs of T_E and T_A into this equation gives the five linear equations

$$60 = c_1 + 15c_2 + 225c_3 + 3375c_4 + 50625c_5,$$

$$\vdots$$

$$81 = c_1 + 27c_2 + 729c_3 + 19683c_4 + 531441c_5.$$

In Matlab/Octave, form these into matrix-vector equation $A\mathbf{c} = \mathbf{t}_A$ for the unknown coefficients $\mathbf{c} = (c_1, c_2, c_3, c_4, c_5)$, the vectors of American temperatures \mathbf{t}_A , and the 5×5 matrix A constructed below (recall from Table 2.3 that $\mathbf{te}.^p$ computes the p th power of each element in the column vector \mathbf{te}).

```
te=[15;26;11;23;27]
ta=[60;80;51;74;81]
plot(te,ta,'o')
A=[ones(5,1) te te.^2 te.^3 te.^4]
```

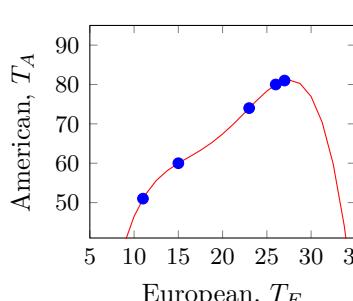
Then solve for the coefficients using $\mathbf{c} = A \setminus \mathbf{t}_A$ to get

```
A =
1      15      225      3375      50625
1      26      676     17576     456976
1      11      121      1331      14641
1      23      529     12167     279841
1      27      729     19683     531441

c =
-163.5469
  46.5194
 -3.6920
  0.1310
 -0.0017
```

Job done—or is it? To check, let's plot the predictions of the quartic polynomial (2.1) with these coefficients. In Matlab/Octave we may plot a graph with the following

```
t=linspace(5,35);
plot(t,c(1)+c(2)*t+c(3)*t.^2+c(4)*t.^3+c(5)*t.^4)
```



and see a graph like the marginal one. Disaster: the quartic polynomial relationship is clearly terrible as it is too wavy and nothing like the straight line we know it should be ($T_A = \frac{9}{5}T_E + 32$).

The problem is we forgot `rcond`. In Matlab/Octave execute `rcond(A)` and discover `rcond` is $3 \cdot 10^{-9}$. This value is in the ‘terrible’ range classified by Procedure 2.2.4. Thus the solution of the linear equations must not be used: here the marginal plot indeed shows the solution coefficients are not acceptable. Always use `rcond` to check for bad systems of linear equations.

■

The previous Example 2.2.10 illustrates one of the ‘rules of thumb’ in science and engineering: *for data fitting, avoid using polynomials of degree higher than cubic.*

Example 2.2.11 (Global Positioning System in space-time). Recall the

Example 2.1.4. Consider the GPS receiver in your smart-phone. The phone’s clock is generally in error, it may only be by a second but the GPS needs micro-second precision. Because of such a timing unknown, five satellites determine our precise position in space *and* time.

Suppose at some time (according to our smart-phone) the phone receives from a GPS satellite that it is at 3D location $(6, 12, 23)$ Mm (Megametres) and that the signal was sent at a true time 0.04 s (seconds) before the phone’s time. But the phone’s time is different to the true time by some unknown amount, say t , then the travel time of the signal from the satellite to the phone is actually $t + 0.04$. Given the speed of light is $c = 300$ Mm/s, this is a distance of $300(t + 0.04) = 300t + 12$ —linear in the discrepancy of the phone’s clock to the GPS clocks. Let (x, y, z) be you and your phone’s position in 3D space, then the distance to the satellite is also $\sqrt{(x - 6)^2 + (y - 12)^2 + (z - 23)^2}$. Equating the squares of these two gives one equation

$$(x - 6)^2 + (y - 12)^2 + (z - 23)^2 = (300t + 12)^2.$$

Similarly other satellites give other equations that help determine our position. But writing “ $300t$ ” all the time is a bit tedious, so replace it with the new unknown $w = 300t$.

Given your phone also detects four other satellites broadcast the following position and time information: $(13, 20, 12)$ time shift 0.04 s before; $(17, 14, 10)$ time shift $0.033\cdots$ s before; $(8, 21, 10)$ time shift $0.033\cdots$ s before; and $(22, 9, 8)$ time shift 0.04 s before. Adapting the approach of Example 2.1.4, use linear algebra to determine your phone’s location in space.

Solution: Let your unknown position be (x, y, z) and the unknown time shift to the phone’s clock t be found from $w = 300t$. Then the five equations from the five satellites are, respectively,

$$(x - 6)^2 + (y - 12)^2 + (z - 23)^2 = (300t + 12)^2 = (w + 12)^2,$$

$$(x - 13)^2 + (y - 20)^2 + (z - 12)^2 = (300t + 12)^2 = (w + 12)^2,$$

$$(x - 17)^2 + (y - 14)^2 + (z - 10)^2 = (300t + 10)^2 = (w + 10)^2,$$

$$(x - 8)^2 + (y - 21)^2 + (z - 10)^2 = (300t + 10)^2 = (w + 10)^2,$$

$$(x - 22)^2 + (y - 9)^2 + (z - 8)^2 = (300t + 12)^2 = (w + 12)^2.$$

Expand all the squares in these equations:

$$\begin{aligned} x^2 - 12x + 36 + y^2 - 24y + 144 + z^2 - 46z + 529 \\ = w^2 + 24w + 144, \\ x^2 - 26x + 169 + y^2 - 40y + 400 + z^2 - 24z + 144 \\ = w^2 + 24w + 144, \\ x^2 - 34x + 289 + y^2 - 28y + 196 + z^2 - 20z + 100 \\ = w^2 + 20w + 100, \\ x^2 - 16x + 64 + y^2 - 42y + 441 + z^2 - 20z + 100 \\ = w^2 + 20w + 100, \\ x^2 - 44x + 484 + y^2 - 18y + 81 + z^2 - 16z + 64 \\ = w^2 + 24w + 144. \end{aligned}$$

These are a system of nonlinear equations and so outside the remit of the course, but a little algebra brings them within. Subtract the last equation, say, from each of the first four equations: then *all* of the nonlinear squares of variables cancel leaving a linear system. Combining the constants on the right-hand side, and moving the w terms to the left gives the system of four linear equations

$$\begin{aligned} 32x - 6y - 30z + 0w &= -80, \\ 18x - 22y - 8z + 0w &= -84, \\ 10x - 10y - 4z + 4w &= 0, \\ 28x - 24y - 4z + 4w &= -20. \end{aligned}$$

Following Procedure 2.2.4, solve this system by forming the corresponding matrix and vector as

```
A=[32 -6 -30 0
   18 -22 -8 0
   10 -10 -4 4
   28 -24 -4 4 ]
b=[-80;-84;0;-20]
```

Check `rcond(A)`: it is acceptably 0.023 so compute the solution via `x=A\b` to find

```
x =
2
4
4
9
```



Hence your phone is at location $(x, y, z) = (2, 4, 4)$ Mm. Further, the time discrepancy between your phone and the GPS satellites' time is proportional to $w = 9$ Mm. Since $w = 300t$, where 300 Mm/s is the speed of light, the time discrepancy is $t = \frac{9}{300} = 0.03$ s.

■

2.2.2 Algebraic manipulation solves systems

A variant of GE [Gaussian Elimination] was used by the Chinese around the first century AD; the *Jiu Zhang Suanshu* (Nine Chapters of the Mathematical Art) contains a worked example for a system of five equations in five unknowns

Higham (1996) [p.195]

To solve linear equations with non-square matrices, or with poorly conditioned matrices we need to know much more details about linear algebra.

This and the next subsection are not essential, but many further courses currently assume knowledge of the content. Theorems 2.2.22 and 2.2.25 are convenient to establish in the next subsection, but could alternatively be established using Procedure 3.3.13.

This subsection systematises the algebraic working of Examples 2.1.3 and 2.1.5. The systematic approach empowers by-hand solution of systems of linear equations, together with two general properties on the number of solutions possible. The algebraic methodology invoked here also reinforces algebraic skills that will help in further courses.

In hand calculations we often want to minimise writing, so the discussion here uses two forms side-by-side for the linear equations: one with all symbols recorded for best clarity; and beside it, one where only coefficients are recorded for quickest writing. Translating from one to the other is crucial even in a computing era as the computer also primarily deals with arrays of numbers, and we must interpret what those arrays of numbers mean in terms of linear equations.

Example 2.2.12. Recall the system of linear equations of Example 2.1.5:

$$\begin{aligned} x_1 + x_2 - x_3 &= -2, \\ x_1 + 3x_2 + 5x_3 &= 8, \\ x_1 + 2x_2 + x_3 &= 1. \end{aligned}$$

The first crucial level of abstraction is to write this in the matrix-vector form, Example 2.2.3,

$$\underbrace{\begin{bmatrix} 1 & 1 & -1 \\ 1 & 3 & 5 \\ 1 & 2 & 1 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_x = \underbrace{\begin{bmatrix} -2 \\ 8 \\ 1 \end{bmatrix}}_b$$

A second step of abstraction omits the symbols “ $\begin{bmatrix} \end{bmatrix}$ ”—often we draw a vertical (dotted) line to show where the symbols “ $\begin{bmatrix} \end{bmatrix}$ ”

were, but this line is not essential and the theoretical statements ignore such a drawn line. Here this second step of abstraction represents this linear system by the so-called augmented matrix

$$\left[\begin{array}{ccc|c} 1 & 1 & -1 & -2 \\ 1 & 3 & 5 & 8 \\ 1 & 2 & 1 & 1 \end{array} \right]$$

■

Definition 2.2.13. *The **augmented matrix** of the system of linear equations $A\mathbf{x} = \mathbf{b}$ is the matrix $[A:\mathbf{b}]$.*

Example 2.2.14. Write down augmented matrices for the two following systems:

$$(a) \begin{aligned} -2r + 3s &= 6, \\ s - 4t &= -\pi, \end{aligned} \quad (b) \begin{aligned} -7y + 3z &= 3, \\ 7y - 5z &= -2, \\ y - 2z &= 1. \end{aligned}$$

Solution:

$$\begin{aligned} \left\{ \begin{array}{l} -2r + 3s = 6 \\ s - 4t = -\pi \end{array} \right. &\iff \left[\begin{array}{ccc|c} -2 & 3 & 0 & 6 \\ 0 & 1 & -4 & -\pi \end{array} \right] \\ \left\{ \begin{array}{l} -7y + 3z = 3 \\ 7y - 5z = -2 \\ y - 2z = 1 \end{array} \right. &\iff \left[\begin{array}{ccc|c} -7 & 3 & 3 \\ 7 & -5 & -2 \\ 1 & -2 & 1 \end{array} \right] \end{aligned}$$

An augmented matrix is not unique: it depends upon the order of the equations, and also the order you choose for the variables in \mathbf{x} . The first example implicitly chose $\mathbf{x} = (r, s, t)$; if instead we choose to order the variables as $\mathbf{x} = (s, t, r)$, then

$$\left\{ \begin{array}{l} 3s - 2r = 6 \\ s - 4t = -\pi \end{array} \right. \iff \left[\begin{array}{ccc|c} 3 & 0 & -2 & 6 \\ 1 & -4 & 0 & -\pi \end{array} \right]$$

Such variations to the augmented matrix are valid, but you must remember the corresponding chosen order of the variables.

■

Recall that Examples 2.1.3 and 2.1.5 manipulated the linear equations to deduce solution(s) to systems of linear equations. The following definition-theorem validates such manipulations in general.

Theorem 2.2.15. *The following **elementary row operations** can be performed on either a system of linear equations or on its corresponding augmented matrix without changing the solutions:*

- (a) interchange two equations/rows; or
- (b) multiply an equation/row by a nonzero constant; or
- (c) add a multiple of an equation/row to another.

Proof. We establish that each of the row operations are reversible. Hence: a solution of the system before a row operation remains a solution after the row operation; and conversely, a solution of the system after a row operation is also a solution before the row operation. When reversible, the possible solutions are the same before and after.

We just address the system of equations form as the augmented matrix form is equivalent but more abstract.

1. Swapping the order of two equations is reversed by the same swap.
2. Multiplying an equation by a non-zero constant is reversed through multiplying by the reciprocal.
3. Adding c times equation i to equation j , for $i \neq j$, is reversed by adding $(-c)$ times equation i to equation j .

□

Example 2.2.16. Use elementary row operations to find the only solution of the following system of linear equations:

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

Confirm with Matlab/Octave.

Solution: In order to know what the row operations should find, let's first solve the system with Matlab/Octave via Procedure 2.2.4. In matrix-vector form the system is

$$\begin{bmatrix} 1 & 2 & 1 \\ 2 & -3 & 0 \\ 0 & -3 & -1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix};$$

hence in Matlab/Octave execute

```
A=[1 2 1;2 -3 0;0 -3 -1]
b=[1;2;2]
rcond(A)
x=A\b
```

`rcond` is just good, 0.0104, so the computed answer $\mathbf{x} = (x, y, z) = (7, 4, -14)$ is the solution.

Second, use elementary row operations. Let's write the working in both full symbolic equations and in augmented matrix form in



order to see the correspondence between the two—you would not have to do both, either one would suffice.

$$\begin{cases} x + 2y + z = 1 \\ 2x - 3y + 0z = 2 \\ 0x - 3y - z = 2 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 2 & 1 & :1 \\ 2 & -3 & 0 & :2 \\ 0 & -3 & -1 & :2 \end{array} \right]$$

Add (-2) times the first equation/row to the second.

$$\begin{cases} x + 2y + z = 1 \\ 0x - 7y - 2z = 0 \\ 0x - 3y - z = 2 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 2 & 1 & :1 \\ 0 & -7 & -2 & :0 \\ 0 & -3 & -1 & :2 \end{array} \right]$$

This makes the first column have a leading one (Definition 2.2.17). Start on the second column by dividing the second equation/row by (-7) .

$$\begin{cases} x + 2y + z = 1 \\ 0x + y + \frac{2}{7}z = 0 \\ 0x - 3y - z = 2 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 2 & 1 & :1 \\ 0 & 1 & \frac{2}{7} & :0 \\ 0 & -3 & -1 & :2 \end{array} \right]$$

Now subtract twice the second equation/row from the first, and add three times the second to the third.

$$\begin{cases} x + 0y + \frac{3}{7}z = 1 \\ 0x + y + \frac{2}{7}z = 0 \\ 0x + 0y - \frac{1}{7}z = 2 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 0 & \frac{3}{7} & :1 \\ 0 & 1 & \frac{2}{7} & :0 \\ 0 & 0 & -\frac{1}{7} & :2 \end{array} \right]$$

This makes the second column have the second leading one (Definition 2.2.17). Start on the third column by multiplying the third equation/row by (-7) .

$$\begin{cases} x + 0y + \frac{3}{7}z = 1 \\ 0x + y + \frac{2}{7}z = 0 \\ 0x + 0y + z = -14 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 0 & \frac{3}{7} & :1 \\ 0 & 1 & \frac{2}{7} & :0 \\ 0 & 0 & 1 & :-14 \end{array} \right]$$

Now subtract $3/7$ of the third equation/row from the first, and $2/7$ from the second.

$$\begin{cases} x + 0y + 0z = 7 \\ 0x + y + 0z = 4 \\ 0x + 0y + z = -14 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 0 & 0 & :7 \\ 0 & 1 & 0 & :4 \\ 0 & 0 & 1 & :-14 \end{array} \right]$$

This completes the transformation of the equations/augmented matrix into a so-called reduced row echelon form (Definition 2.2.17). From this form we read off the solution: the system of equation on the left directly gives $x = 7$, $y = 4$ and $z = -14$, that is, the

solution vector $\mathbf{x} = (x, y, z) = (7, 4, -14)$ (as computed by Matlab/Octave); the transformed augmented matrix on the right tells us exactly the same thing because (Definition 2.2.13) it means the same as the matrix-vector

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 7 \\ 4 \\ -14 \end{bmatrix},$$

which is the same as the system on the left and tells us the solution $\mathbf{x} = (x, y, z) = (7, 4, -14)$. ■

Definition 2.2.17. A system of linear equations or (augmented) matrix is in **reduced row echelon form** if:

- (a) any equations with all zero coefficients, or rows of the matrix consisting entirely of zeros, are at the bottom;
- (b) in each nonzero equation/row, the first nonzero coefficient/entry is a one (called the **leading one**), and is in a variable/column to the left of any leading ones below it; and
- (c) each variable/column containing a leading one has zero coefficients/entries in every other equation/row.

A **free variable** is any variable which is not multiplied by a leading one when the row reduced echelon form is translated to its corresponding algebraic equations.

Example 2.2.18 (reduced row echelon form). Which of the following are in reduced row echelon form (RREF)? For those that are, identify the leading ones, and treating other variables as free variables write down the most general solution of the system of linear equations.

(a)
$$\begin{cases} x_1 + x_2 + 0x_3 - 2x_4 = -2 \\ 0x_1 + 0x_2 + x_3 + 4x_4 = 5 \end{cases}$$

Solution: This is in RREF with leading ones on the variables x_1 and x_3 . Let the other variables be free by say setting $x_2 = s$ and $x_4 = t$ for arbitrary parameters s and t . Then the two equations give $x_1 = -2 - s + 2t$ and $x_3 = 5 - 4t$. Consequently, the most general solution is $\mathbf{x} = (x_1, x_2, x_3, x_4) = (-2 - s + 2t, s, 5 - 4t, t)$ for arbitrary s and t .

(b)
$$\begin{bmatrix} 1 & 0 & -1 & 1 \\ 0 & 1 & -1 & -2 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

Solution: This augmented matrix is in RREF with leading ones in the first and second columns. To find solutions, explicitly write down the corresponding system of linear equations.

But we do not know the variables! If the context does not give variable names, then use the generic x_1, x_2, \dots, x_n . Thus here the corresponding system is

$$x_1 - x_3 = 1, \quad x_2 - x_3 = -2, \quad 0 = 4.$$

The first two equations are valid, but the last is contradictory as $0 \neq 4$. Hence there are no solutions to the system.

$$(c) \left[\begin{array}{ccc|c} 1 & 0 & -1 & 1 \\ 0 & 1 & -1 & -2 \\ 0 & 0 & 0 & 0 \end{array} \right]$$

Solution: This augmented matrix is the same as the previous except for a zero in the bottom right entry. It is in RREF with leading ones in the first and second columns. Explicitly, the corresponding system of linear equations is

$$x_1 - x_3 = 1, \quad x_2 - x_3 = -2, \quad 0 = 0.$$

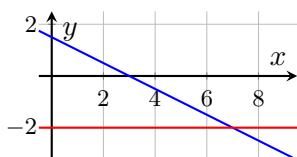
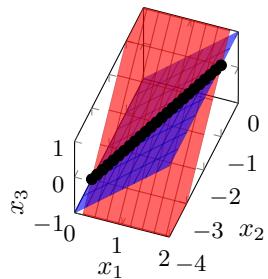
Now the last equation, $0 = 0$, is always satisfied. Hence the first two equations determine solutions to the system: letting free variable $x_3 = s$ for arbitrary s the two equations give solutions $\mathbf{x} = (1 + s, -2 + s, s)$.

$$(d) \begin{cases} x + 2y = 3 \\ 0x + y = -2 \end{cases}$$

Solution: This system is not in RREF: although there are two leading ones multiplying x and y in the first and the second equation respectively, the variable y does not have zero coefficients in the first equation. (A solution to this system exists, shown graphically in the margin, but the question does not ask for it.)

$$(e) \left[\begin{array}{cccc|c} -1 & 4 & 1 & 6 & -1 \\ 3 & 0 & 1 & -2 & -2 \end{array} \right]$$

Solution: This augmented matrix is not in RREF as there are not even any leading ones.



The previous Example 2.2.18 showed that given a system of linear equations in reduced row echelon form we can immediately write down all solutions, or immediately determine if none exists. Generalising Example 2.2.16, the following Gauss–Jordan procedure uses elementary row operations (Theorem 2.2.15) to find an equivalent system of equations in reduced row echelon form so that we can then write down the general solution.

Computers and graphics calculators can perform Gauss–Jordan elimination for you; for example, A\ in Matlab/Octave. However, if you want to use a computer, then the singular value decomposition of Section 3.3 is a far better choice for any system for which A\ is inappropriate.

- Procedure 2.2.19** (Gauss–Jordan elimination).
1. Write down either the full symbolic form of the system of linear equations, or the augmented matrix of the system of linear equations.
 2. Use elementary row operations to reduce the system/augmented matrix to reduced row echelon form.
 3. If the resulting system is consistent, then solve for the leading variables in terms of any remaining free variables.

Example 2.2.20. Use Gauss–Jordan elimination, Procedure 2.2.19, to find all possible solutions to the system

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

depending upon the parameter c .

Solution: Here write both the full symbolic equations and the augmented matrix form—you would only have to do one.

$$\left\{ \begin{array}{l} -x - y = -3 \\ x + 4y = -1 \\ 2x + 4y = c \end{array} \right. \iff \left[\begin{array}{ccc|c} -1 & -1 & -3 \\ 1 & 4 & -1 \\ 2 & 4 & c \end{array} \right]$$

Multiply the first by (-1) .

$$\left\{ \begin{array}{l} x + y = 3 \\ x + 4y = -1 \\ 2x + 4y = c \end{array} \right. \iff \left[\begin{array}{ccc|c} 1 & 1 & 3 \\ 1 & 4 & -1 \\ 2 & 4 & c \end{array} \right]$$

Subtract the first from the second, and twice the first from the third.

$$\left\{ \begin{array}{l} x + y = 3 \\ 0x + 3y = -4 \\ 0x + 2y = c - 6 \end{array} \right. \iff \left[\begin{array}{ccc|c} 1 & 1 & 3 \\ 0 & 3 & -4 \\ 0 & 2 & c - 6 \end{array} \right]$$

Divide the second by three.

$$\left\{ \begin{array}{l} x + y = 3 \\ 0x + y = -\frac{4}{3} \\ 0x + 2y = c - 6 \end{array} \right. \iff \left[\begin{array}{ccc|c} 1 & 1 & 3 \\ 0 & 1 & -\frac{4}{3} \\ 0 & 2 & c - 6 \end{array} \right]$$

Subtract the second from the first, and twice the second from the third.

$$\left\{ \begin{array}{l} x + 0y = \frac{13}{3} \\ 0x + y = -\frac{4}{3} \\ 0x + 0y = c - \frac{10}{3} \end{array} \right. \iff \left[\begin{array}{ccc|c} 1 & 0 & \frac{13}{3} \\ 0 & 1 & -\frac{4}{3} \\ 0 & 0 & c - \frac{10}{3} \end{array} \right]$$

The system is now in reduced row echelon form. The last row immediately tells us that there is no solution for parameter $c \neq \frac{10}{3}$.

as the equation would then be inconsistent. If parameter $c = \frac{10}{3}$, then the system is consistent and the first two rows give that the only solution is $(x, y) = (\frac{13}{3}, -\frac{4}{3})$. ■

Example 2.2.21. Use Gauss–Jordan elimination, Procedure 2.2.19, to find all possible solutions to the system

$$\begin{cases} -2v + 3w = -1, \\ 2u + v + w = -1. \end{cases}$$

Solution: Here write both the full symbolic equations and the augmented matrix form—you would choose one or the other.

$$\begin{cases} 0u - 2v + 3w = -1 \\ 2u + v + w = -1 \end{cases} \iff \left[\begin{array}{ccc|c} 0 & -2 & 3 & -1 \\ 2 & 1 & 1 & -1 \end{array} \right]$$

Swap the two rows to get a non-zero top-left entry.

$$\begin{cases} 2u + v - w = -1 \\ 0u - 2v + 3w = -1 \end{cases} \iff \left[\begin{array}{ccc|c} 2 & 1 & -1 & -1 \\ 0 & -2 & 3 & -1 \end{array} \right]$$

Divide the first row by two.

$$\begin{cases} u + \frac{1}{2}v - \frac{1}{2}w = -\frac{1}{2} \\ 0u - 2v + 3w = -1 \end{cases} \iff \left[\begin{array}{ccc|c} 1 & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ 0 & -2 & 3 & -1 \end{array} \right]$$

Divide the second row by (-2) .

$$\begin{cases} u + \frac{1}{2}v - \frac{1}{2}w = -\frac{1}{2} \\ 0u + v - \frac{3}{2}w = \frac{1}{2} \end{cases} \iff \left[\begin{array}{ccc|c} 1 & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ 0 & 1 & -\frac{3}{2} & \frac{1}{2} \end{array} \right]$$

Subtract half the second row from the first.

$$\begin{cases} u + 0v + \frac{1}{4}w = -\frac{3}{4} \\ 0u + v - \frac{3}{2}w = \frac{1}{2} \end{cases} \iff \left[\begin{array}{ccc|c} 1 & 0 & \frac{1}{4} & -\frac{3}{4} \\ 0 & 1 & -\frac{3}{2} & \frac{1}{2} \end{array} \right]$$

The system is now in reduced row echelon form. The third column is that of a free variable so set the third component $w = t$ for arbitrary t . Then the first row gives $u = -\frac{3}{4} - \frac{1}{4}t$, and the second row gives $v = \frac{1}{2} + \frac{3}{2}t$. That is, the solutions are $(u, v, w) = (-\frac{3}{4} - \frac{1}{4}t, \frac{1}{2} + \frac{3}{2}t, t)$ for arbitrary t . ■

2.2.3 Three possible numbers of solutions

The number of possible solutions to a system of equations is fundamental. We need to know what are the possibilities. As seen in previous examples, the following theorem says there are only three possibilities for linear equations.

Theorem 2.2.22. *For any system of linear equations $A\mathbf{x} = \mathbf{b}$, exactly one of the following is true:*

- there is no solution;
- there is a unique solution;
- there are infinitely many solutions.

Proof. First, if there is exactly none or one solution to $A\mathbf{x} = \mathbf{b}$, then the theorem holds. Second, suppose there are two distinct solutions; let them be \mathbf{y} and \mathbf{z} so $A\mathbf{y} = \mathbf{b}$ and $A\mathbf{z} = \mathbf{b}$. Then consider $\mathbf{x} = t\mathbf{y} + (1 - t)\mathbf{z}$ for all t (a parametric description of the line through \mathbf{y} and \mathbf{z} , section 1.2.2). Consider the first row of $A\mathbf{x}$: by Definition 2.2.2 it is

$$\begin{aligned} & a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ &= a_{11}[t\mathbf{y}_1 + (1 - t)\mathbf{z}_1] + a_{12}[t\mathbf{y}_2 + (1 - t)\mathbf{z}_2] \\ &\quad + \cdots + a_{1n}[t\mathbf{y}_n + (1 - t)\mathbf{z}_n] \\ &\quad (\text{then rearrange this scalar expression}) \\ &= t[a_{11}\mathbf{y}_1 + a_{12}\mathbf{y}_2 + \cdots + a_{1n}\mathbf{y}_n] \\ &\quad + (1 - t)[a_{11}\mathbf{z}_1 + a_{12}\mathbf{z}_2 + \cdots + a_{1n}\mathbf{z}_n] \\ &= t[\text{first row of } A\mathbf{y}] + (1 - t)[\text{first row of } A\mathbf{z}] \\ &= tb_1 + (1 - t)b_1 \quad (\text{as } A\mathbf{y} = \mathbf{b} \text{ and } A\mathbf{z} = \mathbf{b}) \\ &= b_1. \end{aligned}$$

Similarly for all rows of $A\mathbf{x}$: that is, each row in $A\mathbf{x}$ equals the corresponding row of \mathbf{b} . Consequently, $A\mathbf{x} = \mathbf{b}$. Hence, if there are ever two distinct solutions, then there are an infinite number of solutions, $\mathbf{x} = t\mathbf{y} + (1 - t)\mathbf{z}$. \square

An important class of linear equations always has at least one solution, never none. For example, modify Example 2.2.20 to

$$\begin{aligned} -x - y &= 0, \\ x + 4y &= 0, \\ 2x + 4y &= 0, \end{aligned}$$

and then $x = y = 0$ is immediately a solution. The reason is that the right-hand side is all zeros and so $x = y = 0$ makes the left-hand sides also zero.

Definition 2.2.23. *A system of linear equations is called **homogeneous** if the (right-hand side) constant term in each equation is zero; that is, when the system may be written $A\mathbf{x} = \mathbf{0}$. Otherwise the system is termed **non-homogeneous**.*

Example 2.2.24.

(a) $\begin{cases} 3x_1 - 3x_2 = 0 \\ -x_1 - 7x_2 = 0 \end{cases}$ is homogeneous. Solving, the first equa-

tion gives $x_1 = x_2$ and substituting in the second then gives $-x_2 - 7x_2 = 0$ so that $x_1 = x_2 = 0$ is the only solution. It must have $\mathbf{x} = \mathbf{0}$ as a solution as the system is homogeneous.

(b) $\begin{cases} 2r + s - t = 0 \\ r + s + 2t = 0 \\ -2r + s = 3 \\ 2r + 4s - t = 0 \end{cases}$ is not homogeneous because there is a non-zero constant on the right-hand side.

(c) $\begin{cases} -2 + y + 3z = 0 \\ 2x + y + 2z = 0 \end{cases}$ is not homogeneous because there is a non-zero constant in the first equation, the (-2) , even though it is here sneakily written on the left-hand side.

(d) $\begin{cases} x_1 + 2x_2 + 4x_3 - 3x_4 = 0 \\ x_1 + 2x_2 - 3x_3 + 6x_4 = 0 \end{cases}$ is homogeneous. Use Gauss–Jordan elimination, Procedure 2.2.19, to solve:

$$\begin{cases} x_1 + 2x_2 + 4x_3 - 3x_4 = 0 \\ x_1 + 2x_2 - 3x_3 + 6x_4 = 0 \end{cases} \iff \left[\begin{array}{cccc|c} 1 & 2 & 4 & -3 & 0 \\ 1 & 2 & -3 & 6 & 0 \end{array} \right]$$

Subtract the first row from the second.

$$\begin{cases} x_1 + 2x_2 + 4x_3 - 3x_4 = 0 \\ 0x_1 + 0x_2 - 7x_3 + 9x_4 = 0 \end{cases} \iff \left[\begin{array}{cccc|c} 1 & 2 & 4 & -3 & 0 \\ 0 & 0 & -7 & 9 & 0 \end{array} \right]$$

Divide the second row by (-7) .

$$\begin{cases} x_1 + 2x_2 + 4x_3 - 3x_4 = 0 \\ 0x_1 + 0x_2 + x_3 - \frac{9}{7}x_4 = 0 \end{cases} \iff \left[\begin{array}{cccc|c} 1 & 2 & 4 & -3 & 0 \\ 0 & 0 & 1 & -\frac{9}{7} & 0 \end{array} \right]$$

Subtract four times the second row from the first.

$$\begin{cases} x_1 + 2x_2 + 0x_3 + \frac{15}{7}x_4 = 0 \\ 0x_1 + 0x_2 + x_3 - \frac{9}{7}x_4 = 0 \end{cases} \iff \left[\begin{array}{cccc|c} 1 & 2 & 0 & \frac{15}{7} & 0 \\ 0 & 0 & 1 & -\frac{9}{7} & 0 \end{array} \right]$$

The system is now in reduced row echelon form. The second and fourth columns are those of free variables so set the second and fourth component $x_2 = s$ and $x_4 = t$ for arbitrary s and t . Then the first row gives $x_1 = -2s - \frac{15}{7}t$, and the second row gives $x_3 = \frac{9}{7}t$. That is, the solutions are $\mathbf{x} = (-2s - \frac{15}{7}t, s, \frac{9}{7}t, t) = (-2, 1, 0, 0)s + (-\frac{15}{7}, 0, \frac{9}{7}, 1)t$ for arbitrary s and t . These solutions include $\mathbf{x} = \mathbf{0}$ via the choice $s = t = 0$.

■

As this last example illustrates, a further subclass of homogeneous

systems is immediately known to have an infinite number of solutions. Namely, if the number of equations is less than the number of unknowns (two is less than four in the last example), then a homogeneous system always has an infinite number of solutions.

Theorem 2.2.25. *If $Ax = \mathbf{0}$ is a homogeneous system of m linear equations with n variables where $m < n$, then the system has infinitely many solutions.*

Remember that this theorem says nothing about the cases where there are at least as many equations as variables when there may or may not be an infinite number of solutions.

Proof. The zero vector, $\mathbf{x} = \mathbf{0}$ in \mathbb{R}^n , is a solution of $Ax = \mathbf{0}$ so a homogeneous system is always consistent. In the reduced row echelon form there at most m leading variables—one for each row. Here $n > m$ and so the number of free variables is at least $n - m > 0$. Hence there is at least one free variable and consequently an infinite number of solutions. \square

Prefer a matrix/vector level

Working at the element level in this way leads to a profusion of symbols, superscripts, and subscripts that tend to obscure the mathematical structure and hinder insights being drawn into the underlying process. One of the key developments in the last century was the recognition that it is much more profitable to work at the matrix level. *(Higham 2015, §2)*

A large part of this and preceding sections is devoted to arithmetic and algebraic manipulations on the individual coefficients and variables in the system. This is working at the ‘element level’ commented on by Higham. But as Higham also comments, we need to work more at a whole matrix level. This means we need to discuss and manipulate matrices as a whole, not get enmeshed in the intricacies of the element operations. This has close intellectual parallels with computing where abstract data structures empower us to encode complex tasks: here the analogous abstract data structures are matrices and vectors, and working with matrices and vectors as objects in their own right empowers linear algebra. The next chapter proceeds to develop linear algebra at the matrix level. But first, the next Section 2.3 establishes some necessary fundamental aspects at the vector level.

2.2.4 Exercises

Exercise 2.2.1. For each of the following systems, write down two different matrix-vector forms of the equations. For each system: how many

different possible matrix-vector forms could be written down?

$$(a) \begin{aligned} -3x + 6y &= -6 \\ -x - 3y &= 4 \end{aligned}$$

$$(b) \begin{aligned} -2p - q + 1 &= 0 \\ p - 6q &= 2 \end{aligned}$$

$$(c) \begin{aligned} 7.9x - 4.7y &= -1.7 \\ 2.4x - 0.1y &= 1.00 \\ -3.1x + 2.7y &= 2.30 \end{aligned}$$

$$(d) \begin{aligned} 3a + 4b - \frac{5}{2}c &= 0 \\ -\frac{7}{2}a + b - \frac{9}{2}c - \frac{9}{2} &= 0 \end{aligned}$$

$$(e) \begin{aligned} u + v - 2w &= -1 \\ -2u - v + 2w &= 3 \\ u + v + 5w &= 2 \end{aligned}$$

Exercise 2.2.2. Use Procedure 2.2.4 in Matlab/Octave to try to solve each of the systems of Exercise 2.2.1.

Exercise 2.2.3. Use Procedure 2.2.4 in Matlab/Octave to try to solve each of the following systems.

$$\begin{aligned} -5x - 3y + 5z &= -3 \\ (a) \quad 2x + 3y - z &= -5 \\ -2x + 3y + 4z &= -3 \end{aligned}$$

$$\begin{aligned} -p + 2q - r &= -2 \\ (b) \quad -2p + q + 2z &= 1 \\ -3p + 4q &= 4 \end{aligned}$$

$$\begin{aligned} u + 3v + 2w &= -1 \\ (c) \quad 3v + 5w &= 1 \\ -u + 3w &= 2 \end{aligned}$$

$$\begin{aligned} -4a - b - 3c &= -2 \\ (d) \quad 2a - 4c &= 4 \\ a - 7c &= -2 \end{aligned}$$

Exercise 2.2.4. Use elementary row operations (Theorem 2.2.15) to solve the systems in Exercise 2.2.3.

Exercise 2.2.5. Use Procedure 2.2.4 in Matlab/Octave to try to solve each of the following systems.

$$\begin{aligned} 2.2x_1 - 2.2x_2 - 3.5x_3 - 2.2x_4 &= 2.9 \\ (a) \quad 4.8x_1 + 1.8x_2 - 3.1x_3 - 4.8x_4 &= -1.6 \\ -0.8x_1 + 1.9x_2 - 3.2x_3 + 4.1x_4 &= -5.1 \\ -9x_1 + 3.5x_2 - 0.7x_3 + 1.6x_4 &= -3.3 \end{aligned}$$

$$\begin{aligned} 0.7c_1 + 0.7c_2 + 4.1c_3 - 4.2c_4 &= -0.70 \\ (b) \quad c_1 + c_2 + 2.1c_3 - 5.1c_4 &= -2.8 \\ 4.3c_1 + 5.4c_2 + 0.5c_3 + 5.5c_4 &= -6.1 \\ -0.6c_1 + 7.2c_2 + 1.9c_3 - 0.6c_4 &= -0.3 \end{aligned}$$

Exercise 2.2.6. Each of the following show some Matlab/Octave commands and their results. Write down a possible problem that these commands aim to solve, and interpret what the results mean for the problem.

```
(a) >> A=[1.1 2 5.6; 0.4 5.4 0.5; 2 -0.2 -2.8]
A =
    1.1000    2.0000    5.6000
    0.4000    5.4000    0.5000
    2.0000   -0.2000   -2.8000
>> b=[-3;2.9;1]
b =
    -3.0000
    2.9000
    1.0000
>> rcond(A)
ans =
    0.2936
>> x=A\b
x =
    -0.3936
    0.6294
    -0.6832
```

Exercise 2.2.7. Which of the following systems are in reduced row echelon form? For those that are, determine all solutions, if any.

- (a) $x_1 = -194$
 $x_2 = 564$
 $x_3 = -38$
 $x_4 = 275$
- (b) $y_1 - 13.3y_4 = -13.1$
 $y_2 + 6.1y_4 = 5.7$
 $y_3 + 3.3y_4 = 3.1$
- (c) $z_1 - 13.3z_3 = -13.1$
 $z_2 + 6.1z_3 = 5.7$
 $3.3z_3 + z_4 = 3.1$
- (d) $a - d = -4$
 $b - \frac{7}{2}d = -29$
 $c - \frac{1}{4}d = -\frac{7}{2}$
- (e) $x + 0y = 0$
 $0x + y = 0$
 $0x + 0y = 1$
 $0x + 0y = 0$
- (f) $x + 0y = -5$
 $0x + y = 1$
 $0x + 0y = 3$

Exercise 2.2.8. For the following rational expressions, express the task of finding the partial fraction decomposition as a system of linear equations. Solve the system to find the decomposition. Record your working.

(a)

$$\frac{-x^2 + 2x - 5}{x^2(x - 1)}$$

(b)

$$\frac{-4x^3 + 2x^2 - x + 2}{(x + 1)^2(x - 1)^2}$$

(c)

$$\frac{5x^4 - x^3 + 3x^2 + 10x - 1}{x^2(x + 2)(x - 1)^2}$$

(d)

$$\frac{4x^4 + 2x^3 - x^2 - 7x - 2}{(x + 1)^3(x - 1)^2}$$

Exercise 2.2.9. For each of the following tables of data, use a system of linear equations to determine the nominated polynomial that finds the second column as a function of the first column. Sketch a graph of your fitted polynomial and the data points. Record your working.

(a) linear

x	y
2	-4
3	4

(b) quadratic

x	y
-2	-1
1	0
2	5

(c) quadratic

p	q
0	-1
2	3
3	4

(d) cubic

r	t
-3	-4
-2	0
-1	-3
0	-6

Exercise 2.2.10. In three consecutive years a company sells goods to the value of \$51M, \$81M and \$92M (in millions of dollars). Find a quadratic that fits this data, and use the quadratic to predict the value of sales in the fourth year.

Exercise 2.2.11. In 2011 there were 98 wolves in Yellowstone National Park; in 2012 there were 83 wolves; and in 2013 there were 95 wolves. Find a quadratic that fits this data, and use the quadratic to predict the number of wolves in 2014. To keep the coefficients manageable, write the quadratic in terms of the number of years from the starting year of 2011.

Table 2.4: orbital periods for four planets of the solar system: the periods are in (Earth) days; the distance is the length of the semi-major axis of the orbits [Wikipedia, 2014].

planet	distance (Gigametres)	period (days)
Mercury	57.91	87.97
Venus	108.21	224.70
Earth	149.60	365.26
Mars	227.94	686.97

Exercise 2.2.12. Table 2.4 lists the time taken by a planet to orbit the Sun and a typical distance of the planet from the Sun. Analogous to Example 2.2.10, fit a quadratic polynomial $T = c_1 + c_2R + c_3R^2$ for the period T as a function of distance R . Use the data for Mercury, Venus and Earth. Then use the quadratic to predict the period of Mars: what is the error in your prediction? (Example 3.5.7 shows a power law fit is better, and that the power law agrees with Kepler's law.)

Exercise 2.2.13 (Global Positioning System in space-time). For each case below, and in space-time, suppose you know from five GPS satellites that you and your GPS receiver are the given measured time shift away from the given locations of each of the three satellites (locations are in Mm). Following Example 2.2.11, determine both your position and the discrepancy in time between your GPS receiver and the satellites GPS time. Which case needs another satellite?

- (a) $\begin{cases} 0.03 \text{ s} & \text{time shift before } (17, 11, 17) \\ 0.03 \text{ s} & \text{time shift before } (11, 20, 14) \\ 0.0233\cdots \text{ s} & \text{time shift before } (20, 10, 9) \\ 0.03 \text{ s} & \text{time shift before } (9, 13, 21) \\ 0.03 \text{ s} & \text{time shift before } (7, 24, 8) \end{cases}$
- (b) $\begin{cases} 0.1 \text{ s} & \text{time shift before } (11, 12, 18) \\ 0.1066\cdots \text{ s} & \text{time shift before } (18, 6, 19) \\ 0.1 \text{ s} & \text{time shift before } (11, 19, 9) \\ 0.1066\cdots \text{ s} & \text{time shift before } (9, 10, 22) \\ 0.1 \text{ s} & \text{time shift before } (23, 3, 9) \end{cases}$
- (c) $\begin{cases} 0.03 \text{ s} & \text{time shift before } (17, 11, 17) \\ 0.03 \text{ s} & \text{time shift before } (19, 12, 14) \\ 0.0233\cdots \text{ s} & \text{time shift before } (20, 10, 9) \\ 0.03 \text{ s} & \text{time shift before } (9, 13, 21) \\ 0.03 \text{ s} & \text{time shift before } (7, 24, 8) \end{cases}$

Exercise 2.2.14. Formulate the following two thousand year old Chinese puzzle as a system of linear equations. Use algebraic manipulation to solve the system.

There are three classes of grain, of which three bundles of the first class, two of the second, and one of the third make 39 measure. Two of the first, three of the second, and one of the third make 34 measures. And one of the first, two of the second, and three of the third make 26 measures. How many measures of grain are contained in one bundle of each class?

Jiuzhang Suanshu, 200BC (Chartier 2015, p.3)

Exercise 2.2.15. Suppose you are given data at n points, equi-spaced in x . Say the known data points are $(1, y_1), (2, y_2), \dots, (n, y_n)$ for some given y_1, y_2, \dots, y_n . Seek a polynomial fit to the data of degree $(n - 1)$; that is, seek the fit $y = c_1 + c_2x + c_3x^2 + \dots + c_nx^{n-1}$. In Matlab/Octave, form the matrix of the linear equations that need to be solved for the coefficients c_1, c_2, \dots, c_n . According to Procedure 2.2.4, for what number n of data points is `rcond` good? poor? bad? terrible?

Exercise 2.2.16 (rational functions). Sometimes we wish to fit rational functions to data. This fit also reduces to solving linear equations for the coefficients of the rational function. For example, to fit the rational function $y = a/(1 + bx)$ to data points $(x, y) = (-7, -\frac{1}{9})$ and $(x, y) = (2, \frac{1}{3})$ we need to satisfy the two equations

$$-\frac{1}{9} = \frac{a}{1 - 7b} \quad \text{and} \quad \frac{1}{3} = \frac{a}{1 + 2b}.$$

Multiply both sides of each by their denominator to require

$$\begin{aligned} -\frac{1}{9}(1 - 7b) &= a \quad \text{and} \quad \frac{1}{3}(1 + 2b) = a \\ \iff a - \frac{7}{9}b &= -\frac{1}{9} \quad \text{and} \quad a - \frac{2}{3}b = \frac{1}{3}. \end{aligned}$$

By hand (or Matlab/Octave) solve this pair of linear equations to find $a = 3$ and $b = 4$. Hence the required rational function is $y = 3/(1 + 4x)$.

- (a) Similarly, use linear equations to fit the rational function $y = (a_0 + a_1x)/(1 + b_1x)$ to the three data points $(1, \frac{7}{2}), (3, \frac{19}{4})$ and $(4, 5)$.
- (b) Similarly, use linear equations to fit the rational function $y = (a_0 + a_1x + a_2x^2)/(1 + b_1x + b_2x^2)$ to the five data points $(-\frac{5}{2}, \frac{69}{44}), (-1, \frac{3}{2}), (\frac{1}{2}, \frac{9}{8}), (1, \frac{5}{4})$ and $(2, \frac{15}{11})$.

2.3 Linear combinations span sets

Section Contents

2.3.1 Exercises	127
---------------------------	-----

A common feature in the solution to linear equations is the appearance of combinations of several vectors. For example, the general solution to Example 2.2.24d is

$$\begin{aligned}\mathbf{x} &= \left(-2s - \frac{15}{7}t, s, \frac{9}{7}t, t\right) \\ &= \underbrace{s(-2, 1, 0, 0) + t\left(-\frac{15}{7}, 0, \frac{9}{7}, 1\right)}_{\text{linear combination}}.\end{aligned}$$

The general solution to Example 2.2.18a is

$$\begin{aligned}\mathbf{x} &= (-2 - s + 2t, s, 5 - 4t, t) \\ &= \underbrace{1 \cdot (-2, 0, 5, 0) + s(-1, 1, 0, 0) + t(2, 0, -4, 1)}_{\text{linear combination}}.\end{aligned}$$

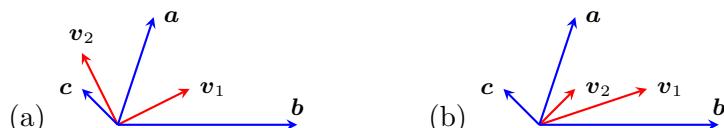
Such so-called linear combinations occur in many other contexts. Recall the standard unit vectors in \mathbb{R}^3 are $\mathbf{e}_1 = (1, 0, 0)$, $\mathbf{e}_2 = (0, 1, 0)$ and $\mathbf{e}_3 = (0, 0, 1)$ (Definition 1.2.5): any other vector in \mathbb{R}^3 may be written as

$$\begin{aligned}\mathbf{x} &= (x_1, x_2, x_3) \\ &= x_1(1, 0, 0) + x_2(0, 1, 0) + x_3(0, 0, 1) \\ &= \underbrace{x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + x_3\mathbf{e}_3}_{\text{linear combination}}.\end{aligned}$$

The wide-spread appearance of such ‘linear combinations’ calls for the following definition.

Definition 2.3.1. A vector \mathbf{v} is a **linear combination** of vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ if there are scalars c_1, c_2, \dots, c_k (called the **coefficients**) such that $\mathbf{v} = c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k$.

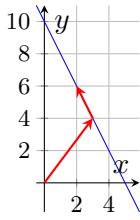
Example 2.3.2. Estimate roughly each of the blue vectors as a linear combination of the given red vectors in the following graphs (estimate coefficients to say 10% error).



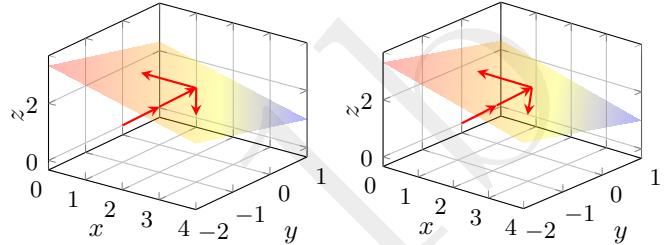
Solution: By visualising various combinations: vector $\mathbf{a} \approx -1\mathbf{v}_1 + 4\mathbf{v}_2$; vector $\mathbf{b} \approx 2\mathbf{v}_1 - 2\mathbf{v}_2$; vector $\mathbf{c} \approx -1\mathbf{v}_1 + 2\mathbf{v}_2$.

Solution: By visualising various combinations: vector $\mathbf{a} \approx -1\mathbf{v}_1 + 4\mathbf{v}_2$; vector $\mathbf{b} \approx 2\mathbf{v}_1 - 2\mathbf{v}_2$; vector $\mathbf{c} \approx -1\mathbf{v}_1 + 2\mathbf{v}_2$.

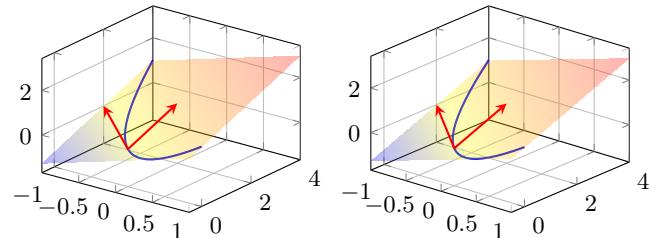
Example 2.3.3. Parametric descriptions of lines and planes involve linear combinations (Sections 1.2–1.3).



- (a) For each value of t , $(3, 4) + t(-1, 2)$ is a linear combination of the two vectors $(3, 4)$ and $(-1, 2)$. Over all values of parameter t it describes the line illustrated in the margin. (The line is alternatively described as $2x + y = 10$.)
- (b) For each value of s and t , $2(1, 0, 1) + s(-1, -\frac{1}{2}, \frac{1}{2}) + t(1, -1, 0)$ is a linear combination of the three vectors $(1, 0, 1)$, $(-1, -\frac{1}{2}, \frac{1}{2})$ and $(1, -1, 0)$. Over all values of the parameters s and t it describes the plane illustrated below. (Alternatively the plane could be described as $x + y + 3z = 8$).



- (c) $t(-1, 2, 0) + t^2(0, 2, 1)$ is a linear combination of the two vectors $(-1, 2, 0)$ and $(0, 2, 1)$ as the vectors are multiplied by scalars and then added. That a coefficient is a nonlinear function of some parameter is irrelevant to the property of linear combination. This expression is the parametric description of a parabola in \mathbb{R}^3 , as illustrated below, and very soon we will be able to say it is a parabola in the plane spanned by $(-1, 2, 0)$ and $(0, 2, 1)$.



Example 2.3.4. The matrix-vector form $A\mathbf{x} = \mathbf{b}$ of a system of linear equations involves a linear combination on the left-hand side. Recall from Definition 2.2.2 that $A\mathbf{x} = \mathbf{b}$ is our abstract abbreviation for the system of m equations

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1, \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2, \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n &= b_m. \end{aligned}$$

Form both sides into a vector so that

$$\begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}.$$

Then use addition and scalar multiplication of vectors (Definition 1.2.3) to rewrite the left-hand side vector as

$$\begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{m1} \end{bmatrix} x_1 + \begin{bmatrix} a_{12} \\ a_{22} \\ \vdots \\ a_{m2} \end{bmatrix} x_2 + \cdots + \begin{bmatrix} a_{1n} \\ a_{2n} \\ \vdots \\ a_{mn} \end{bmatrix} x_n = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}.$$

This left-hand side is a linear combination of the columns of matrix A : define from the columns of A the n vectors, $\mathbf{a}_1 = (a_{11}, a_{21}, \dots, a_{m1})$, $\mathbf{a}_2 = (a_{12}, a_{22}, \dots, a_{m2})$, \dots , $\mathbf{a}_n = (a_{1n}, a_{2n}, \dots, a_{mn})$, then the left-hand side is a linear combination of these vectors, with the coefficients of the linear combination being x_1, x_2, \dots, x_n . ■

Be aware of a subtle twist going on here: this theorem turns a question about the existence of n variable solution \mathbf{x} , into a question about vectors with m components; and vice-versa.

Proof. Example 2.3.4 establishes that if a solution \mathbf{x} exists, then \mathbf{b} is a linear combination of the columns. Conversely, if \mathbf{b} is a linear combination of the columns, then a solution \mathbf{x} exists with components of \mathbf{x} set to the coefficients in the linear combination. □

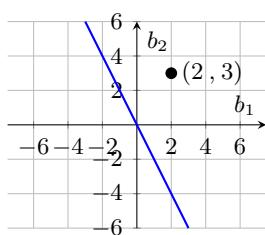
Example 2.3.6. This first example considers the simplest case when the matrix has only one column, and so any linear combination is only a scalar multiple of that column. Compare the consistency of the equations with the right-hand side being a linear combination of the column of the matrix.

$$(a) \begin{bmatrix} -1 \\ 2 \end{bmatrix} x = \begin{bmatrix} -2 \\ 4 \end{bmatrix}.$$

Solution: The system is consistent because $x = 2$ is a solution (Procedure 2.2.19). Also, the right-hand side $\mathbf{b} = (-2, 4)$ is the linear combination $2(-1, 2)$ of the column of the matrix.

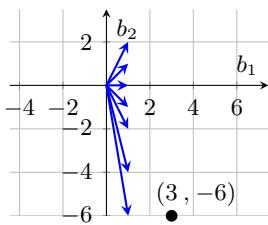
$$(b) \begin{bmatrix} -1 \\ 2 \end{bmatrix} x = \begin{bmatrix} 2 \\ 3 \end{bmatrix}.$$

Solution: The system is inconsistent as the first equation requires $x = -2$ whereas the second requires $x = \frac{3}{2}$ and these cannot hold simultaneously (Procedure 2.2.19). Also, there is no multiple of $(-1, 2)$ that gives the right-hand side $\mathbf{b} = (2, 3)$ so the right-hand side cannot be a linear combination of the column of the matrix—as illustrated in the margin.



(c) $\begin{bmatrix} 1 \\ a \end{bmatrix} x = \begin{bmatrix} 3 \\ -6 \end{bmatrix}$ depending upon parameter a .

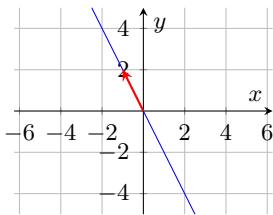
Solution: The first equation requires $x = 3$ whereas the second equation requires $ax = -6$; that is, $a \cdot 3 = -6$, that is, $a = -2$. Thus it is only for $a = -2$ that the system is consistent; for $a \neq -2$ the system is inconsistent. Also, plotted in the margin are vectors $(1, a)$ for various a . It is only for $a = -2$ that the vector is aligned towards the given $(3, -6)$. Hence it is only for $a = -2$ that a linear combination of $(1, a)$ can give the required $(3, -6)$. ■



In the examples of linear combination, the coefficients mostly involve a variable parameter or unknown. Consequently, mostly we are interested in the range of possibilities encompassed by a given set of vectors.

Definition 2.3.7. Let a set of vectors in \mathbb{R}^n be $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$, then the set of all linear combinations of $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ is called the **span** of $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$, and is denoted by $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ or $\text{span } S$.⁵

Example 2.3.8. (a) Let the set $S = \{(-1, 2)\}$ with just one vector. Then $\text{span } S = \text{span}\{(-1, 2)\}$ is the set of all vectors encompassed by the form $t(-1, 2)$: from the parametric equation of a line (Definition 1.2.9), $\text{span } S$ is all vectors in the line $y = -2x$ as shown in the margin.



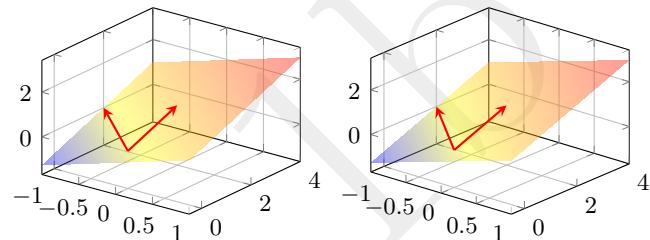
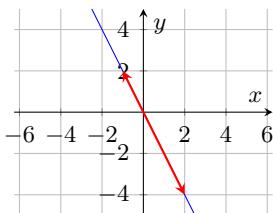
(b) With two vectors in the set, $\text{span}\{(-1, 2), (3, 4)\} = \mathbb{R}^2$ is the entire 2D plane. To see this, recall that any point in the span must be of the form $s(-1, 2) + t(3, 4)$. Given any vector (x_1, x_2) in \mathbb{R}^2 we choose $s = (-4x_1 + 3x_2)/10$ and $t = (2x_1 + x_2)/10$ and then the linear combination

$$\begin{aligned} s \begin{bmatrix} -1 \\ 2 \end{bmatrix} + t \begin{bmatrix} 3 \\ 4 \end{bmatrix} &= \frac{-4x_1 + 3x_2}{10} \begin{bmatrix} -1 \\ 2 \end{bmatrix} + \frac{2x_1 + x_2}{10} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \\ &= \left(\frac{-4}{10} \begin{bmatrix} -1 \\ 2 \end{bmatrix} + \frac{2}{10} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \right) x_1 \\ &\quad + \left(\frac{3}{10} \begin{bmatrix} -1 \\ 2 \end{bmatrix} + \frac{1}{10} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \right) x_2 \\ &= \begin{bmatrix} 1 \\ 0 \end{bmatrix} x_1 + \begin{bmatrix} 0 \\ 1 \end{bmatrix} x_2 \\ &= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}. \end{aligned}$$

⁵ In the degenerate case of the set S being the empty set, we take its span to be just the zero vector; that is, by convention $\text{span}\{\} = \{\mathbf{0}\}$. But we rarely need this degenerate case.

Since every vector in \mathbb{R}^2 can be expressed as $s(-1, 2) + t(3, 4)$, then $\mathbb{R}^2 = \text{span}\{(-1, 2), (3, 4)\}$

- (c) But if two vectors are proportional to each other then their span is a line: for example, $\text{span}\{(-1, 2), (2, -4)\}$ is the set of all vectors of the form $r(-1, 2) + s(2, -4) = r(-1, 2) + (-2s)(-1, 2) = (r - 2s)(-1, 2) = t(-1, 2)$ for $t = r - 2s$. That is, $\text{span}\{(-1, 2), (2, -4)\} = \text{span}\{(-1, 2)\}$ as illustrated in the margin.
- (d) In 3D, $\text{span}\{(-1, 2, 0), (0, 2, 1)\}$ is the set of all linear combinations $s(-1, 2, 0) + t(0, 2, 1)$ which here is a parametric form of the plane illustrated below (Definition 1.3.25). The plane passes through the origin $\mathbf{0}$, obtained when $s = t = 0$.



One could also check that the vector $(2, 1, -2)$ is orthogonal to these two vectors, hence is a normal to the plane, and so the plane may be also expressed as $2x + y - 2z = 0$.

- (e) For the complete set of n standard unit vectors in \mathbb{R}^n , $\text{span}\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} = \mathbb{R}^n$. This is because any vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ in \mathbb{R}^n may be written as the linear combination $\mathbf{x} = x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + \dots + x_n\mathbf{e}_n$ and hence is in $\text{span}\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$.
- (f) The homogeneous system (Definition 2.2.23) of linear equations from Example 2.2.24d has solutions $\mathbf{x} = (-2s - \frac{15}{7}t, s, \frac{9}{7}t, t) = (-2, 1, 0, 0)s + (-\frac{15}{7}, 0, \frac{9}{7}, 1)t$ for arbitrary s and t . That is, the set of solutions is $\text{span}\{(-2, 1, 0, 0), (-\frac{15}{7}, 0, \frac{9}{7}, 1)\}$, a subset of \mathbb{R}^4 .

Generally, the set of solutions to a homogeneous system is the span of some set.

- (g) However, the set of solutions to a non-homogeneous system is generally not equal to the span of some set. For example, the solutions to Example 2.2.21 are all of the form $(u, v, w) = (-\frac{3}{4} - \frac{1}{4}t, \frac{1}{2} + \frac{3}{2}t, t) = (-\frac{3}{4}, \frac{1}{2}, 0) + t(-\frac{1}{4}, \frac{3}{2}, 1)$ for arbitrary t . True, each of these solutions is a linear combination of vectors $(-\frac{3}{4}, \frac{1}{2}, 0)$ and $(-\frac{1}{4}, \frac{3}{2}, 1)$. But the multiple of $(-\frac{3}{4}, \frac{1}{2}, 0)$ is always fixed, whereas the span invokes *all* multiples. Consequently, all the possible solutions cannot be the same as the span of some vectors.



Example 2.3.9. Describe in other words $\text{span}\{\mathbf{i}, \mathbf{k}\}$ in \mathbb{R}^3 .

Solution: All vectors in $\text{span}\{\mathbf{i}, \mathbf{k}\}$ are of the form $c_1\mathbf{i} + c_2\mathbf{k} = c_1(1, 0, 0) + c_2(0, 0, 1) = (c_1, 0, c_2)$. Hence the span is all vectors with second component zero—the plane $y = 0$ in (x, y, z) coordinates. ■

Example 2.3.10. Find a set S such that $\text{span } S = \{(3b, a+b, -2a-4b) : a, b \text{ scalars}\}$. Similarly, find a set T such that $\text{span } T = \{(-a-2b-2, -b+1, 3b-1) : a, b \text{ scalars}\}$.

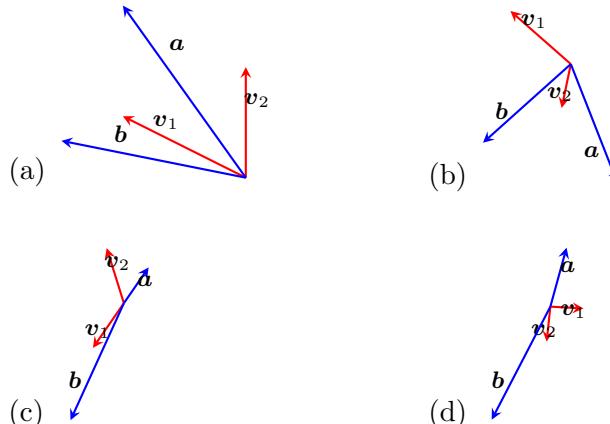
Solution: Because vectors $(3b, a+b, -2a-4b) = a(0, 1, -2) + b(3, 1, -4)$ for all scalars a and b , a suitable set is $S = \{(0, 1, -2), (3, 1, -4)\}$.

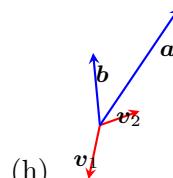
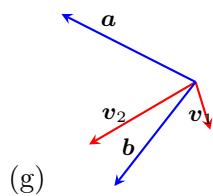
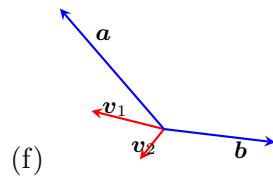
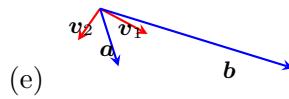
Second, vectors $(-a-2b-2, -b+1, 3b-1) = a(-1, 0, 0) + b(-2, -1, 3) + (-2, 1, -1)$ which are linear combinations for all a and b . However, the vectors cannot form a span due to the constant vector $(-2, 1, -1)$ because a span requires *all* linear combinations of its component vectors. The given set cannot be expressed as a span. ■

Geometrically, the span of a set of vectors is always all vectors lying in either a line, a plane, or a higher dimensional hyper-plane, that passes *through the origin* (discussed further by section 3.4).

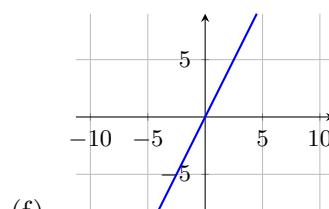
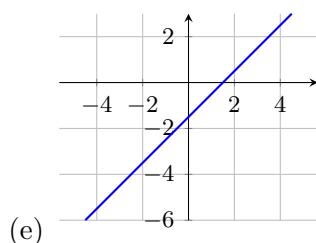
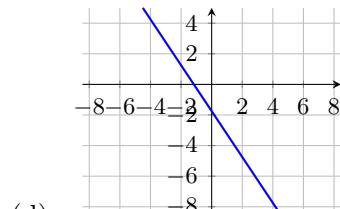
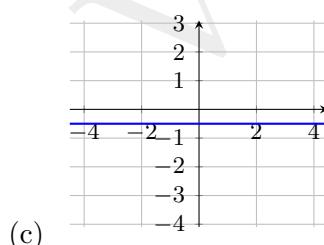
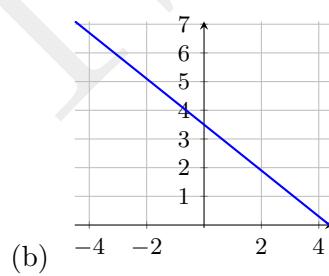
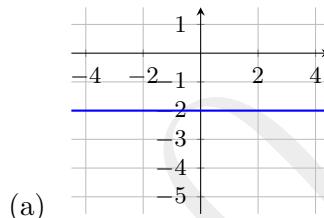
2.3.1 Exercises

Exercise 2.3.1. For each of the following, express vectors \mathbf{a} and \mathbf{b} as a linear combination of vectors \mathbf{v}_1 and \mathbf{v}_2 . Estimate the coefficients roughly (to say 10%).





Exercise 2.3.2. For each of the following lines in 2D, write down a parametric equation of the line as a linear combination of two vectors, one of which is multiplied by the parameter.



Exercise 2.3.3. Write each of the following systems of linear equations as one vector equation involving a linear combination of vectors.

$$(a) \begin{aligned} -2x + y - 2z &= -2 \\ -4x + 2y - z &= 2 \end{aligned}$$

$$(b) \begin{array}{l} -3x + 2y - 3z = 0 \\ y - z = 0 \\ x - 3y = 0 \end{array}$$

$$(c) \begin{array}{l} x_1 + 3x_2 + x_3 - 2x_4 = 2 \\ 2x_1 + x_2 + 4x_3 - 2x_4 = -1 \\ -x_1 + 2x_2 - 2x_3 - x_4 = 3 \end{array} \quad (d) \begin{array}{l} -2p - 2q = -1 \\ q = 2 \\ 3p - q = 1 \end{array}$$

Exercise 2.3.4. For each of the cases in Exercise 2.3.3, by attempting to solve the system, determine if the right-hand side vector is in the span of the vectors on the left-hand side.

Exercise 2.3.5. For each of the following sets, write the set as a span, if possible. Give reasons.

- (a) $\{(p - 4q, p + 2q, p + 2q) : p, q \text{ scalars}\}$
- (b) $\{(-p + 2r, 2p - 2q, p + 2q + r, -q - 3r) : p, q, r \text{ scalars}\}$
- (c) The line $y = 2x + 1$ in \mathbb{R}^2 .
- (d) The line $x = y = z$ in \mathbb{R}^3 .
- (e) The set of vectors \mathbf{x} in \mathbb{R}^4 with component $x_3 = 0$.

Exercise 2.3.6. Show the following identities hold for any given vectors \mathbf{u}, \mathbf{v} and \mathbf{w} :

- (a) $\text{span}\{\mathbf{u}, \mathbf{v}\} = \text{span}\{\mathbf{u} - \mathbf{v}, \mathbf{u} + \mathbf{v}\};$
- (b) $\text{span}\{\mathbf{u}, \mathbf{v}, \mathbf{w}\} = \text{span}\{\mathbf{u}, \mathbf{u} - \mathbf{v}, \mathbf{u} + \mathbf{v} + \mathbf{w}\}.$

Exercise 2.3.7. Suppose $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_s$ are any s vectors in \mathbb{R}^n . Let set $R = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ for some $r < s$, and set $S = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r, \mathbf{u}_{r+1}, \dots, \mathbf{u}_s\}$.

- (a) Prove that $\text{span } R \subseteq \text{span } S$.
- (b) Hence deduce that if $\text{span } R = \mathbb{R}^n$, then $\text{span } S = \mathbb{R}^n$.

Exercise 2.3.8. Suppose $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r$ and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s$ are all vectors in \mathbb{R}^n .

- (a) Prove that if every vector \mathbf{u}_j is a linear combination of $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s$, then $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\} \subseteq \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s\}$.
- (b) Prove that if, additionally, every vector \mathbf{v}_j is a linear combination of $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r$, then $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\} = \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s\}$.

Exercise 2.3.9. Let $S = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s\}$ be a set of vectors in \mathbb{R}^n such that vector \mathbf{v}_1 is a linear combination of $\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_s$. Prove that $\text{span } S = \text{span}\{\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_s\}$.

Answers to selected exercises

2.1.1b : $(x, y) = (1, 2)$

2.1.1d : line $y = \frac{3}{2}x - 1$

2.1.1f : $(x, y) = (-1, -1)$

2.1.1h : $(p, q) = (\frac{1}{2}, -\frac{1}{2})$

2.1.1j : no solution

2.1.1l : line $t = 4s - 2$

2.1.2b : $(1.8, 1.04)$

2.1.2d : $(1.25, 0.41)$

2.1.2f : $(2.04, 0.88)$

2.1.2h : $(2.03, 0.34)$

2.1.3b : no solution

2.1.3d : no solution

2.1.3f : no solution

2.1.4b : $(4, 5)$

2.1.4d : $(-1, 11)/9$ surely indicates an error.

2.2.1b : e.g. $\begin{bmatrix} -2 & -1 \\ 1 & -6 \end{bmatrix} \begin{bmatrix} p \\ q \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$, $\begin{bmatrix} 1 & -6 \\ -2 & -1 \end{bmatrix} \begin{bmatrix} p \\ q \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$. Four:
two orderings of rows, and two orderings of the variables.

2.2.1d : Twelve possibilities: two orderings of rows, and $3! = 6$ orderings of the variables.

2.2.2 :

1. $(x, y) = (-0.4, -1.2)$

2.

$(p, q) = (0.6154, -0.2308)$

3. No solution as `rcond` requires a square matrix.

4. No solution as `rcond` requires a square matrix.

5. $(u, v, w) = (-2, 1.8571, 0.4286)$

2.2.3b : $(p, q, r) = (32, 25, 20)$

2.2.3d : $(a, b, c) = (3.6, -14.8, 0.8)$

2.2.5a : $\mathbf{x} = (-0.26, -1.33, 0.18, -0.54)$ (2 d.p.)

2.2.6a : Solve the system $1.1x_1 + 2x_2 + 5.6x_3 = -3$, $0.4x_1 + 5.4x_2 + 0.5x_3 = 2.9$, $2x_1 - 0.2x_2 - 2.8x_3 = 1$. Since `rcond` is good, the solution is $x_1 = -0.39$, $x_2 = 0.63$ and $x_3 = -0.68$ (2 d.p.).

2.2.7b : Yes, $\mathbf{y} = (-13.1 + 13.3t, 5.7 - 6.1t, 3.1 - 3.3t, t)$ for all t

2.2.7d : Yes. $(a, b, c, d) = (-4 + t, -29 + \frac{7}{2}t, -\frac{7}{2} + \frac{1}{4}t, t)$ for all t

2.2.7f : Not in RREF as the last does not have a leading one.

$$2.2.8b : -\frac{2}{x+1} + \frac{9/4}{(x+1)^2} - \frac{2}{x-1} - \frac{1}{(x-1)^2}$$

$$2.2.8d : \frac{19/8}{x-1} - \frac{1/2}{(x-1)^2} + \frac{13/8}{x+1} - \frac{9/4}{(x+1)^2} + \frac{3/2}{(x+1)^3}$$

$$2.2.9b : y = -\frac{8}{3} + \frac{3}{2}x + \frac{7}{6}x^2$$

$$2.2.9d : t = -6 - \frac{2}{3}t + \frac{7}{2}t^2 + \frac{7}{6}t^3$$

2.2.11 : 134 wolves

2.2.13a : $\text{rcond} = 0.014$, $(3, 4, 3)$, shift = $12/300 = 0.04$ s

2.2.13c : $\text{rcond} = 0$, singular, needs another satellite.

2.2.15 : Good, $\{1, 2\}$; poor, $\{3, 4\}$; bad, $\{5, 6, 7\}$; terrible, $\{8, 9, \dots\}$.

$$2.2.16b : y = (1 + x + 3x^2)/(1 + x + 2x^2)$$

$$2.3.1b : \mathbf{a} = -1\mathbf{v}_1 + 1.5\mathbf{v}_2, \mathbf{b} = 1\mathbf{v}_1 + 3\mathbf{v}_2$$

$$2.3.1d : \mathbf{a} = 0.3\mathbf{v}_1 - 1.7\mathbf{v}_2, \mathbf{b} = -1.4\mathbf{v}_1 + 3.3\mathbf{v}_2$$

$$2.3.1f : \mathbf{a} = 2.3\mathbf{v}_1 - 2.6\mathbf{v}_2, \mathbf{b} = -1.4\mathbf{v}_1 - 0.4\mathbf{v}_2$$

$$2.3.1h : \mathbf{a} = -1.8\mathbf{v}_1 + 1.5\mathbf{v}_2, \mathbf{b} = -1.5\mathbf{v}_1 - 0.6\mathbf{v}_2$$

$$2.3.2b : \text{e.g. } (0, 3.5) + t(-2.5, 2)$$

$$2.3.2d : \text{e.g. } (-1.5, 0.5) + t(1, -1.5)$$

$$2.3.2f : \text{e.g. } (0.5, 1) + t(0.5, 1)$$

$$2.3.3b : \begin{bmatrix} -3 \\ 0 \\ 1 \end{bmatrix} x + \begin{bmatrix} 2 \\ 1 \\ -3 \end{bmatrix} y + \begin{bmatrix} -3 \\ 1 \\ 0 \end{bmatrix} z = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$2.3.3d : \begin{bmatrix} -2 \\ 0 \\ 3 \end{bmatrix} p + \begin{bmatrix} -2 \\ 1 \\ -1 \end{bmatrix} q = \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix}$$

$$2.3.5a : \text{e.g. } \text{span}\{(1, 1, 1), (-4, 2, 2)\}$$

2.3.5c : Not a span.

2.3.5e : e.g. $\text{span}\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_4\}$.

3 Matrices encode system interactions

Chapter Contents

3.1	Matrix operations and algebra	134
3.1.1	Basic matrix terminology	134
3.1.2	Addition, subtraction and multiplication with matrices	136
3.1.3	Familiar algebraic properties of matrix oper- ations	153
3.1.4	Exercises	158
3.2	The inverse of a matrix	164
3.2.1	Introducing the unique inverse	164
3.2.2	Diagonal matrices stretch and shrink	174
3.2.3	Orthogonal matrices rotate	182
3.2.4	Exercises	190
3.3	Factorise to the singular value decomposition	205
3.3.1	Introductory examples	205
3.3.2	The SVD solves general systems	209
3.3.3	Prove the SVD Theorem 3.3.5	225
3.3.4	Exercises	230
3.4	Subspaces, basis and dimension	240
3.4.1	Subspaces are lines, planes, and so on	240
3.4.2	Orthonormal bases form a foundation	249
3.4.3	Is it a line? a plane? The dimension answers .	259
3.4.4	Exercises	268
3.5	Project to solve inconsistent equations	276
3.5.1	Make a minimal change to the problem	276
3.5.2	Compute the smallest appropriate solution	290
3.5.3	Orthogonal projection resolves vector compo- nents	296
3.5.4	Exercises	323

3.6	Introducing linear transformations	336
3.6.1	Matrices characterise linear transforms	341
3.6.2	The pseudo-inverse of a matrix	345
3.6.3	Function composition connects to matrix in- verse	353
3.6.4	Exercises	359

Section 2.2 introduced matrices in the matrix-vector form $A\mathbf{x} = \mathbf{b}$ of a system of linear equations. This chapter starts with Sections 3.1 and 3.2 developing the basic operations on matrices that make them so useful in applications and theory—including making sense of the ‘product’ $A\mathbf{x}$. Section 3.3 then explores how the so-called “singular value decomposition (SVD)” of a matrix empowers us to understand how to solve general linear systems of equations, and a graphical meaning of a matrix in terms of rotations and stretching. The structures discovered by an SVD lead to further conceptual development (Section 3.4) that underlies the at first paradoxical solution of inconsistent equations (Section 3.5). Finally, Section 3.6 unifies the geometric views invoked.

the language of mathematics reveals itself unreasonably effective in the natural sciences . . . a wonderful gift which we neither understand nor deserve. We should be grateful for it and hope that it will remain valid in future research and that it will extend, for better or for worse, to our pleasure even though perhaps also to our bafflement, to wide branches of learning

Wigner, 1960 ([Mandelbrot 1982](#), p.3)

3.1 Matrix operations and algebra

Section Contents

3.1.1	Basic matrix terminology	134
3.1.2	Addition, subtraction and multiplication with matrices	136
	Matrix addition and subtraction	137
	Scalar multiplication of matrices	138
	Matrix-vector multiplication transforms	139
	Matrix-matrix multiplication	142
	The transpose of a matrix	145
	Compute in Matlab/Octave	147
3.1.3	Familiar algebraic properties of matrix operations	153
3.1.4	Exercises	158

This section introduces some basic matrix concepts, operations and algebra. Many of you will have met much of it in previous study. Consequently, this introduction is fairly quick.

3.1.1 Basic matrix terminology

Let's start with some basic definitions of terminology.

- A **matrix** is a rectangular array of real numbers, written inside **brackets** $[]$, such as the following six:¹

$$\begin{bmatrix} -2 & -5 & 4 \\ 1 & -3 & 0 \\ 2 & 4 & 0 \end{bmatrix}, \quad \begin{bmatrix} -2.33 & 3.66 \\ -4.17 & -0.36 \end{bmatrix}, \quad \begin{bmatrix} 0.56 \\ 3.99 \\ -5.22 \end{bmatrix}, \\ \begin{bmatrix} 1 & -\sqrt{3} & \pi \\ -5/3 & \sqrt{5} & -1 \end{bmatrix}, \quad \begin{bmatrix} 1 & \frac{10}{3} & \frac{\pi^2}{4} \end{bmatrix}, \quad [0.35]. \quad (3.1)$$

- The **size** of a matrix is its number of rows and columns—written $m \times n$ where m is the number of rows and n is the number of columns. The six example matrices of (3.1) are of size, respectively, 3×3 , 2×2 , 3×1 , 2×3 , 1×3 , and 1×1 .

Recall from Definition 2.2.2 that if for a matrix the number of rows equals the number of columns, $m = n$, then it is called a square matrix. For example, the first, second and last matrices in (3.1) are square; the others are not.

¹ Chapter 7 will start using complex numbers in a matrix, but until then we stay within the realm of real numbers.

- To correspond with vectors, we often invoke the term **column vector** which means a matrix with only one column; that is, a matrix of size $m \times 1$ for some m . For convenience and compatibility with vectors, we often write a column matrix horizontally within **parentheses** (). The third matrix in (3.1) is an example, and may also be written as $(0.56, 3.99, -5.22)$.

Occasionally we refer to a **row vector** to mean a matrix with one row; that is, a $1 \times n$ matrix for some n , such as the fifth matrix in (3.1).

- The numbers appearing in a matrix are called the **entries**, **elements** or **components** of the matrix. For example, the first matrix in (3.1) has entries/elements/components of the numbers $-5, -3, -2, 0, 1, 2$ and 4 .
- But it is important to identify where the numbers appear in a matrix: the **double subscript** notation identifies the location of an entry. For a matrix A , the entry in row i and column j is denoted by a_{ij} : by convention we use capital (uppercase) letters for a matrix, and the corresponding lowercase letter subscripted for its entries.² For example, let matrix

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

then entries $a_{12} = -5$, $a_{22} = -3$ and $a_{31} = 2$.

- The first of two special matrices is a **zero matrix** of all zeros and of any size: $O_{m \times n}$ denotes the $m \times n$ zero matrix, such as

$$O_{2 \times 4} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

The symbol O_n denotes the square zero matrix of size $n \times n$, whereas the plain symbol O denotes a zero matrix whose size is apparent from the context.

- Arising from the nature of matrix multiplication (section 3.1.2), the second special matrix is the **identity matrix**: I_n denotes a $n \times n$ square matrix which has zero entries except for the diagonal from the top-left to the bottom-right which are all ones. Occasionally we invoke non-square ‘identity’ matrices denoted $I_{m \times n}$. For examples,

$$I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad I_{2 \times 3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad I_{4 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$

² Some people use the capital letter subscripted for its entries: that is, some use A_{ij} to denote the entry in the i th row and j th column of matrix A .

The plain symbol I denotes an identity matrix whose size is apparent from the context.

- Using the double subscript notation, and as already used in Definition 2.2.2, a general $m \times n$ matrix

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}.$$

Often, as already seen in Example 2.3.4, it is useful to write a matrix A in terms of its n column vectors \mathbf{a}_j , $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$. For example, matrix

$$B = \begin{bmatrix} 1 & -\sqrt{3} & \pi \\ -5/3 & \sqrt{5} & -1 \end{bmatrix} = [\mathbf{b}_1 \ \mathbf{b}_2 \ \mathbf{b}_3]$$

for the three column vectors

$$\mathbf{b}_1 = \begin{bmatrix} 1 \\ -5/3 \end{bmatrix}, \quad \mathbf{b}_2 = \begin{bmatrix} -\sqrt{3} \\ \sqrt{5} \end{bmatrix}, \quad \mathbf{b}_3 = \begin{bmatrix} \pi \\ -1 \end{bmatrix}.$$

Alternatively we could write these column vectors as $\mathbf{b}_1 = (1, -5/3)$, $\mathbf{b}_2 = (-\sqrt{3}, \sqrt{5})$ and $\mathbf{b}_3 = (\pi, -1)$.

- Lastly, two matrices are **equal** ($=$) if they both have the same size *and* their corresponding entries are equal. Otherwise the matrices are not equal. For example, consider matrices

$$A = \begin{bmatrix} 2 & \pi \\ 3 & 9 \end{bmatrix}, \quad B = \begin{bmatrix} \sqrt{4} & \pi \\ 2+1 & 3^2 \end{bmatrix},$$

$$C = [2 \ \pi], \quad D = \begin{bmatrix} 2 \\ \pi \end{bmatrix} = (2, \pi).$$

The matrices $A = B$ because they are the same size and their corresponding entries are equal, such as $a_{11} = 2 = \sqrt{4} = b_{11}$. Matrix A cannot be equal to C because their sizes are different. Matrices C and D are not equal, despite having the same elements in the same order, because they have different sizes: 1×2 and 2×1 respectively.

3.1.2 Addition, subtraction and multiplication with matrices

A matrix is not just an array of numbers: associated with a matrix is a suite of operations that empower a matrix in applications. We start with addition and multiplication: ‘division’ is addressed by Section 3.2 and others.

An analogue in computing science is the concept of object orientated programming. In object oriented programming one defines not just data structures, but also the functions that operate on those

structures. Analogously, an array may be just a group of numbers, but a matrix is an array together with many operations explicitly available. The power and beauty of matrices results from the implications of its associated operations.

Matrix addition and subtraction

Corresponding to vector addition and subtraction (Definition 1.2.3), matrix addition and subtraction is done component wise, but only between matrices of the same size.

Example 3.1.1. Let matrices

$$A = \begin{bmatrix} 4 & 0 \\ -5 & -4 \\ 0 & -3 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 2 \\ -3 & 0 & 3 \end{bmatrix}, \quad C = \begin{bmatrix} -4 & -1 \\ -4 & -1 \\ 1 & 4 \end{bmatrix},$$

$$D = \begin{bmatrix} -2 & -1 & -3 \\ 1 & 3 & 0 \end{bmatrix}, \quad E = \begin{bmatrix} 5 & -2 & -2 \\ 0 & -3 & 2 \\ -4 & 7 & -1 \end{bmatrix}.$$

Then the addition and subtraction

$$\begin{aligned} A + C &= \begin{bmatrix} 4 & 0 \\ -5 & -4 \\ 0 & -3 \end{bmatrix} + \begin{bmatrix} -4 & -1 \\ -4 & -1 \\ 1 & 4 \end{bmatrix} \\ &= \begin{bmatrix} 4 + (-4) & 0 + (-1) \\ -5 + (-4) & -4 + (-1) \\ 0 + 1 & -3 + 4 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ -9 & -5 \\ 1 & 1 \end{bmatrix}, \\ B - D &= \begin{bmatrix} 1 & 0 & 2 \\ -3 & 0 & 3 \end{bmatrix} - \begin{bmatrix} -2 & -1 & -3 \\ 1 & 3 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 1 - (-2) & 0 - (-1) & 2 - (-3) \\ -3 - 1 & 0 - 3 & 3 - 0 \end{bmatrix} = \begin{bmatrix} 3 & 1 & 5 \\ -4 & -3 & 3 \end{bmatrix}. \end{aligned}$$

But because the matrices are of different sizes, the following are not defined and must not be attempted: $A + B$, $A - D$, $E - A$, $B + C$, $E - C$, for example. ■

In general, suppose A and B are both $m \times n$ matrices, with entries a_{ij} and b_{ij} respectively, then we define their **sum** or **addition**, $A + B$, as the $m \times n$ matrix whose (i, j) th entry is $a_{ij} + b_{ij}$. Similarly, define the **difference** or **subtraction** $A - B$ as the $m \times n$ matrix whose (i, j) th entry is $a_{ij} - b_{ij}$. That is,

$$A + B = \begin{bmatrix} a_{11} + b_{11} & a_{12} + b_{12} & \cdots & a_{1n} + b_{1n} \\ a_{21} + b_{21} & a_{22} + b_{22} & \cdots & a_{2n} + b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} + b_{m1} & a_{m2} + b_{m2} & \cdots & a_{mn} + b_{mn} \end{bmatrix},$$

$$A - B = \begin{bmatrix} a_{11} - b_{11} & a_{12} - b_{12} & \cdots & a_{1n} - b_{1n} \\ a_{21} - b_{21} & a_{22} - b_{22} & \cdots & a_{2n} - b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} - b_{m1} & a_{m2} - b_{m2} & \cdots & a_{mn} - b_{mn} \end{bmatrix}.$$

For example, letting O denote the zero matrix of the appropriate size, then

$$A \pm O = A, \quad O + A = A, \quad \text{and} \quad A - A = O.$$

Scalar multiplication of matrices

Corresponding to multiplication of a vector by a scalar (Definition 1.2.3), multiplication of a matrix by a scalar means that every entry of the matrix is multiplied by the scalar.

Example 3.1.2. Let matrices

$$A = \begin{bmatrix} 5 & 2 \\ -2 & 3 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \\ -6 \end{bmatrix}, \quad C = \begin{bmatrix} 5 & -6 & 4 \\ -1 & -3 & -3 \end{bmatrix}.$$

Then the scalar multiplications

$$\begin{aligned} 3A &= \begin{bmatrix} 3 \cdot 5 & 3 \cdot 2 \\ 3 \cdot (-2) & 3 \cdot 3 \end{bmatrix} = \begin{bmatrix} 15 & 6 \\ -6 & 9 \end{bmatrix}, \\ -B &= (-1)B = \begin{bmatrix} (-1) \cdot 1 \\ (-1) \cdot 0 \\ (-1) \cdot (-6) \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 6 \end{bmatrix}, \\ -\pi C &= (-\pi)C = \begin{bmatrix} 5\pi & -6\pi & 4\pi \\ -\pi & -3\pi & -3\pi \end{bmatrix}. \end{aligned}$$

■

In general, suppose A is an $m \times n$ matrix, with entries a_{ij} , then we define the **scalar product** by c , either cA or Ac , as the $m \times n$ matrix whose (i, j) th entry is ca_{ij} .³ That is,

$$cA = Ac = \begin{bmatrix} ca_{11} & ca_{12} & \cdots & ca_{1n} \\ ca_{21} & ca_{22} & \cdots & ca_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ ca_{m1} & ca_{m2} & \cdots & ca_{mn} \end{bmatrix}.$$

³ Strictly speaking the product ‘[0.35] A ’ is not defined because strictly speaking [0.35] is not a scalar but is a 1×1 matrix. However, Matlab/Octave reasonably treats multiplication by a ‘ 1×1 matrix’ as a scalar multiplication.

Matrix-vector multiplication transforms

Recall that the matrix-vector form of a system of linear equations, Definition 2.2.2, wrote $A\mathbf{x} = \mathbf{b}$. In this form, $A\mathbf{x}$ denotes a matrix-vector product. As implied by Definition 2.2.2, we define the general **matrix-vector product**

$$A\mathbf{x} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \cdots + a_{mn}x_n \end{bmatrix}$$

for $m \times n$ matrix A and vector \mathbf{x} in \mathbb{R}^n with entries/components

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \quad \text{and} \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}.$$

This product is only defined when the number of columns of matrix A are the same as the number of components of vector \mathbf{x} .
⁴

Example 3.1.3. Let matrices

$$A = \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 5 & -6 & 4 \\ -1 & -3 & -3 \end{bmatrix},$$

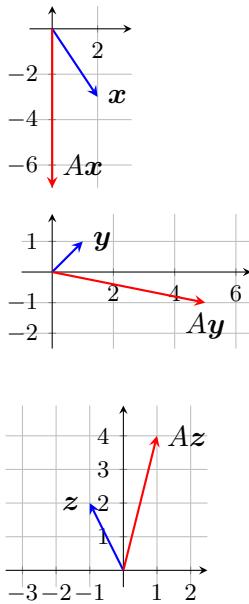
and vectors $\mathbf{x} = (2, -3)$ and $\mathbf{b} = (1, 0, 4)$. Then the matrix-vector products

$$\begin{aligned} A\mathbf{x} &= \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \end{bmatrix} = \begin{bmatrix} 3 \cdot 2 + 2 \cdot (-3) \\ (-2) \cdot 2 + 1 \cdot (-3) \end{bmatrix} = \begin{bmatrix} 0 \\ -7 \end{bmatrix}, \\ B\mathbf{b} &= \begin{bmatrix} 5 & -6 & 4 \\ -1 & -3 & -3 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 4 \end{bmatrix} \\ &= \begin{bmatrix} 5 \cdot 1 + (-6) \cdot 0 + 4 \cdot 4 \\ (-1) \cdot 1 + (-3) \cdot 0 + (-3) \cdot 4 \end{bmatrix} = \begin{bmatrix} 9 \\ -13 \end{bmatrix}. \end{aligned}$$

The combinations $A\mathbf{b}$ and $B\mathbf{x}$ are not defined as the number of columns of the matrices are not equal to the number of components in the vectors.

Further, we do not here define products such as $\mathbf{x}A$ or $\mathbf{b}B$: the order of multiplication matters with matrices and so these are not in the scope of the definition. ■

⁴ Some of you who have studied calculus may wonder about what might be called ‘continuous matrices’ $A(x, y)$ which multiply a function $f(x)$ according to the integral $\int_a^b A(x, y)f(y) dy$. Then you might wonder about solving problems such as find the unknown $f(x)$ such that $\int_0^1 A(x, y)f(y) dy = \sin \pi x$ for given $A(x, y) := \min(x, y)[1 - \max(x, y)]$; you may check that here the solution is $f = \frac{1}{\pi^2} \sin \pi x$. Such notions are a useful generalisation of our linear algebra: they are called integral equations; the main structures and patterns developed by this course also apply.



Geometric interpretation Multiplication of a vector by a square matrix transforms the vector into another in the same space as shown in the margin for $A\mathbf{x}$ from Example 3.1.3. For another vector $\mathbf{y} = (1, 1)$ the product

$$A\mathbf{y} = \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \cdot 1 + 2 \cdot 1 \\ (-2) \cdot 1 + 1 \cdot 1 \end{bmatrix} = \begin{bmatrix} 5 \\ -1 \end{bmatrix},$$

as illustrated in the second marginal picture. Similarly, for the vector $\mathbf{z} = (-1, 2)$ the product

$$A\mathbf{z} = \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 2 \end{bmatrix} = \begin{bmatrix} 3 \cdot (-1) + 2 \cdot 2 \\ (-2) \cdot (-1) + 1 \cdot 2 \end{bmatrix} = \begin{bmatrix} 1 \\ 4 \end{bmatrix},$$

as illustrated in the third marginal picture. Such a geometric interpretation underlies the use of matrix multiplication in video and picture processing, for example. Such processing employs stretching and shrinking (section 3.2.2), rotations (section 3.2.3), among more general transformations (Section 3.6).

Example 3.1.4. Recall I_n is the $n \times n$ identity matrix. Then the products

$$\begin{aligned} I_2\mathbf{x} &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -3 \end{bmatrix} = \begin{bmatrix} 1 \cdot 2 + 0 \cdot (-3) \\ 0 \cdot 2 + 1 \cdot (-3) \end{bmatrix} = \begin{bmatrix} 2 \\ -3 \end{bmatrix}, \\ I_3\mathbf{b} &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \cdot 1 + 0 \cdot 0 + 0 \cdot 4 \\ 0 \cdot 1 + 1 \cdot 0 + 0 \cdot 4 \\ 0 \cdot 1 + 0 \cdot 0 + 1 \cdot 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 4 \end{bmatrix}. \end{aligned}$$

That is, and justifying its name of “identity”, the products with an identity matrix give the result that is the vector itself: $I_2\mathbf{x} = \mathbf{x}$ and $I_3\mathbf{b} = \mathbf{b}$. Multiplication by the identity matrix leaves the vector unchanged (Theorem 3.1.18e). ■

Example 3.1.5 (age structured population). An ecologist studies an isolated population of a species of animal. The growth of the population depends primarily upon the females so it is only these that are counted. The females are grouped into three ages: female pups (in their first year), juvenile females (one year old), and mature females (two years or older). During the study, the ecologist observes the following happens over the period of a year:

- half of the female pups survive and become juvenile females;
- one-third of the juvenile females survive and become mature females;
- each mature female breeds and produces four female pups;
- one-third of the mature females survive to breed in the following year;
- female pups and juvenile females do not breed.

- (a) Let x_1 , x_2 and x_3 be the number of females at the start of a year, of ages zero, one and two+ respectively, and let x'_1 , x'_2 and x'_3 be their number at the start of the next year. Use the observations to write x'_1 , x'_2 and x'_3 as a function of x_1 , x_2 and x_3 (this is called a Markov chain).
- (b) Letting vectors $\mathbf{x} = (x_1, x_2, x_3)$ and $\mathbf{x}' = (x'_1, x'_2, x'_3)$ write down your function as the matrix-vector product $\mathbf{x}' = L\mathbf{x}$ for some matrix L (called a Leslie matrix).
- (c) Suppose the ecologist observes the numbers of females at the start of a given year is $\mathbf{x} = (60, 70, 20)$, use your matrix to predict the numbers \mathbf{x}' at the start of the next year. Continue similarly to predict the numbers after two years (\mathbf{x}'')? and three years (\mathbf{x}''')?

Solution: (a) Since mature females breed and produce four female pups, $x'_1 = 4x_3$. Since half of the female pups survive and become juvenile females, $x'_2 = \frac{1}{2}x_1$. Since one-third of the juvenile females survive and become mature females, $\frac{1}{3}x_2$ contribute to x'_3 , but additionally one-third of the mature females survive to breed in the following year, so $x'_3 = \frac{1}{3}x_2 + \frac{1}{3}x_3$.

- (b) Writing these equations into vector form

$$\mathbf{x}' = \begin{bmatrix} x'_1 \\ x'_2 \\ x'_3 \end{bmatrix} = \begin{bmatrix} 4x_3 \\ \frac{1}{2}x_1 \\ \frac{1}{3}x_2 + \frac{1}{3}x_3 \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix}}_L \mathbf{x}.$$

- (c) Given the initial numbers of female animals is $\mathbf{x} = (60, 70, 20)$, the number of females after one year is then predicted by the matrix-vector product

$$\mathbf{x}' = L\mathbf{x} = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 60 \\ 70 \\ 20 \end{bmatrix} = \begin{bmatrix} 80 \\ 30 \\ 30 \end{bmatrix}.$$

That is, the predicted numbers of females are 80 pups, 30 juveniles, and 30 mature.

After a second year the number of females is then predicted by the matrix-vector product $\mathbf{x}'' = L\mathbf{x}'$. Here

$$\mathbf{x}'' = L\mathbf{x}' = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 80 \\ 30 \\ 30 \end{bmatrix} = \begin{bmatrix} 120 \\ 40 \\ 20 \end{bmatrix}.$$

After a third year the number of females is predicted by the

matrix-vector product $\mathbf{x}''' = L\mathbf{x}''$. Here

$$\mathbf{x}''' = L\mathbf{x}'' = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 120 \\ 40 \\ 20 \end{bmatrix} = \begin{bmatrix} 80 \\ 60 \\ 20 \end{bmatrix}.$$

■

Matrix-matrix multiplication

Matrix-vector multiplication explicitly writes the vector in its equivalent form as an $n \times 1$ matrix—a matrix with one column. Such multiplication immediately generalises to the case of a right-hand matrix with multiple columns.

Example 3.1.6. Let matrices

$$A = \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 5 & -6 & 4 \\ -1 & -3 & -3 \end{bmatrix},$$

then the matrix multiplication AB may be done as the matrix A multiplying each of the three columns in B . That is, in detail write

$$\begin{aligned} AB &= A \begin{bmatrix} 5 & -6 & 4 \\ -1 & -3 & -3 \end{bmatrix} \\ &= A \begin{bmatrix} 5 \\ -1 \\ -3 \end{bmatrix} : \begin{bmatrix} -6 \\ -3 \\ -3 \end{bmatrix} : \begin{bmatrix} 4 \\ -3 \\ -3 \end{bmatrix} \\ &= \left[A \begin{bmatrix} 5 \\ -1 \end{bmatrix} : A \begin{bmatrix} -6 \\ -3 \end{bmatrix} : A \begin{bmatrix} 4 \\ -3 \end{bmatrix} \right] \\ &= \left[\begin{bmatrix} 13 \\ -11 \end{bmatrix} : \begin{bmatrix} -24 \\ 9 \end{bmatrix} : \begin{bmatrix} 6 \\ -11 \end{bmatrix} \right] \\ &= \begin{bmatrix} 13 & -24 & 6 \\ -11 & 9 & -11 \end{bmatrix}. \end{aligned}$$

Conversely, the product BA cannot be done because if we follow the same procedure

$$\begin{aligned} BA &= B \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} \\ &= B \left[\begin{bmatrix} 3 \\ -2 \end{bmatrix} : \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right] \\ &= \left[B \begin{bmatrix} 3 \\ -2 \end{bmatrix} : B \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right], \end{aligned}$$

and neither of these matrix-vector products can be done as, for example,

$$B \begin{bmatrix} 3 \\ -2 \end{bmatrix} = \begin{bmatrix} 5 & -6 & 4 \\ -1 & -3 & -3 \end{bmatrix} \begin{bmatrix} 3 \\ -2 \end{bmatrix}$$

the number of columns of the left matrix is not equal to the number of elements of the vector on the right. Hence the product BA is not defined. \blacksquare

Example 3.1.7. Let matrices

$$C = \begin{bmatrix} -4 & -1 \\ -4 & -1 \\ 1 & 4 \end{bmatrix}, \quad D = \begin{bmatrix} -2 & -1 & -3 \\ 1 & 3 & 0 \end{bmatrix}.$$

Compute, if possible, CD and DC ; compare these products.

Solution: • On the one hand,

$$\begin{aligned} CD &= C \begin{bmatrix} -2 & -1 & -3 \\ 1 & 3 & 0 \end{bmatrix} \\ &= \left[C \begin{bmatrix} -2 \\ 1 \end{bmatrix} : C \begin{bmatrix} -1 \\ 3 \end{bmatrix} : C \begin{bmatrix} -3 \\ 0 \end{bmatrix} \right] \\ &= \begin{bmatrix} 7 & 1 & 12 \\ 7 & 1 & 12 \\ 2 & 11 & -3 \end{bmatrix}. \end{aligned}$$

• Conversely,

$$\begin{aligned} DC &= D \begin{bmatrix} -4 & -1 \\ -4 & -1 \\ 1 & 4 \end{bmatrix} \\ &= \left[D \begin{bmatrix} -4 \\ -4 \\ 1 \end{bmatrix} : D \begin{bmatrix} -1 \\ -1 \\ 4 \end{bmatrix} \right] \\ &= \begin{bmatrix} 9 & -9 \\ -16 & -4 \end{bmatrix}. \end{aligned}$$

Interestingly, $CD \neq DC$. \blacksquare

Definition 3.1.8 (matrix product). *Let matrix A be $m \times n$, and matrix B be $n \times p$, then the **matrix product** $C = AB$ is the $m \times p$ matrix whose (i, j) th entry is*

$$c_{ij} = a_{i1}b_{1j} + a_{i2}b_{2j} + \cdots + a_{in}b_{nj}.$$

This formula looks like a dot product of two vectors (1.3.2): indeed we do use that the expression for the (i, j) th entry is the dot product

of the i th row of A and the j th column of B as illustrated by

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ \vdots & \vdots & & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{in} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix} \begin{bmatrix} b_{11} & \cdots & b_{1j} & \cdots & b_{1p} \\ b_{21} & \cdots & b_{2j} & \cdots & b_{2p} \\ \vdots & & \vdots & & \vdots \\ b_{n1} & \cdots & b_{nj} & \cdots & b_{np} \end{bmatrix}.$$

As seen in the examples, although the two matrices A and B may be of different sizes, the number of columns of A must equal the number of rows of B in order for the product AB to be defined.

Example 3.1.9. Matrix multiplication leads to powers of a square matrix.

Let matrix

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

then by A^2 we mean the product

$$AA = \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 8 \\ -8 & -3 \end{bmatrix},$$

and by A^3 we mean the product

$$AAA = AA^2 = \begin{bmatrix} 5 & 8 \\ -8 & -3 \end{bmatrix} \begin{bmatrix} 3 & 2 \\ -2 & 1 \end{bmatrix} = \begin{bmatrix} -1 & 18 \\ -18 & -19 \end{bmatrix},$$

and so on. ■

In general, for an $n \times n$ square matrix A and a positive integer exponent p we define the **matrix power**

$$A^p = \underbrace{AA \cdots A}_{p \text{ factors}}.$$

The matrix powers A^p are also $n \times n$ square matrices.

Example 3.1.10 (age structured ecology). Matrix powers occur naturally in modelling populations by ecologists such as the animals of Example 3.1.5. Recall that given the numbers of female pups, juveniles and mature aged formed into a vector $\mathbf{x} = (x_1, x_2, x_3)$, the number in each age one year later (indicated here by a dash) is $\mathbf{x}' = L\mathbf{x}$ for Leslie matrix

$$L = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix}.$$

Hence the number in each age category two years later (indicated here by two dashes) is

$$\mathbf{x}'' = L\mathbf{x}' = L(L\mathbf{x}) = (LL)\mathbf{x} = L^2\mathbf{x},$$

provided that matrix multiplication is associative (Theorem 3.1.18c) to enable us to write $L(L\mathbf{x}) = (LL)\mathbf{x}$. Then the matrix square

$$L^2 = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} = \begin{bmatrix} 0 & \frac{4}{3} & \frac{4}{3} \\ 0 & 0 & 2 \\ \frac{1}{6} & \frac{1}{9} & \frac{1}{9} \end{bmatrix}.$$

Continuing to use such associativity, the number in each age category three years later (indicated here by threes dashes) is

$$\mathbf{x}''' = L\mathbf{x}'' = L(L^2\mathbf{x}) = (LL^2)\mathbf{x} = L^3\mathbf{x},$$

where the matrix cube

$$L^3 = LL^2 = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 0 & \frac{4}{3} & \frac{4}{3} \\ 0 & 0 & 2 \\ \frac{1}{6} & \frac{1}{9} & \frac{1}{9} \end{bmatrix} = \begin{bmatrix} \frac{2}{3} & \frac{4}{9} & \frac{4}{9} \\ 0 & \frac{2}{3} & \frac{2}{3} \\ \frac{1}{18} & \frac{1}{27} & \frac{19}{27} \end{bmatrix}.$$

That is, the powers of the Leslie matrix help predict what happens two, three, or more years into the future. ■

The transpose of a matrix

The operations so far defined for matrices correspond directly to analogous operations for real numbers. The transpose has no corresponding analogue. At first mysterious, the transpose occurs frequently—often due to it linking the dot product of vectors with matrix multiplication. The transpose also reflects symmetry in applications (Chapter 4), such as Newton’s law that every action has an equal and opposite reaction.

Example 3.1.11. Let matrices

$$A = \begin{bmatrix} -4 & 2 \\ -3 & 4 \\ -1 & -7 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & 0 & -1 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 1 & 1 \\ -1 & -3 & 0 \\ 2 & 3 & 2 \end{bmatrix}.$$

Then obtain the transpose of each of these three matrices by writing each of their rows as columns, in order:

$$A^T = \begin{bmatrix} -4 & -3 & -1 \\ 2 & 4 & -7 \end{bmatrix}, \quad B^T = \begin{bmatrix} 2 \\ 0 \\ -1 \end{bmatrix}, \quad C^T = \begin{bmatrix} 1 & -1 & 2 \\ 1 & -3 & 3 \\ 1 & 0 & 2 \end{bmatrix}.$$

■

These examples illustrate the following definition.

Definition 3.1.12 (transpose). *The transpose of an $m \times n$ matrix A is the $n \times m$ matrix, denoted A^T , obtained by writing the i th row of A as the i th column of A^T , or equivalently by writing the j th column of A to be the j th row of A^T . That is, if $B = A^T$, then $b_{ij} = a_{ji}$.*

Example 3.1.13 (transpose and dot product). Consider two vectors in \mathbb{R}^n , say $\mathbf{u} = (u_1, u_2, \dots, u_n)$ and $\mathbf{v} = (v_1, v_2, \dots, v_n)$; that is,

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}.$$

Then the dot product between the two vectors

$$\mathbf{u} \cdot \mathbf{v} = u_1 v_1 + u_2 v_2 + \cdots + u_n v_n \quad (\text{Defn. 1.3.2})$$

$$= [u_1 \ u_2 \ \cdots \ u_n] \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \quad (\text{matrix mult.})$$

$$= \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}^T \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \quad (\text{matrix transpose})$$

$$= \mathbf{u}^T \mathbf{v}.$$

Subsequent sections and chapters often use this identity, that $\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v}$. ■

Definition 3.1.14 (symmetry). A matrix A is a **symmetric matrix** if $A^T = A$; that is, if the matrix is equal to its transpose.

A symmetric matrix must be a square matrix—as otherwise the sizes of A and A^T would be different and so the matrices could not be equal.

Example 3.1.15. None of the three matrices in Example 3.1.11 are symmetric: the first two matrices are not square so cannot be symmetric, and the third $C \neq C^T$. The following matrix is symmetric:

$$D = \begin{bmatrix} 2 & 0 & 1 \\ 0 & -6 & 3 \\ 1 & 3 & 4 \end{bmatrix} = D^T.$$

When is the following general 2×2 matrix symmetric?

$$E = \begin{bmatrix} a & b \\ c & d \end{bmatrix}.$$

Solution: Consider the transpose

$$E^T = \begin{bmatrix} a & c \\ b & d \end{bmatrix} \quad \text{compared with } E = \begin{bmatrix} a & b \\ c & d \end{bmatrix}.$$

The top-left and bottom-right elements are always the same. The top-right and bottom-left elements will be the same if and only if (iff) $b = c$. That is, the 2×2 matrix E is symmetric iff $b = c$.

Table 3.1: As well as the basics of Matlab/Octave listed in Table 1.2 and 2.3, we need these matrix operations.

- `size(A)` returns the number of rows and columns of matrix A : if A is $m \times n$, then `size(A)` returns $[m \ n]$.
- $A(i,j)$ is the (i, j) th entry of a matrix A , $A(:,j)$ is the j th column, $A(i,:)$ is the i th row; either to use the value(s) or to assign value(s).
- $+, -, *$ is matrix/vector/scalar addition, subtraction, and multiplication, but only provided the sizes of the two operands are compatible.
- A^p for scalar p computes the p th power of square matrix A (in contrast to $A.^p$ which computes the p th power of each element of A , Table 2.3).
- The character single quote, A' , transposes the matrix A .
- Predefined matrices include:
 - `zeros(m,n)` is the zero matrix $O_{m \times n}$;
 - `eye(m,n)` is $m \times n$ ‘identity matrix’ $I_{m \times n}$;
 - `ones(m,n)` is the $m \times n$ matrix where all entries are one;
 - `randn(m,n)` is a $m \times n$ matrix with random entries (distributed Normally, mean zero, standard deviation one).

A single argument gives the square matrix version:

- `zeros(n)` is $O_n = O_{n \times n}$;
- `eye(n)` is the $n \times n$ identity matrix $I_n = I_{n \times n}$;
- `ones(n)` is the $n \times n$ matrix of all ones;
- `randn(n)` is an $n \times n$ matrix with random entries.

With no argument, these functions return the corresponding scalar: for example, `randn` computes a single random number.

- Very large and small magnitude numbers are printed in Matlab/Octave like the following:
 - `4.852e+08` denotes the large $4.852 \cdot 10^8$; whereas
 - `3.469e-16` denotes the small $3.469 \cdot 10^{-16}$.

Symmetric matrices of note are the $n \times n$ identity matrix and $n \times n$ zero matrix, I_n and O_n .

Compute in Matlab/Octave

Matlab/Octave empowers us to compute all these operations quickly, especially for the large problems found in applications: after all, Matlab is an abbreviation of *Matrix Laboratory*. Table 3.1 summarises the Matlab/Octave version of the operations introduced so far, and used in subsequent sections.

Matrix size and elements Let the matrix

$$A = \begin{bmatrix} 0 & 0 & -2 & -11 & 5 \\ 0 & 1 & -1 & 11 & -8 \\ -4 & 2 & 10 & 2 & -3 \end{bmatrix}.$$

We readily see this is a 3×5 matrix, but to check that Matlab/Octave agrees, execute the following in Matlab/Octave:

```
A=[0 0 -2 -11 5
   0 1 -1 11 -8
   -4 2 10 2 -3]
size(A)
```



The answer, “3 5”, confirms A is 3×5 . Matlab/Octave accesses individual elements, rows and columns. For example, execute each of the following:

- $A(2,4)$ gives a_{24} which here is 11;
- $A(:,5)$ is the the fifth column vector, here $\begin{bmatrix} 5 \\ -8 \\ -3 \end{bmatrix}$;
- $A(1,:)$ is the first row, here $[0 \ 0 \ -2 \ -11 \ 5]$.

One may also use these constructs to change the elements in matrix A : for example, executing $A(2,4)=9$ changes matrix A to

$$\begin{aligned} A = \\ \begin{array}{ccccc} 0 & 0 & -2 & -11 & 5 \\ 0 & 1 & -1 & 9 & -8 \\ -4 & 2 & 10 & 2 & -3 \end{array} \end{aligned}$$

then $A(:,5)=[2;-3;1]$ changes matrix A to

$$\begin{aligned} A = \\ \begin{array}{ccccc} 0 & 0 & -2 & -11 & 2 \\ 0 & 1 & -1 & 9 & -3 \\ -4 & 2 & 10 & 2 & 1 \end{array} \end{aligned}$$

whereas $A(1,:)=[1 \ 2 \ 3 \ 4 \ 5]$ changes matrix A to

$$\begin{aligned} A = \\ \begin{array}{ccccc} 1 & 2 & 3 & 4 & 5 \\ 0 & 1 & -1 & 9 & -3 \\ -4 & 2 & 10 & 2 & 1 \end{array} \end{aligned}$$

Matrix addition and subtraction To illustrate further operations let's use some random matrices generated by Matlab/Octave: you will generate different matrices to the following, but the operations will work the same. Table 3.1 mentions that `randn(m,n)` generates a random matrix so execute say

```
A=randn(4)
B=randn(4)
C=randn(4,2)
```



and obtain matrices such as (2 d.p.)

```
A =
-1.31   2.07   0.08   2.05
 1.25  -1.35  -1.00   1.94
 1.08   1.79  -0.99   0.93
 1.34  -0.99  -0.23  -0.22

B =
 1.21  -0.46   0.09   0.58
 1.67  -1.96   1.26   1.93
 0.24  -0.46   2.77  -0.59
 0.03  -0.28  -0.76   0.13

C =
 1.14   0.85
-0.48   0.17
 0.37  -0.64
 0.62  -1.17
```

Then $A+B$ gives here the sum

```
ans =
-0.10   1.62   0.17   2.63
 2.92  -3.31   0.26   3.87
 1.31   1.33   1.78   0.34
 1.37  -1.27  -0.99  -0.09
```

and $A-B$ the difference

```
ans =
-2.52   2.53  -0.01   1.46
-0.41   0.62  -2.25   0.01
 0.84   2.26  -3.76   1.52
 1.31  -0.71   0.53  -0.35
```

You could check that $B+A$ gives the same matrix as $A+B$ (Theorem 3.1.16a) by seeing that their difference is the 3×5 zero matrix: execute $(A+B)-(B+A)$ (the parentheses control the order of evaluation). However, expressions such as $B+C$ and $A-C$ give an error, because the matrices are of incompatible sizes, reported by Matlab as

```
Error using +
Matrix dimensions must agree.
```

or reported by Octave as

```
error: operator +: nonconformant arguments
```



Scalar multiplication of matrices In Matlab/Octave the asterisk indicates multiplication. Scalar multiplication can be done either way around. For example, generate a random 2×4 matrix A and compute $2A$ and $A\frac{1}{10}$. These commands

```
A=randn(4,3)
```

```
2*A
```

```
A*0.1
```

might give the following (2 d.p.)

```
A =
```

0.82	2.54	-0.98
2.30	0.05	2.63
-1.45	2.15	0.89
-2.58	-0.09	-0.55

```
>> 2*A
```

```
ans =
```

1.64	5.07	-1.97
4.61	0.10	5.25
-2.90	4.30	1.77
-5.16	-0.18	-1.11

```
>> A*0.1
```

```
ans =
```

0.08	0.25	-0.10
0.23	0.00	0.26
-0.15	0.21	0.09
-0.26	-0.01	-0.06

Division by a scalar is also defined in Matlab/Octave and means multiplication by the reciprocal; for example, the product $A*0.1$ could equally well be computed as $A/10$.

In mathematical algebra we would not normally accept statements such as $A + 3$ or $2A - 5$ because addition and subtraction with matrices has only been defined between matrices of the same size.⁵ However, Matlab/Octave usefully extends addition and subtraction so that $A+3$ and $2*A-5$ mean add three to *every* element of A and subtract five from *every* element of $2A$. For example, with the above random 4×3 matrix A ,

```
>> A+3
```

```
ans =
```

3.82	5.54	2.02
5.30	3.05	5.63
1.55	5.15	3.89
0.42	2.91	2.45

⁵ Although in some contexts such mathematical expressions are routinely accepted, be careful of their meaning.

```
>> 2*A-5
ans =
-3.36    0.07   -6.97
-0.39   -4.90    0.25
-7.90   -0.70   -3.23
-10.16   -5.18   -6.11
```

This last computation illustrates that in any expression the operations of multiplication and division are performed before additions and subtractions—as normal in mathematics.

Matrix multiplication In Matlab/Octave the asterisk also invokes matrix-matrix and matrix-vector multiplication. For example, generate and multiply two random matrices say of size 3×4 and 4×3 (2 d.p.) with

```
A=randn(3,4)
B=randn(4,2)
C=A*B
```

might give the following result

```
A =
-0.02    1.31   -0.74   -0.49
-0.36   -1.30   -0.23    0.41
-0.88   -0.34    0.28   -0.99

B =
-1.32   -0.79
  0.71    1.48
 -0.48    2.79
  1.40   -0.41

>> C=A*B
C =
  0.62    0.10
  0.24   -2.44
 -0.60    1.38
```

Without going into excruciating arithmetic detail this product is hard to check. However, we can check several things such as the first row of A times the first column of B gives c_{11} by computing $A(1,:)*B(:,1)$ and seeing it does give 0.62 as required. Also check that the two columns of C may be viewed as the two matrix-vector products Ab_1 and Ab_2 by comparing C with $[A*B(:,1) A*B(:,2)]$ and seeing they are the same.

Recall that in a matrix product the number of columns of the left matrix have to be the same as the number of rows of the right matrix. Matlab/Octave gives an error message if this is not the



case, such as occurs upon asking it to compute $B*A$ when Matlab reports

```
Error using *
Inner matrix dimensions must agree.
```

and Octave reports

```
error: operator *: nonconformant arguments
```

The caret symbol computes matrix powers in Matlab/Octave, such as the cube A^3 , but it only makes sense and works for square matrices A .⁶ For example, if matrix A was 3×4 , then $A^2 = AA$ would involve multiplying a 3×4 matrix by a 3×4 matrix: since the number of columns of the left A is not the same as the number of rows of the right A such a multiplication is not allowed.

The transpose and symmetry In Matlab/Octave the single apostrophe denotes matrix transpose. For example, see it transpose a couple of random matrices with

```
A=randn(3,4)
B=randn(4,2)
A'
B'
```

giving here for example (2 d.p.)

```
A =
0.80 0.30 -0.12 -0.57
0.07 -0.51 -0.81 1.95
0.29 -0.10 0.17 0.70
```

```
B =
-0.71 -0.34
-0.33 -0.73
1.11 -0.21
0.41 0.33
```

```
>> A'
ans =
0.80 0.07 0.29
0.30 -0.51 -0.10
-0.12 -0.81 0.17
-0.57 1.95 0.70
```

```
>> B'
ans =
-0.71 -0.33 1.11 0.41
-0.34 -0.73 -0.21 0.33
```

⁶ We define matrix powers for only integer power. Matlab/Octave will compute the power of a square matrix for any real/complex exponent, but its meaning involves matrix exponentials and logarithms that we do not explore here.

One can do further operations after the transposition, such as checking the multiplication rule that $(AB)^T = B^T A^T$ (Theorem 3.1.21d) by verifying $(A*B)' - B'*A'$ is the zero matrix, here $O_{2 \times 3}$.

You can generate a symmetric matrix by adding a square matrix to its transpose (Theorem 3.1.21f): for example, generate a random square matrix by first `C=randn(3)` then `C=C+C'` makes a random symmetric matrix such as the following (2 d.p.)

```
>> C=randn(3)
C =
-0.33    0.65   -0.62
-0.43   -2.18   -0.28
 1.86   -1.00   -0.52

>> C=C+C'
C =
-0.65    0.22    1.24
 0.22   -4.36   -1.28
 1.24   -1.28   -1.04

>> C-C'
ans =
 0.00    0.00    0.00
 0.00    0.00    0.00
 0.00    0.00    0.00
```

That the resulting matrix C is symmetric is checked by this last step which computes the difference between C and C^T and confirming the difference is zero.

3.1.3 Familiar algebraic properties of matrix operations

Almost all of the familiar algebraic properties of scalar addition, subtraction and multiplication—namely commutativity, associativity and distributivity—hold for matrix addition, subtraction and multiplication.

The one outstanding exception is that matrix multiplication is *not* commutative: for matrices A and B the products AB and BA are usually not equal. We are used to non-commutativity in life. For example, when you go home, to enter your house you first open the door, second walk in, and third close the door. You cannot swap the order and try to walk in before opening the door—the operations do not commute. Similarly, for another example, I often teach classes on the third floor of a building next to my office: after finishing classes, first I walk downstairs to ground level, and second I cross the road to my office. If I try to cross the road before going downstairs, then the force of gravity has something very painful to say about the outcome—the operations do not commute. Similar

to these analogues, the order of matrix multiplication makes a difference to the result.

Theorem 3.1.16 (Properties of addition and scalar multiplication). *Let matrices A , B and C be of the same size, and let c and d be scalars. Then:*

- (a) $A + B = B + A$ (*commutativity of addition*);
- (b) $(A + B) + C = A + (B + C)$ (*associativity of addition*);
- (c) $A \pm O = A = O + A$;
- (d) $c(A \pm B) = cA \pm cB$ (*distributivity over matrix addition*);
- (e) $(c \pm d)A = cA \pm dA$ (*distributivity over scalar addition*);
- (f) $c(dA) = (cd)A$ (*associativity of scalar multiplication*);
- (g) $1A = A$; and
- (h) $0A = O$.

Proof. The proofs directly match those of the corresponding vector properties and are set as exercises. \square

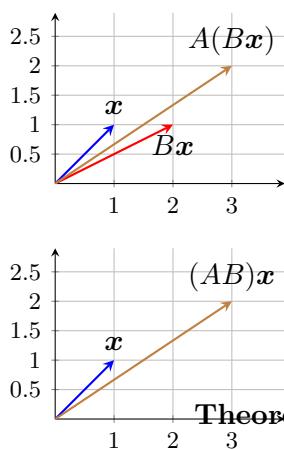
Example 3.1.17 (geometry of associativity). Many properties of matrix multiplication have a useful geometric interpretation such as that discussed for matrix-vector products. Recall the earlier Example 3.1.10 invoked the associativity Theorem 3.1.18c. For another example, consider the two matrices and vector

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & 0 \\ 2 & -1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}.$$

Now the transform $\mathbf{x}' = B\mathbf{x} = (2, 1)$, and then transforming with A gives $\mathbf{x}'' = A\mathbf{x}' = A(B\mathbf{x}) = (3, 2)$, as illustrated in the margin. This is the same results as forming the product

$$AB = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 2 & -1 \end{bmatrix} = \begin{bmatrix} 4 & -1 \\ 2 & 0 \end{bmatrix}$$

and then computing $(AB)\mathbf{x} = (3, 2)$ as also illustrated in the margin. Such associativity asserts that $A(B\mathbf{x}) = (AB)\mathbf{x}$: that is, the geometric transform of \mathbf{x} by matrix B followed by the transform of matrix A is the same result as just transforming by the matrix formed from the product AB —as assured by Theorem 3.1.18c. ■



Theorem 3.1.18 (properties of matrix multiplication). *Let matrices A , B and C be of sizes such that the following expressions are defined, and let c be a scalar, then:*

- (a) $A(B \pm C) = AB \pm AC$ (*distributivity of matrix multiplication*);
- (b) $(A \pm B)C = AC \pm BC$ (*distributivity of matrix multiplication*);

- (c) $A(BC) = (AB)C$ (associativity of matrix multiplication);
- (d) $c(AB) = (cA)B = A(cB)$;
- (e) $I_m A = A = AI_n$ for $m \times n$ matrix A (multiplicative identity);
- (f) $O_m A = O_{m \times n} = AO_n$ for $m \times n$ matrix A ;
- (g) $A^p A^q = A^{p+q}$, $(A^p)^q = A^{pq}$ and $(cA)^p = c^p A^p$ for square A and for positive integers p and q .⁷

Proof. Let's document a few proofs, others are exercises.

3.1.18a : The direct proof involves some long expressions involving the entries of $m \times n$ matrix A , and $n \times p$ matrices B and C . Let $(\cdot)_{ij}$ denote the (i,j) th entry of whatever matrix expression is inside the parentheses. By Definition 3.1.8 of matrix multiplication

$$\begin{aligned}
 & (A(B \pm C))_{ij} \\
 &= a_{i1}(B \pm C)_{1j} + a_{i2}(B \pm C)_{2j} + \cdots + a_{in}(B \pm C)_{nj} \\
 &\quad (\text{by definition of matrix addition}) \\
 &= a_{i1}(b_{1j} \pm c_{1j}) + a_{i2}(b_{2j} \pm c_{2j}) + \cdots + a_{in}(b_{nj} \pm c_{nj}) \\
 &\quad (\text{distributing the scalar multiplications}) \\
 &= a_{i1}b_{1j} \pm a_{i1}c_{1j} + a_{i2}b_{2j} \pm a_{i2}c_{2j} + \cdots + a_{in}b_{nj} \pm a_{in}c_{nj} \\
 &\quad (\text{upon reordering terms in the sum}) \\
 &= a_{i1}b_{1j} + a_{i2}b_{2j} + \cdots + a_{in}b_{nj} \\
 &\quad \pm (a_{i1}c_{1j} + a_{i2}c_{2j} + \cdots + a_{in}c_{nj}) \\
 &\quad (\text{using Definition 3.1.8 for matrix products}) \\
 &= (AB)_{ij} \pm (AC)_{ij}.
 \end{aligned}$$

Since this identity holds for all indices i and j , the matrix identity $A(B \pm C) = AB \pm AC$ holds, proving Theorem 3.1.18a.

3.1.18c : Associativity involves some longer expressions involving the entries of $m \times n$ matrix A , $n \times p$ matrix B , and $p \times q$ matrix C . By Definition 3.1.8 of matrix multiplication

$$\begin{aligned}
 (A(BC))_{ij} &= a_{i1}(BC)_{1j} + a_{i2}(BC)_{2j} + \cdots + a_{in}(BC)_{nj} \\
 &\quad (\text{then using Definition 3.1.8 for } BC) \\
 &= a_{i1}(b_{11}c_{1j} + b_{12}c_{2j} + \cdots + b_{1p}c_{pj}) \\
 &\quad + a_{i2}(b_{21}c_{1j} + b_{22}c_{2j} + \cdots + b_{2p}c_{pj}) \\
 &\quad + \cdots \\
 &\quad + a_{in}(b_{n1}c_{1j} + b_{n2}c_{2j} + \cdots + b_{np}c_{pj}) \\
 &\quad (\text{distributing the scalar multiplications}) \\
 &= a_{i1}b_{11}c_{1j} + a_{i1}b_{12}c_{2j} + \cdots + a_{i1}b_{1p}c_{pj} \\
 &\quad + a_{i2}b_{21}c_{1j} + a_{i2}b_{22}c_{2j} + \cdots + a_{i2}b_{2p}c_{pj}
 \end{aligned}$$

⁷ Generically, these exponent properties hold for all scalar p and q , although we have to be very careful with non-integer exponents.

$$\begin{aligned}
& + \dots \\
& + a_{in}b_{n1}c_{1j} + a_{in}b_{n2}c_{2j} + \dots + a_{in}b_{np}c_{pj} \\
& \quad (\text{reordering the terms—transpose}) \\
= & \quad a_{i1}b_{11}c_{1j} + a_{i2}b_{21}c_{1j} + \dots + a_{in}b_{n1}c_{1j} \\
& + a_{i1}b_{12}c_{2j} + a_{i2}b_{22}c_{2j} + \dots + a_{in}b_{n2}c_{2j} \\
& + \dots \\
& + a_{i1}b_{1p}c_{pj} + a_{i2}b_{2p}c_{pj} + \dots + a_{in}b_{np}c_{pj} \\
& \quad (\text{factoring } c_{1j}, c_{2j}, \dots, c_{pj}) \\
= & \quad (a_{i1}b_{11} + a_{i2}b_{21} + \dots + a_{in}b_{n1})c_{1j} \\
& + (a_{i1}b_{12} + a_{i2}b_{22} + \dots + a_{in}b_{n2})c_{2j} \\
& + \dots \\
& + (a_{i1}b_{1p} + a_{i2}b_{2p} + \dots + a_{in}b_{np})c_{pj} \\
& \quad (\text{recognising the entries for } (AB)_{ik}) \\
= & \quad (AB)_{i1}c_{1j} + (AB)_{i2}c_{2j} + \dots + (AB)_{ip}c_{pj} \\
& \quad (\text{again using Definition 3.1.8}) \\
= & \quad ((AB)C)_{ij}.
\end{aligned}$$

Since this identity holds for all indices i and j , the matrix identity $A(BC) = (AB)C$ holds, proving Theorem 3.1.18c.

3.1.18g : Other proofs develop from previous parts of the theorem. For example, to establish $A^p A^q = A^{p+q}$ start from the definition of matrix powers:

$$\begin{aligned}
A^p A^q &= \underbrace{AA \cdots A}_{p \text{ times}} \underbrace{AA \cdots A}_{q \text{ times}} \\
&\quad (\text{using associativity, Thm. 3.1.18c}) \\
&= \underbrace{AA \cdots A}_{p+q \text{ times}} \\
&= A^{p+q}.
\end{aligned}$$

□

Example 3.1.19. Show that $(A + B)^2 \neq A^2 + 2AB + B^2$ in general.

Solution: Consider

$$\begin{aligned}
(A + B)^2 &= (A + B)(A + B) \quad (\text{matrix power}) \\
&= A(A + B) + B(A + B) \quad (\text{Thm. 3.1.18b}) \\
&= AA + AB + BA + BB \quad (\text{Thm. 3.1.18a}) \\
&= A^2 + AB + BA + B^2 \quad (\text{matrix power}).
\end{aligned}$$

This expression is only equal to $A^2 + 2AB + B^2$ if we can replace BA by AB . But this requires $BA = AB$ which is generally not true. That is, $(A + B)^2 = A^2 + 2AB + B^2$ only if $BA = AB$.

■

Example 3.1.20. Show that the matrix $J = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$ is not a multiplicative identity (despite having ones down a diagonal, this diagonal is the wrong one for an identity).

Solution: Among many other ways to show J is not a multiplicative identity, let's invoke a general 3×3 matrix

$$A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix},$$

and evaluate the product

$$JA = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} = \dots = \begin{bmatrix} g & h & i \\ d & e & f \\ a & b & c \end{bmatrix} \neq A.$$

Since $JA \neq A$ then matrix J cannot be a multiplicative identity (the multiplicative identity is only when the ones are along the diagonal from top-left to bottom-right). ■

Theorem 3.1.21 (properties of transpose). *Let matrices A and B be of sizes such that the following expressions are defined, then:*

- (a) $(A^T)^T = A$;
- (b) $(A \pm B)^T = A^T \pm B^T$;
- (c) $(cA)^T = c(A^T)$ for any scalar c ;
- (d) $(AB)^T = B^T A^T$;
- (e) $(A^p)^T = (A^T)^p$ for all positive integer exponents p ;⁸
- (f) $A + A^T$, $A^T A$ and AA^T are symmetric matrices.

Proof. Let's document a few proofs, others are exercises. Some proofs use primitive definitions—usually using $(\cdot)_{ij}$ to denote the (i,j) th entry of whatever matrix expression is inside the parentheses—others invoke earlier proved parts.

3.1.21b : Recall from Definition 3.1.12 of the transpose that

$$\begin{aligned} ((A \pm B)^T)_{ij} &= (A \pm B)_{ji} \\ &\quad (\text{by definition of addition}) \\ &= a_{ji} \pm b_{ji} \\ &\quad (\text{by Defn. 3.1.12 of transpose}) \end{aligned}$$

⁸ With care, this property also holds for all scalar exponents p .

$$= (A^T)_{ij} \pm (B^T)_{ij}.$$

Since this identity holds for all indices i and j , then $(A \pm B)^T = A^T \pm B^T$.

3.1.21d : The transpose of matrix multiplication is more involved. Let matrices A and B be of sizes $m \times n$ and $n \times p$ respectively. Then from Definition 3.1.12 of the transpose

$$\begin{aligned} ((AB)^T)_{ij} &= (AB)_{ji} \\ &\quad (\text{by Defn. 3.1.8 of multiplication}) \\ &= a_{j1}b_{1i} + a_{j2}b_{2i} + \cdots + a_{jn}b_{ni} \\ &\quad (\text{commuting the products}) \\ &= b_{1i}a_{j1} + b_{2i}a_{j2} + \cdots + b_{ni}a_{jn} \\ &\quad (\text{by Defn. 3.1.12 of transpose}) \\ &= (B^T)_{i1}(A^T)_{1j} + (B^T)_{i2}(A^T)_{2j} + \cdots + (B^T)_{in}(A^T)_{nj} \\ &\quad (\text{by Defn. 3.1.8 of multiplication}) \\ &= (B^T A^T)_{ij}. \end{aligned}$$

Since this identity holds for all indices i and j , then $(AB)^T = B^T A^T$.

3.1.21f : To prove the second, that $A^T A$ is equal to its transpose, we invoke earlier parts. Consider the transpose

$$\begin{aligned} (A^T A)^T &= (A)^T (A^T)^T \quad (\text{by 3.1.21d}) \\ &= A^T A \quad (\text{by 3.1.21a}). \end{aligned}$$

Since $A^T A$ equals its transpose, it is symmetric.

□

3.1.4 Exercises

Exercise 3.1.1. Consider the following six matrices: $A = \begin{bmatrix} -1 & 3 \\ 0 & -5 \\ 0 & -7 \end{bmatrix}$; $B = \begin{bmatrix} -4 & -3 & -3 & 1 \\ -3 & -2 & 0 & -1 \end{bmatrix}$; $C = \begin{bmatrix} -3 & 1 \end{bmatrix}$; $D = \begin{bmatrix} 0 & 6 & 6 & 3 \\ 2 & 2 & 0 & -5 \end{bmatrix}$; $E = \begin{bmatrix} 0 & 1 & 1 & -2 \\ -1 & 5 & 4 & -1 \\ 1 & -3 & 7 & 3 \\ -6 & -3 & 0 & 2 \end{bmatrix}$; $F = \begin{bmatrix} 4 & 1 & 0 \\ -1 & 1 & 6 \\ -4 & 5 & -2 \end{bmatrix}$.

- (a) What is the size of each of these matrices?
- (b) Which pairs of matrices may be added or subtracted?
- (c) Which matrix multiplications can be performed between two of the matrices?

Exercise 3.1.2. Consider the following six matrices: $A = \begin{bmatrix} 3 & \frac{17}{6} \\ -\frac{5}{3} & \frac{1}{2} \\ -\frac{1}{6} & -\frac{1}{6} \\ \frac{5}{3} & 1 \end{bmatrix}$; $B = \begin{bmatrix} \frac{7}{6} & \frac{1}{3} & \frac{17}{3} \end{bmatrix}$; $C = \begin{bmatrix} -\frac{11}{3} & -\frac{7}{3} \\ \frac{2}{3} & \frac{4}{3} \\ \frac{3}{2} & -\frac{17}{6} \end{bmatrix}$; $D = \begin{bmatrix} 0 & -\frac{13}{6} & 0 & \frac{13}{3} \\ \frac{20}{3} & 2 & -\frac{8}{3} & -\frac{7}{2} \\ \frac{5}{6} & \frac{1}{3} & \frac{13}{6} & -\frac{16}{3} \end{bmatrix}$; $E = \begin{bmatrix} \frac{13}{6} & -\frac{1}{6} \\ -\frac{7}{3} & -5 \end{bmatrix}$; $F = \begin{bmatrix} -\frac{1}{3} \\ \frac{13}{3} \end{bmatrix}$.

- (a) What is the size of each of these matrices?
- (b) Which pairs of matrices may be added or subtracted?
- (c) Which matrix multiplications can be performed between two of the matrices?

Exercise 3.1.3. Given the matrix

$$A = \begin{bmatrix} -0.3 & 2.1 & -4.8 \\ -5.9 & 3.6 & -1.3 \end{bmatrix} :$$

write down its column vectors; what are the values of elements a_{13} and a_{21} ?

Exercise 3.1.4. Given the matrix

$$B = \begin{bmatrix} 7.6 & -1.1 & -0.7 & -4.5 \\ -1.1 & -9.3 & 0.1 & 8.2 \\ 2.6 & 6.9 & 1.2 & -3.6 \\ -1.5 & -7.5 & 3.7 & 2.6 \\ -0.2 & 5.5 & -0.9 & 2.4 \end{bmatrix} :$$

write down its column vectors; what are the values of entries b_{13} , b_{31} , b_{42} ?

Exercise 3.1.5. Write down the column vectors of the identity I_4 . What do we call these column vectors?

Exercise 3.1.6. For the following pairs of matrices, calculate their sum and difference.

$$(a) A = \begin{bmatrix} 2 & 1 & -1 \\ -4 & 1 & -3 \\ -2 & 2 & -1 \end{bmatrix}, B = \begin{bmatrix} 1 & 1 & 0 \\ 4 & -6 & -6 \\ -6 & 4 & 0 \end{bmatrix}$$

$$(b) C = \begin{bmatrix} -2 & -2 & -7 \end{bmatrix}, D = \begin{bmatrix} 4 & 2 & -2 \end{bmatrix}$$

$$(c) P = \begin{bmatrix} -2 & 5 & 1 \\ 3 & -3 & 2 \\ -3 & 3 & -3 \end{bmatrix}, Q = \begin{bmatrix} -1 & -3 & -1 \\ 6 & -4 & -2 \\ 3 & -3 & 1 \end{bmatrix}$$

$$(d) R = \begin{bmatrix} -2.5 & -0.4 \\ -1.0 & -3.5 \\ -3.3 & 1.8 \end{bmatrix}, S = \begin{bmatrix} -0.9 & 4.9 \\ -1.2 & -0.7 \\ -4.0 & -5.4 \end{bmatrix}$$

Exercise 3.1.7. For the given matrix, evaluate the following matrix-scalar products.

$$(a) A = \begin{bmatrix} -3 & -2 \\ 4 & -2 \\ 2 & -4 \end{bmatrix}: -2A, 2A, \text{ and } 3A.$$

$$(b) B = \begin{bmatrix} 4 & 0 \\ -1 & -1 \end{bmatrix}: 1.9B, 2.6B, \text{ and } -6.9B.$$

$$(c) U = \begin{bmatrix} -3.9 & -0.3 & -2.9 \\ 3.1 & -3.9 & -1. \\ 3.1 & -6.5 & 0.9 \end{bmatrix}: -4U, 2U, \text{ and } 4U.$$

$$(d) V = \begin{bmatrix} -2.6 & -3.2 \\ 3.3 & -0.8 \\ -0.3 & 0.3 \end{bmatrix}: 1.3V, -3.7V, \text{ and } 2.5V.$$

Exercise 3.1.8. Use Matlab/Octave to generate some random matrices of a suitable size of your choice, and some random scalars (see Table 3.1). Then confirm the addition and scalar multiplication properties of Theorem 3.1.16. Record all your commands and the output from Matlab/Octave.

Exercise 3.1.9. Use the definition of matrix addition and scalar multiplication to prove the basic properties of Theorem 3.1.16.

Exercise 3.1.10. For each of the given matrices, calculate the specified matrix-vector products.

$$(a) \text{ For } A = \begin{bmatrix} 4 & -3 \\ -2 & 5 \end{bmatrix} \text{ and vectors } \mathbf{p} = \begin{bmatrix} -6 \\ -5 \end{bmatrix}, \mathbf{q} = \begin{bmatrix} -2 \\ -4 \end{bmatrix}, \text{ and } \mathbf{r} = \begin{bmatrix} -3 \\ 1 \end{bmatrix}, \text{ calculate } A\mathbf{p}, A\mathbf{q} \text{ and } A\mathbf{r}.$$

$$(b) \text{ For } B = \begin{bmatrix} 1 & 6 \\ 4 & -5 \end{bmatrix} \text{ and vectors } \mathbf{p} = \begin{bmatrix} -3 \\ -3 \end{bmatrix}, \mathbf{q} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \text{ and } \mathbf{r} = \begin{bmatrix} -5 \\ 2 \end{bmatrix}, \text{ calculate } B\mathbf{p}, B\mathbf{q} \text{ and } B\mathbf{r}.$$

$$(c) \text{ For } C = \begin{bmatrix} -3 & 0 & -3 \\ -1 & -1 & 1 \end{bmatrix} \text{ and vectors } \mathbf{u} = \begin{bmatrix} -4 \\ 3 \\ 2 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} -3 \\ 1 \\ 2 \end{bmatrix}, \text{ and } \mathbf{w} = \begin{bmatrix} -4 \\ 5 \\ -4 \end{bmatrix}, \text{ calculate } C\mathbf{u}, C\mathbf{v} \text{ and } C\mathbf{w}.$$

$$(d) \text{ For } D = \begin{bmatrix} 0 & 4 \\ 1 & 2 \\ -1 & 1 \end{bmatrix} \text{ and vectors } \mathbf{u} = \begin{bmatrix} 3 \\ -0.9 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} 0.9 \\ 6.8 \end{bmatrix}, \text{ and } \mathbf{w} = \begin{bmatrix} 0.3 \\ 7.3 \end{bmatrix}, \text{ calculate } D\mathbf{u}, D\mathbf{v} \text{ and } D\mathbf{w}.$$

Exercise 3.1.11. For each of the given matrices and vectors, calculate the matrix-vector products. Plot in 2D, and label, the vectors and the specified matrix-vector products.

- (a) $A = \begin{bmatrix} 3 & 2 \\ -3 & -1 \end{bmatrix}$, $\mathbf{u} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$, $\mathbf{v} = \begin{bmatrix} 0 \\ -3 \end{bmatrix}$, and $\mathbf{w} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$.
- (b) $B = \begin{bmatrix} 3 & -2 \\ 3 & 2 \end{bmatrix}$, $\mathbf{p} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$, $\mathbf{q} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$, and $\mathbf{r} = \begin{bmatrix} -2 \\ 1 \end{bmatrix}$.
- (c) $C = \begin{bmatrix} -2.1 & 1.1 \\ 4.6 & -1 \end{bmatrix}$, $\mathbf{x}_1 = \begin{bmatrix} 2.1 \\ 0 \end{bmatrix}$, $\mathbf{x}_2 = \begin{bmatrix} -0.1 \\ 1.1 \end{bmatrix}$, and $\mathbf{x}_3 = \begin{bmatrix} -0.3 \\ -1 \end{bmatrix}$.
- (d) $D = \begin{bmatrix} 0.1 & 3.4 \\ 3.9 & 5.1 \end{bmatrix}$, $\mathbf{a} = \begin{bmatrix} 0.2 \\ 0.5 \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} -0.3 \\ 0.3 \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} -0.2 \\ -0.6 \end{bmatrix}$.

Exercise 3.1.12. For each of the given matrices and vectors, calculate the matrix-vector products. Plot in 2D, and label, the vectors and the specified matrix-vector products. For each of the matrices, interpret the matrix multiplication of the vectors as either a rotation, a reflection, a stretch, or none of these.

- (a) $P = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$, $\mathbf{u} = \begin{bmatrix} 1 \\ -1.4 \end{bmatrix}$, $\mathbf{v} = \begin{bmatrix} -3.6 \\ -1.7 \end{bmatrix}$, and $\mathbf{w} = \begin{bmatrix} 0.1 \\ 2.3 \end{bmatrix}$.
- (b) $Q = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$, $\mathbf{p} = \begin{bmatrix} 2.1 \\ 1.9 \end{bmatrix}$, $\mathbf{q} = \begin{bmatrix} 2.8 \\ -1.1 \end{bmatrix}$, and $\mathbf{r} = \begin{bmatrix} 0.8 \\ 3.3 \end{bmatrix}$.
- (c) $R = \begin{bmatrix} 0.8 & -0.6 \\ 0.6 & 0.8 \end{bmatrix}$, $\mathbf{x}_1 = \begin{bmatrix} -4 \\ 2 \end{bmatrix}$, $\mathbf{x}_2 = \begin{bmatrix} 4 \\ -3 \end{bmatrix}$, and $\mathbf{x}_3 = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$.
- (d) $S = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$, $\mathbf{a} = \begin{bmatrix} -1.1 \\ 0 \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} -4.6 \\ -1.5 \end{bmatrix}$, and $\mathbf{c} = \begin{bmatrix} -3.1 \\ 0.9 \end{bmatrix}$.

Exercise 3.1.13. Using the matrix-vector products you calculated for Exercise 3.1.10, write down the results of the following matrix-matrix products.

- (a) For $A = \begin{bmatrix} 4 & -3 \\ -2 & 5 \end{bmatrix}$, write down the matrix products

i. $A \begin{bmatrix} -6 & -2 \\ -5 & -4 \end{bmatrix}$,	ii. $A \begin{bmatrix} -6 & -3 \\ -5 & 1 \end{bmatrix}$,
iii. $A \begin{bmatrix} -2 & -3 \\ -4 & 1 \end{bmatrix}$,	iv. $A \begin{bmatrix} -6 & -2 & -3 \\ -5 & -4 & 1 \end{bmatrix}$.

- (b) For $B = \begin{bmatrix} 1 & 6 \\ 4 & -5 \end{bmatrix}$, write down the matrix products

i. $B \begin{bmatrix} -3 & 2 \\ -3 & 1 \end{bmatrix}$,	ii. $B \begin{bmatrix} -5 & 2 \\ 2 & 1 \end{bmatrix}$,
---	---

$$\text{iii. } B \begin{bmatrix} -5 & -3 \\ 2 & -3 \end{bmatrix}, \quad \text{iv. } B \begin{bmatrix} -5 & 2 & -3 \\ 2 & 1 & -3 \end{bmatrix}.$$

(c) For $C = \begin{bmatrix} -3 & 0 & -3 \\ -1 & -1 & 1 \end{bmatrix}$, write down the matrix products

$$\text{i. } C \begin{bmatrix} -4 & -3 \\ 3 & 1 \\ 2 & 2 \end{bmatrix}, \quad \text{ii. } C \begin{bmatrix} -4 & -3 \\ 5 & 1 \\ -4 & 2 \end{bmatrix},$$

$$\text{iii. } C \begin{bmatrix} -4 & -4 \\ 5 & 3 \\ -4 & 2 \end{bmatrix}, \quad \text{iv. } C \begin{bmatrix} -4 & -3 & -4 \\ 5 & 1 & 3 \\ -4 & 2 & 2 \end{bmatrix}.$$

(d) For $D = \begin{bmatrix} 0 & 4 \\ 1 & 2 \\ -1 & 1 \end{bmatrix}$, write down the matrix products

$$\text{i. } D \begin{bmatrix} 0.9 & 0.3 \\ 6.8 & 7.3 \end{bmatrix}, \quad \text{ii. } D \begin{bmatrix} 0.9 & 3 \\ 6.8 & -0.9 \end{bmatrix},$$

$$\text{iii. } D \begin{bmatrix} 0.3 & 3 \\ 7.3 & -0.9 \end{bmatrix}, \quad \text{iv. } D \begin{bmatrix} 0.9 & 0.3 & 3 \\ 6.8 & 7.3 & -0.9 \end{bmatrix}.$$

Exercise 3.1.14. Use Matlab/Octave to generate some random matrices of a suitable size of your choice, and some random scalars (see Table 3.1). Choose some suitable exponents. Then confirm the matrix multiplication properties of Theorem 3.1.18. Record all your commands and the output from Matlab/Octave.

In checking some properties you may get matrices with elements such as $2.2204\text{e-}16$: recall from Table 3.1 that this denotes the very small number $2.2204 \cdot 10^{-16}$. When adding and subtracting numbers of size one or so, the result $2.2204\text{e-}16$ is effectively zero (due to the sixteen digit precision of Matlab/Octave, Table 1.2).

Exercise 3.1.15. Use Definition 3.1.8 of matrix-matrix multiplication to prove multiplication properties of Theorem 3.1.18. Prove parts: 3.1.18b, distributivity; 3.1.18d, scalar associativity; 3.1.18e, identity; 3.1.18f, zeros.

Exercise 3.1.16. Use the other parts of Theorem 3.1.18 to prove part 3.1.18g that $(A^p)^q = A^{pq}$ and $(cA)^p = c^p A^p$ for square matrix A , scalar c , and for positive integer exponents p and q .

Exercise 3.1.17 (Tasmanian Devils). Ecologists studying a colony of Tasmanian Devils, an Australian marsupial, observed the following: two-thirds of the female newborns survive to be one year old; two-thirds of female one year olds survive to be two years old; one-half of female two year olds survive to be three years old; each year, each female aged two or three years gives birth to two female offspring; female Tasmanian Devils survive for four years, at most.

Analogous to Example 3.1.5 define a vector \mathbf{x} in \mathbb{R}^4 to be the number of females of specified ages. Use the above information to write down the Leslie matrix L that predicts the number in the next year, \mathbf{x}' , from the number in any year, \mathbf{x} . Given the observed initial female numbers of 18 newborns, 9 one year olds, 18 two year olds, and 18 three year olds, use matrix multiplication to predict the numbers of female Tasmanian Devils one, two and three years later. Does the population appear to be increasing? or decreasing?

Exercise 3.1.18. Write down the transpose of each of the following matrices. Which of the following matrices are a symmetric matrix?

$$(a) \begin{bmatrix} -2 & 3 \\ 3 & 0 \\ -8 & 2 \\ -2 & -4 \end{bmatrix}$$

$$(b) \begin{bmatrix} 3 & -4 & -2 & 2 \\ -5 & 2 & -3 & 3 \end{bmatrix}$$

$$(c) \begin{bmatrix} 14 & 5 & 3 & 2 \\ 5 & 0 & -1 & 1 \\ 3 & -1 & -6 & -4 \\ 2 & 1 & -4 & 4 \end{bmatrix}$$

$$(d) [3 \ 1 \ -2 \ -3]$$

$$(e) \begin{bmatrix} 5 & -1 & -2 & 2 \\ 1 & -2 & -2 & 0 \\ -1 & -5 & 4 & -1 \\ 5 & 2 & -1 & -2 \end{bmatrix}$$

$$(f) \begin{bmatrix} -4 & -5.1 & 0.3 \\ -5.1 & -7.4 & -3 \\ 0.3 & -3 & 2.6 \end{bmatrix}$$

$$(g) \begin{bmatrix} -1.5 & -0.6 & -1.7 \\ -1 & -0.4 & -5.6 \end{bmatrix}$$

$$(h) \begin{bmatrix} 1.7 & -0.2 & -0.4 \\ 0.7 & -0.3 & -0.4 \\ 0.6 & 3 & -2.2 \end{bmatrix}$$

Exercise 3.1.19. Are the following matrices symmetric? I_4 , $I_{3 \times 4}$, O_3 , and $O_{3 \times 1}$.

Exercise 3.1.20. Use Matlab/Octave to generate some random matrices of a suitable size of your choice, and some random scalars (see Table 3.1). Choose some suitable exponents. Recalling that in Matlab/Octave the dash ' performs the transpose, confirm the matrix transpose properties of Theorem 3.1.21. Record all your commands and the output from Matlab/Octave.

Exercise 3.1.21. Use Definition 3.1.12 of the matrix transpose to prove properties 3.1.21a and 3.1.21c of Theorem 3.1.21.

Exercise 3.1.22. Use the other parts of Theorem 3.1.21 to prove parts 3.1.21e and 3.1.21f.

3.2 The inverse of a matrix

Section Contents

3.2.1	Introducing the unique inverse	164
3.2.2	Diagonal matrices stretch and shrink	174
	Solve systems whose matrix is diagonal	175
	But do not divide by zero	178
	Stretch or squash the unit square	179
	Sketch convenient coordinates	181
3.2.3	Orthogonal matrices rotate	182
	Orthogonal set of vectors	183
	Orthogonal matrices	184
3.2.4	Exercises	190

The previous Section 3.1 introduced addition, subtraction, multiplication, and other operations of matrices. Conspicuously missing from the list is ‘division’ by a matrix: this section addresses ‘division’ by a matrix as multiplication by the inverse of a matrix. The analogue in ordinary arithmetic is that division by ten is the same as multiplying by its reciprocal, one-tenth. But the inverse of a matrix looks nothing like a reciprocal.

3.2.1 Introducing the unique inverse

Let’s start with an example that illustrates an analogy with the reciprocal/inverse of a scalar number.

Example 3.2.1. Recall that a crucial property is that a number multiplied by its reciprocal/inverse is one: for example, $2 \times 0.5 = 1$ so 0.5 is the reciprocal/inverse of 2. Similarly, show that matrix

$$B = \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} \text{ is an inverse of } A = \begin{bmatrix} 1 & -1 \\ 4 & -3 \end{bmatrix}$$

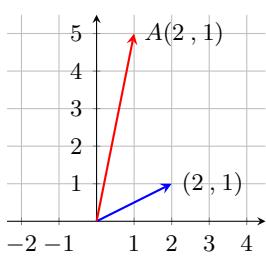
by showing their product is the 2×2 identity I_2 .

Solution: Multiply

$$AB = \begin{bmatrix} 1 & -1 \\ 4 & -3 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2$$

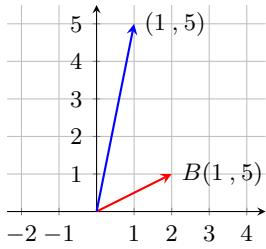
the multiplicative identity. But matrix multiplication is not commutative (section 3.1.3), so also consider

$$BA = \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ 4 & -3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2.$$



That these products are the identity—analogous to the number one in scalar arithmetic—means that the matrix A has the same relation to the matrix B as a number has to its reciprocal/inverse.

Being the inverse, matrix B ‘undoes’ the action of matrix A —as illustrated in the margin. The first picture shows multiplication by A transforms vector $(2, 1)$ to vector $(1, 5)$: $A \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 5 \end{bmatrix}$. The second picture shows that multiplication by B undoes the transform because $B \begin{bmatrix} 1 \\ 5 \end{bmatrix} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ the original vector. ■



The previous example shows at least one case when we can do some sort of matrix ‘division’: that is, multiplying by B is equivalent to ‘dividing’ by A . One restriction is that a clearly defined ‘division’ only works for square matrices because we need to be able to compute both AB and BA .

Definition 3.2.2 (inverse). *If A is an $n \times n$ square matrix, an **inverse** of A is an $n \times n$ matrix B such that $AB = I_n$ and $BA = I_n$. If such a matrix B exists, then A is called **invertible**.*

Example 3.2.3. Show that matrix

$$B = \begin{bmatrix} 0 & -\frac{1}{4} & -\frac{1}{8} \\ \frac{3}{2} & 1 & \frac{7}{8} \\ \frac{1}{2} & \frac{1}{4} & \frac{3}{8} \end{bmatrix} \text{ is an inverse of } A = \begin{bmatrix} 1 & -1 & 5 \\ -5 & -1 & 3 \\ 2 & 2 & -6 \end{bmatrix}.$$

Solution: First compute

$$\begin{aligned} AB &= \begin{bmatrix} 1 & -1 & 5 \\ -5 & -1 & 3 \\ 2 & 2 & -6 \end{bmatrix} \begin{bmatrix} 0 & -\frac{1}{4} & -\frac{1}{8} \\ \frac{3}{2} & 1 & \frac{7}{8} \\ \frac{1}{2} & \frac{1}{4} & \frac{3}{8} \end{bmatrix} \\ &= \begin{bmatrix} 1 \cdot 0 - 1 \cdot \frac{3}{2} + 5 \cdot \frac{1}{2} & 1 \cdot (-\frac{1}{4}) - 1 \cdot 1 + 5 \cdot \frac{1}{4} & 1 \cdot (-\frac{1}{8}) - 1 \cdot \frac{7}{8} + 5 \cdot \frac{3}{8} \\ -5 \cdot 0 - 1 \cdot \frac{3}{2} + 3 \cdot \frac{1}{2} & -5 \cdot (-\frac{1}{4}) - 1 \cdot 1 + 3 \cdot \frac{1}{4} & -5 \cdot (-\frac{1}{8}) - 1 \cdot \frac{7}{8} + 3 \cdot \frac{3}{8} \\ 2 \cdot 0 + 2 \cdot \frac{3}{2} - 6 \cdot \frac{1}{2} & 2 \cdot (-\frac{1}{4}) + 2 \cdot 1 - 6 \cdot \frac{1}{4} & 2 \cdot (-\frac{1}{8}) + 2 \cdot \frac{7}{8} - 6 \cdot \frac{3}{8} \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = I_3. \end{aligned}$$

Second compute

$$BA = \begin{bmatrix} 0 & -\frac{1}{4} & -\frac{1}{8} \\ \frac{3}{2} & 1 & \frac{7}{8} \\ \frac{1}{2} & \frac{1}{4} & \frac{3}{8} \end{bmatrix} \begin{bmatrix} 1 & -1 & 5 \\ -5 & -1 & 3 \\ 2 & 2 & -6 \end{bmatrix}$$

$$\begin{aligned}
&= \begin{bmatrix} 0 \cdot 1 - \frac{1}{4} \cdot (-5) - \frac{1}{8} \cdot 2 & 0 \cdot (-1) - \frac{1}{4} \cdot (-1) - \frac{1}{8} \cdot 2 & 0 \cdot 5 - \frac{1}{4} \cdot 3 - \frac{1}{8} \cdot (-6) \\ \frac{3}{2} \cdot 1 + 1 \cdot (-5) + \frac{7}{8} \cdot 2 & \frac{3}{2} \cdot (-1) + 1 \cdot (-1) + \frac{7}{8} \cdot 2 & \frac{3}{2} \cdot 5 + 1 \cdot 3 + \frac{7}{8} \cdot (-6) \\ \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot (-5) + \frac{3}{8} \cdot 2 & \frac{1}{2} \cdot (-1) + \frac{1}{4} \cdot (-1) + \frac{3}{8} \cdot 2 & \frac{1}{2} \cdot 5 + \frac{1}{4} \cdot 3 + \frac{3}{8} \cdot (-6) \end{bmatrix} \\
&= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = I_3.
\end{aligned}$$

Since both of these products are the identity, then matrix A is invertible, and B is an inverse of A . ■

But even among square matrices, there are many non-zero matrices which do not have inverses! The next Section 3.3 characterises why some matrices do not have an inverse: the reason is associated with both `rcond = 0` (Procedure 2.2.4) and the so-called determinant being zero (Chapter 6).

Example 3.2.4 (no inverse). Prove that the matrix

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

does not have an inverse.

Solution: Assume there is an inverse matrix

$$B = \begin{bmatrix} a & b \\ c & d \end{bmatrix}.$$

Then by Definition 3.2.2 the product $AB = I_2$; that is,

$$\begin{aligned}
AB &= \begin{bmatrix} 1 & -2 \\ -3 & 6 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \\
&= \begin{bmatrix} a - 2c & b - 2d \\ -3a + 6c & -3b + 6d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.
\end{aligned}$$

The bottom-left entry in this matrix equality asserts $-3a + 6c = 0$ which is $-3(a - 2c) = 0$, that is, $a - 2c = 0$. But the top-left entry in the matrix equality asserts $a - 2c = 1$. Both of these equations involving a and c cannot be true simultaneously; therefore the assumption of an inverse must be incorrect. The matrix A does not have an inverse. ■

Theorem 3.2.5 (unique inverse). *If A is an invertible matrix, then its inverse is unique (and denoted by A^{-1}).*

Proof. We suppose there are two inverses, say B_1 and B_2 , and proceed to show they must be the same. Since they are inverses, by Definition 3.2.2 both $AB_1 = B_1A = I_n$ and $AB_2 = B_2A = I_n$. Consequently, using associativity of matrix multiplication,

$$B_1 = B_1I_n = B_1(AB_2) = (B_1A)B_2 = I_nB_2 = B_2.$$

That is, $B_1 = B_2$ and so the inverse is unique. \square

In the elementary case of 1×1 matrices, that is $A = [a_{11}]$, the inverse is simply the reciprocal of the entry, that is $A^{-1} = [1/a_{11}]$ provided a_{11} is non-zero. The reason is that $AA^{-1} = [a_{11} \cdot \frac{1}{a_{11}}] = [1] = I_1$ and $A^{-1}A = [\frac{1}{a_{11}} \cdot a_{11}] = [1] = I_1$.

In the case of 2×2 matrices the inverse is a little more complicated, but should be remembered.

Theorem 3.2.6 (2×2 inverse). *Let 2×2 matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$. Then A is invertible if the **determinant** $ad - bc \neq 0$, in which case*

$$A^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}. \quad (3.2)$$

If the determinant $ad - bc = 0$, then A is not invertible.

Example 3.2.7. (a) Recall Example 3.2.1 verified that

$$B = \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} \text{ is an inverse of } A = \begin{bmatrix} 1 & -1 \\ 4 & -3 \end{bmatrix}.$$

Formula (3.2) gives this inverse from the matrix A : its elements are $a = 1$, $b = -1$, $c = 4$ and $d = -3$ so the determinant $ad - bc = 1 \cdot (-3) - (-1) \cdot 4 = 1$ and hence formula (3.2) derives the inverse

$$A^{-1} = \frac{1}{1} \begin{bmatrix} -3 & -(-1) \\ -4 & 1 \end{bmatrix} = \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} = B.$$

(b) Further, recall Example 3.2.4 proved there is no inverse for matrix

$$A = \begin{bmatrix} 1 & -2 \\ -3 & 6 \end{bmatrix}.$$

Theorem 3.2.6 also establishes this matrix is not invertible because the matrix determinant $ad - bc = 1 \cdot 6 - (-2) \cdot (-3) = 6 - 6 = 0$. \blacksquare

Proof. To prove Theorem 3.2.6, first show the given A^{-1} satisfies Definition 3.2.2 when the determinant $ad - bc \neq 0$ (and using associativity of scalar-matrix multiplication, Theorem 3.1.18d). For the proposed A^{-1} , on the one hand,

$$\begin{aligned} A^{-1}A &= \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \\ &= \frac{1}{ad - bc} \begin{bmatrix} da - bc & db - bd \\ -ca + ac & -cb + ad \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2. \end{aligned}$$

On the other hand,

$$\begin{aligned} AA^{-1} &= \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \frac{1}{ad - bc} \\ &= \begin{bmatrix} ad - bc & -ab + ba \\ cd - dc & -cb + da \end{bmatrix} \frac{1}{ad - bc} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2. \end{aligned}$$

By uniqueness (Theorem 3.2.5), equation (3.2) is the only inverse when $ad - bc \neq 0$.

Now eliminate the case when $ad - bc = 0$. If an inverse exists, say X , then it must satisfy $AX = I_2$. The top-left entry of this matrix equality requires $ax_{11} + bx_{21} = 1$, whereas the bottom-left equality requires $cx_{11} + dx_{21} = 0$. Regard these as a system of two linear equations for the as yet unknowns x_{11} and x_{21} : from $d \times$ the first subtract $b \times$ the second to deduce that an inverse requires $dax_{11} + dbx_{21} - bcx_{11} - bdx_{21} = d \cdot 1 - b \cdot 0$. By cancellation and factorising x_{11} , this equation then requires $(ad - bc)x_{11} = d$. But the determinant is zero, so this equation requires $0 \cdot x_{11} = d$.

- If element d is non-zero, then this equation cannot be satisfied and hence no inverse can be found.
- Otherwise, if d is zero, then any x_{11} satisfies this equation so if there is an inverse then there would be an infinite number of inverses (through the free variable x_{11}). But this contradicts the uniqueness Theorem 3.2.5 so an inverse cannot exist in this case either.

That is, if the determinant $ad - bc = 0$, then the 2×2 matrix is not invertible. \square

Almost anything you can do with A^{-1} can be done without it.

G. E. Forsythe and C. B. Moler, 1967 (Higham 1996, p.261)

Computer considerations Except for easy cases such as 2×2 matrices, we rarely explicitly compute the inverse of a matrix. Computationally there are (almost) always better ways such as the Matlab/Octave operation $A \setminus b$ of Procedure 2.2.4. The inverse is a crucial theoretical device, rarely a practical tool.

The following theorem is an example: for a system of linear equations the theorem connects the existence of a unique solution to the invertibility of the matrix of coefficients. Further, section 3.3.2 connects solutions to the `rcond` invoked by Procedure 2.2.4. Although in theoretical statements we write expressions like $\mathbf{x} = A^{-1}\mathbf{b}$, practically, once we know a solution exists (`rcond` is acceptable), we generally compute the solution without ever constructing A^{-1} .

Theorem 3.2.8. *If A is an invertible $n \times n$ matrix, then the system of linear equations $A\mathbf{x} = \mathbf{b}$ has the unique solution $\mathbf{x} = A^{-1}\mathbf{b}$ for any \mathbf{b} in \mathbb{R}^n .*

Consequently, if a system of linear equations has no solution or an infinite number of solutions (Theorem 2.2.22), then this theorem establishes that the matrix of the system is not invertible.

Proof. The proof has two parts: first showing $\mathbf{x} = A^{-1}\mathbf{b}$ is a solution, and second showing that there are no others. First, try $\mathbf{x} = A^{-1}\mathbf{b}$ and use associativity (Theorem 3.1.18c) and the inverse Definition 3.2.2:

$$A\mathbf{x} = A(A^{-1}\mathbf{b}) = (AA^{-1})\mathbf{b} = I_n\mathbf{b} = \mathbf{b}.$$

Second, suppose \mathbf{y} is any solution, that is, $A\mathbf{y} = \mathbf{b}$. Multiply both sides by the inverse A^{-1} , and again use associativity and the definition of the inverse, to deduce

$$\begin{aligned} A^{-1}(A\mathbf{y}) &= A^{-1}\mathbf{b} \implies (A^{-1}A)\mathbf{y} = \mathbf{x} \\ &\implies I_n\mathbf{y} = \mathbf{x} \\ &\implies \mathbf{y} = \mathbf{x}. \end{aligned}$$

That is, $\mathbf{x} = A^{-1}\mathbf{b}$ is the unique solution. □

Example 3.2.9. Use the matrices of Examples 3.2.1, 3.2.3 and 3.2.4 to decide whether the following systems have a unique solution, or not.

$$(a) \begin{cases} x - y = 4, \\ 4x - 3y = 3. \end{cases} \quad (b) \begin{cases} u - v + 5w = 2, \\ -5u - v + 3w = 5, \\ 2u + 2v - 6w = 1. \end{cases}$$

$$(c) \begin{cases} r - 2s = -1, \\ -3r + 6s = 3. \end{cases}$$

Solution: (a) A matrix for this system is $\begin{bmatrix} 1 & -1 \\ 4 & -3 \end{bmatrix}$ which Example 3.2.1 shows has an inverse. Theorem 3.2.8 then assures us the system has a unique solution.

(b) A matrix for this system is $\begin{bmatrix} 1 & -1 & 5 \\ -5 & -1 & 3 \\ 2 & 2 & -6 \end{bmatrix}$ which Example 3.2.3 shows has an inverse. Theorem 3.2.8 then assures us the system has a unique solution.

(c) A matrix for this system is $\begin{bmatrix} 1 & -2 \\ -3 & 6 \end{bmatrix}$ which Example 3.2.4 shows is not invertible. Theorem 3.2.8 then assures us the system does not have a unique solution. By Theorem 2.2.22 there may be either no solution or an infinite number of solutions—the matrix alone does not tell us which.

■

Example 3.2.10. Given the following information about solutions of systems of linear equations, write down if the matrix associated with each system is invertible, or not, or there is not enough given information to decide. Give reasons.

- | | |
|--|---|
| (a) The general solution is
$(1, -5, 0, 3)$. | (b) The general solution is
$(3, -5 + 3t, 3 - t, -1)$. |
| (c) A solution of a system is
$(-3/2, -2, -\pi, 2, -4)$. | (d) A solution of a homogeneous system is
$(1, 2, -8)$. |

Solution: (a) Since the solution is unique, the matrix in the system must be invertible.
 (b) This solution has an apparent free parameter, t , and so there are many solutions which implies the matrix is not invertible.
 (c) Not enough information is given as we do not know whether there are any more solutions.
 (d) Since a homogeneous system always has $\mathbf{0}$ as a solution (section 2.2.3), then we know that there are at least two solutions to the system, and hence the matrix is not invertible.

■

Recall from Section 3.1 the properties of scalar multiplication, matrix powers, transpose, and their computation Table 3.1. The next theorem incorporates the inverse into this suite of properties.

Theorem 3.2.11 (properties of the inverse). *Let A and B be invertible matrices of the same size, then:*

- (a) matrix A^{-1} is invertible and $(A^{-1})^{-1} = A$;
- (b) if scalar $c \neq 0$, then cA is invertible and $(cA)^{-1} = \frac{1}{c}A^{-1}$;
- (c) matrix AB is invertible and $(AB)^{-1} = B^{-1}A^{-1}$;
- (d) matrix A^T is invertible and $(A^T)^{-1} = (A^{-1})^T$;
- (e) matrices A^p are invertible for all $p = 1, 2, 3, \dots$ and $(A^p)^{-1} = (A^{-1})^p$.

Proof. Three parts are proved, and two are left as exercises.

3.2.11a : By Definition 3.2.2 the matrix A^{-1} satisfies $A^{-1}A = AA^{-1} = I$. But also by Definition 3.2.2 this is exactly the identities we need to assert that matrix A is the inverse of matrix A^{-1} . Hence $A = (A^{-1})^{-1}$.

3.2.11c : Test that $B^{-1}A^{-1}$ has the required properties for the inverse of AB . First, by associativity (Theorem 3.1.18c) and multiplication by the identity (Theorem 3.1.18e)

$$(B^{-1}A^{-1})(AB) = B^{-1}(A^{-1}A)B = B^{-1}IB = B^{-1}B = I.$$

Second, and similarly

$$(AB)(B^{-1}A^{-1}) = A(BB^{-1})A^{-1} = AIA^{-1} = AA^{-1} = I.$$

Hence by Definition 3.2.2 and the uniqueness Theorem 3.2.5, matrix AB is invertible and $B^{-1}A^{-1}$ is the inverse.

3.2.11e : Prove by induction and use 3.2.11c.

- For the case of exponent $p = 1$, $(A^1)^{-1} = (A)^{-1} = A^{-1} = (A^{-1})^1$ and so the identity holds.
- For any integer exponent $p \geq 2$, assume the identity $(A^{p-1})^{-1} = (A^{-1})^{p-1}$. Consider

$$\begin{aligned} (A^p)^{-1} &= (AA^{p-1})^{-1} \quad (\text{by power law Thm. 3.1.18g}) \\ &= (A^{p-1})^{-1}A^{-1} \quad (\text{by Thm. 3.2.11c, } B = A^{p-1}) \\ &= (A^{-1})^{p-1}A^{-1} \quad (\text{by inductive assumption}) \\ &= (A^{-1})^p \quad (\text{by power law Thm. 3.1.18g}). \end{aligned}$$

- By induction, the identity $(A^p)^{-1} = (A^{-1})^p$ holds for exponents $p = 1, 2, 3, \dots$

□

Definition 3.2.12 (non-positive powers). *If A is an invertible matrix, then define $A^0 = I$ and for any positive integer p define $A^{-p} = (A^{-1})^p$ (or by Theorem 3.2.11e equivalently as $(A^p)^{-1}$).*

Example 3.2.13. Recall from Example 3.2.1 that matrix

$$A = \begin{bmatrix} 1 & -1 \\ 4 & -3 \end{bmatrix} \text{ has inverse } A^{-1} = \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix}.$$

Compute A^{-2} and A^{-4} .

Solution: From Definition 3.2.12,

$$\begin{aligned} A^{-2} &= (A^{-1})^2 = \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} \begin{bmatrix} -3 & 1 \\ -4 & 1 \end{bmatrix} = \begin{bmatrix} 5 & -2 \\ 8 & -3 \end{bmatrix}, \\ A^{-4} &= (A^{-1})^4 = [(A^{-1})^2]^2 \\ &= \begin{bmatrix} 5 & -2 \\ 8 & -3 \end{bmatrix} \begin{bmatrix} 5 & -2 \\ 8 & -3 \end{bmatrix} = \begin{bmatrix} 9 & -4 \\ 16 & -7 \end{bmatrix}, \end{aligned}$$

upon using one of the power laws of Theorem 3.1.18g. ■

Example 3.2.14 (predict the past). Recall Example 3.1.5 introduced how to use a Leslie matrix to predict the future population of an animal. If $\mathbf{x} = (60, 70, 20)$ is the current number of pups, juveniles, and mature females respectively, then, for the Leslie matrix

$$L = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix}, \quad \text{which has inverse } L^{-1} = \begin{bmatrix} 0 & 2 & 0 \\ -\frac{1}{4} & 0 & 3 \\ \frac{1}{4} & 0 & 0 \end{bmatrix},$$

by the modelling the predicted population numbers after a year is $\mathbf{x}' = L\mathbf{x}$, after two years is $\mathbf{x}'' = L\mathbf{x}' = L^2\mathbf{x}$, and so on. Assume the same rule applies for earlier years.

- Letting the population numbers a year ago be denoted by \mathbf{x}^- then by the modelling the current population $\mathbf{x} = L\mathbf{x}^-$. Multiply by the inverse: $L^{-1}\mathbf{x} = L^{-1}L\mathbf{x}^- = \mathbf{x}^-$; that is, the population a year before the current is $\mathbf{x}^- = L^{-1}\mathbf{x}$.
- Similarly, letting the population numbers two years ago be denoted by $\mathbf{x}^=$ then by the modelling $\mathbf{x}^- = L\mathbf{x}^=$ and multiplication by L^{-1} gives $\mathbf{x}^= = L^{-1}\mathbf{x}^- = L^{-1}L^{-1}\mathbf{x} = L^{-2}\mathbf{x}$.
- One more year earlier, letting the population numbers two years ago be denoted by \mathbf{x}^{\equiv} then by the modelling $\mathbf{x}^= = L\mathbf{x}^{\equiv}$ and multiplication by L^{-1} gives $\mathbf{x}^{\equiv} = L^{-1}\mathbf{x}^= = L^{-1}L^{-2}\mathbf{x} = L^{-3}\mathbf{x}$.

Hence use the inverse powers of L to predict the earlier history of the population of female animals in the given example: but first verify the given inverse is correct.

Solution: Verify the given inverse by evaluating (showing only non-zero terms in a sum)

$$\begin{aligned}
 LL^{-1} &= \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 0 & 2 & 0 \\ -\frac{1}{4} & 0 & 3 \\ \frac{1}{4} & 0 & 0 \end{bmatrix} \\
 &= \begin{bmatrix} 4 \cdot \frac{1}{4} & 0 & 0 \\ 0 & \frac{1}{2} \cdot 2 & 0 \\ \frac{1}{3} \cdot (-\frac{1}{4}) + \frac{1}{3} \cdot \frac{1}{4} & 0 & \frac{1}{3} \cdot 3 \end{bmatrix} \\
 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = I_3, \\
 L^{-1}L &= \begin{bmatrix} 0 & 2 & 0 \\ -\frac{1}{4} & 0 & 3 \\ \frac{1}{4} & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix} \\
 &= \begin{bmatrix} 2\frac{1}{2} & 0 & 0 \\ 0 & 3 \cdot \frac{1}{3} & -\frac{1}{4} \cdot 4 + 3 \cdot \frac{1}{3} \\ 0 & 0 & \frac{1}{4} \cdot 4 \end{bmatrix} \\
 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = I_3.
 \end{aligned}$$

Hence the given L^{-1} is indeed the inverse. For the current population $\mathbf{x} = (60, 70, 20)$, now use the inverse to compute earlier populations.

- The population of females one year ago was

$$\mathbf{x}^- = L^{-1}\mathbf{x} = \begin{bmatrix} 0 & 2 & 0 \\ -\frac{1}{4} & 0 & 3 \\ \frac{1}{4} & 0 & 0 \end{bmatrix} \begin{bmatrix} 60 \\ 70 \\ 20 \end{bmatrix} = \begin{bmatrix} 140 \\ 45 \\ 15 \end{bmatrix}.$$

That is, there were 140 pups, 45 juveniles, and 15 mature females.

- Computing the square of the inverse

$$L^{-2} = (L^{-1})^2 = \begin{bmatrix} 0 & 2 & 0 \\ -\frac{1}{4} & 0 & 3 \\ \frac{1}{4} & 0 & 0 \end{bmatrix}^2 = \begin{bmatrix} -\frac{1}{2} & 0 & 6 \\ \frac{3}{4} & -\frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 0 \end{bmatrix},$$

we predict the population of females two years ago was

$$\mathbf{x}^{\pm} = L^{-2}\mathbf{x} = \begin{bmatrix} -\frac{1}{2} & 0 & 6 \\ \frac{3}{4} & -\frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} 60 \\ 70 \\ 20 \end{bmatrix} = \begin{bmatrix} 90 \\ 10 \\ 35 \end{bmatrix}.$$

- Similarly, computing the cube of the inverse

$$L^{-3} = L^{-2}L^{-1} = \dots = \begin{bmatrix} \frac{3}{2} & -1 & 0 \\ \frac{1}{8} & \frac{3}{2} & -\frac{3}{2} \\ -\frac{1}{8} & 0 & \frac{3}{2} \end{bmatrix},$$

we predict the population of females three years ago was

$$\mathbf{x}^{\equiv} = L^{-3}\mathbf{x} = \begin{bmatrix} \frac{3}{2} & -1 & 0 \\ \frac{1}{8} & \frac{3}{2} & -\frac{3}{2} \\ -\frac{1}{8} & 0 & \frac{3}{2} \end{bmatrix} \begin{bmatrix} 60 \\ 70 \\ 20 \end{bmatrix} = \begin{bmatrix} 20 \\ 82.5 \\ 22.5 \end{bmatrix}.$$

(Predicting half animals in this last calculation is because the modelling only deals with average numbers, not exact numbers.)

■

Example 3.2.15. As an alternative to the hand calculations of Example 3.2.14, predict earlier populations by computing in Matlab/Octave without ever explicitly finding the inverse or powers of the inverse. The procedure is to solve the linear system $L\mathbf{x}^- = \mathbf{x}$ for the population \mathbf{x}^- a year ago, and then similarly solve $L\mathbf{x}^{\equiv} = \mathbf{x}^-$, $L\mathbf{x}^{\equiv} = \mathbf{x}^{\equiv}$, and so on.

Solution: Execute



```
L=[0 0 4;1/2 0 0;0 1/3 1/3]
x=[60;70;20]
rcond(L)
xm=L\x
xmm=L\xm
xmmp=L\xmm
```

Since `rcond` is 0.08 (good), this code uses `L\` to solve the linear systems and confirm the population of females in previous years is as determined by Example 3.2.14, namely

$$\mathbf{x}_m = \begin{bmatrix} 140 \\ 45 \\ 15 \end{bmatrix}, \quad \mathbf{x}_{mm} = \begin{bmatrix} 90 \\ 10 \\ 35 \end{bmatrix}, \quad \mathbf{x}_{mmp} = \begin{bmatrix} 20 \\ 82.5 \\ 22.5 \end{bmatrix}.$$

■

3.2.2 Diagonal matrices stretch and shrink

Recall that the identity matrices are zero except for a diagonal of ones from the top-left to the bottom-right of the matrix. Because of the nature of matrix multiplication it is this diagonal that is special. Because of the special nature of this diagonal, this section explores matrices which are zero except for the numbers (not generally ones) in the top-left to bottom-right diagonal.

Example 3.2.16. That is, this section explores the nature of so-called diagonal matrices such as

$$\begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix}, \quad \begin{bmatrix} 0.58 & 0 & 0 \\ 0 & -1.61 & 0 \\ 0 & 0 & 2.17 \end{bmatrix}.$$

We use the term diagonal matrix to also include non-square matrices such as

$$\begin{bmatrix} -\sqrt{2} & 0 \\ 0 & \frac{1}{2} \\ 0 & 0 \end{bmatrix}, \quad \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & \pi & 0 & 0 & 0 \\ 0 & 0 & e & 0 & 0 \end{bmatrix}.$$

The term diagonal matrix does *not* apply to say

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 2 & 0 \\ -\frac{1}{2} & 0 & 0 \end{bmatrix}, \quad \text{and not } \begin{bmatrix} -0.17 & 0 & 0 & 0 \\ 0 & -4.22 & 0 & 0 \\ 0 & 0 & 0 & 3.05 \end{bmatrix}.$$

■

Amazingly, the singular value decomposition of Section 3.3 proves that diagonal matrices lie at the very heart of the action of *all* matrices.

Definition 3.2.17 (diagonal matrix). *Given an $m \times n$ matrix A , the **diagonal entries** of A are $a_{11}, a_{22}, \dots, a_{pp}$ where $p = \min(m, n)$. A matrix whose non-diagonal entries are all zero is called a **diagonal matrix**. For brevity we may write $\text{diag}(v_1, v_2, \dots, v_n)$ to denote the $n \times n$ square matrix with diagonal entries v_1, v_2, \dots, v_n , or $\text{diag}_{m \times n}(v_1, v_2, \dots, v_p)$ for an $m \times n$ matrix with diagonal entries v_1, v_2, \dots, v_p .*

Example 3.2.18. The four diagonal matrices of Example 3.2.16 could equivalently be written as $\text{diag}(3, 2)$, $\text{diag}(2, \frac{2}{3}, -1)$, $\text{diag}_{3 \times 2}(-\sqrt{2}, \frac{1}{2})$ and $\text{diag}_{3 \times 5}(1, \pi, e)$, respectively. ■

Solve systems whose matrix is diagonal

Solving a system of linear equations (Definition 2.1.2) is particularly straightforward when the matrix of the system is diagonal. Indeed much mathematics in both theory and applications is devoted to transforming a given problem so that the matrix appearing in the system is diagonal (e.g., sections 2.2.2 and 3.3.2, and Chapters 4 and 7).

Example 3.2.19. Solve

$$\begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 2 \\ -5 \end{bmatrix}$$

Table 3.2: As well as the basics of Matlab/Octave listed in Tables 1.2, 2.3 and 3.1, we need these matrix operations.

- `diag(v)` where v is a row/column vector of length p generates the $p \times p$ matrix

$$\text{diag}(v_1, v_2, \dots, v_p) = \begin{bmatrix} v_1 & 0 & \cdots & 0 \\ 0 & v_2 & & \vdots \\ \vdots & & \ddots & \\ 0 & \cdots & & v_p \end{bmatrix}.$$

- In Matlab/Octave (but not usually in algebra), `diag` also does the opposite: for an $m \times n$ matrix A such that both $m, n \geq 2$, `diag(A)` returns the (column) vector $(a_{11}, a_{22}, \dots, a_{pp})$ of diagonal entries where the result vector length $p = \min(m, n)$.
- The dot operators `./` and `.*` do element-by-element division and multiplication of two matrices/vectors of the same size. For example,
 $[5 \ 14 \ 33] ./ [5 \ 7 \ 3] = [1 \ 2 \ 11]$
- Section 3.5 also needs to compute the logarithm of data: `log10(v)` finds the logarithm to base 10 of each component of v and returns the results in a vector of the same size; `log(v)` does the same but for the natural logarithm (not `ln(v)`).

Solution: Algebraically this matrix-vector equation means

$$\begin{array}{l} 3x_1 + 0x_2 = 2 \\ 0x_1 + 2x_2 = -5 \end{array} \iff \begin{array}{l} 3x_1 = 2 \\ 2x_2 = -5 \end{array} \iff \begin{array}{l} x_1 = 2/3 \\ x_2 = -5/2 \end{array}.$$

The solution is $\mathbf{x} = (2/3, -5/2)$. ■

Example 3.2.20. Solve

$$\begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{2}{3} & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

Solution: Algebraically this equation means

$$\begin{array}{l} 2x_1 + 0x_2 + 0x_3 = b_1 \\ 0x_1 + \frac{2}{3}x_2 + 0x_3 = b_2 \\ 0x_1 + 0x_2 - 1x_3 = b_3 \end{array} \iff \begin{array}{l} 2x_1 = b_1 \\ \frac{2}{3}x_2 = b_2 \\ -x_3 = b_3 \end{array} \iff \begin{array}{l} x_1 = \frac{1}{2}b_1 \\ x_2 = \frac{3}{2}b_2 \\ x_3 = -b_3 \end{array}.$$

The solution is

$$\mathbf{x} = \begin{bmatrix} \frac{1}{2}b_1 \\ \frac{3}{2}b_2 \\ -b_3 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & \frac{3}{2} & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}.$$

Consequently, by its uniqueness (Theorem 3.2.5), the inverse of the given diagonal matrix must be

$$\begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{2}{3} & 0 \\ 0 & 0 & -1 \end{bmatrix}^{-1} = \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & \frac{3}{2} & 0 \\ 0 & 0 & -1 \end{bmatrix},$$

which interestingly is the diagonal of reciprocals of the given matrix. ■

Theorem 3.2.21 (inverse of diagonal matrix). *Let D be the $n \times n$ diagonal matrix $D = \text{diag}(d_1, d_2, \dots, d_n)$. If all the diagonal entries are nonzero, $d_i \neq 0$ for $i = 1, 2, \dots, n$, then D is invertible and the inverse $D^{-1} = \text{diag}(1/d_1, 1/d_2, \dots, 1/d_n)$.*

Proof. Consider the matrix product

$$\begin{aligned} & \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} \begin{bmatrix} \frac{1}{d_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{d_2} & & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{d_n} \end{bmatrix} \\ &= \begin{bmatrix} d_1 \frac{1}{d_1} + 0 + \cdots + 0 & d_1 0 + 0 \frac{1}{d_1} + \cdots + 0 & \cdots & d_1 0 + 0 + \cdots + 0 \frac{1}{d_n} \\ 0 \frac{1}{d_1} + d_2 0 + \cdots + 0 & 0 + d_2 \frac{1}{d_2} + \cdots + 0 & \cdots & 0 + d_2 0 + \cdots + 0 \frac{1}{d_n} \\ \vdots & \ddots & \ddots & \vdots \\ 0 \frac{1}{d_1} + 0 + \cdots + d_n 0 & 0 + 0 \frac{1}{d_2} + \cdots + d_n 0 & \cdots & 0 + 0 + \cdots + d_n \frac{1}{d_n} \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = I_n. \end{aligned}$$

Similarly for the reverse product. By Definition 3.2.2, D is invertible with the given inverse. □

Example 3.2.22. The previous example gave the inverse of a 3×3 . For the 2×2 matrix $D = \text{diag}(3, 2) = \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix}$ the inverse is $D^{-1} = \text{diag}(\frac{1}{3}, \frac{1}{2}) = \begin{bmatrix} \frac{1}{3} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$. Then the solution to

$$\begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 2 \\ -5 \end{bmatrix} \quad \text{is } \mathbf{x} = \begin{bmatrix} \frac{1}{3} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 \\ -5 \end{bmatrix} = \begin{bmatrix} 2/3 \\ -5/2 \end{bmatrix}.$$
■

Compute in Matlab/Octave. To solve $D\mathbf{x} = \mathbf{b}$ recognise this matrix-vector equation means

$$\begin{aligned} \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} &= \begin{bmatrix} d_1 x_1 \\ d_2 x_2 \\ \vdots \\ d_n x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} \\ \iff d_1 x_1 &= b_1 \quad x_1 = b_1/d_1 \\ \iff d_2 x_2 &= b_2 \quad x_2 = b_2/d_2 \\ \vdots & \quad \vdots \\ d_n x_n &= b_n \quad x_n = b_n/d_n \end{aligned} \tag{3.3}$$

- Suppose you have a column vector \mathbf{d} of the diagonal entries of D and a column vector \mathbf{b} of the RHS; then compute a solution by, for example,

```
d=[2;2/3;-1]
b=[1;2;3]
x=b./d
```

where $.$ / does element-by-element division in Matlab/Octave (Table 3.2) to here find the answer $[0.5; 3; -3]$.

- When you have the diagonal matrix in full: extract the diagonal elements into a column vector with `diag()` (Table 3.2); then execute the element-by-element division; for example,

```
D=[2 0 0;0 2/3 0;0 0 -1]
b=[1;2;3]
x=b./diag(D)
```

But do not divide by zero

Dividing by zero is almost always nonsense. Instead use reasoning. Consider solving $D\mathbf{x} = \mathbf{b}$ for diagonal $D = \text{diag}(d_1, d_2, \dots, d_n)$ where $d_n = 0$ (and similarly for any others that are zero). From (3.3) we need to solve $d_n x_n = 0 \cdot x_n = b_n$, that is, $0 = b_n$. There are two cases:

- if $b_n \neq 0$, then there is no solution; conversely
- if $b_n = 0$, then there is an infinite number of solutions as any x_n satisfies $0 \cdot x_n = 0$.

Example 3.2.23. Solve the two systems (the only difference is the last component on the RHS)

$$\begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{2}{3} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{2}{3} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$$

Solution: Algebraically, the first system means

$$\begin{array}{l} 2x_1 + 0x_2 + 0x_3 = 1 \\ 0x_1 + \frac{2}{3}x_2 + 0x_3 = 2 \\ 0x_1 + 0x_2 + 0x_3 = 3 \end{array} \iff \begin{array}{l} 2x_1 = 1 \\ \frac{2}{3}x_2 = 2 \\ 0x_3 = 3 \end{array}$$

There is no solution in the first case as there is no choice of x_3 such that $0x_3 = 3$.

Algebraically, the second system means

$$\begin{array}{l} 2x_1 + 0x_2 + 0x_3 = 1 \\ 0x_1 + \frac{2}{3}x_2 + 0x_3 = 2 \\ 0x_1 + 0x_2 + 0x_3 = 0 \end{array} \iff \begin{array}{l} 2x_1 = 1 \\ \frac{2}{3}x_2 = 2 \\ 0x_3 = 0 \end{array}$$

In this second case we satisfy the equation $0x_3 = 0$ with any x_3 . Hence there are an infinite number of solutions, namely $\mathbf{x} = (\frac{1}{2}, 3, t)$ for all t —a free variable just as in Gauss–Jordan elimination (Procedure 2.2.19). ■

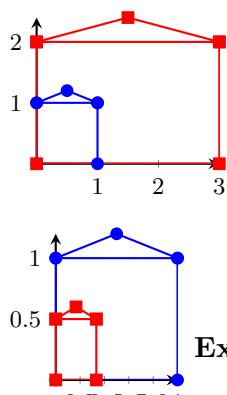
Stretch or squash the unit square

Equations are just the boring part of mathematics. I attempt to see things in terms of geometry.

Stephen Hawking, 2005

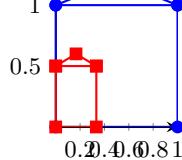
This geometric interpretation prepares for the geometry of rotation by orthogonal matrices.

Multiplication by matrices transforms shapes: multiplication by diagonal matrices just stretches or squashes and/or reflects in the direction of the coordinate axes.



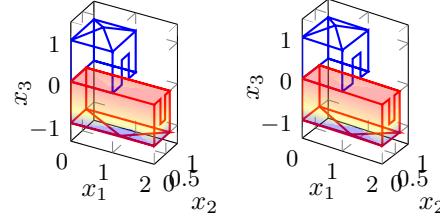
Example 3.2.24. Consider $A = \text{diag}(3, 2) = \begin{bmatrix} 3 & 0 \\ 0 & 2 \end{bmatrix}$. The marginal pictures shows this matrix stretches the (blue) unit square (with ‘roof’) by a factor of three horizontally and two vertically (to the red). Recall that (x_1, x_2) denotes the corresponding column vector. As seen in the corner points of the graphic in the margin, $A \times (1, 0) = (3, 0)$, $A \times (0, 1) = (0, 2)$, $A \times (0, 0) = (0, 0)$, and $A \times (1, 1) = (3, 2)$. The ‘roof’ just helps to track which corner goes where.

The inverse $A^{-1} = \text{diag}(\frac{1}{3}, \frac{1}{2}) = \begin{bmatrix} \frac{1}{3} & 0 \\ 0 & \frac{1}{2} \end{bmatrix}$ undoes the stretching of the matrix A by squashing both horizontally and vertically (from blue to red). ■



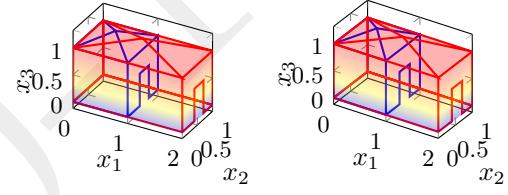
Example 3.2.25. Consider $\text{diag}(2, \frac{2}{3}, -1) = \begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{2}{3} & 0 \\ 0 & 0 & -1 \end{bmatrix}$: the stereo pair below illustrates how this diagonal matrix stretches in one direction,

squashes in another, and reflects in the vertical. By multiplying the matrix by corner vectors $(1, 0, 0)$, $(0, 1, 0)$, $(0, 0, 1)$ and so on, we see that the blue unit cube (with roof and door) maps to the red.

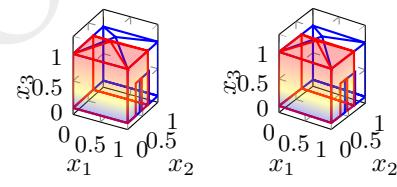


One great aspect of a diagonal matrix is that it is easy to separate its effects into separate effects in each coordinate direction. For example, the above 3×3 matrix is the same as the combined effects of the following three.

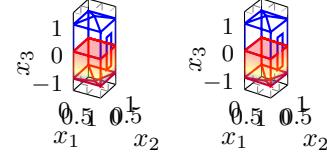
$\begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$. Stretch by a factor of two in the x_1 direction.



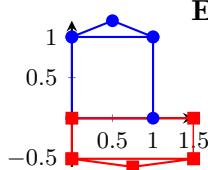
$\begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{2}{3} & 0 \\ 0 & 0 & 1 \end{bmatrix}$. Squash by a factor of $2/3$ in the x_2 direction.



$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$. Reflect in the vertical x_3 direction.



Example 3.2.26. What diagonal matrix transforms the blue unit square to the red in the illustration in the margin?

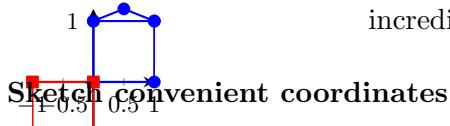


Solution: In the illustration, the horizontal is stretched by a factor of $3/2$, whereas the vertical is squashed by a factor of $1/2$, and reflected (minus sign). Hence the matrix is $\text{diag}(\frac{3}{2}, -\frac{1}{2}) = \begin{bmatrix} \frac{3}{2} & 0 \\ 0 & -\frac{1}{2} \end{bmatrix}$.

■

Some diagonal matrices rotate Now consider the transformation of multiplying by matrix $\begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}$: the two reflections of this

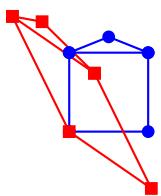
diagonal matrix, the two -1 s, have the same effect as one rotation, here by 180° , as shown to the left. Matrices that rotate are incredibly useful and is the topic of the next section 3.2.3.



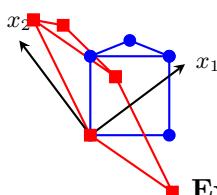
Optional preliminary to diagonalisation.

One of the fundamental principles of applying mathematics in science and engineering is that the real world, nature, does its thing irrespective of our mathematical description. Hence we often simplify our mathematical description of real world applications by choosing a coordinate system to suit its nature. That is, although this book (almost) always draws the x or x_1 axis horizontally, and the y or x_2 axis vertically, in applications it is often better to draw the axes in some other directions which are convenient for the application. This example illustrates the principle.

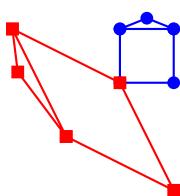
Example 3.2.27. Consider the transformation shown in the margin (it might arise from the deformation of some material and we need to know the internal stretching and shrinking to predict failure): it has no coordinate axes shown because it supposed to be some transformation in nature; now we wish to impose on nature our mathematical description. Draw approximate coordinate axes, with origin at the common point at the lower-left corner, so the transformation becomes that of the diagonal matrix $\text{diag}(\frac{1}{2}, 2) = \begin{bmatrix} 0.5 & 0 \\ 0 & 2 \end{bmatrix}$.



Solution: From the diagonal matrix we first look for a direction in which the transformation squashes by a factor of $1/2$: from the marginal graph, this direction must be towards the top-right. Second, from the diagonal matrix we look for a direction in which the transformation stretches by a factor of two: from the marginal picture this direction must be aligned to the top-left. Because the top-right corner of the square is stretched a little in this second direction, the first direction must be aimed a little lower than this corner. Hence, coordinate axes that make the transformation the given diagonal matrix are as shown in the margin. ■

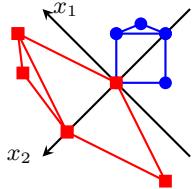


Example 3.2.28. Consider the transformation shown in the margin: it has no coordinate axes shown because it supposed to be some transformation in nature; now impose on nature our mathematical description. Draw approximate coordinate axes, with origin at the common corner point, so the transformation becomes that of the diagonal matrix $\text{diag}(3, -1) = \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix}$.



Solution: From the diagonal matrix we first look for a direction in which the transformation stretches by a factor of three: from the marginal graph, this direction must be aligned along the diagonal

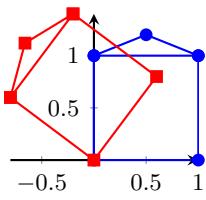
top-left to bottom-right. Second, from the diagonal matrix we look for a direction in which the transformation reflects: from the marginal picture this direction must be aligned along the top-right to bottom-left. Hence, coordinate axes that make the transformation the given diagonal matrix are as shown in the margin.



Finding such coordinate systems in which a given real world transformation is diagonal is important in science, engineering, and computer science. Systematic methods for such diagonalisation are developed in Section 3.3, and Chapters 4 and 7. These rely on understanding the algebra and geometry of rotations, which is our next topic. ■

3.2.3 Orthogonal matrices rotate

Whereas diagonal matrices stretch and squash, the so-called ‘orthogonal matrices’ represent just rotations (and/or reflection). For example, this section shows that multiplying by the ‘orthogonal matrix’ $\begin{bmatrix} 3/5 & -4/5 \\ 4/5 & 3/5 \end{bmatrix}$ rotates by 53.13° as shown in the marginal picture. Orthogonal matrices are the best to compute with, such as to solve linear equations, since they all have `rcond = 1`. To see these and related marvellous properties, we must invoke the geometry of lengths and angles.



Recall the dot product determines lengths and angles Section 1.3 introduced the dot product between two vectors (Definition 1.3.2). For any two vectors in \mathbb{R}^n , $\mathbf{u} = (u_1, \dots, u_n)$ and $\mathbf{v} = (v_1, \dots, v_n)$, define the **dot product**

$$\begin{aligned}\mathbf{u} \cdot \mathbf{v} &= (u_1, \dots, u_n) \cdot (v_1, \dots, v_n) \\ &= u_1 v_1 + u_2 v_2 + \cdots + u_n v_n.\end{aligned}$$

Considering the two vectors as column matrices, the dot product is the same as the matrix product (Example 3.1.13)

$$\mathbf{u} \cdot \mathbf{v} = \mathbf{u}^T \mathbf{v} = \mathbf{v}^T \mathbf{u} = \mathbf{v} \cdot \mathbf{u}.$$

Also (Theorem 1.3.13a), the length of a vector $\mathbf{v} = (v_1, v_2, \dots, v_n)$ in \mathbb{R}^n is the real number

$$|\mathbf{v}| = \sqrt{\mathbf{v} \cdot \mathbf{v}} = \sqrt{v_1^2 + v_2^2 + \cdots + v_n^2},$$

and that unit vectors are those of length one. For two non-zero vectors \mathbf{u}, \mathbf{v} in \mathbb{R}^n , Theorem 1.3.4 defines the angle θ between the vectors via

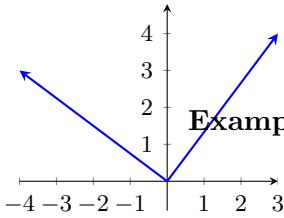
$$\cos \theta = \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}| |\mathbf{v}|}, \quad 0 \leq \theta \leq \pi.$$

If the two vectors are at right-angles, then the dot product is zero and the two vectors are termed orthogonal (Definition 1.3.15).

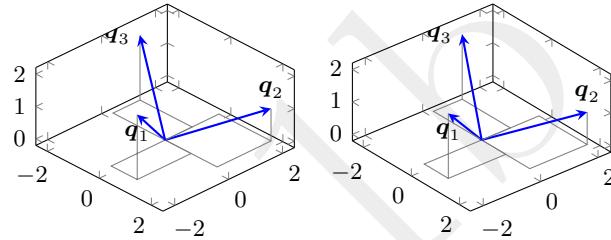
Orthogonal set of vectors

We need sets of orthogonal vectors (vectors which are all at right-angles to each other). One example is the set of standard unit vectors $\{e_1, e_2, \dots, e_n\}$ aligned with the coordinate axes in \mathbb{R}^n .

Example 3.2.29. The set of two vectors $\{(3, 4), (-4, 3)\}$ shown in the margin is an orthogonal set as the two vectors have dot product $= 3 \cdot (-4) + 4 \cdot 3 = -12 + 12 = 0$. ■



Example 3.2.30. Let vectors $q_1 = (1, -2, 2)$, $q_2 = (2, 2, 1)$ and $q_3 = (-2, 1, 2)$, illustrated in stereo below. Is $\{q_1, q_2, q_3\}$ an orthogonal set?



Solution: Yes, because all the pairwise dot products are zero:
 $q_1 \cdot q_2 = 2 - 4 + 2 = 0$; $q_1 \cdot q_3 = -2 - 2 + 4 = 0$; $q_2 \cdot q_3 = -4 + 2 + 2 = 0$. ■

Definition 3.2.31. A set of non-zero vectors $\{q_1, q_2, \dots, q_k\}$ in \mathbb{R}^n is called an **orthogonal set** if all pairs of distinct vectors in the set are orthogonal: that is, $q_i \cdot q_j = 0$ whenever $i \neq j$ for $i, j = 1, 2, \dots, k$. A set of vectors in \mathbb{R}^n is called an **orthonormal set** if it is an orthogonal set of unit vectors.⁹

Example 3.2.32. Any set, or subset, of standard unit vectors in \mathbb{R}^n (Definition 1.2.5) are an orthonormal set as they are all at right-angles (orthogonal), and all of length one. ■

Example 3.2.33. Let vectors $q_1 = (\frac{1}{3}, -\frac{2}{3}, \frac{2}{3})$, $q_2 = (\frac{2}{3}, \frac{2}{3}, \frac{1}{3})$, $q_3 = (-\frac{2}{3}, \frac{1}{3}, \frac{2}{3})$. Show the set $\{q_1, q_2, q_3\}$ is an orthonormal set.

Solution: These vectors are all 1/3 of the vectors in Example 3.2.30 and so are orthogonal. They all have length one: $|q_1|^2 = \frac{1}{9} + \frac{4}{9} + \frac{4}{9} = 1$; $|q_2|^2 = \frac{4}{9} + \frac{4}{9} + \frac{1}{9} = 1$; $|q_3|^2 = \frac{4}{9} + \frac{1}{9} + \frac{4}{9} = 1$. Hence $\{q_1, q_2, q_3\}$ is an orthonormal set in \mathbb{R}^3 . ■

⁹ A single non-zero vector always forms an orthogonal set. A single unit vector always forms an orthonormal set.

Orthogonal matrices

Example 3.2.34. Example 3.2.29 showed $\{(3, 4), (-4, 3)\}$ is an orthogonal set. Each vector has length five so dividing by its length means $\{(\frac{3}{5}, \frac{4}{5}), (-\frac{4}{5}, \frac{3}{5})\}$ is an orthonormal set. Form the matrix Q with these two vectors as its columns:

$$Q = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix}.$$

Then consider

$$Q^T Q = \begin{bmatrix} \frac{3}{5} & \frac{4}{5} \\ -\frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} = \begin{bmatrix} \frac{9+16}{25} & \frac{-12+12}{25} \\ \frac{-12+12}{25} & \frac{16+9}{25} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Similarly $QQ^T = I_2$. Consequently, Q^T is the inverse of Q (Definition 3.2.2). The transpose being the inverse is no accident.

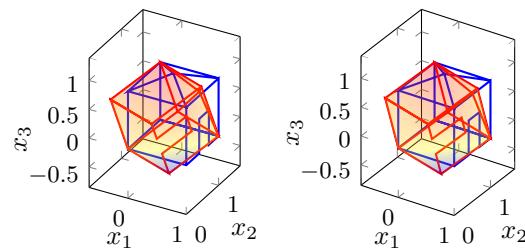
Also no accident is that multiplication by this Q gives the rotation illustrated at the start of this section, (§3.2.3). ■

Definition 3.2.35 (orthogonal matrices). A square $n \times n$ matrix Q is called an **orthogonal matrix** if $Q^T Q = I_n$. Because of its special properties (Theorem 3.2.39), multiplication by an orthogonal matrix is called a **rotation and/or reflection**; ¹⁰ for brevity and depending upon the circumstances it may be called just a **rotation** or just a **reflection**.

Example 3.2.36. In the following equation, check that the matrix is orthogonal, and hence solve the equation $Qx = b$:

$$\begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} x = \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix}.$$

The stereo pair below illustrates the rotation of the unit cube under multiplication by this matrix: every point x in the (blue) unit cube, is mapped to the point Qx to form the (red) result.



¹⁰ Although herein we term multiplication by an orthogonal matrix as a ‘rotation’ it generally is not a rotation about a single axis. Instead, generally the ‘rotation’ characterised by any one orthogonal matrix may be composed of a sequence of ‘rotations about different axes’—each axis with a different orientation.

Solution: In Matlab/Octave (recall the single quote (prime) gives the transpose, Table 3.1),

```
Q=[1,-2,2;2,2,1;-2,1,2]/3
Q'*Q
```

Since the product $Q^T Q$ is I_3 , the matrix is orthogonal. Multiplying by Q^T both sides of the equation $Qx = b$ gives $Q^T Qx = Q^T b$; that is, $I_3 x = Q^T b$, equivalently, $x = Q^T b$. Here,

```
x=Q'*[0;2;-1]
```

gives the solution $x = (2, 1, 0)$.



Example 3.2.37. Given the matrix is orthogonal, solve the linear equation

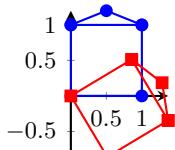
$$\begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix} x = \begin{bmatrix} 1 \\ -1 \\ 1 \\ 3 \end{bmatrix}.$$

Solution: Given the matrix is orthogonal, calling the matrix Q we know $Q^T Q = I_4$, just multiply the equation $Qx = b$ by the transpose Q^T to deduce $Q^T Qx = Q^T b$, that is, $x = Q^T b$. Here this gives the solution to be

$$x = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \\ -2 \\ 0 \end{bmatrix}.$$



Example 3.2.38. The marginal graph shows a rotation of the unit square.



From the graph estimate roughly the matrix Q that performs the rotation. Confirm that your estimated matrix is orthogonal (approximately).

Solution: Consider what happens to the standard unit vectors: $(1, 0) \mapsto (0.5, -0.9)$ roughly (one decimal place); and $(0, 1) \mapsto (0.9, 0.5)$ roughly. To do this the matrix must have these two vectors as its two columns, that is,

$$Q \approx \begin{bmatrix} 0.5 & 0.9 \\ -0.9 & 0.5 \end{bmatrix}.$$

We may check by what happens to the corner point $(1, 1)$: $Q(1, 1) \approx (1.4, -0.4)$ which looks approximately correct. To confirm orthogonality of Q , find

$$Q^T Q = \begin{bmatrix} 0.5 & -0.9 \\ 0.9 & 0.5 \end{bmatrix} \begin{bmatrix} 0.5 & 0.9 \\ -0.9 & 0.5 \end{bmatrix} = \begin{bmatrix} 1.06 & 0 \\ 0 & 1.06 \end{bmatrix} \approx I_2,$$

and so is approximately orthogonal.

■

Because orthogonal matrices represent rotations, they arise frequently in engineering and scientific mechanics of bodies. Also, the ease in solving equations with orthogonal matrices puts orthogonal matrices at the heart of coding and decoding photographs (jpeg), videos (mpeg), signals (Fourier transforms), and so on. Furthermore, an extension of orthogonal matrices to complex valued matrices, the so-called unitary matrices, is at the core of quantum physics and quantum computing. Moreover, the next Section 3.3 establishes that orthogonal matrices express the orientation of the action of *every* matrix and hence are a vital component of solving linear equations in general. But to utilise orthogonal matrices across the wide range of applications we need to establish the following properties.

Theorem 3.2.39. *Let Q be an $n \times n$ matrix, then the following statements are equivalent:*

- (a) Q is an orthogonal matrix;
- (b) the column vectors of Q form an orthonormal set;
- (c) Q is invertible and $Q^{-1} = Q^T$;
- (d) Q^T is an orthogonal matrix;
- (e) the row vectors of Q form an orthonormal set;
- (f) multiplication by Q preserves all lengths and angles (and hence corresponds to our intuition of a rotation and/or reflection).

Proof. Write matrix $Q = [\mathbf{q}_1 \ \mathbf{q}_2 \ \cdots \ \mathbf{q}_n]$ in terms of its n columns $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$.

3.2.39a \iff 3.2.39b : Consider

$$Q^T Q = \begin{bmatrix} \mathbf{q}_1^T \\ \mathbf{q}_2^T \\ \vdots \\ \mathbf{q}_n^T \end{bmatrix} [\mathbf{q}_1 \ \mathbf{q}_2 \ \cdots \ \mathbf{q}_n]$$

$$\begin{aligned}
&= \begin{bmatrix} \mathbf{q}_1^T \mathbf{q}_1 & \mathbf{q}_1^T \mathbf{q}_2 & \cdots & \mathbf{q}_1^T \mathbf{q}_n \\ \mathbf{q}_2^T \mathbf{q}_1 & \mathbf{q}_2^T \mathbf{q}_2 & \cdots & \mathbf{q}_2^T \mathbf{q}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_n^T \mathbf{q}_1 & \mathbf{q}_n^T \mathbf{q}_2 & \cdots & \mathbf{q}_n^T \mathbf{q}_n \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{q}_1 \cdot \mathbf{q}_1 & \mathbf{q}_1 \cdot \mathbf{q}_2 & \cdots & \mathbf{q}_1 \cdot \mathbf{q}_n \\ \mathbf{q}_2 \cdot \mathbf{q}_1 & \mathbf{q}_2 \cdot \mathbf{q}_2 & \cdots & \mathbf{q}_2 \cdot \mathbf{q}_n \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{q}_n \cdot \mathbf{q}_1 & \mathbf{q}_n \cdot \mathbf{q}_2 & \cdots & \mathbf{q}_n \cdot \mathbf{q}_n \end{bmatrix}
\end{aligned}$$

Matrix Q is orthogonal iff this product is the identity (Definition 3.2.35) which is iff $\mathbf{q}_i \cdot \mathbf{q}_j = 0$ for $i \neq j$ and $|\mathbf{q}_i|^2 = \mathbf{q}_i \cdot \mathbf{q}_i = 1$, that is, iff the columns are orthonormal (Definition 3.2.31).

3.2.39b \implies 3.2.39c : First, consider the homogeneous system $Q^T \mathbf{x} = \mathbf{0}$ for \mathbf{x} in \mathbb{R}^n . We establish $\mathbf{x} = \mathbf{0}$ is the only solution. The system $Q^T \mathbf{x} = \mathbf{0}$, written in terms of the orthonormal columns of Q , is

$$\begin{bmatrix} \mathbf{q}_1^T \\ \mathbf{q}_2^T \\ \vdots \\ \mathbf{q}_n^T \end{bmatrix} \mathbf{x} = \begin{bmatrix} \mathbf{q}_1^T \mathbf{x} \\ \mathbf{q}_2^T \mathbf{x} \\ \vdots \\ \mathbf{q}_n^T \mathbf{x} \end{bmatrix} = \begin{bmatrix} \mathbf{q}_1 \cdot \mathbf{x} \\ \mathbf{q}_2 \cdot \mathbf{x} \\ \vdots \\ \mathbf{q}_n \cdot \mathbf{x} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Since the dot products are all zero, either $\mathbf{x} = \mathbf{0}$ or \mathbf{x} is orthogonal (at right angles) to all of the n orthonormal vectors $\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n\}$. In \mathbb{R}^n we cannot have $(n+1)$ non-zero vectors all at right angles to each other (Theorem 1.3.20), consequently $\mathbf{x} = \mathbf{0}$ is the only possibility as the solution of $Q^T \mathbf{x} = \mathbf{0}$.

Second, let $n \times n$ matrix $X = I_n - QQ^T$. Pre-multiply by Q^T : $Q^T X = Q^T(I_n - QQ^T) = Q^T I_n - (Q^T Q)Q^T = Q^T I_n - I_n Q^T = Q^T - Q^T = O_n$. That is, $Q^T X = O_n$ but each column is of the form $Q^T \mathbf{x} = \mathbf{0}$ which we know requires $\mathbf{x} = \mathbf{0}$, hence $X = O_n$. Then $X = I_n - QQ^T = O_n$ which rearranged gives $QQ^T = I_n$. Put $QQ^T = I_n$ together with $Q^T Q = I_n$ (Definition 3.2.35), then by Definition 3.2.2 Q is invertible with inverse Q^T .

3.2.39c \implies 3.2.39a, 3.2.39d : Part 3.2.39c asserts Q is invertible with inverse Q^T , that is, $Q^T Q = QQ^T = I_n$. Since $Q^T Q = I_n$, matrix Q is orthogonal.

Since $I_n = QQ^T = (Q^T)^T Q^T$, by Definition 3.2.35 Q^T is orthogonal.

3.2.39d \iff 3.2.39e : The proof is similar to that for 3.2.39a \iff 3.2.39b, but for the rows of Q and $QQ^T = I_n$.

3.2.39e \implies 3.2.39c : Similar to that for 3.2.39b \implies 3.2.39c, but for the rows of Q , $Q\mathbf{x} = \mathbf{0}$ and $X = I_n - Q^T Q$.

3.2.39a \implies 3.2.39f : We prove that multiplication by orthogonal Q preserves all lengths and angles, as illustrated in the picture in Examples 3.2.36 and 3.2.38, by comparing the properties of

transformed vectors Qu with the properties of the original u . For any u, v in \mathbb{R}^n , consider the dot product $(Qu) \cdot (Qv) = (Qu)^T Qv = u^T Q^T Qv = u^T I_n v = u^T v = u \cdot v$. We use this identity $(Qu) \cdot (Qv) = u \cdot v$ below.

- Firstly, the length $|Qu| = \sqrt{(Qu) \cdot (Qu)} = \sqrt{u \cdot u} = |u|$ is preserved, and correspondingly for v .
- Secondly, let θ be the angle between u and v and θ' be the angle between Qu and Qv (recall $0 \leq \text{angle} \leq \pi$), then

$$\cos \theta' = \frac{(Qu) \cdot (Qv)}{|Qu||Qv|} = \frac{u \cdot v}{|u||v|} = \cos \theta,$$

and so all angles are preserved.

3.2.39f \implies 3.2.39b : Look at the consequences of matrix Q preserving all lengths and angles when applied to the standard unit vectors e_1, e_2, \dots, e_n . Observe $Qe_j = q_j$, the j th column of matrix Q . Then for all j , the length of the j th column $|q_j| = |Qe_j| = |e_j| = 1$ by the preservation of the length of the standard unit vector. Also, for all $i \neq j$ the dot product of columns $q_i \cdot q_j = |q_i||q_j| \cos \theta' = 1 \cdot 1 \cdot \cos \frac{\pi}{2} = 0$ where θ' is the angle between q_i and q_j which is the angle between e_i and e_j by preservation, namely the angle $\frac{\pi}{2}$. That is, the columns of Q form an orthonormal set.

□

Another important property, proved by Exercise 3.2.19, is that the product of orthogonal matrices is also an orthogonal matrix.

Example 3.2.40. Show that these matrices are orthogonal and hence write down their inverses:

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}.$$

Solution: For the first matrix each column is a unit vector, and each column is orthogonal to each other: since the matrix has orthonormal columns, then the matrix is orthogonal (Theorem 3.2.39b). Its inverse is the transpose (Theorem 3.2.39c)

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

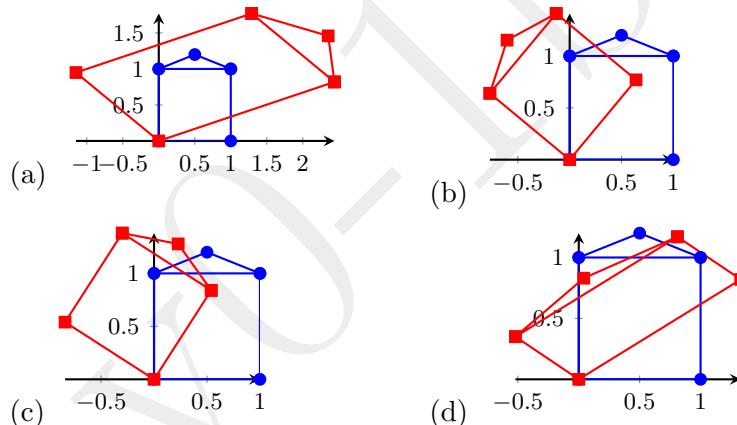
For the second matrix the two columns are unit vectors as $|\cos \theta, \sin \theta|^2 = \cos^2 \theta + \sin^2 \theta = 1$ and $|-\sin \theta, \cos \theta|^2 = \sin^2 \theta + \cos^2 \theta = 1$. The two columns are orthogonal as the dot product $(\cos \theta, \sin \theta) \cdot (-\sin \theta, \cos \theta) = -\cos \theta \sin \theta + \sin \theta \cos \theta = 0$.

$(-\sin \theta, \cos \theta) = -\cos \theta \sin \theta + \sin \theta \cos \theta = 0$. Since the matrix has orthonormal columns, then the matrix is orthogonal (Theorem 3.2.39b). Its inverse is the transpose (Theorem 3.2.39c)

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

■

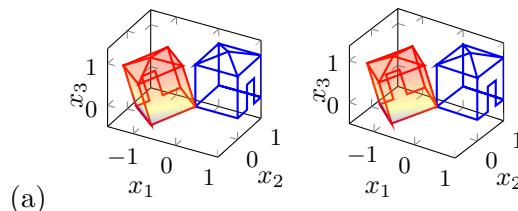
Example 3.2.41. The following graphs illustrate the transformation of the unit square through multiplying by some different matrices: using Theorem 3.2.39f, which transformations appear to be that of multiplying by an orthogonal matrix?

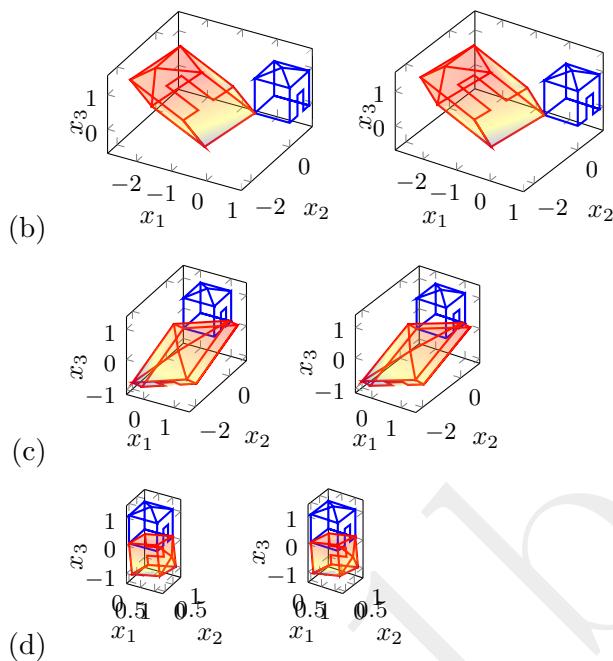


Solution: (a) No—the square is stretched and angles changed.
 (b) Yes—the square is just rotated.
 (c) Yes—the square is rotated and reflected.
 (d) No—the square is squashed and angles changed.

■

Example 3.2.42. The following stereo pairs illustrate the transformation of the unit cube through multiplying by some different matrices: using Theorem 3.2.39f, which transformations appear to be that of multiplying by an orthogonal matrix?





- Solution:*
- (a) Yes—the cube is just rotated.
 - (b) No—the cube is stretched and angles changed.
 - (c) No—the cube is stretched and angles changed.
 - (d) Yes—the cube is just rotated and reflected.

■

3.2.4 Exercises

Exercise 3.2.1. By direct multiplication, both ways, confirm that for each of the following pairs, matrix B is an inverse of matrix A , or not.

$$(a) A = \begin{bmatrix} 0 & -4 \\ 4 & 4 \end{bmatrix}, B = \begin{bmatrix} 1/4 & 1/4 \\ -1/4 & 0 \end{bmatrix}$$

$$(b) A = \begin{bmatrix} -3 & 3 \\ -3 & 1 \end{bmatrix}, B = \begin{bmatrix} 1/6 & -1/2 \\ 1/2 & -1/2 \end{bmatrix}$$

$$(c) A = \begin{bmatrix} 5 & -1 \\ 3 & -2 \end{bmatrix}, B = \begin{bmatrix} 3/7 & -1/7 \\ 3/7 & -5/7 \end{bmatrix}$$

$$(d) A = \begin{bmatrix} -1 & 1 \\ -5 & -1 \end{bmatrix}, B = \begin{bmatrix} -1/6 & -1/6 \\ 5/6 & -1/6 \end{bmatrix}$$

$$(e) A = \begin{bmatrix} -2 & 4 & 2 \\ 1 & 1 & -1 \\ 1 & -1 & 1 \end{bmatrix}, B = \begin{bmatrix} 0 & 1/2 & 1/2 \\ 1/6 & 1/3 & 0 \\ 1/6 & -1/6 & 1/2 \end{bmatrix}$$

$$(f) A = \begin{bmatrix} -3 & 0 & -1 \\ 1 & 4 & 2 \\ 3 & -4 & -1 \end{bmatrix}, B = \begin{bmatrix} 1 & 1 & 1 \\ 7/4 & 3/2 & 5/4 \\ -4 & -3 & -3 \end{bmatrix}$$

(g) $A = \begin{bmatrix} -3 & -3 & 3 \\ 4 & 3 & -3 \\ -2 & -1 & 4 \end{bmatrix}$, $B = \begin{bmatrix} 1 & 1 & 0 \\ -1 & -2/3 & 1/3 \\ 2/9 & 1/3 & 1/3 \end{bmatrix}$

(h) $A = \begin{bmatrix} -1 & 3 & 4 & -1 \\ 2 & 2 & -2 & 0 \\ 1 & 2 & -2 & 0 \\ 4 & 2 & 4 & -1 \end{bmatrix}$, $B = \begin{bmatrix} 0 & 1 & -1 & 0 \\ 1 & 5 & -5 & -1 \\ 1 & 11/2 & -6 & -1 \\ 6 & 36 & -38 & -7 \end{bmatrix}$, use Matlab/Octave

(i) $A = \begin{bmatrix} 0 & -2 & -2 & -3 \\ -4 & -5 & 0 & 1 \\ -2 & 0 & 4 & 5 \\ -1 & 1 & 0 & -2 \end{bmatrix}$, $B = \begin{bmatrix} 2/3 & -1/3 & 1/3 & -1/3 \\ -2/3 & 1/9 & -1/3 & 2/9 \\ 7/6 & -4/9 & 5/6 & 1/9 \\ -2/3 & 2/9 & -1/3 & -2/9 \end{bmatrix}$, use Matlab/Octave

(j) $A = \begin{bmatrix} 0 & 0 & 7 & -1 \\ 2 & 4 & 0 & 4 \\ -2 & -2 & 1 & -2 \\ -3 & -3 & 1 & -3 \end{bmatrix}$, $B = \begin{bmatrix} 0 & -1/2 & 2 & -2 \\ 1 & 1/2 & -21 & 15 \\ 0 & 0 & 3 & -2 \\ -1 & 0 & 21 & -14 \end{bmatrix}$, use Matlab/Octave

(k) $A = \begin{bmatrix} 1 & -7 & -4 & 3 \\ 0 & 2 & 1 & -1 \\ 0 & -2 & -3 & 2 \\ 3 & -2 & -4 & 1 \end{bmatrix}$, $B = \begin{bmatrix} -2 & -7 & -1 & 1 \\ -1 & -8/3 & 0 & 1/3 \\ -2 & -22/3 & -1 & 2/3 \\ -4 & -41/3 & -1 & 4/3 \end{bmatrix}$, use Matlab/Octave

Exercise 3.2.2. Use Theorem 3.2.6 to calculate the inverse, when it exists, of the following 2×2 matrices.

(a) $\begin{bmatrix} -2 & 2 \\ -1 & 4 \end{bmatrix}$

(b) $\begin{bmatrix} -5 & -10 \\ -1 & -2 \end{bmatrix}$

(c) $\begin{bmatrix} -2 & -4 \\ 5 & 2 \end{bmatrix}$

(d) $\begin{bmatrix} -3 & 2 \\ -1 & -2 \end{bmatrix}$

(e) $\begin{bmatrix} 2 & -4 \\ 3 & 0 \end{bmatrix}$

(f) $\begin{bmatrix} -0.6 & -0.9 \\ 0.8 & -1.4 \end{bmatrix}$

(g) $\begin{bmatrix} 0.3 & 0 \\ 0.9 & 1.9 \end{bmatrix}$

(h) $\begin{bmatrix} 0.6 & 0.5 \\ -0.3 & 0.5 \end{bmatrix}$

Exercise 3.2.3. Given the inverses of Exercises 3.2.1, solve each of the following systems of linear equations with a matrix-vector multiply (Theorem 3.2.8).

(a) $\begin{cases} -4y = 1 \\ 4x + 4y = -5 \end{cases}$

(b) $\begin{cases} -3p + 3q = 3 \\ -3p + q = -1 \end{cases}$

$$(c) \begin{cases} m - x = 1 \\ -m - 5x = -1 \end{cases} \quad (d) \begin{cases} -3x - z = 3 \\ x + 4y + 2z = -3 \\ 3x - 4y - z = 2 \end{cases}$$

$$(e) \begin{cases} 2p - 2q + 4r = -1 \\ -p + q + r = -2 \\ p + q - r = 1 \end{cases} \quad (f) \begin{cases} -x_1 + 3x_2 + 4x_3 - x_4 = 0 \\ 2x_1 + 2x_2 - 2x_3 = -1 \\ x_1 + 2x_2 - 2x_3 = 3 \\ 4x_1 + 2x_2 + 4x_3 - x_4 = -5 \end{cases}$$

$$(g) \begin{cases} p - 7q - 4r + 3s = -1 \\ 2q + r - s = -5 \\ -2q - 3r + 2s = 3 \\ 3p - 2q - 4r + s = -1 \end{cases} \quad (h) \begin{cases} -3b - 2c - 2d = 4 \\ -4a + b - 5d = -3 \\ -2a + 5b + 4c = -2 \\ -a - 2b + d = 0 \end{cases}$$

Exercise 3.2.4. Given the following information about solutions of systems of linear equations, write down if the matrix associated with each system is invertible, or not, or there is not enough given information to decide. Give reasons.

- (a) The general solution is $(-2, 1, 2, 0, 2)$.
- (b) A solution of a system is $(2.4, -2.8, -3.6, -2.2, -3.8)$.
- (c) A solution of a homogeneous system is $(0.8, 0.4, -2.3, 2.5)$.
- (d) The general solution of a system is $(4, 1, 0, 2)t$ for all t .
- (e) The general solution of a homogeneous system is $(0, 0, 0, 0)$.
- (f) A solution of a homogeneous system is $(0, 0, 0, 0, 0)$.

Exercise 3.2.5. Use Matlab/Octave to generate some random matrices of a suitable size of your choice, and some random scalar exponents (see Table 3.1). Then confirm the properties of inverse matrices given by Theorem 3.1.16. For the purposes of this exercise, use the Matlab/Octave function `inv(A)` that computes the inverse of the matrix A if it exists (as commented, computing the inverse of a matrix is generally not desirable—the inverse is primarily a useful theoretical device: this exercise only computes the inverse for educational purposes). Record all your commands and the output from Matlab/Octave.

Exercise 3.2.6. Consider Theorem 3.2.11 on the properties of the inverse. Invoking properties of matrix operations from section 3.1.3,

- (a) prove Part 3.2.11b using associativity, and
- (b) prove Part 3.2.11d using the transpose.

Exercise 3.2.7. Using the inverses identified in Exercise 3.2.1, and matrix multiplication, calculate the following matrix powers.

$$(a) \begin{bmatrix} 0 & -4 \\ 4 & 4 \end{bmatrix}^{-2}$$

$$(b) \begin{bmatrix} 0 & -4 \\ 4 & 4 \end{bmatrix}^{-3}$$

$$(c) \begin{bmatrix} -3 & 3 \\ -3 & 1 \end{bmatrix}^{-2}$$

$$(d) \begin{bmatrix} -1/6 & -1/6 \\ 5/6 & -1/6 \end{bmatrix}^{-4}$$

$$(e) \begin{bmatrix} -3 & 0 & -1 \\ 1 & 4 & 2 \\ 3 & -4 & -1 \end{bmatrix}^2$$

$$(f) \begin{bmatrix} 0 & 1/2 & 1/2 \\ 1/6 & 1/3 & 0 \\ 1/6 & -1/6 & 1/2 \end{bmatrix}^{-2}$$

$$(g) \begin{bmatrix} -1 & 3 & 4 & -1 \\ 2 & 2 & -2 & 0 \\ 1 & 2 & -2 & 0 \\ 4 & 2 & 4 & -1 \end{bmatrix}^{-2}$$

Matlab/Octave

$$(h) \begin{bmatrix} -2 & -7 & -1 & 1 \\ -1 & -8/3 & 0 & 1/3 \\ -2 & -22/3 & -1 & 2/3 \\ -4 & -41/3 & -1 & 4/3 \end{bmatrix}^{-3}$$

use Matlab/Octave.

Exercise 3.2.8. Which of the following matrices are diagonal? For those that are diagonal, write down how they may be represented with the diag function (algebraic, not Matlab/Octave).

$$(a) \begin{bmatrix} 9 & 0 & 0 \\ 0 & -5 & 0 \\ 0 & 0 & 4 \end{bmatrix}$$

$$(b) \begin{bmatrix} 0 & 0 & 1 \\ 0 & 2 & 0 \\ -2 & 0 & 0 \end{bmatrix}$$

$$(c) \begin{bmatrix} -5 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 9 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$(d) \begin{bmatrix} 6 & 0 & 0 & 0 \\ 0 & 1 & 0 & -9 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$(e) \begin{bmatrix} 0 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -5 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$(f) \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 2 & 0 \\ 0 & 0 \end{bmatrix}$$

$$(g) \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$(h) \begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 0 \end{bmatrix}$$

$$(i) \begin{bmatrix} -3 & 0 & c & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$(j) \begin{bmatrix} -3c & 0 & 0 & -4d \\ 0 & 5b & 0 & 0 \\ a & 0 & 0 & 0 \end{bmatrix}$$

Exercise 3.2.9. Write down the individual algebraic equations represented

by each of the following diagonal matrix-vector equations. Hence, where possible solve each system.

$$(a) \begin{bmatrix} -4 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 6 \\ -1 \\ -4 \\ -2 \\ 4 \end{bmatrix}$$

$$(b) \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = \begin{bmatrix} -8 \\ 3 \\ -1 \end{bmatrix}$$

$$(c) \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$$

$$(d) \begin{bmatrix} -6 & 0 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ 2 \end{bmatrix}$$

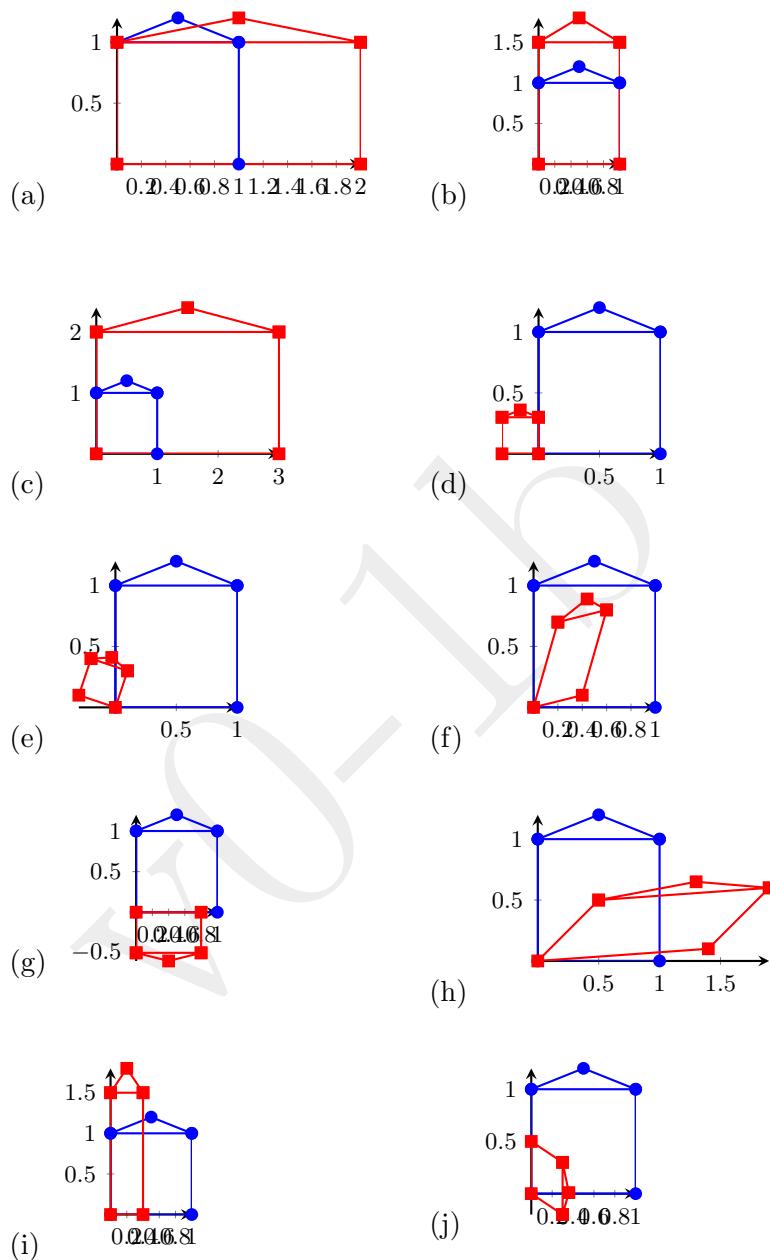
$$(e) \begin{bmatrix} -3 & 0 & 0 & 0 & 0 \\ 0 & -4 & 0 & 0 & 0 \\ 0 & 0 & -6 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix} = \begin{bmatrix} -5 \\ -5 \\ 1 \\ -1 \end{bmatrix}$$

$$(f) \begin{bmatrix} 0.5 & 0 & 0 & 0 \\ 0 & 0.2 & 0 & 0 \end{bmatrix} \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 3 \\ -1 \end{bmatrix}$$

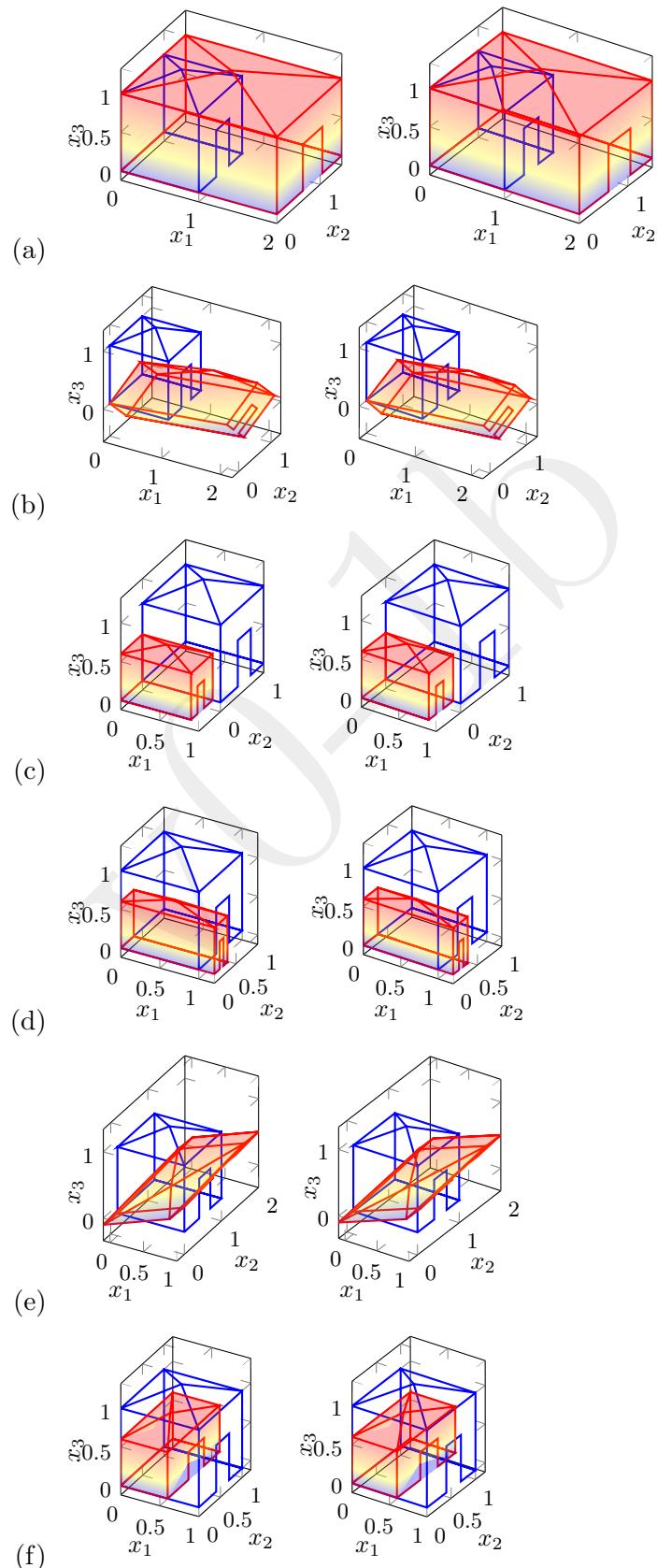
$$(g) \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -3 & 0 & 0 \\ 0 & 0 & -3 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} p \\ q \\ r \\ s \end{bmatrix} = \begin{bmatrix} -2 \\ -6 \\ 8 \\ 0 \end{bmatrix}$$

$$(h) \begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \quad \\ \quad \end{bmatrix} = \begin{bmatrix} -3 \\ 0 \\ 3 \end{bmatrix}$$

Exercise 3.2.10. In each of the following illustrations, the unit square (blue) is transformed by a matrix multiplication to some shape (red). Which of these transformations correspond to multiplication by a diagonal matrix? For those that are, estimate the elements of the diagonal matrix.

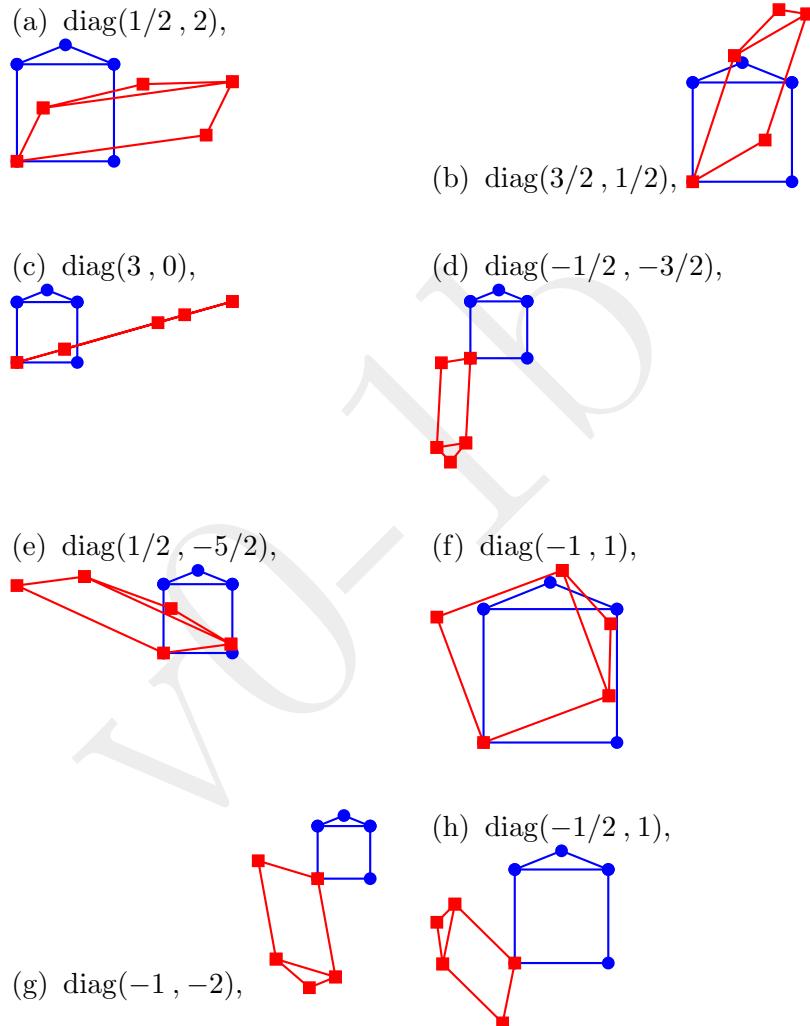


Exercise 3.2.11. In each of the following stereo illustrations, the unit cube (blue) is transformed by a matrix multiplication to some shape (red). Which of these transformations correspond to multiplication by a diagonal matrix? For those that are, estimate the elements of the diagonal matrix.

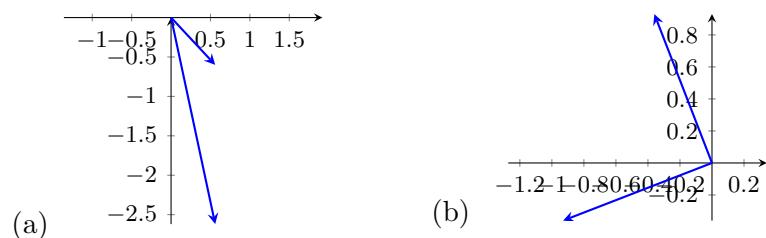


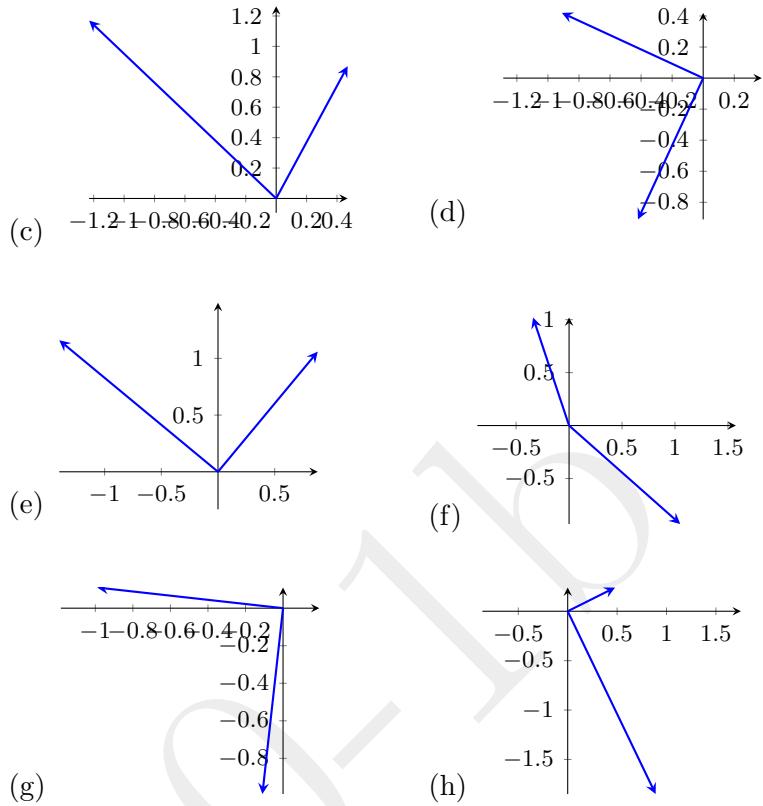
Exercise 3.2.12. Consider each of the transformations shown below that transform the from blue unit square to the red parallelogram.

They each have no coordinate axes shown because it is supposed to be some transformation in nature. Now impose on nature our mathematical description. Draw approximate orthogonal coordinate axes, with origin at the common corner point, so the transformation becomes that of multiplication by the specified diagonal matrix.

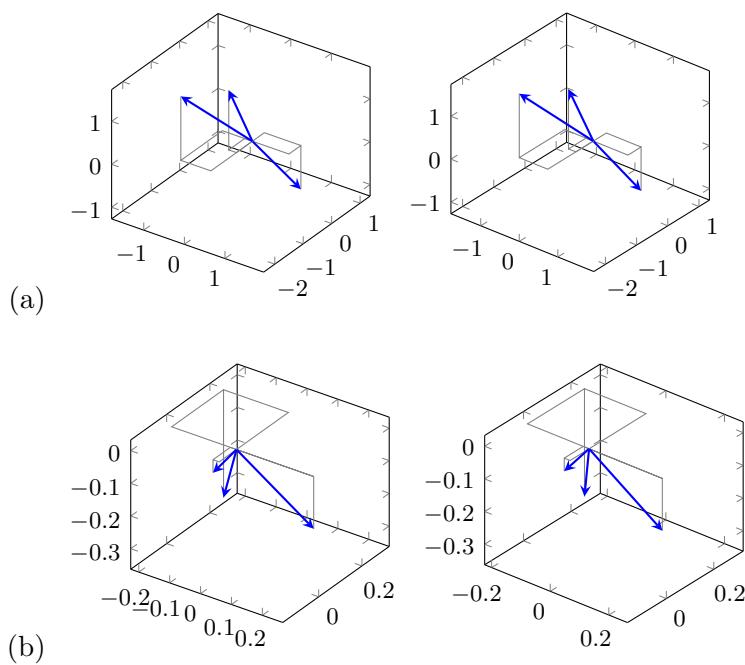


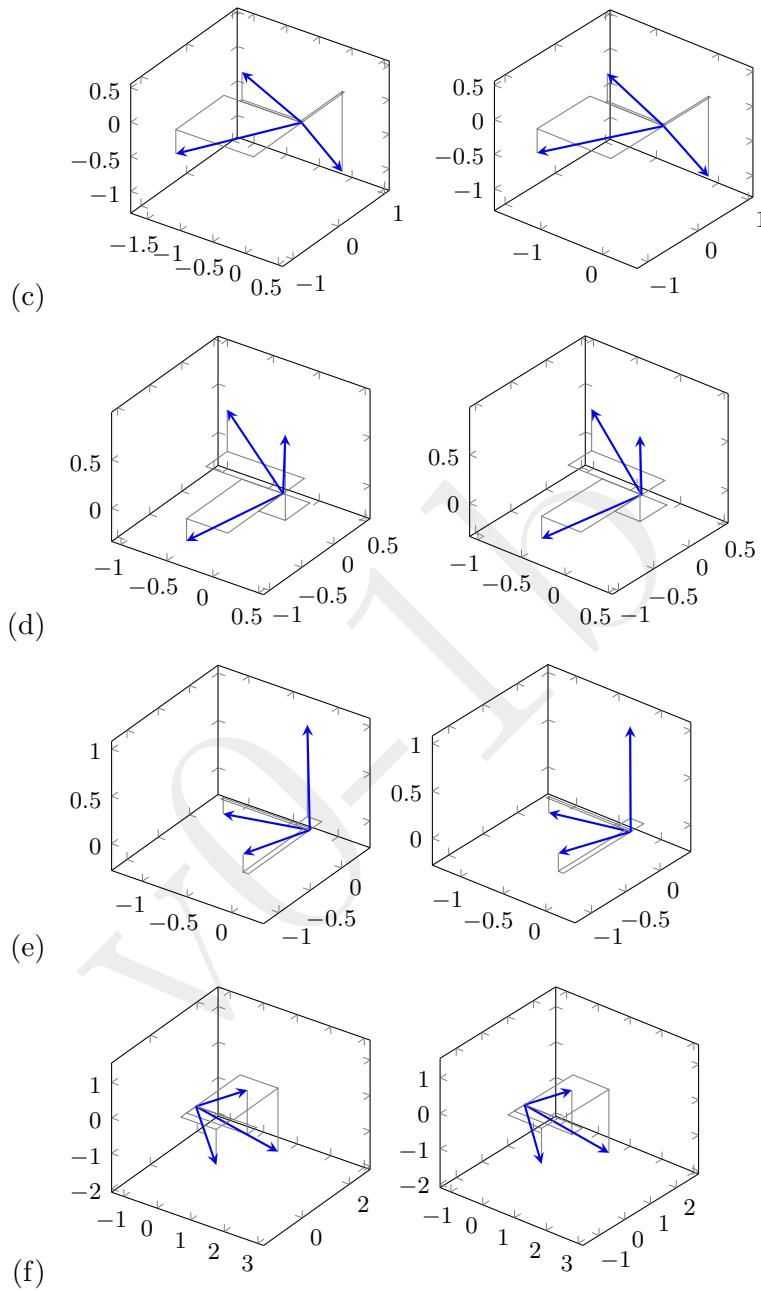
Exercise 3.2.13. Which of the following pairs of vectors appear to form an orthogonal set of two vectors? Which appear to form an orthonormal set of two vectors?





Exercise 3.2.14. Which of the following sets of vectors, drawn as stereo pairs, appear to form an orthogonal set? Which appear to form an orthonormal set?





Exercise 3.2.15. Use the dot product to determine which of the following sets of vectors are orthogonal sets. For the orthogonal sets, scale the vectors to form an orthonormal set.

- $\{(2, 3, 6), (3, -6, 2), (6, 2, -3)\}$
- $\{(4, 4, 7), (1, -8, 4), (-8, 1, 4)\}$
- $\{(2, 6, 9), (9, -6, 2)\}$
- $\{(6, 3, 2), (3, -6, 2), (2, -3, 6)\}$
- $\{(1, 1, 1, 1), (1, 1, -1, -1), (1, -1, -1, 1), (1, -1, 1, -1)\}$
- $\{(1, 2, 2, 4), (2, -1, -4, 2)\}$

$$(g) \{(1, 2, 2, 4), (2, -1, -4, 2), (-4, 2, 2, -1)\}$$

$$(h) \{(5, 6, 2, 4), (-2, 6, -5, -4), (6, -5, -4, 2)\}$$

Exercise 3.2.16. Using Definition 3.2.35 determine which of the following matrices are orthogonal matrices. For those matrices which are orthogonal, confirm Theorem Part 3.2.39c.

$$(a) \begin{bmatrix} \frac{5}{13} & \frac{12}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{bmatrix}$$

$$(b) \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$(c) \begin{bmatrix} -3 & 4 \\ 4 & 3 \end{bmatrix}$$

$$(d) \begin{bmatrix} \frac{2}{7} & \frac{3}{7} & \frac{6}{7} \\ \frac{3}{7} & -\frac{6}{7} & \frac{2}{7} \\ \frac{6}{7} & \frac{2}{7} & -\frac{3}{7} \end{bmatrix}$$

$$(e) \begin{bmatrix} \frac{2}{11} & \frac{9}{11} \\ \frac{6}{11} & -\frac{6}{11} \\ \frac{9}{11} & \frac{2}{11} \end{bmatrix}$$

$$(f) \begin{bmatrix} \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \end{bmatrix}$$

$$(g) \frac{1}{9} \begin{bmatrix} 4 & 1 & 8 \\ 4 & -8 & -1 \\ 7 & 4 & -4 \end{bmatrix}$$

$$(h) \frac{1}{7} \begin{bmatrix} 1 & 4 & 4 & 4 \\ 4 & -5 & 2 & 2 \\ 4 & 2 & -5 & 2 \\ 4 & 2 & 2 & -5 \end{bmatrix}$$

$$(i) \begin{bmatrix} 0.2 & 0.4 & 0.4 \\ 0.4 & -0.2 & 0.8 \\ 0.4 & -0.8 & -0.2 \\ 0.8 & 0.4 & -0.4 \end{bmatrix}$$

$$(j) \frac{1}{6} \begin{bmatrix} 1 & 1 & 3 & 5 \\ -5 & -3 & 1 & 1 \\ 3 & -5 & -1 & 1 \\ 1 & -1 & 5 & 3 \end{bmatrix}$$

$$(k) \begin{bmatrix} 0.1 & 0.5 & 0.5 & 0.7 \\ 0.5 & -0.1 & -0.7 & 0.5 \\ 0.5 & 0.7 & -0.1 & -0.5 \\ 0.7 & -0.5 & 0.5 & -0.1 \end{bmatrix}$$

Exercise 3.2.17. Each part gives an orthogonal matrix Q and two vectors \mathbf{u} and \mathbf{v} . For each part calculate the lengths of \mathbf{u} and \mathbf{v} , and the angle between \mathbf{u} and \mathbf{v} . Confirm these are the same as the lengths of $Q\mathbf{u}$ and $Q\mathbf{v}$, and the angle between $Q\mathbf{u}$ and $Q\mathbf{v}$, respectively.

$$(a) Q = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} 12 \\ 5 \end{bmatrix}$$

$$(b) Q = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$(c) Q = \begin{bmatrix} -\frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 2 \\ -3 \end{bmatrix}, \mathbf{v} = [0 \ -4]$$

$$(d) Q = \begin{bmatrix} \frac{3}{5} & \frac{4}{5} \\ -\frac{4}{5} & \frac{3}{5} \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 7 \\ -4 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$(e) Q = \frac{1}{11} \begin{bmatrix} 7 & 6 & 6 \\ 6 & 2 & -9 \\ 6 & -9 & 2 \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$(f) Q = \frac{1}{9} \begin{bmatrix} 7 & 4 & 4 \\ 4 & -8 & 1 \\ 4 & 1 & -8 \end{bmatrix}, \mathbf{u} = \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} -1 \\ 2 \\ -3 \end{bmatrix}$$

$$(g) Q = \begin{bmatrix} 0.1 & 0.1 & 0.7 & 0.7 \\ 0.1 & -0.1 & -0.7 & 0.7 \\ 0.7 & 0.7 & -0.1 & -0.1 \\ 0.7 & -0.7 & 0.1 & -0.1 \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 0 \\ 0 \\ -2 \\ -3 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} -1 \\ -1 \\ -2 \\ -1 \end{bmatrix}$$

$$(h) Q = \begin{bmatrix} 0.1 & 0.3 & 0.3 & 0.9 \\ 0.3 & -0.1 & -0.9 & 0.3 \\ 0.3 & 0.9 & -0.1 & -0.3 \\ 0.9 & -0.3 & 0.3 & -0.1 \end{bmatrix}, \mathbf{u} = \begin{bmatrix} 3 \\ 1 \\ 2 \\ 1 \end{bmatrix}, \mathbf{v} = \begin{bmatrix} -2 \\ -1 \\ 0 \\ 2 \end{bmatrix}$$

Exercise 3.2.18. Using one or other of the orthogonal matrices appearing in Exercise 3.2.17, solve each of the following systems of linear equations by a matrix-vector multiplication.

$$(a) \begin{cases} \frac{3}{5}x + \frac{4}{5}y = 5 \\ -\frac{4}{5}x + \frac{3}{5}y = 2 \end{cases}$$

$$(b) \begin{cases} -\frac{3}{5}x + \frac{4}{5}y = 1.6 \\ \frac{4}{5}x + \frac{3}{5}y = -3.5 \end{cases}$$

$$(c) \begin{cases} \frac{1}{\sqrt{2}}(x+y) = 3 \\ \frac{1}{\sqrt{2}}(x-y) = 2 \end{cases}$$

$$(d) \begin{cases} 3x + 4y = 20 \\ -4x + 3y = 5 \end{cases}$$

$$(e) \begin{cases} \frac{7}{9}p + \frac{4}{9}q + \frac{4}{9}r = 2 \\ \frac{4}{9}p - \frac{8}{9}q + \frac{1}{9}r = 3 \\ \frac{4}{9}p + \frac{1}{9}q - \frac{8}{9}r = 7 \end{cases}$$

$$(f) \begin{cases} \frac{7}{9}u + \frac{4}{9}v + \frac{4}{9}w = 1 \\ \frac{4}{9}u + \frac{1}{9}v - \frac{8}{9}w = 2 \\ \frac{4}{9}u - \frac{8}{9}v + \frac{1}{9}w = 0 \end{cases}$$

$$(g) \begin{cases} 7a + 6b + 6c = 22 \\ 6a + 2b - 9c = 11 \\ 6a - 9b + 2c = -22 \end{cases}$$

$$(h) \begin{cases} 0.1x_1 + 0.1x_2 + 0.7x_3 + 0.7x_4 = 1 \\ 0.1x_1 - 0.1x_2 - 0.7x_3 + 0.7x_4 = -1 \\ 0.7x_1 + 0.7x_2 - 0.1x_3 - 0.1x_4 = 0 \\ 0.7x_1 - 0.7x_2 + 0.1x_3 - 0.1x_4 = -2 \end{cases}$$

$$(i) \begin{cases} 0.1y_1 + 0.1y_2 + 0.7y_3 + 0.7y_4 = 1 \\ 0.7y_1 + 0.7y_2 - 0.1y_3 - 0.1y_4 = -2.5 \\ 0.7y_1 - 0.7y_2 + 0.1y_3 - 0.1y_4 = 2 \\ 0.1y_1 - 0.1y_2 - 0.7y_3 + 0.7y_4 = 2.5 \end{cases}$$

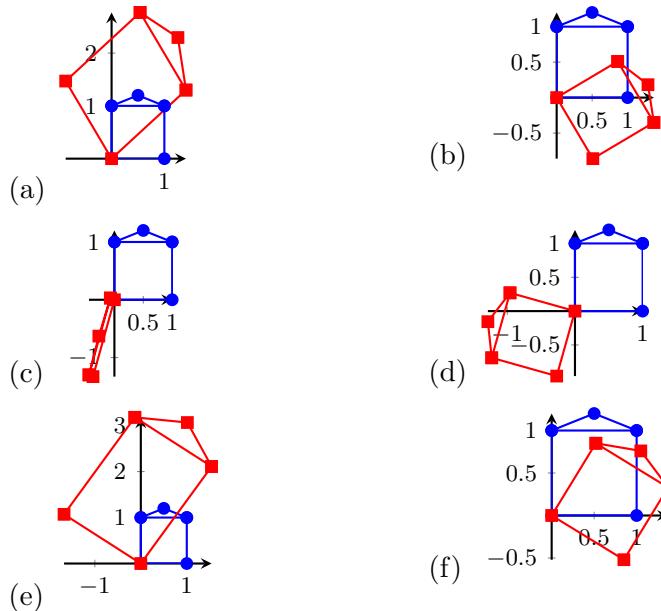
$$(j) \begin{cases} z_1 + 3z_2 + 3z_3 + 9z_4 = 5 \\ 3z_1 - z_2 - 9z_3 + 3z_4 = 0 \\ 3z_1 + 9z_2 - z_3 - 3z_4 = -1 \\ 9z_1 - 3z_2 + 3z_3 - z_4 = -3 \end{cases}$$

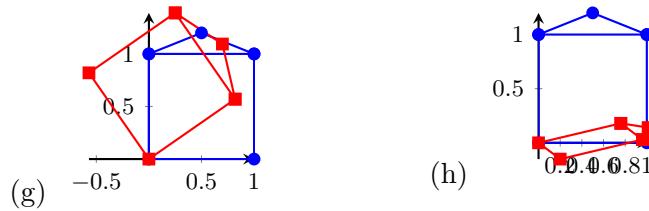
Exercise 3.2.19 (product of orthogonal matrices). Use Definition 3.2.35 to prove that if Q_1 and Q_2 are orthogonal matrices of the same size, then so is the product Q_1Q_2 . Consider $(Q_1Q_2)^T(Q_1Q_2)$.

Exercise 3.2.20. Fill in details of the proof for Theorem 3.2.39 to establish that a matrix Q^T is orthogonal iff the row vectors of Q form an orthonormal set.

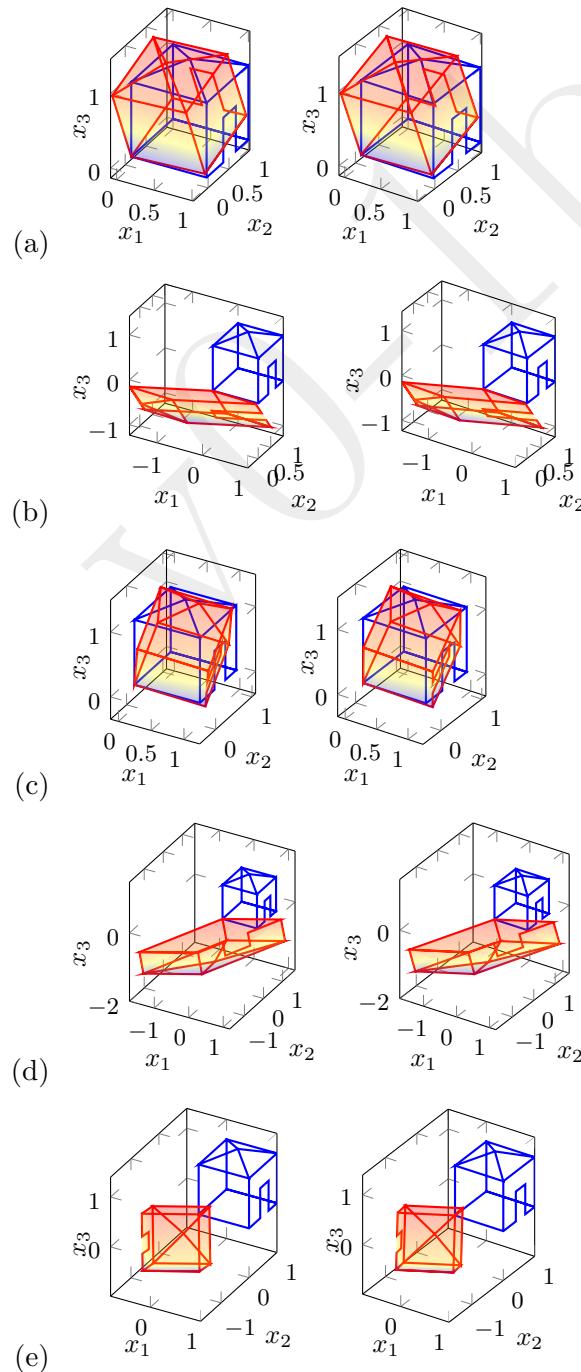
Exercise 3.2.21. Fill in details of the proof for Theorem 3.2.39 to establish that if the row vectors of Q form an orthonormal set, then Q is invertible and $Q^{-1} = Q^T$.

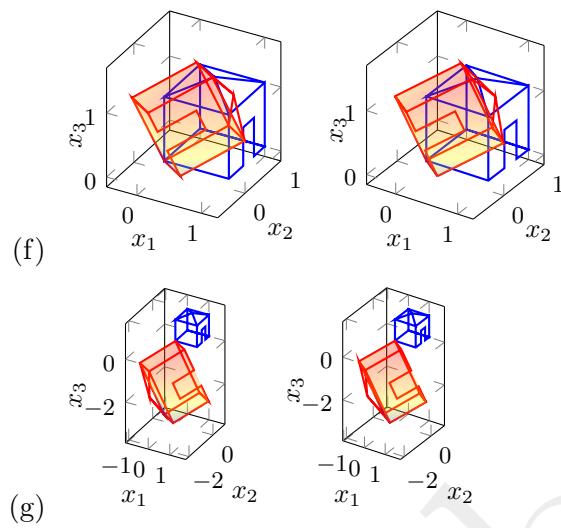
Exercise 3.2.22. The following graphs illustrate the transformation of the unit square through multiplying by some different matrices. Using Theorem 3.2.39f, which transformations appear to be that of multiplying by an orthogonal matrix?





Exercise 3.2.23. The following stereo pairs illustrate the transformation of the unit cube through multiplying by some different matrices. Using Theorem 3.2.39f, which transformations appear to be that of multiplying by an orthogonal matrix?





3.3 Factorise to the singular value decomposition

Section Contents

3.3.1	Introductory examples	205
3.3.2	The SVD solves general systems	209
	Computers empower use of the SVD	211
	Condition number and rank determine possibilities	217
3.3.3	Prove the SVD Theorem 3.3.5	225
	Prelude to the proof	225
	Detailed proof of the SVD Theorem 3.3.5 . .	228
3.3.4	Exercises	230

Beltrami first derived the SVD in 1873. The first reliable method for computing an SVD was developed by Golub and Kahan in 1965, and only thereafter did applications proliferate.

The singular value decomposition (SVD) is sometimes called the jewel in the crown of linear algebra. Its importance is certified by the many names by which it is invoked in scientific and engineering applications: principal component analysis, singular spectrum analysis, principal orthogonal decomposition, latent semantic indexing, Schmidt decomposition, correspondence analysis, Lanczos methods, dimension reduction, and so on. Let's start seeing what it can do for us.

3.3.1 Introductory examples

Introduce an analogous procedure so the SVD version follows more easily.

You are a contestant in a quiz show. The final million dollar question is:

in your head, without a calculator, solve $42x = 1554$
within twenty seconds,

your time starts now

Solution: Long division is hopeless in the time available. However, recognise $42 = 2 \cdot 3 \cdot 7$ and so divide 1554 by 2 to get 777, divide 777 by 3 to get 259, and divide 259 by 7 to get 37, and win the prize.

Example 3.3.1. Given $154 = 2 \cdot 7 \cdot 11$, solve in your head $154x = 8008$ or 9856 or 12628 or 13090 or 14322 (teacher to choose): first to answer wins. ■

Such examples show factorisation can turn a hard problem into a sequence of easy problems. We adopt a matrix factorisation for general linear equations.

To illustrate the procedure to come, let's write the above solution steps in detail: we solve $42x = 1554$.

- Factorise the coefficient $42 = 2 \cdot 3 \cdot 7$ so the equation becomes

$$2 \cdot 3 \cdot \underbrace{7x}_{=z} = 1554,$$

and introduce two intermediate unknowns y and z as shown above.

- Solve $2z = 1554$ to get $z = 777$.
- Solve $3y = z = 777$ to get $y = 259$.
- Solve $7x = y = 259$ to get $x = 37$ —the answer.

Now let's proceed to small matrix examples—albeit easier to solve directly—to introduce the general matrix procedure.

Example 3.3.2. Solve the 2×2 system

$$\begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 18 \\ -1 \end{bmatrix}$$

given the matrix factorisation

$$\begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix} = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T$$

(note the transpose on the last matrix).

Solution: Optionally check the factorisation if you like:

$$\begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} = \begin{bmatrix} 6\sqrt{2} & -4\sqrt{2} \\ 8\sqrt{2} & 3\sqrt{2} \end{bmatrix};$$

then $\begin{bmatrix} 6\sqrt{2} & -4\sqrt{2} \\ 8\sqrt{2} & 3\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix}.$

The following four steps forms the general procedure.

- Write the system using the factorisation, and with two intermediate unknowns \mathbf{y} and \mathbf{z} :

$$\underbrace{\begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix}}_{=z} \underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T}_{=y} \mathbf{x} = \begin{bmatrix} 18 \\ -1 \end{bmatrix}.$$

- Solve $\begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \mathbf{z} = \begin{bmatrix} 18 \\ -1 \end{bmatrix}$: recall that the matrix appearing here is orthogonal (and this is no accident), so multiplying by the transpose gives the intermediary

$$\mathbf{z} = \begin{bmatrix} \frac{3}{5} & \frac{4}{5} \\ -\frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 18 \\ -1 \end{bmatrix} = \begin{bmatrix} 10 \\ -15 \end{bmatrix}.$$

(c) Now solve $\begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \mathbf{y} = \mathbf{z} = \begin{bmatrix} 10 \\ -15 \end{bmatrix}$: the matrix appearing here is diagonal (and this is no accident), so dividing by the respective diagonal elements gives the intermediary

$$\mathbf{y} = \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix}^{-1} \begin{bmatrix} 10 \\ -15 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} \\ -3/\sqrt{2} \end{bmatrix}.$$

(d) Finally solve $\begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T \mathbf{x} = \mathbf{y} = \begin{bmatrix} 1/\sqrt{2} \\ -3/\sqrt{2} \end{bmatrix}$: now the matrix appearing here is also orthogonal (its orthogonality is also no accident), so multiplying by itself (the transpose of the transpose) gives the solution

$$\mathbf{x} = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} \\ -3/\sqrt{2} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} + \frac{3}{2} \\ \frac{1}{2} - \frac{3}{2} \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}.$$

■

Example 3.3.3. Solve the 3×3 system

$$A\mathbf{x} = \begin{bmatrix} 10 \\ 2 \\ -2 \end{bmatrix} \quad \text{for matrix } A = \begin{bmatrix} -4 & -2 & 4 \\ -8 & -1 & -4 \\ 6 & 6 & 0 \end{bmatrix}$$

using the following given matrix factorisation (note the last is transposed)

$$A = \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} \underbrace{\begin{bmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 3 \end{bmatrix}}_{=S} \begin{bmatrix} -\frac{8}{9} & -\frac{1}{9} & -\frac{4}{9} \\ -\frac{4}{9} & \frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & \frac{4}{9} \end{bmatrix}^T.$$

Solution: Use Matlab/Octave. Enter the matrices and the right-hand side, and check the factorisation (and the typing):

```
U=[1,-2,2;2,2,1;-2,1,2]/3
S=[12,0,0;0,6,0;0,0,3]
V=[-8,-1,-4;-4,4,7;-1,-8,4]/9
b=[10;2;-2]
A=U*S*V'
```



(a) Write the system $A\mathbf{x} = \mathbf{b}$ using the factorisation, and with two intermediate unknowns \mathbf{y} and \mathbf{z} :

$$\underbrace{\begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 3 \end{bmatrix}}_{=\mathbf{z}} \underbrace{\begin{bmatrix} -\frac{8}{9} & -\frac{1}{9} & -\frac{4}{9} \\ -\frac{4}{9} & \frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & \frac{4}{9} \end{bmatrix}^T}_{=\mathbf{y}} \mathbf{x} = \begin{bmatrix} 10 \\ 2 \\ -2 \end{bmatrix}.$$

(b) Solve $\begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} \mathbf{z} = \begin{bmatrix} 10 \\ 2 \\ -2 \end{bmatrix}$. Now this matrix, called \mathbf{U} , is orthogonal—check by computing $\mathbf{U}' * \mathbf{U}$ —so multiplying by the transpose gives the intermediary: $\mathbf{z} = \mathbf{U}' * \mathbf{b} = (6, -6, 6)$.

(c) Then solve $\begin{bmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 3 \end{bmatrix} \mathbf{y} = \mathbf{z} = \begin{bmatrix} 6 \\ -6 \\ 6 \end{bmatrix}$: this matrix, called \mathbf{S} , is diagonal, so dividing by the respective diagonal elements gives the intermediary $\mathbf{y} = \mathbf{z} ./ \text{diag}(\mathbf{S}) = (\frac{1}{2}, -1, 2)$.

(d) Finally solve $\begin{bmatrix} -\frac{8}{9} & -\frac{1}{9} & -\frac{4}{9} \\ -\frac{4}{9} & \frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & \frac{4}{9} \end{bmatrix}^T \mathbf{x} = \mathbf{y} = \begin{bmatrix} \frac{1}{2} \\ -1 \\ 2 \end{bmatrix}$. This matrix, called \mathbf{V} , is also orthogonal—check by computing $\mathbf{V}' * \mathbf{V}$ —so multiplying by itself (the transpose of the transpose) gives the final solution $\mathbf{x} = \mathbf{V} * \mathbf{y} = (-\frac{11}{9}, \frac{8}{9}, \frac{31}{18})$.

■

Warning: do *not* solve in reverse order

Example 3.3.4. Reconsider Example 3.3.2 wrongly.

(a) After writing the system using the SVD as

$$\underbrace{\begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix}}_{=\mathbf{z}} \underbrace{\begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix}}_{=\mathbf{y}} \underbrace{\begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T}_{=\mathbf{x}} = \begin{bmatrix} 18 \\ -1 \end{bmatrix},$$

one might be inadvertently tempted to ‘solve’ the system by using the matrices in reverse order as in the following: *do not do this*.

(b) First solve $\begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T \mathbf{x} = \begin{bmatrix} 18 \\ -1 \end{bmatrix}$: this matrix is orthogonal, so multiplying by itself (the transpose of the transpose) gives

$$\mathbf{x} = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 18 \\ -1 \end{bmatrix} = \begin{bmatrix} 19/\sqrt{2} \\ 17/\sqrt{2} \end{bmatrix}.$$

(c) Inappropriately ‘solve’, $\begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \mathbf{y} = \begin{bmatrix} 19/\sqrt{2} \\ 17/\sqrt{2} \end{bmatrix}$: this matrix is diagonal, so dividing by the diagonal elements gives

$$\mathbf{y} = \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix}^{-1} \begin{bmatrix} 19/\sqrt{2} \\ 17/\sqrt{2} \end{bmatrix} = \begin{bmatrix} \frac{19}{20} \\ \frac{17}{10} \end{bmatrix}.$$

(d) Inappropriately ‘solve’ $\begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \mathbf{z} = \begin{bmatrix} \frac{19}{20} \\ \frac{17}{10} \end{bmatrix}$: this matrix is orthogonal, so multiplying by the transpose gives

$$\mathbf{z} = \begin{bmatrix} \frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix}^{-1} \begin{bmatrix} \frac{19}{20} \\ \frac{17}{10} \end{bmatrix} = \begin{bmatrix} 1.93 \\ 0.26 \end{bmatrix}.$$

And then, since the solution is to be called \mathbf{x} , we might inappropriately call what we just calculated as the solution $\mathbf{x} = (1.93, 0.26)$. ■

Avoid this reverse process as it is wrong. Matrix multiplicative is *not* commutative (section 3.1.3). We must use an SVD factorisation in the correct order: to solve linear equations use the matrices in an SVD from left to right.

3.3.2 The SVD solves general systems

<http://www.youtube.com/watch?v=JEYLfIVvR9I> is an entertaining prelude

Theorem 3.3.5 (SVD factorisation). *Every $m \times n$ real matrix A can be factored into a product of three matrices*

$$A = USV^T, \quad (3.4)$$

called a **singular value decomposition** (SVD), where

- $m \times m$ matrix $U = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_m]$ is orthogonal,
- $n \times n$ matrix $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$ is orthogonal, and
- $m \times n$ diagonal matrix S is zero except for non-negative diagonal elements called **singular values** (denoted by greek letter sigma) $\sigma_1, \sigma_2, \dots, \sigma_{\min(m,n)}$, which are unique when ordered from largest to smallest so that $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{\min(m,n)} \geq 0$.

The orthonormal vectors \mathbf{u}_j and \mathbf{v}_j are called **singular vectors**. 11

Proof. Detailed in Section 3.3.3. Importantly, the singular values are unique (when ordered), but the orthogonal matrices U and V are not unique (e.g., one may always change the sign of corresponding columns in U and V). Nonetheless, although there are many SVDs of a matrix, all SVDS are equivalent in application. □

¹¹ This enormously useful theorem also generalises from matrices of finite dimension to analogues in ‘infinite’ dimensions: an SVD exists for all compact linear operators (Kress 2015, §7).

Some may be disturbed by the non-uniqueness of an SVD. But the non-uniqueness is analogous to the non-uniqueness of row reduction upon arbitrary re-ordering of equations, and/or re-ordering the variables in the equations.

Example 3.3.6. Example 3.3.2 invoked the SVD

$$\begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix} = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T,$$

where the two outer matrices are orthogonal (check), so the singular values of this matrix are $\sigma_1 = 10\sqrt{2}$ and $\sigma_2 = 5\sqrt{2}$.

Example 3.3.3 invoked the SVD

$$\begin{bmatrix} -4 & -2 & 4 \\ -8 & -1 & -4 \\ 6 & 6 & 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} -\frac{8}{9} & -\frac{1}{9} & -\frac{4}{9} \\ -\frac{4}{9} & \frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & \frac{4}{9} \end{bmatrix}^T,$$

where the two outer matrices are orthogonal (check), so the singular values of this matrix are $\sigma_1 = 12$, $\sigma_2 = 6$ and $\sigma_3 = 3$. ■

Example 3.3.7. Any orthogonal matrix Q , say $n \times n$, has an SVD $Q = QI_nI_n^T$; that is, $U = Q$, $S = V = I_n$. Hence every $n \times n$ orthogonal matrix has singular values $\sigma_1 = \sigma_2 = \dots = \sigma_n = 1$. ■

Example 3.3.8 (some non-uniqueness). has an SVD $I_n = I_nI_nI_n^T$.

- An identity matrix, say I_n ,
- But, for *any* $n \times n$ orthogonal matrix Q , the identity I_n also has the SVD $I_n = QI_nQ^T$ (as this right-hand side $QI_nQ^T = QQ^T = I_n$).
- Lastly, any constant multiple of an identity, say $sI_n = \text{diag}(s, s, \dots, s)$, has the same non-uniqueness: an SVD is $sI_n = USV^T$ for matrices $U = Q$, $S = sI_n$ and $V = Q$ for any $n \times n$ orthogonal Q (provided $s \geq 0$).

The matrices in this example are characterised by all their singular values having an identical value. In general, analogous non-uniqueness occurs whenever two or more singular values are identical in value. ■

Example 3.3.9 (positive ordering). Find an SVD of the diagonal matrix

$$D = \begin{bmatrix} 2.7 & 0 & 0 \\ 0 & -3.9 & 0 \\ 0 & 0 & -0.9 \end{bmatrix}.$$

Table 3.3: As well as the Matlab/Octave commands and operations listed in Tables 1.2, 2.3, 3.1 and 3.2, we need these matrix operations.

- $[U, S, V] = \text{svd}(A)$ computes the three matrices U , S and V in a singular value decomposition (SVD) of the $m \times n$ matrix: $A = USV^T$ for $m \times m$ orthogonal matrix U , $n \times n$ orthogonal matrix V , and $m \times n$ non-negative diagonal matrix S (Theorem 3.3.5).
- $\text{svd}(A)$ just reports the singular values in a vector.
- Complementing information of Table 3.1, to extract and compute with a subset of rows/columns of a matrix, specify the vector of indices. For examples:
 - $V(:, 1:r)$ selects the first r columns of V ;
 - $A([2 3 5], :)$ selects the second, third and fifth row of matrix A ;
 - $B(4:6, 1:3)$ selects the 3×3 submatrix of the first three columns of the fourth, fifth and sixth rows.

Solution: Singular values cannot be negative so a factorisation is

$$D = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} 2.7 & 0 & 0 \\ 0 & 3.9 & 0 \\ 0 & 0 & 0.9 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}^T,$$

where the (-1) s in the first matrix encode the signs of the corresponding diagonal elements. However, Theorem 3.3.5 requires that singular values be ordered in decreasing magnitude, so sort the diagonal of the middle matrix into order and correspondingly permute the columns of the outer two matrices to obtain the following SVD:

$$D = \begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & -1 \end{bmatrix} \begin{bmatrix} 3.9 & 0 & 0 \\ 0 & 2.7 & 0 \\ 0 & 0 & 0.9 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}^T.$$

Alternatively, one could use the rightmost matrix to contain the pattern of signs.

■

Computers empower use of the SVD

Except for simple cases such as 2×2 matrices (Example 3.3.26), constructing an SVD is usually far too laborious by hand.¹² Typically, this book either gives an SVD (as in the earlier two examples) or computes an SVD in Matlab/Octave with $[U, S, V] = \text{svd}(A)$ (Table 3.3).

¹² For those interested advanced students, Trefethen & Bau (1997) [p.234] discusses how the standard method of numerically computing an SVD is based upon first transforming to bidiagonal form, and then using an iteration based upon a so-called QR factorisation.

The following examples illustrate no or infinite solutions, to follow from the unique solutions of the first two examples.

Example 3.3.10 (rate sport teams/players). Consider three table tennis players, Anne, Bob and Chris: Anne beat Bob 3 games to 2 games; Anne beat Chris 3-1; Bob beat Chris 3-2. How good are they? What is their rating?

Solution: Denote Anne's rating by x_1 , Bob's rating by x_2 , and Chris' rating by x_3 . The ratings should predict the results of matches, so from the above three match results, surely

- Anne beat Bob 3 games to 2 $\leftrightarrow x_1 - x_2 = 3 - 2 = 1$;
- Anne beat Chris 3-1 $\leftrightarrow x_1 - x_3 = 3 - 1 = 2$; and
- Bob beat Chris 3-2 $\leftrightarrow x_2 - x_3 = 3 - 2 = 1$.

In matrix-vector form, $A\mathbf{x} = \mathbf{b}$,

$$\begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \\ 0 & 1 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}.$$

In Matlab/Octave, we might try Procedure 2.2.4:

```
A=[1,-1,0;1,0,-1;0,1,-1]
b=[1;2;1]
rcond(A)
```

but find $rcond=0$ which is extremely terrible so we cannot use $\mathbf{A}\backslash\mathbf{b}$ to solve the system $A\mathbf{x} = \mathbf{b}$. Whenever, difficulties arise, use an SVD.

- (a) Compute an SVD $A = USV^T$ with $[U, S, V] = \text{svd}(A)$ (Table 3.3): here

```
U =
    0.4082   -0.7071    0.5774
   -0.4082   -0.7071   -0.5774
   -0.8165   -0.0000    0.5774
S =
    1.7321         0         0
        0    1.7321         0
        0         0    0.0000
V =
    0.0000   -0.8165    0.5774
   -0.7071    0.4082    0.5774
    0.7071    0.4082    0.5774
```

so the singular values are $\sigma_1 = \sigma_2 = 1.7321 = \sqrt{3}$ and $\sigma_3 = 0$ (different computers may give different U and V , but any



deductions will be equivalent). The system of equations for the ratings becomes

$$A\mathbf{x} = U \underbrace{S V^T \mathbf{x}}_{=\mathbf{z}} = \mathbf{b} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}.$$

- (b) As U is orthogonal, $U\mathbf{z} = \mathbf{b}$ has unique solution $\mathbf{z} = U^T \mathbf{b}$ computed by $\mathbf{z}=U'*\mathbf{b}$:

$$\begin{aligned}\mathbf{z} = & \\ -1.2247 & \\ -2.1213 & \\ 0 & \end{aligned}$$

- (c) Now solve $S\mathbf{y} = \mathbf{z}$. But S has a troublesome zero on the diagonal. So interpret the equation $S\mathbf{y} = \mathbf{z}$ in detail as

$$\begin{bmatrix} 1.7321 & 0 & 0 \\ 0 & 1.7321 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{y} = \begin{bmatrix} -1.2247 \\ -2.1213 \\ 0 \end{bmatrix}:$$

- i. the first line implies $y_1 = -1.2247/1.7321$;
- ii. the second line implies $y_2 = -2.1213/1.7321$;
- iii. the third line is $0y_3 = 0$ which is satisfied for all y_3 .

In using Matlab/Octave you must notice $\sigma_3 = 0$, check that the corresponding $z_3 = 0$, and then compute a *particular solution* from the first two components to give the first two components of \mathbf{y} :

$$\begin{aligned}\mathbf{y} = & \mathbf{z}(1:2) ./ \text{diag}(S(1:2,1:2)) \\ \mathbf{y} = & \\ -0.7071 & \\ -1.2247 & \end{aligned}$$

The third component, involving the free variable y_3 , we omit from this numerical computation.

- (d) Finally, as V is orthogonal, $V^T \mathbf{x} = \mathbf{y}$ has the solution $\mathbf{x} = V\mathbf{y}$ (unique for each valid \mathbf{y}): in Matlab/Octave, compute a particular solution with $\mathbf{x}=V(:,1:2)*\mathbf{y}$

$$\begin{aligned}\mathbf{x} = & \\ 1.0000 & \\ 0.0000 & \\ -1.0000 & \end{aligned}$$

Then for a general solution remember to add an arbitrary multiple, y_3 , of $V(:,3) = (0.5774, 0.5774, 0.5774) = (1, 1, 1)/\sqrt{3}$.

Thus the three player ratings may be any one the infinity of possible solutions from the general solution

$$(x_1, x_2, x_3) = (1, 0, -1) + y_3(1, 1, 1)/\sqrt{3}.$$

In this application we only care about relative ratings, not absolute, so here adding any multiple of $(1, 1, 1)$ is immaterial. This solution for the ratings indicates Anne is the best player, and Chris the worst.

■

Compute in Matlab/Octave. As seen in the previous example, often we need to compute with a subset of the components of matrices (Table 3.3):

- $\mathbf{b}(1:r)$ selects the first r entries of vector \mathbf{b}
- $S(1:r, 1:r)$ selects the top-left $r \times r$ submatrix of S ;
- $V(:, 1:r)$ selects the first r columns of matrix V .

Example 3.3.11. But what if Bob beat Chris 3-1?

Solution: The only change to the problem is the new right-hand side $\mathbf{b} = (1, 2, 2)$.

- (a) An SVD of matrix A remains the same.
- (b) $U\mathbf{z} = \mathbf{b}$ has unique solution $\mathbf{z}=U'*\mathbf{b}$ of

$$\begin{aligned}\mathbf{z} &= \\ &-2.0412 \\ &-2.1213 \\ &0.5774\end{aligned}$$

- (c) We need to interpret $S\mathbf{y} = \mathbf{z}$,

$$\begin{bmatrix} 1.7321 & 0 & 0 \\ 0 & 1.7321 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{y} = \begin{bmatrix} -2.0412 \\ -2.1213 \\ 0.5774 \end{bmatrix}.$$

The third line of this system says $0y_3 = 0.5774$ which is impossible for any y_3 .

In this case there is no solution of the system of equations. It would appear we cannot assign ratings to the players!

■

Section 3.5 further explores systems with no solution and uses the SVD to determine a good approximate solution (Example 3.5.2).

Example 3.3.12. Find the value(s) of the parameter c such that the following system has a solution, and find a general solution for that (those) parameter value(s):

$$\begin{bmatrix} -9 & -15 & -9 & -15 \\ -10 & 2 & -10 & 2 \\ 8 & 4 & 8 & 4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} c \\ 8 \\ -5 \end{bmatrix}.$$

Solution: (a) In Matlab/Octave, compute an SVD of this 3×4 matrix with

```
A=[-9 -15 -9 -15; -10 2 -10 2; 8 4 8 4]
[U,S,V]=svd(A)
U =
    0.8571    0.4286    0.2857
    0.2857   -0.8571    0.4286
   -0.4286    0.2857    0.8571
S =
    28.0000         0         0         0
         0    14.0000         0         0
         0         0    0.0000         0
V =
   -0.5000    0.5000   -0.1900   -0.6811
   -0.5000   -0.5000    0.6811   -0.1900
   -0.5000    0.5000    0.1900    0.6811
   -0.5000   -0.5000   -0.6811    0.1900
```

Depending upon Matlab/Octave you may get different alternatives for the last two columns for V —adjust accordingly.

The singular values are $\sigma_1 = 28$, $\sigma_2 = 14$ and the problematic $\sigma_3 = 0$ (it is computed as the negligible 10^{-15}).

- (b) We want to solve $U\mathbf{z} = \mathbf{b}$ via $\mathbf{z} = U^T\mathbf{b}$. But for the next step we must have the third component of \mathbf{z} to be zero as otherwise there is no solution. Now $z_3 = \mathbf{u}_3^T\mathbf{b}$ (where \mathbf{u}_3 is the third column of U); that is, $z_3 = 0.2857 \times c + 0.4286 \times 8 + 0.8571 \times (-5)$ needs to be zero, which requires $c = -(0.4286 \times 8 + 0.8571 \times (-5))/0.2857$. Recognise this expression is equivalent to $c = -(0.2857 \times 0 + 0.4286 \times 8 + 0.8571 \times (-5))/0.2857 = \mathbf{u}_3 \cdot (0, 8, -5)/0.2857$ and so compute

```
c=-U(:,3)'*[0;8;-5]/U(1,3)
```

Having found $c = 3$, compute \mathbf{z} from $\mathbf{z}=U'*[3;8;-5]$ to find $\mathbf{z} = (7, -7, 0)$.

- (c) Find a general solution of the diagonal system $S\mathbf{y} = \mathbf{z}$:

$$\begin{bmatrix} 28 & 0 & 0 & 0 \\ 0 & 14 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \mathbf{y} = \begin{bmatrix} 7 \\ -7 \\ 0 \end{bmatrix}.$$

The first line gives $y_1 = 7/28 = 1/4$, the second line gives $y_2 = -7/14 = -1/2$, and the third line is $0y_3 + 0y_4 = 0$ which is satisfied for all y_3 (because we chose c correctly) and for

all y_4 . Thus $\mathbf{y} = (\frac{1}{4}, -\frac{1}{2}, y_3, y_4)$ is a general solution for this intermediary. Obtain the particular solution with $y_3 = y_4 = 0$ via

```
y=z(1:2)./diag(S(1:2,1:2))
```

(d) Finally solve $V^T \mathbf{x} = \mathbf{y}$ as $\mathbf{x} = V\mathbf{y}$, namely

$$\mathbf{x} = \begin{bmatrix} -0.5 & 0.5 & -0.1900 & -0.6811 \\ -0.5 & -0.5 & 0.6811 & -0.1900 \\ -0.5 & 0.5 & 0.1900 & 0.6811 \\ -0.5 & -0.5 & -0.6811 & 0.1900 \end{bmatrix} \begin{bmatrix} 1/4 \\ -1/2 \\ y_3 \\ y_4 \end{bmatrix}$$

Obtain a particular solution with $\mathbf{x}=V(:,1:2)*\mathbf{y}$ of $\mathbf{x} = (-3, 1, -3, 1)/8$, and then add the free components:

$$\mathbf{x} = \begin{bmatrix} -\frac{3}{8} \\ \frac{1}{8} \\ -\frac{3}{8} \\ \frac{1}{8} \end{bmatrix} + \begin{bmatrix} -0.1900 \\ 0.6811 \\ 0.1900 \\ -0.6811 \end{bmatrix} y_3 + \begin{bmatrix} -0.6811 \\ -0.1900 \\ 0.6811 \\ 0.1900 \end{bmatrix} y_4.$$

■

Procedure 3.3.13 (general solution). *Obtain a general solution of the system $A\mathbf{x} = \mathbf{b}$ using an SVD and via intermediate unknowns.*

1. Obtain an SVD factorisation $A = USV^T$.
2. Solve $U\mathbf{z} = \mathbf{b}$ by $\mathbf{z} = U^T\mathbf{b}$ (unique given U).
3. When possible, solve $S\mathbf{y} = \mathbf{z}$ as follows.¹³ Identify the non-zero and the zero singular values: suppose $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ and $\sigma_{r+1} = \dots = \sigma_{\min(m,n)} = 0$:
 - if $z_i \neq 0$ for any $i = r+1, \dots, m$, then there is no solution (the equations are **inconsistent**);
 - otherwise determine the i th component of \mathbf{y} by $y_i = z_i/\sigma_i$ for $i = 1, \dots, r$ (for which $\sigma_i > 0$), and let y_i be a free variable for $i = r+1, \dots, n$.
4. Solve $V^T \mathbf{x} = \mathbf{y}$ (unique given V and for each \mathbf{y}) to obtain a general solution as $\mathbf{x} = V\mathbf{y}$.

Proof. Consider each and every solution of $A\mathbf{x} = \mathbf{b}$:

$$\begin{aligned} A\mathbf{x} = \mathbf{b} &\iff USV^T\mathbf{x} = \mathbf{b} \quad (\text{by step 1}) \\ &\iff S(V^T\mathbf{x}) = U^T\mathbf{b} \\ &\iff S\mathbf{y} = \mathbf{z} \quad (\text{by steps 2 and 4}), \end{aligned}$$

¹³ Being diagonal, S is in a modification of row reduced echelon form (Definition 2.2.17).

and step 3 determines all possible \mathbf{y} satisfying $S\mathbf{y} = \mathbf{z}$. Hence Procedure 3.3.13 determines all possible solutions of $A\mathbf{x} = \mathbf{b}$.¹⁴ \square

This Procedure 3.3.13 determines for us that there is either none, one or an infinite number of solutions, as Theorem 2.2.22 requires.

However, Matlab/Octave's “ $\mathbf{A}\backslash$ ” gives one ‘answer’ for all of these cases. The function `rcond(A)` indicates whether the ‘answer’ is a good unique solution of $A\mathbf{x} = \mathbf{b}$ (Procedure 2.2.4). Section 3.5 addresses what the ‘answer’ by Matlab/Octave means in the other cases.

Condition number and rank determine possibilities

The expression ‘ill-conditioned’ is sometimes used merely as a term of abuse . . . It is characteristic of ill-conditioned sets of equations that small percentage errors in the coefficients given may lead to large percentage errors in the solution. *Alan Turing, 1934 (Higham 1996, p.131)*

The Matlab/Octave function `rcond()` roughly estimates the reciprocal of what is called the condition number (estimates to within a factor of two or three).

Definition 3.3.14. For any $m \times n$ matrix A , the **condition number** is the ratio of the largest to smallest of its singular values: $\text{cond } A = \sigma_1/\sigma_{\min(m,n)}$. By convention we write $\text{cond } A = \infty$ if $\sigma_{\min(m,n)} = 0$; also, for zero matrices $\text{cond } O_{m \times n} = \infty$.

Example 3.3.15. Example 3.3.6 gives the singular values of two matrices: for the 2×2 matrix the condition number $\sigma_1/\sigma_2 = (10\sqrt{2})/(5\sqrt{2}) = 2$ (`rcond = 0.5`); for the 3×3 matrix the condition number $\sigma_1/\sigma_3 = 12/3 = 4$ (`rcond = 0.25`). Example 3.3.7 comments that every $n \times n$ orthogonal matrix has singular values $\sigma_1 = \dots = \sigma_n = 1$; hence an orthogonal matrix has condition number one (`rcond = 1`). Such small condition numbers (non-small `rcond`) indicate all orthogonal matrices are “good” matrices (as classified by Procedure 2.2.4).

However, the matrix in the sports ranking Example 3.3.10 has singular values $\sigma_1 = \sigma_2 = \sqrt{3}$ and $\sigma_3 = 0$ so its condition number $\sigma_1/\sigma_3 = \sqrt{3}/0 = \infty$ which suggests the equations are likely to be unsolvable. (In Matlab/Octave, see that $\sigma_3 = 2 \cdot 10^{-17}$ so a numerical calculation would give condition number $1.7321/\sigma_3 = 7 \cdot 10^{16}$ which is effectively infinite.) ■

In practice, a condition number $> 10^8$ is effectively infinite (equivalently `rcond < 10^-8` is effectively zero, and hence called “terrible”

¹⁴ Any non-uniqueness in the orthogonal U and V just gives rise to equivalent different algebraic expressions for the set of possibilities.

by Procedure 2.2.4). The closely related important property of a matrix is the *number* of singular values that are nonzero. When applying the following definition in practical computation (e.g., Matlab/Octave), any singular values $< 10^{-8}\sigma_1$ are effectively zero.

Definition 3.3.16. *The rank of a matrix A is the number of nonzero singular values in an SVD, $A = USV^T$: letting $r = \text{rank } A$,*

$$S = \begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots & O_{r \times (n-r)} \\ 0 & \cdots & \sigma_r \\ O_{(m-r) \times r} & O_{(m-r) \times (n-r)} \end{bmatrix},$$

equivalently $S = \text{diag}_{m \times n}(\sigma_1, \sigma_2, \dots, \sigma_r, 0, \dots, 0)$.

Example 3.3.17. In the four matrices of Example 3.3.15, the respective ranks are 2, 3, n and 2. ■

Theorem 3.3.5 asserts the singular values are unique for a given matrix, so the rank of a matrix is independent of its different SVDs.

Example 3.3.18. Use Matlab/Octave to find the ranks of the two matrices

$$(a) \begin{bmatrix} 0 & 1 & 0 \\ 1 & 1 & -1 \\ 1 & 0 & -1 \\ 2 & 0 & -2 \end{bmatrix}$$

$$(b) \begin{bmatrix} 1 & -2 & -1 & 2 & 1 \\ -2 & -2 & -0 & 2 & -0 \\ -2 & -3 & 1 & -1 & 1 \\ -3 & 0 & 1 & -0 & -1 \\ 2 & 1 & 1 & 2 & -1 \end{bmatrix}$$

Solution: (a) Enter the matrix into Matlab/Octave and compute its singular values with `svd(A)`:¹⁵

```
A=[0 1 0
   1 1 -1
   1 0 -1
   2 0 -2]
svd(A)
```

The singular values are 3.49, 1.34 and $1.55 \cdot 10^{-16} \approx 0$ (2 d.p.). Since two singular values are nonzero, the rank of the matrix is two.

(b) Enter the matrix into Matlab/Octave and compute its singular values with `svd(A)`:

¹⁵ Some advanced students will know that Matlab/Octave provides the `rank()` function to directly compute the rank. However, this example is to reinforce its meaning in terms of singular values.

```
A=[1 -2 -1 2 1
-2 -2 -0 2 -0
-2 -3 1 -1 1
-3 0 1 -0 -1
2 1 1 2 -1 ]
svd(A)
```

The singular values are 5.58, 4.17, 3.13, 1.63 and $2.99 \cdot 10^{-16} \approx 0$ (2 d.p.). Since four singular values are nonzero, the rank of the matrix is four.

■

Theorem 3.3.19. *For any matrix A , let an SVD of A be USV^T , then the transpose A^T has an SVD of $V(S^T)U^T$. Further, $\text{rank}(A^T) = \text{rank } A$.*

Proof. Let $m \times n$ matrix A have SVD USV^T . Using the properties of the matrix transpose (Theorem 3.1.21),

$$A^T = (USV^T)^T = (V^T)^T S^T U^T = V(S^T)U^T$$

which is an SVD for A^T since U and V are orthogonal, and S^T has the necessary diagonal structure. Since the number of non-zero values along the diagonal of S^T is precisely the same as that of the diagonal of S , $\text{rank}(A^T) = \text{rank } A$. □

Example 3.3.20. From earlier examples, write down an SVD of the matrices

$$\begin{bmatrix} 10 & 5 \\ 2 & 11 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} -4 & -8 & 6 \\ -2 & -1 & 6 \\ 4 & -4 & 0 \end{bmatrix}.$$

Solution: These matrices are the transpose of the two matrices whose SVDs are given in Example 3.3.6. Hence their SVDs are the transpose of the SVDs in that example (remembering that the transpose of a product is the product of the transpose but in reverse order):

$$\begin{bmatrix} 10 & 5 \\ 2 & 11 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix}^T,$$

and

$$\begin{bmatrix} -4 & -8 & 6 \\ -2 & -1 & 6 \\ 4 & -4 & 0 \end{bmatrix} = \begin{bmatrix} -\frac{8}{9} & -\frac{1}{9} & -\frac{4}{9} \\ -\frac{4}{9} & \frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & \frac{4}{9} \end{bmatrix} \begin{bmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix}^T.$$

■

Let's now return to the topic of linear equations and connect new concepts to the task of solving linear equations. In particular, the following theorem addresses when a unique solution exists to a system of linear equations. Concepts developed in subsequent sections extend this theorem further.

Theorem 3.3.21 (Unique Solutions: version 1). *Let A be an $n \times n$ square matrix. The following statements are equivalent:*

- (a) A is invertible;
- (b) $A\mathbf{x} = \mathbf{b}$ has a unique solution for every \mathbf{b} in \mathbb{R}^n ;
- (c) $A\mathbf{x} = \mathbf{0}$ has only the zero solution;
- (d) all n singular values of A are nonzero;
- (e) $\text{rank } A = n$.

Proof. Prove a circular chain of implications.

3.3.21a \implies 3.3.21b : Established by Theorem 3.2.8.

3.3.21b \implies 3.3.21c : Now $\mathbf{x} = \mathbf{0}$ is always a solution of $A\mathbf{x} = \mathbf{0}$. If property 3.3.21b holds, then this is the only solution.

3.3.21c \implies 3.3.21d : Use contradiction. If any singular value is zero, then Procedure 3.3.13 finds an infinite number of solutions to the homogeneous system $A\mathbf{x} = \mathbf{0}$, which contradicts 3.3.21c.

3.3.21d \implies 3.3.21e : Property 3.3.21e is direct from Definition 3.3.16.

3.3.21e \implies 3.3.21a : Find an SVD $A = USV^T$. Then form $B = VS^{-1}U^T$ which is well defined as all singular values in diagonal S are non-zero (Theorem 3.2.21). Then $AB = USV^TVS^{-1}U^T = USS^{-1}U^T = UU^T = I_n$. Similarly $BA = I_n$. From Definition 3.2.2, A is invertible (with $B = VS^{-1}U^T$ as its inverse).

□

Practical shades of grey The Unique Solution Theorem 3.3.21 is ‘black-and-white’: either a solution exists, or it does not. But in applications, problems arise in all shades of grey. Practical issues in applications are better phrased in terms of reliability, uncertainty, and error estimates. For example, suppose in an experiment you measure quantities \mathbf{b} to three significant digits, then solve the linear equations $A\mathbf{x} = \mathbf{b}$ to estimate quantities of interest \mathbf{x} : how accurate are your estimates of the interesting quantities \mathbf{x} ? or are your estimates complete nonsense?

Optional: this discussion and theorem reinforces why we must check condition numbers in computation.

Example 3.3.22. Consider the following innocuous looking system of linear equations

$$\begin{cases} -2q + r = 3 \\ p - 5q + r = 8 \\ -3p + 2q + 3r = -5 \end{cases}$$

Solve by hand (Procedure 2.2.19) to find the unique solution is $(p, q, r) = (2, -1, 1)$.

But, and it is a big but in practical applications, what happens if the right-hand side comes from experimental measurements with a relative error of 1%? Let's explore by writing the system in matrix-vector form and using Matlab/Octave to solve with various example errors.

- (a) First solve the system as stated. Denoting the unknowns by vector $\mathbf{x} = (p, q, r)$, write the system as $A\mathbf{x} = \mathbf{b}$ for matrix

$$A = \begin{bmatrix} 0 & -2 & 1 \\ 1 & -5 & 1 \\ -3 & 2 & 3 \end{bmatrix}, \quad \text{and right-hand side } \mathbf{b} = \begin{bmatrix} 3 \\ 8 \\ -5 \end{bmatrix}.$$



Use Procedure 2.2.4 to solve the system in Matlab/Octave:

- enter the matrix and vector with

```
A=[0 -2 1; 1 -5 1; -3 2 3]
b=[3;8;-5]
```

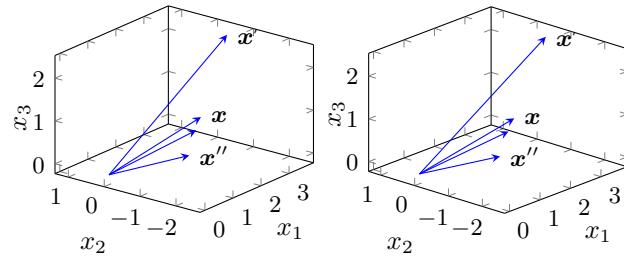
- find `rcond(A)` is 0.0031 which is poor;
- then `x=A\b` gives the solution $\mathbf{x} = (2, -1, 1)$ as before.



- (b) Now recognise that the right-hand side comes from experimental measurements with a 1% error. In Matlab/Octave, `norm(b)` computes the length $|\mathbf{b}| = 9.90$ (2 d.p.). Thus a 1% error corresponds to changing \mathbf{b} by $0.01 \times 9.90 \approx 0.1$. Let's say the first component of \mathbf{b} is in error by this amount and see what the new solution would be:

- executing `x1=A\b+[0.1;0;0]` adds the 1% error (0.1, 0, 0) to \mathbf{b} and then solves the new system to find $\mathbf{x}' = (3.7, -0.4, 2.3)$ —this solution is very different to the original solution $\mathbf{x} = (2, -1, 1)$.
- `relerr1=norm(x-x1)/norm(x)` computes its relative error $|\mathbf{x} - \mathbf{x}'|/|\mathbf{x}|$ to be 0.91.

As illustrated below, the large difference between \mathbf{x} and \mathbf{x}' indicates ‘the solution’ \mathbf{x} is almost complete nonsense. How can a 1% error in \mathbf{b} turn into the astonishingly large 91% error in solution \mathbf{x} ? Theorem 3.3.24 shows it is no accident that the magnification of the error by a factor of 91 is of the same order of magnitude as the condition number = 152.27 computed by `s=svd(A)` then `condA=s(1)/s(3)`.



- (c) To explore further, let's say the second component of \mathbf{b} is in error by 1% of \mathbf{b} , that is, by 0.1. As in the previous case, add $(0, 0.1, 0)$ to the right-hand side and solve to find now $\mathbf{x}'' = (1.2, -1.3, 0.4)$ which is quite different to both \mathbf{x} and \mathbf{x}' . Compute its relative error $|\mathbf{x} - \mathbf{x}''|/|\mathbf{x}| = 0.43$. At 43%, the relative error in solution \mathbf{x}'' is still much larger than the 1% error in \mathbf{b} .
- (d) Lastly, let's say the third component of \mathbf{b} is in error by 1% of \mathbf{b} , that is, by 0.1. As in the previous cases, add $(0, 0, 0.1)$ to the right-hand side and solve to find now $\mathbf{x}''' = (1.7, -1.1, 0.8)$ which at least is roughly \mathbf{x} . Compute its relative error $|\mathbf{x} - \mathbf{x}'''|/|\mathbf{x}| = 0.15$. At 15%, the relative error in solution \mathbf{x}''' is still significantly larger than the 1% error in \mathbf{b} .

This example shows that the apparently innocuous matrix A variously multiples measurement errors in \mathbf{b} by factors of 91, 41 or 15 when finding ‘the solution’ \mathbf{x} to $A\mathbf{x} = \mathbf{b}$. The matrix A must, after all, be a bad matrix. Theorem 3.3.24 shows this badness is quantified by its condition number 152.27, and its estimated reciprocal `rcond`. ■

Example 3.3.23. Consider solving the system of linear equations

$$\begin{bmatrix} 0.4 & 0.4 & -0.2 & 0.8 \\ -0.2 & 0.8 & -0.4 & -0.4 \\ 0.4 & -0.4 & -0.8 & -0.2 \\ -0.8 & -0.2 & -0.4 & 0.4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -3 \\ 3 \\ -9 \\ -1 \end{bmatrix}.$$

Use Matlab/Octave to explore the effect on the solution \mathbf{x} of 1% errors in the right-hand side vector.

Solution: Enter the matrix and right-hand side vector into Matlab/Octave, then solve with Procedure 2.2.4:

```
Q=[0.4 0.4 -0.2 0.8
   -0.2 0.8 -0.4 -0.4
   0.4 -0.4 -0.8 -0.2
   -0.8 -0.2 -0.4 0.4]
b=[-3;3;-9;-1]
rcond(Q)
x=Q\b
```



to find the solution $\mathbf{x} = (-4.6, 5, 7, -2.2)$.

Now see the effect on this solution of 1% errors in \mathbf{b} . Since the length $|\mathbf{b}| = \text{norm}(\mathbf{b}) = 10$ we find the solution for various changes to \mathbf{b} of size 0.1.

- For example, adding the 1% error $(0.1, 0, 0, 0)$ to \mathbf{b} , the Matlab/Octave commands

```
x1=Q\(\mathbf{b}+[0.1;0;0;0])
relerr1=norm(x-x1)/norm(x)
```

show the changed solution is $\mathbf{x}' = (-4.56, 5.04, 6.98, -2.12)$ which here is reasonably close to \mathbf{x} . Indeed its relative error $|\mathbf{x} - \mathbf{x}'|/|\mathbf{x}|$ is computed to be $0.0100 = 1\%$. Here the relative error in solution \mathbf{x} is exactly the same as the relative error in \mathbf{b} .

- Exploring further, upon adding the 1% error $(0, 0.1, 0, 0)$ to \mathbf{b} , analogous commands show the changed solution is $\mathbf{x}'' = (-4.62, 5.08, 6.96, -2.24)$ which has relative error $|\mathbf{x} - \mathbf{x}''|/|\mathbf{x}| = 0.0100 = 1\%$ again.
- Whereas, upon adding 1% error $(0, 0, 0.1, 0)$ to \mathbf{b} , analogous commands show the changed solution is $\mathbf{x}''' = (-4.56, 4.96, 6.92, -2.22)$ which has relative error $|\mathbf{x} - \mathbf{x}'''|/|\mathbf{x}| = 0.0100 = 1\%$ again.
- Lastly, upon adding 1% error $(0, 0, 0, 0.1)$ to \mathbf{b} , analogous commands show the changed solution is $\mathbf{x}'''' = (-4.68, 4.98, 6.96, -2.16)$ which has relative error $|\mathbf{x} - \mathbf{x}''''|/|\mathbf{x}| = 0.0100 = 1\%$ yet again.

In this example, and in contrast to the previous example, throughout the relative error in solution \mathbf{x} is exactly the same as the relative error in \mathbf{b} . The reason is that here the matrix Q is an orthogonal matrix—check by computing $Q'*Q$ (Definition 3.2.35). Being orthogonal, its action is to only rotate and reflect, and so it never stretches (Theorem 3.2.39c). Consequently errors remain the same size when multiplied by such orthogonal matrices, as seen in this example, and as reflected in the condition number of Q being one (as computed by $s=\text{svd}(Q)$ and then $\text{condQ}=s(1)/s(4)$). ■

The condition number determines the reliability of the solution of a system of linear equations. This is why we should always precede the computation of a solution with an estimate of the condition number such as that provided by the reciprocal `rcond()` (Procedure 2.2.4). The next theorem establishes that the condition number characterises the amplification of errors that occurs in solving a linear system. Hence solving a system of linear equations with a large condition number (small `rcond`) means that errors are amplified by a large factor as happens in Example 3.3.22.

Theorem 3.3.24 (error magnification). *Consider solving $A\mathbf{x} = \mathbf{b}$ for $n \times n$ matrix A with full rank $A = n$. Suppose the right-hand side \mathbf{b} has relative error of size ϵ , then the solution \mathbf{x} has relative error $\leq \epsilon \operatorname{cond} A$, with equality in the worst case.*

Proof. Let the length of the right-hand side vector be $b = |\mathbf{b}|$. Then the error in \mathbf{b} has size ϵb since ϵ is the relative error. Following Procedure 3.3.13, let $A = USV^T$ be an SVD for matrix A . Compute $\mathbf{z} = U^T\mathbf{b}$: recall that multiplication by orthogonal U preserves lengths (Theorem 3.2.39), so not only is $|\mathbf{z}| = b$, but also \mathbf{z} will be in error by an amount ϵb since \mathbf{b} has this error. Consider solving $S\mathbf{y} = \mathbf{z}$: the diagonals of S stretch and shrink both ‘the signal and the noise’. The *worst case* is when $\mathbf{z} = (b, 0, \dots, 0, \epsilon b)$; that is, when all the ‘signal’ happens to be in the first component of \mathbf{z} , and all the ‘noise’, the error, is in the last component. Then the intermediary $\mathbf{y} = (b/\sigma_1, 0, \dots, \epsilon b/\sigma_n)$. Consequently, the intermediary has relative error $(\epsilon b/\sigma_n)/(b/\sigma_1) = \epsilon(\sigma_1/\sigma_n) = \epsilon \operatorname{cond} A$. Again because multiplication by orthogonal V preserves lengths, the solution $\mathbf{x} = V\mathbf{y}$ has the same relative error: in the worst case of $\epsilon \operatorname{cond} A$. \square

Example 3.3.25. Each of the following cases involves solving a linear system $A\mathbf{x} = \mathbf{b}$ to determine quantities of interest \mathbf{x} from some measured quantities \mathbf{b} . From the given information estimate the maximum relative error in \mathbf{x} , if possible, otherwise say so.

- (a) Quantities \mathbf{b} are measured to a relative error 0.001, and matrix A has condition number of ten.
- (b) Quantities \mathbf{b} are measured to three significant digits and $\operatorname{rcond}(A) = 0.025$.
- (c) Measurements are accurate to two decimal places, and matrix A has condition number of twenty.
- (d) Measurements are correct to two significant digits and $\operatorname{rcond}(A) = 0.002$.

Solution: (a) The relative error in \mathbf{x} could be as big as $0.001 \times 10 = 0.01$.

- (b) Measuring to three significant digits means the relative error is 0.0005, while with $\operatorname{rcond}(A) = 0.025$, matrix A has condition number of roughly 40, so the relative error of \mathbf{x} is less than $0.0005 \times 40 = 0.02$; that is, up to 2%.
- (c) There is not enough information as we cannot determine the relative error in measurements \mathbf{b} .
- (d) Two significant digits means the relative error is 0.005, while matrix A has condition number of roughly $1/0.002 = 500$ so the relative error of \mathbf{x} could be as big as $0.005 \times 500 = 2.5$; that is, the estimate \mathbf{x} is possibly complete rubbish.



This issue of the amplification of errors occurs in other contexts. The eminent mathematician Henri Poincaré (1854–1912) was the first to detect possible chaos in the orbits of the planets.

If we knew exactly the laws of nature and the situation of the universe at the initial moment, we could predict exactly the situation of that same universe at a succeeding moment. But even if it were the case that the natural laws had no longer any secret for us, we could still only know the initial situation approximately. If that enabled us to predict the succeeding situation with the same approximation, that is all we require, and we should say that the phenomenon had been predicted, that it is governed by laws. But it is not always so; it may happen that small differences in the initial conditions produce very great ones in the final phenomena. A small error in the former will produce an enormous error in the latter. Prediction becomes impossible, and we have the fortuitous phenomenon. *Poincaré, 1903*

The analogue for us in solving linear equations such as $A\mathbf{x} = \mathbf{b}$ is the following: it may happen that a small error in the elements of \mathbf{b} will produce an enormous error in the final \mathbf{x} . The condition number warns when this happens by characterising the amplification.

3.3.3 Prove the SVD Theorem 3.3.5

When doing maths there's this great feeling. You start with a problem that just mystifies you. You can't understand it, it's so complicated, you just can't make head nor tail of it. But then when you finally resolve it, you have this incredible feeling of how beautiful it is, how it all fits together so elegantly. *Andrew Wiles, C1993*

This proof may be delayed until the last week of a semester. It may be given together with the closely related classic proof of Theorem 4.2.15 on the eigenvectors of symmetric matrices.

Two preliminary examples introduce the structure of the general proof that an SVD exists. As in this example prelude, the proof of a general singular value decomposition is similarly constructive.

Prelude to the proof

These first two examples are optional: their purpose is to introduce key parts of the general proof in a definite setting.

Example 3.3.26 (a 2×2 case). Recall Example 3.3.2 factorised the matrix

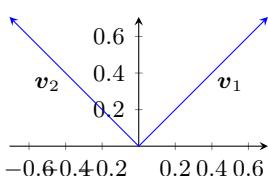
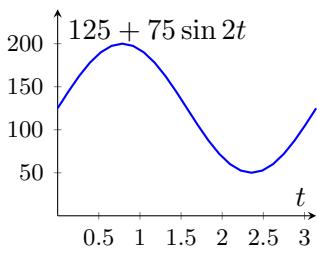
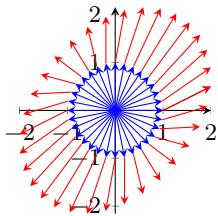
$$A = \begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix} = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T.$$

Find this factorisation, $A = USV^T$, by maximising $|Av|$ over all unit vectors \mathbf{v} .

Solution: In 2D, all unit vectors are of the form $\mathbf{v} = (\cos t, \sin t)$ for $-\pi < t \leq \pi$. The marginal picture plots these unit vectors \mathbf{v} in blue for a selection of angles t . Plotted in red from the end of each \mathbf{v} is the vector Av (scaled down by a factor of ten for clarity). Our aim is to find the \mathbf{v} that maximises the length of the corresponding adjoined Av . By inspection, the longest red vectors Av occur towards the top-right or the bottom-left, either of these directions \mathbf{v} are what we first find.

Maximising $|Av|$ is the same as maximising $|Av|^2$ which is what the following considers: since

$$\begin{aligned} Av &= \begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix} \begin{bmatrix} \cos t \\ \sin t \end{bmatrix} = \begin{bmatrix} 10 \cos t + 2 \sin t \\ 5 \cos t + 11 \sin t \end{bmatrix}, \\ |Av|^2 &= (10 \cos t + 2 \sin t)^2 + (5 \cos t + 11 \sin t)^2 \\ &= 100 \cos^2 t + 40 \cos t \sin t + 4 \sin^2 t \\ &\quad + 25 \cos^2 t + 110 \cos t \sin t + 121 \sin^2 t \\ &= 125(\cos^2 t + \sin^2 t) + 150 \sin t \cos t \\ &= 125 + 75 \sin 2t \quad (\text{shown in the margin}). \end{aligned}$$



Since $\sin \theta$ has maximum of one at angle $\theta = \frac{\pi}{2}$, the maximum of $|Av|^2$ is $125+75 = 200$ for $2t = \frac{\pi}{2}$, that is, for $t = \frac{\pi}{4}$ corresponding to unit vector $\mathbf{v}_1 = (\cos \frac{\pi}{4}, \sin \frac{\pi}{4}) = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ —this vector points to the top-right as identified from the marginal figure. This vector is the first column of V .

Now multiply to find $Av_1 = (6\sqrt{2}, 8\sqrt{2})$. The length of this vector is $\sqrt{72 + 128} = \sqrt{200} = 10\sqrt{2} = \sigma_1$ the leading singular value. Normalise the vector Av_1 by $Av_1/\sigma_1 = (6\sqrt{2}, 8\sqrt{2})/(10\sqrt{2}) = (\frac{3}{5}, \frac{4}{5}) = \mathbf{u}_1$, the first column of U .

The other column of V must be orthogonal (at right-angles) to \mathbf{v}_1 in order for matrix V to be orthogonal. Thus set $\mathbf{v}_2 = (-\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$ as shown in the marginal graph. Now multiply to find $Av_2 = (-4\sqrt{2}, 3\sqrt{2})$: magically, and a crucial part of the general proof, this vector is orthogonal to \mathbf{u}_1 . The length of $Av_2 = (-4\sqrt{2}, 3\sqrt{2})$ is $\sqrt{32 + 18} = \sqrt{50} = 5\sqrt{2} = \sigma_2$ the other singular value. Normalise the vector to $Av_2/\sigma_2 = (-4\sqrt{2}, 3\sqrt{2})/(2\sqrt{2}) = (-\frac{4}{5}, \frac{3}{5}) = \mathbf{u}_2$, the second column of U .

This construction establishes that here $AV = US$; that is, an SVD is $A = USV^T$.

In this example we could have chosen the negative of \mathbf{v}_1 (angle $t = -\frac{3\pi}{4}$), and/or chosen the negative of \mathbf{v}_2 . The result would still be a valid SVD of the matrix A . The orthogonal matrices in an SVD are not unique, and need not be; but the singular values are unique.

■

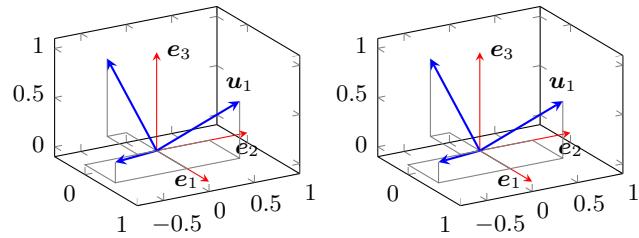
Example 3.3.27 (a 3×1 case). Find the following SVD:

$$A = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{3}} & \cdot & \cdot \\ \frac{1}{\sqrt{3}} & \cdot & \cdot \\ \frac{1}{\sqrt{3}} & \cdot & \cdot \end{bmatrix} \begin{bmatrix} \sqrt{3} \\ 0 \\ 0 \end{bmatrix} [1]^T = USV^T,$$

where we do not worry about the elements denoted by dots as they are multiplied by the zeros in $S = (\sqrt{3}, 0, 0)$.

Solution: We seek to maximise $|A\mathbf{v}|^2$ but here vector \mathbf{v} is in \mathbb{R}^1 . Being of unit magnitude, there are two alternatives: $\mathbf{v} = (\pm 1)$. Each alternative gives the same $|A\mathbf{v}|^2 = |(\pm 1, \pm 1, \pm 1)| = 3$. Choose one alternative, say $\mathbf{v}_1 = (1)$, fixes the matrix $V = [1]$.

Then $A\mathbf{v}_1 = (1, 1, 1)$ which is of length $\sqrt{3}$. This length is the singular value $\sigma_1 = \sqrt{3}$. Dividing $A\mathbf{v}_1$ by its length gives the unit vector $\mathbf{u}_1 = (\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}})$, the first column of U . To find the other columns of U , consider the three standard unit vectors in \mathbb{R}^3 (red in the illustration below), rotate them all together so that one lines up with \mathbf{u}_1 , and then the other two rotated unit vectors form the other two columns of U (blue vectors below). Since the columns of U are then orthonormal, U is an orthogonal matrix (Theorem 3.2.39).



Outline of the general proof We use induction on the size $m \times n$ of the matrix.

- First zero matrices have trivial SVD, and $m \times 1$ and $1 \times n$ matrices have straightforward SVD.
- Choose \mathbf{v}_1 to maximise $|A\mathbf{v}|^2$ for unit vectors \mathbf{v} in \mathbb{R}^n .
- Crucially, we establish that $\mathbf{v} \perp \mathbf{v}_1 \implies A\mathbf{v} \perp A\mathbf{v}_1$.

- Then rotate the standard unit vectors to align one with \mathbf{v}_1 . Similarly for $A\mathbf{v}_1$.
- This rotation transforms the matrix A to strip off the leading singular value, and effectively leave an $(m-1) \times (n-1)$ matrix.
- By induction on the size, an SVD exists for all sizes.

This proof corresponds closely to the proof of the spectral theorem 4.2.15 for symmetric matrices of section 4.2.

Detailed proof of the SVD Theorem 3.3.5

Use induction on the size $m \times n$ of the matrix A : we assume an SVD exists for all $(m-1) \times (n-1)$ matrices, and prove that consequently an SVD must exist for all $m \times n$ matrices. There are three base cases to establish: one for $m \leq n$, one for $m \geq n$, and one for matrix $A = O$; then the induction extends to all sized matrices.

Case $A = O_{m \times n}$: When $m \times n$ matrix $A = O_{m \times n}$ then choose $U = I_m$ (orthogonal), $S = O_{m \times n}$ (diagonal), and $V = I_n$ (orthogonal) so then $USV^T = I_m O_{m \times n} I_n^T = O_{m \times n} = A$.

Consequently, the rest of the proof only considers the non-trivial cases when the matrix A is not all zero.

Case $m \times 1$ ($n = 1$): Here the $m \times 1$ nonzero matrix $A = [\mathbf{a}_1]$ for $\mathbf{a}_1 = (a_{11}, a_{21}, \dots, a_{m1})$. Set the singular value $\sigma_1 = |\mathbf{a}_1| = \sqrt{a_{11}^2 + a_{21}^2 + \dots + a_{m1}^2}$ and unit vector $\mathbf{u}_1 = \mathbf{a}_1 / \sigma_1$. Set 1×1 orthogonal matrix $V = [1]$; $m \times 1$ diagonal matrix $S = (\sigma_1, 0, \dots, 0)$; and $m \times m$ orthogonal matrix $U = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_m]$. Matrix U exists because we can take the orthonormal set of standard unit vectors in \mathbb{R}^m and rotate them all together so that the first lines up with \mathbf{u}_1 : the other $(m-1)$ unit vectors then become the other \mathbf{u}_j . Then an SVD for the $m \times 1$ matrix A is

$$\begin{aligned} USV^T &= [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} 1^T \\ &= \sigma_1 \mathbf{u}_1 = [\mathbf{a}_1] = A. \end{aligned}$$

Case $1 \times n$ ($m = 1$): use an exactly complementary argument to the case $m \geq n = 1$.

Induction Assume an SVD exists for all $(m-1) \times (n-1)$ matrices: prove that consequently an SVD must exist for all $m \times n$ matrices. Consider any $m \times n$ nonzero matrix A for $m, n \geq 2$. Set vector \mathbf{v}_1 in \mathbb{R}^n to be a *unit vector* that maximises $|A\mathbf{v}|^2$ for unit vectors \mathbf{v} in \mathbb{R}^n ; that is, vector \mathbf{v}_1 achieves the maximum in $\max_{|\mathbf{v}|=1} |A\mathbf{v}|^2$.

- Such a maximum exists by the Extreme Value Theorem in Calculus. Proved in higher level analysis.

As matrix A is nonzero, there exists \mathbf{v} such that $|A\mathbf{v}| > 0$. Since \mathbf{v}_1 maximises $|A\mathbf{v}|$ it follows that $|A\mathbf{v}_1| > 0$.

The vector \mathbf{v}_1 is not unique: for example, the negative $-\mathbf{v}_1$ is another unit vector that achieves the maximum value. Sometimes there are other unit vectors that achieve the maximum value. Choose any one of them.

Nonetheless, the maximum value of $|A\mathbf{v}|^2$ is unique, and so the following singular value σ_1 is unique.

- Set the singular value $\sigma_1 := |A\mathbf{v}_1| > 0$ and unit vector $\mathbf{u}_1 := (A\mathbf{v}_1)/\sigma_1$ in \mathbb{R}^m . Let \mathbf{v} be any unit vector orthogonal to \mathbf{v}_1 : we prove that then the vector $A\mathbf{v}$ is orthogonal to \mathbf{u}_1 . Let $\mathbf{u} := A\mathbf{v}$ in \mathbb{R}^m and consider $f(t) := |A(\mathbf{v}_1 \cos t + \mathbf{v} \sin t)|^2$. Since \mathbf{v}_1 achieves the maximum, and $\mathbf{v}_1 \cos t + \mathbf{v} \sin t$ is a unit vector for all t (Exercise 3.3.16), then $f(t)$ must have a maximum at $t = 0$ (maybe at other t as well), and so $f'(0) = 0$ (from the Calculus of a maximum). On the other hand,

$$\begin{aligned} f(t) &= |A\mathbf{v}_1 \cos t + A\mathbf{v} \sin t|^2 \\ &= |\sigma_1 \mathbf{u}_1 \cos t + \mathbf{u} \sin t|^2 \\ &= (\sigma_1 \mathbf{u}_1 \cos t + \mathbf{u} \sin t) \cdot (\sigma_1 \mathbf{u}_1 \cos t + \mathbf{u} \sin t) \\ &= \sigma_1^2 \cos^2 t + \sigma_1 \mathbf{u} \cdot \mathbf{u}_1 2 \sin t \cos t + |\mathbf{u}|^2 \sin^2 t; \end{aligned}$$

differentiating $f(t)$ and evaluating at zero gives $0 = f'(0) = \sigma_1 \mathbf{u} \cdot \mathbf{u}_1$. Since the singular value $\sigma_1 > 0$, we must have $\mathbf{u} \cdot \mathbf{u}_1 = 0$ and so \mathbf{u}_1 and \mathbf{u} are orthogonal (Definition 1.3.15).

- Consider the orthonormal set of standard unit vectors in \mathbb{R}^n : rotate them so that the first unit vector lines up with \mathbf{v}_1 , and let the other $(n - 1)$ rotated unit vectors become the columns of the $n \times (n - 1)$ matrix \bar{V} . Then set the $n \times n$ matrix $V_1 := [\mathbf{v}_1 \ \bar{V}]$ which is orthogonal as its columns are orthonormal (Theorem 3.2.39b). Similarly set an $m \times m$ orthogonal matrix $U_1 := [\mathbf{u}_1 \ \bar{U}]$. Compute the $m \times n$ matrix

$$\begin{aligned} A_1 := U_1^T A V_1 &= \begin{bmatrix} \mathbf{u}_1^T \\ \bar{U}^T \end{bmatrix} A [\mathbf{v}_1 \ \bar{V}] \\ &= \begin{bmatrix} \mathbf{u}_1^T A \mathbf{v}_1 & \mathbf{u}_1^T A \bar{V} \\ \bar{U}^T A \mathbf{v}_1 & \bar{U}^T A \bar{V} \end{bmatrix} \end{aligned}$$

where

- the top-left entry $\mathbf{u}_1^T A \mathbf{v}_1 = \mathbf{u}_1^T \sigma_1 \mathbf{u}_1 = \sigma_1 |\mathbf{u}_1|^2 = \sigma_1$,
- the bottom-left column $\bar{U}^T A \mathbf{v}_1 = \bar{U}^T \sigma_1 \mathbf{u}_1 = O_{m-1 \times 1}$ as the columns of \bar{U} are orthogonal to \mathbf{u}_1 ,

- the top-right row $\mathbf{u}_1^T A \bar{V} = O_{1 \times n-1}$ as each column of \bar{V} is orthogonal to \mathbf{v}_1 and hence each column of $A \bar{V}$ is orthogonal to \mathbf{u}_1 ,
- and set the bottom-right block $B := \bar{U}^T A \bar{V}$ which is an $(m-1) \times (n-1)$ matrix as \bar{U}^T is $(m-1) \times m$ and \bar{V} is $n \times (n-1)$.

Consequently,

$$A_1 = \begin{bmatrix} \sigma_1 & O_{1 \times n-1} \\ O_{m-1 \times 1} & B \end{bmatrix}.$$

Note: rearranging $A_1 := U_1^T A V_1$ gives $A V_1 = U_1 A_1$.

4. *By induction assumption, $(m-1) \times (n-1)$ matrix B has an SVD, and so we now construct an SVD for $m \times n$ matrix A . Let $B = \hat{U} \hat{S} \hat{V}^T$ be an SVD for B . Then construct*

$$U := U_1 \begin{bmatrix} 1 & 0 \\ 0 & \hat{U} \end{bmatrix}, \quad V := V_1 \begin{bmatrix} 1 & 0 \\ 0 & \hat{V} \end{bmatrix}, \quad S := \begin{bmatrix} \sigma_1 & 0 \\ 0 & \hat{S} \end{bmatrix}.$$

Matrices U and V are orthogonal as each are the product of two orthogonal matrices (Exercise 3.2.19), and matrix S is diagonal. These form an SVD for matrix A since

$$\begin{aligned} AV &= AV_1 \begin{bmatrix} 1 & 0 \\ 0 & \hat{V} \end{bmatrix} = U_1 A_1 \begin{bmatrix} 1 & 0 \\ 0 & \hat{V} \end{bmatrix} \\ &= U_1 \begin{bmatrix} \sigma_1 & 0 \\ 0 & B \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \hat{V} \end{bmatrix} = U_1 \begin{bmatrix} \sigma_1 & 0 \\ 0 & B \hat{V} \end{bmatrix} \\ &= U_1 \begin{bmatrix} \sigma_1 & 0 \\ 0 & \hat{U} \hat{S} \end{bmatrix} = U_1 \begin{bmatrix} 1 & 0 \\ 0 & \hat{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 \\ 0 & \hat{S} \end{bmatrix} \\ &= US. \end{aligned}$$

Hence $A = USV^T$. By induction, an SVD exists for all $m \times n$ matrices.

This argument establishes the SVD Theorem 3.3.5.

3.3.4 Exercises

Exercise 3.3.1. Using a factorisation of the left-hand side coefficient, quickly solve by hand the following equations.

- | | |
|--------------------|--------------------|
| (a) $18x = 1134$ | (b) $42x = 2226$ |
| (c) $66x = 3234$ | (d) $70x = 3150$ |
| (e) $99x = 8118$ | (f) $154x = 7854$ |
| (g) $175x = 14350$ | (h) $242x = 20086$ |
| (i) $245x = 12495$ | (j) $363x = 25047$ |
| (k) $385x = 15785$ | (l) $539x = 28028$ |

Exercise 3.3.2. Find a general solution, if a solution exists, of each of the following systems of linear equations using Procedure 3.3.13. Calculate by hand using the given SVD factorisation; record your working.

$$(a) \underbrace{\begin{bmatrix} -\frac{9}{5} & \frac{12}{5} \\ -4 & -3 \end{bmatrix}}_{=A} \mathbf{x} = \begin{bmatrix} -\frac{9}{5} \\ \frac{17}{2} \end{bmatrix} \text{ given the SVD}$$

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 5 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} -\frac{4}{5} & -\frac{3}{5} \\ -\frac{3}{5} & \frac{4}{5} \end{bmatrix}^T$$

$$(b) \underbrace{\begin{bmatrix} \frac{15}{13} & \frac{36}{13} \\ \frac{36}{13} & -\frac{15}{13} \end{bmatrix}}_{=B} \mathbf{x} = \begin{bmatrix} \frac{54}{13} \\ -\frac{45}{26} \end{bmatrix} \text{ given the SVD}$$

$$B = \begin{bmatrix} -\frac{12}{13} & \frac{5}{13} \\ \frac{5}{13} & \frac{12}{13} \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}^T$$

$$(c) \underbrace{\begin{bmatrix} -0.96 & 1.28 \\ -0.72 & 0.96 \end{bmatrix}}_{=C} \mathbf{x} = \begin{bmatrix} 2.88 \\ 2.16 \end{bmatrix} \text{ given the SVD}$$

$$C = \begin{bmatrix} \frac{4}{5} & \frac{3}{5} \\ \frac{3}{5} & -\frac{4}{5} \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} -\frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & -\frac{3}{5} \end{bmatrix}^T$$

$$(d) \underbrace{\begin{bmatrix} -\frac{5}{26} & -\frac{6}{13} \\ -\frac{12}{13} & \frac{5}{13} \end{bmatrix}}_{=D} \mathbf{x} = \begin{bmatrix} -\frac{7}{13} \\ \frac{34}{13} \end{bmatrix} \text{ given the SVD}$$

$$D = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} -\frac{12}{13} & -\frac{5}{13} \\ \frac{5}{13} & -\frac{12}{13} \end{bmatrix}^T$$

$$(e) \underbrace{\begin{bmatrix} -\frac{2}{3} & \frac{23}{51} & \frac{22}{51} \\ \frac{1}{6} & \frac{7}{51} & -\frac{31}{51} \end{bmatrix}}_{=E} \mathbf{x} = \begin{bmatrix} -\frac{115}{102} \\ -\frac{35}{102} \end{bmatrix} \text{ given the SVD}$$

$$E = \begin{bmatrix} \frac{15}{17} & -\frac{8}{17} \\ -\frac{8}{17} & -\frac{15}{17} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} -\frac{2}{3} & \frac{1}{3} & -\frac{2}{3} \\ \frac{1}{3} & -\frac{2}{3} & -\frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & -\frac{1}{3} \end{bmatrix}^T$$

$$(f) \underbrace{\begin{bmatrix} \frac{3}{35} & \frac{9}{35} & \frac{9}{70} \\ -\frac{4}{35} & -\frac{12}{35} & -\frac{6}{35} \end{bmatrix}}_{=F} \mathbf{x} = \begin{bmatrix} \frac{3}{8} \\ -\frac{1}{2} \end{bmatrix} \text{ given the SVD}$$

$$F = \begin{bmatrix} -\frac{3}{5} & \frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -\frac{2}{7} & -\frac{6}{7} & \frac{3}{7} \\ -\frac{6}{7} & \frac{3}{7} & \frac{2}{7} \\ -\frac{3}{7} & -\frac{2}{7} & -\frac{6}{7} \end{bmatrix}^T$$

$$(g) \underbrace{\begin{bmatrix} \frac{7}{39} & -\frac{17}{39} \\ -\frac{22}{39} & -\frac{19}{39} \\ -\frac{4}{39} & -\frac{53}{78} \end{bmatrix}}_{=G} \mathbf{x} = \begin{bmatrix} -\frac{1}{3} \\ -\frac{2}{3} \\ -\frac{2}{3} \end{bmatrix} \text{ given the SVD}$$

$$G = \begin{bmatrix} -\frac{1}{3} & \frac{2}{3} & \frac{2}{3} \\ -\frac{2}{3} & -\frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & -\frac{2}{3} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{2} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{5}{13} & \frac{12}{13} \\ \frac{12}{13} & -\frac{5}{13} \end{bmatrix}^T$$

$$(h) \underbrace{\begin{bmatrix} \frac{36}{119} & -\frac{11}{17} \\ \frac{164}{119} & -\frac{18}{17} \\ -\frac{138}{119} & -\frac{6}{17} \end{bmatrix}}_{=H} \mathbf{x} = \begin{bmatrix} \frac{11}{17} \\ \frac{9}{17} \\ \frac{3}{17} \end{bmatrix} \text{ given the SVD}$$

$$H = \begin{bmatrix} \frac{2}{7} & -\frac{3}{7} & -\frac{6}{7} \\ \frac{6}{7} & -\frac{2}{7} & \frac{3}{7} \\ -\frac{3}{7} & -\frac{6}{7} & \frac{2}{7} \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{15}{17} & \frac{8}{17} \\ -\frac{8}{17} & \frac{15}{17} \end{bmatrix}^T$$

$$(i) \underbrace{\begin{bmatrix} -\frac{17}{18} & -\frac{8}{9} & -\frac{8}{9} \\ 1 & \frac{2}{3} & -\frac{2}{3} \\ -\frac{11}{9} & \frac{8}{9} & -\frac{7}{9} \end{bmatrix}}_{=\mathcal{I}} \mathbf{x} = \begin{bmatrix} -\frac{17}{18} \\ \frac{5}{3} \\ -\frac{7}{18} \end{bmatrix} \text{ given the SVD}$$

$$\mathcal{I} = \begin{bmatrix} -\frac{2}{3} & -\frac{1}{3} & -\frac{2}{3} \\ \frac{1}{3} & \frac{2}{3} & -\frac{2}{3} \\ -\frac{2}{3} & \frac{2}{3} & \frac{1}{3} \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{3}{2} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{8}{9} & \frac{1}{9} & -\frac{4}{9} \\ \frac{1}{9} & \frac{8}{9} & \frac{4}{9} \\ \frac{4}{9} & -\frac{4}{9} & \frac{7}{9} \end{bmatrix}^T$$

$$(j) \underbrace{\begin{bmatrix} -\frac{10}{27} & -\frac{2}{27} & \frac{31}{54} \\ -\frac{4}{27} & \frac{10}{27} & \frac{17}{27} \\ -\frac{8}{27} & -\frac{7}{27} & \frac{7}{27} \end{bmatrix}}_{=J} \mathbf{x} = \begin{bmatrix} \frac{83}{54} \\ \frac{49}{54} \\ \frac{17}{54} \end{bmatrix} \text{ given the SVD}$$

$$J = \begin{bmatrix} -\frac{2}{3} & \frac{1}{3} & -\frac{2}{3} \\ -\frac{2}{3} & -\frac{2}{3} & \frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{4}{9} & -\frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & -\frac{4}{9} \\ -\frac{8}{9} & -\frac{1}{9} & \frac{4}{9} \end{bmatrix}^T$$

$$(k) \underbrace{\begin{bmatrix} \frac{4}{33} & \frac{4}{11} & \frac{6}{11} \\ \frac{4}{33} & \frac{4}{11} & \frac{6}{11} \\ \frac{2}{33} & \frac{2}{11} & \frac{3}{11} \end{bmatrix}}_{=K} \mathbf{x} = \begin{bmatrix} -\frac{7}{3} \\ -\frac{7}{3} \\ -\frac{7}{6} \end{bmatrix} \text{ given the SVD}$$

$$K = \begin{bmatrix} \frac{2}{3} & -\frac{2}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{1}{3} & -\frac{2}{3} \\ \frac{1}{3} & \frac{2}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{2}{11} & -\frac{9}{11} & \frac{6}{11} \\ \frac{6}{11} & \frac{6}{11} & \frac{7}{11} \\ \frac{9}{11} & -\frac{2}{11} & -\frac{6}{11} \end{bmatrix}^T$$

$$(l) \underbrace{\begin{bmatrix} -\frac{6}{11} & -\frac{1}{11} & \frac{81}{22} \\ \frac{7}{11} & \frac{3}{11} & \frac{27}{11} \\ -\frac{6}{11} & \frac{9}{22} & -\frac{9}{11} \end{bmatrix}}_{=L} \mathbf{x} = \begin{bmatrix} -\frac{35}{2} \\ -\frac{41}{4} \\ \frac{15}{8} \end{bmatrix} \text{ given the SVD}$$

$$L = \begin{bmatrix} \frac{9}{11} & \frac{6}{11} & \frac{2}{11} \\ \frac{6}{11} & -\frac{7}{11} & -\frac{6}{11} \\ -\frac{2}{11} & \frac{6}{11} & -\frac{9}{11} \end{bmatrix} \begin{bmatrix} \frac{9}{2} & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & -1 \\ 1 & 0 & 0 \end{bmatrix}^T$$

Exercise 3.3.3. Find a general solution, if a solution exists, of each of the following systems of linear equations. Calculate by hand using the given SVD factorisation; check the SVD and confirm your calculations with Matlab/Octave (Procedure 3.3.13); then compare and contrast the two methods.



$$(a) \underbrace{\begin{bmatrix} \frac{7}{180} & \frac{8}{45} & \frac{41}{180} \\ \frac{19}{180} & -\frac{22}{45} & \frac{101}{180} \\ -\frac{19}{180} & \frac{4}{45} & \frac{133}{180} \\ \frac{59}{60} & -\frac{2}{15} & \frac{91}{60} \end{bmatrix}}_{=A} \mathbf{x} = \begin{bmatrix} -\frac{13}{40} \\ -\frac{9}{8} \\ -\frac{17}{20} \\ -\frac{15}{4} \end{bmatrix} \text{ given the SVD}$$

$$A = \begin{bmatrix} -\frac{1}{10} & \frac{3}{10} & \frac{3}{10} & -\frac{9}{10} \\ -\frac{3}{10} & -\frac{9}{10} & -\frac{1}{10} & -\frac{3}{10} \\ -\frac{3}{10} & -\frac{1}{10} & \frac{9}{10} & \frac{3}{10} \\ -\frac{9}{10} & \frac{3}{10} & -\frac{3}{10} & \frac{1}{10} \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -\frac{4}{9} & \frac{4}{9} & -\frac{7}{9} \\ \frac{1}{9} & \frac{8}{9} & \frac{4}{9} \\ -\frac{8}{9} & -\frac{1}{9} & \frac{4}{9} \end{bmatrix}^T$$



$$(b) \underbrace{\begin{bmatrix} \frac{79}{66} & \frac{7}{33} & -\frac{29}{33} \\ -\frac{65}{66} & -\frac{13}{33} & \frac{35}{33} \\ \frac{31}{66} & -\frac{29}{33} & -\frac{37}{33} \\ \frac{17}{66} & -\frac{23}{33} & -\frac{43}{33} \end{bmatrix}}_{=B} \mathbf{x} = \begin{bmatrix} -\frac{22}{6} \\ \frac{20}{6} \\ -\frac{1}{6} \\ \frac{1}{6} \end{bmatrix} \text{ given the SVD}$$

$$B = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} \frac{8}{3} & 0 & 0 \\ 0 & \frac{4}{3} & 0 \\ 0 & 0 & \frac{1}{3} \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -\frac{6}{11} & -\frac{6}{11} & -\frac{7}{11} \\ \frac{2}{11} & -\frac{9}{11} & \frac{6}{11} \\ \frac{9}{11} & -\frac{2}{11} & -\frac{6}{11} \end{bmatrix}^T$$



(c) $\underbrace{\begin{bmatrix} \frac{14}{15} & -\frac{14}{15} & \frac{7}{15} & -\frac{28}{15} \\ \frac{2}{5} & -\frac{8}{5} & -\frac{4}{5} & \frac{4}{5} \\ -\frac{6}{5} & -\frac{6}{5} & \frac{12}{5} & \frac{3}{5} \end{bmatrix}}_{=C} \mathbf{x} = \begin{bmatrix} 0 \\ 4 \\ 27 \end{bmatrix}$ given the SVD

$$C = \begin{bmatrix} 0 & -1 & 0 \\ 0 & 0 & -1 \\ -1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 & 0 \\ 0 & \frac{7}{3} & 0 & 0 \\ 0 & 0 & 2 & 0 \end{bmatrix} \begin{bmatrix} \frac{2}{5} & -\frac{2}{5} & -\frac{1}{5} & \frac{4}{5} \\ \frac{2}{5} & \frac{2}{5} & \frac{4}{5} & \frac{1}{5} \\ -\frac{4}{5} & -\frac{1}{5} & \frac{2}{5} & \frac{2}{5} \\ -\frac{1}{5} & \frac{4}{5} & -\frac{2}{5} & \frac{2}{5} \end{bmatrix}^T$$



(d) $\underbrace{\begin{bmatrix} \frac{57}{22} & -\frac{3}{22} & -\frac{45}{22} & -\frac{9}{22} \\ -\frac{14}{11} & \frac{32}{11} & -\frac{4}{11} & -\frac{14}{11} \\ -\frac{9}{22} & -\frac{3}{22} & -\frac{45}{22} & \frac{57}{22} \end{bmatrix}}_{=D} \mathbf{x} = \begin{bmatrix} 117 \\ -72 \\ 63 \end{bmatrix}$ given the SVD

$$D = \begin{bmatrix} -\frac{6}{11} & \frac{9}{11} & \frac{2}{11} \\ \frac{7}{11} & \frac{6}{11} & -\frac{6}{11} \\ -\frac{6}{11} & -\frac{2}{11} & -\frac{9}{11} \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 3 & 0 \end{bmatrix} \begin{bmatrix} -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \end{bmatrix}^T$$



(e) $\underbrace{\begin{bmatrix} -\frac{2}{5} & -\frac{2}{5} & -\frac{26}{45} & \frac{26}{45} \\ \frac{11}{9} & \frac{11}{9} & -\frac{1}{3} & \frac{1}{3} \\ \frac{31}{90} & \frac{31}{90} & \frac{17}{90} & -\frac{17}{90} \\ \frac{4}{9} & \frac{4}{9} & -\frac{2}{9} & \frac{2}{9} \end{bmatrix}}_{=E} \mathbf{x} = \begin{bmatrix} 3 \\ -6 \\ -3 \\ -2 \end{bmatrix}$ given the SVD

$$E = \begin{bmatrix} \frac{2}{9} & \frac{8}{9} & \frac{1}{3} & \frac{2}{9} \\ -\frac{8}{9} & \frac{2}{9} & -\frac{2}{9} & \frac{1}{3} \\ -\frac{2}{9} & -\frac{1}{3} & \frac{8}{9} & \frac{2}{9} \\ -\frac{1}{3} & \frac{2}{9} & \frac{2}{9} & -\frac{8}{9} \end{bmatrix} \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -\frac{7}{10} & -\frac{1}{10} & \frac{1}{10} & -\frac{7}{10} \\ -\frac{7}{10} & -\frac{1}{10} & -\frac{1}{10} & \frac{7}{10} \\ \frac{1}{10} & -\frac{7}{10} & -\frac{7}{10} & -\frac{1}{10} \\ -\frac{1}{10} & \frac{7}{10} & -\frac{7}{10} & -\frac{1}{10} \end{bmatrix}^T$$



(f) $\underbrace{\begin{bmatrix} \frac{5}{14} & -\frac{41}{14} & \frac{1}{2} & -\frac{3}{2} \\ \frac{22}{7} & -\frac{4}{7} & 0 & 2 \\ -\frac{9}{14} & -\frac{13}{14} & \frac{5}{2} & -\frac{1}{2} \\ -\frac{2}{7} & -\frac{6}{7} & 2 & 2 \end{bmatrix}}_{=F} \mathbf{x} = \begin{bmatrix} -45 \\ -50 \\ 18 \\ 20 \end{bmatrix}$ given the SVD

$$F = \begin{bmatrix} \frac{4}{7} & \frac{2}{7} & -\frac{5}{7} & \frac{2}{7} \\ -\frac{4}{7} & \frac{5}{7} & -\frac{2}{7} & -\frac{2}{7} \\ \frac{4}{7} & \frac{2}{7} & \frac{2}{7} & -\frac{5}{7} \\ \frac{1}{7} & \frac{4}{7} & \frac{4}{7} & \frac{4}{7} \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{bmatrix}^T$$

Exercise 3.3.4. Find a general solution, if possible, of each of the following systems of linear equations with Matlab/Octave and using Procedure 3.3.13.



$$(a) \begin{bmatrix} 2.4 & 1.6 & 1 & -0.8 \\ -1.2 & 3.2 & -2 & -0.4 \\ -1.2 & -0.8 & 2 & -1.6 \\ 0.6 & -1.6 & -4 & -0.8 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -29.4 \\ -12.4 \\ 13.2 \\ -0.8 \end{bmatrix}$$



$$(b) \begin{bmatrix} -0.7 & -0.7 & -2.5 & -0.7 \\ 1 & -2.2 & 0.2 & -0.2 \\ -1 & 1.4 & -1.4 & -2.6 \\ 2.6 & -1.4 & -1 & -1.4 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -4 \\ 2.4 \\ 3.2 \\ 0 \end{bmatrix}$$



$$(c) \begin{bmatrix} -3.14 & -1.18 & 0.46 & -0.58 \\ 0.66 & 0.18 & -0.06 & 2.22 \\ -1.78 & -2.54 & -1.82 & -5.26 \\ 0.58 & 1.06 & -0.82 & 0.26 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -17.38 \\ -1.14 \\ 5.22 \\ 12.26 \end{bmatrix}$$



$$(d) \begin{bmatrix} 1.38 & 0.50 & 3.30 & 0.34 \\ -0.66 & -0.70 & 1.50 & -2.38 \\ -0.90 & 2.78 & -0.54 & 0.10 \\ 0.00 & 1.04 & -0.72 & -1.60 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -7.64 \\ -7.72 \\ -20.72 \\ -20.56 \end{bmatrix}$$



$$(e) \begin{bmatrix} 1.32 & 1.40 & 1.24 & -0.20 \\ 1.24 & 3.00 & 2.68 & 1.00 \\ 1.90 & -1.06 & -1.70 & 2.58 \\ -1.30 & 0.58 & 0.90 & -0.94 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -5.28 \\ 2.04 \\ 6.30 \\ 2.50 \end{bmatrix}$$



$$(f) \begin{bmatrix} 2.16 & 0.82 & -2.06 & 0.72 \\ -0.18 & -0.56 & 1.84 & -0.78 \\ 1.68 & -0.14 & 0.02 & -0.24 \\ -1.14 & -0.88 & -2.48 & 0.66 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -12.6 \\ 13.8 \\ 0.2 \\ -32.6 \end{bmatrix}$$



$$(g) \begin{bmatrix} 0.00 & -0.54 & -0.72 & 0.90 \\ 0.40 & 0.74 & 0.32 & -0.10 \\ 1.20 & 2.22 & 0.96 & -0.30 \\ -0.00 & -0.18 & -0.24 & 0.30 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -1.8 \\ -3.2 \\ -9.6 \\ -0.6 \end{bmatrix}$$



$$(h) \begin{bmatrix} 7 & 1 & -1 & 4 \\ 2 & 4 & -4 & 0 \\ 0 & 4 & 0 & -1 \\ -4 & 1 & 1 & -1 \\ -1 & 0 & -1 & 3 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 22.4 \\ 11.2 \\ -6.1 \\ -8.3 \\ 17.8 \end{bmatrix}$$



$$(i) \begin{bmatrix} 7 & 1 & -1 & 4 \\ 2 & 4 & -4 & 0 \\ 0 & 4 & 0 & -1 \\ -4 & 1 & 1 & -1 \\ -1 & 0 & -1 & 3 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2.1 \\ 2.2 \\ 4.6 \\ -0.7 \\ 5.5 \end{bmatrix}$$



$$(j) \begin{bmatrix} -1 & 0 & -6 & 0 & 5 \\ 0 & -3 & 2 & 1 & 7 \\ 0 & 2 & -3 & -2 & 2 \\ 0 & -3 & 7 & -5 & 0 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 30.7 \\ -17.0 \\ 21.3 \\ -45.7 \end{bmatrix}$$



$$(k) \begin{bmatrix} 1 & 6 & 1 & 1 & -4 \\ 3 & -2 & 0 & -4 & 7 \\ 1 & -3 & -1 & -5 & -2 \\ -1 & 4 & -2 & -1 & -2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 4 \\ -7 \\ 2 \\ -3 \end{bmatrix}$$

Exercise 3.3.5. Recall Theorems 2.2.22 and 2.2.25 on the existence of none, one, or an infinite number of solutions to linear equations. Use procedure 3.3.13 to provide an alternative proof to each of these two theorems.

Exercise 3.3.6. Write down the condition number and the rank of each of the matrices A, \dots, L in Exercise 3.3.2 using the given SVDS.

Exercise 3.3.7. Write down the condition number and the rank of each of the matrices A, \dots, F in Exercise 3.3.3 using the given SVDS. For each square matrix, compute `rcond` and comment on its relation to the condition number.

Exercise 3.3.8. In Matlab/Octave, use `randn()` to generate some random matrices A of chosen sizes, and some correspondingly sized random right-hand side vectors \mathbf{b} . For each, find a general solution, if possible, of the system $A\mathbf{x} = \mathbf{b}$ with Matlab/Octave and using Procedure 3.3.13. Record each step, the condition number and rank of A , and comment on what is interesting about the sizes you choose.

Exercise 3.3.9. Let $m \times n$ matrix A have the svd $A = USV^T$. Derive that the matrix $A^T A$ has an svd $A^T A = V \tilde{S} V^T$, for what matrix \tilde{S} ? Derive that the matrix AA^T has an svd $AA^T = U \tilde{S} U^T$, for what matrix \tilde{S} ?

Exercise 3.3.10. Consider the problems (a)–(l) in Exercise 3.3.2 and problems (a)–(f) in Exercise 3.3.3. For each of these problems comment on the applicability of the Unique Solution Theorem 3.3.21, and comment on how the solution(s) illustrate the theorem.

Exercise 3.3.11. Recall Definition 3.2.2 says that a square matrix A is invertible if there exists a matrix B such that *both* $AB = I$ and $BA = I$. We now see that we need only one of these to ensure the matrix is invertible.

- (a) Use Theorem 3.3.21c to now prove that a square matrix A is invertible if there exists a matrix B such that $BA = I$.
- (b) Use the transpose and Theorems 3.3.21e and 3.3.19 to then prove that a square matrix A is invertible if there exists a matrix B such that $AB = I$.

Exercise 3.3.12. For each of the following systems, explore the effect on the solution of 1% errors in the right-hand side, and comment on the relation to the given condition number of the matrix.

$$(a) \begin{bmatrix} 1 & 0 \\ -4 & 1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 10 \end{bmatrix}, \text{cond} = 17.94$$

$$(b) \begin{bmatrix} 2 & -4 \\ -2 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 10 \\ 0 \end{bmatrix}, \text{cond} = 2$$

$$(c) \begin{bmatrix} -3 & 1 \\ -4 & 2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 6 \\ 8 \end{bmatrix}, \text{cond} = 14.93$$

$$(d) \begin{bmatrix} -1 & 1 \\ 4 & -5 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -2 \\ 10 \end{bmatrix}, \text{cond} = 42.98$$

$$(e) \begin{bmatrix} -1 & -2 \\ 3 & 1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -5 \\ 9 \end{bmatrix}, \text{cond} = 2.618$$

Exercise 3.3.13. For each of the following systems, use Matlab/Octave to explore the effect on the solution of 0.1% errors in the right-hand side. Record your commands and output, and comment on the relation to the condition number of the matrix.



$$(a) \begin{bmatrix} 1 & 2 & 2 \\ -1 & -1 & 0 \\ 0 & 3 & 1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -7 \\ 2 \\ -7 \end{bmatrix}$$



$$(b) \begin{bmatrix} -1 & 6 & -1 \\ 0 & 1 & 3 \\ -1 & 7 & 3 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 6 \\ 2 \\ 8 \end{bmatrix}$$



$$(c) \begin{bmatrix} 1 & 3 & 4 & 0 \\ 0 & 0 & -5 & 5 \\ 3 & 1 & 0 & 8 \\ 1 & 2 & 1 & 5 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 0 \\ 5 \\ 7 \\ 4 \end{bmatrix}$$



$$(d) \begin{bmatrix} -3 & -2 & -2 & -2 \\ 2 & 1 & -5 & -7 \\ 2 & 4 & 3 & 3 \\ 2 & 1 & 1 & 1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -3 \\ -8 \\ 5 \\ 2 \end{bmatrix}$$



$$(e) \begin{bmatrix} -1 & 6 & -6 & 2 & 7 \\ -7 & 4 & 3 & 1 & -8 \\ 7 & 6 & 4 & 0 & 5 \\ -8 & 3 & 3 & 2 & 4 \\ 2 & 0 & -3 & 1 & 0 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 5 \\ -7 \\ 5 \\ -2 \\ 1 \end{bmatrix}$$

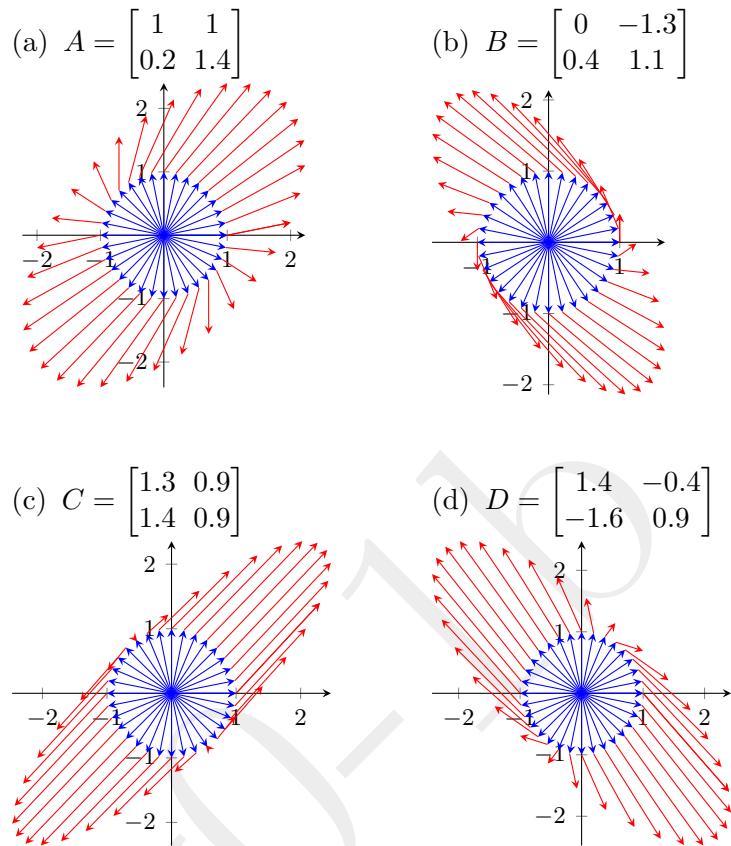


$$(f) \begin{bmatrix} 9 & 0 & -10 & -8 & -1 \\ 9 & 3 & -5 & -4 & 4 \\ -1 & 0 & -3 & -6 & -6 \\ 4 & 6 & 0 & -5 & -14 \\ -2 & -1 & -4 & -7 & 5 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 7 \\ 4 \\ 3 \\ 5 \\ 1 \end{bmatrix}$$

Exercise 3.3.14. For any $m \times n$ matrix A , use an SVD $A = USV^T$ to prove that $\text{rank}(A^T A) = \text{rank } A$ and that $\text{cond}(A^T A) = \text{cond}(A)^2$ (see Exercise 3.3.9).

Exercise 3.3.15. Recall Example 3.3.26 introduced that finding a singular vector and singular value of a matrix A came from maximising $|A\mathbf{v}|$. Each of the following matrices, say A for discussion, has plotted $A\mathbf{v}$ (red) adjoined the corresponding unit vector \mathbf{v} (blue). For each case:

- (i) by inspection of the plot, estimate a singular vector \mathbf{v}_1 that appears to maximise $|A\mathbf{v}_1|$ (to one decimal place say);
- (ii) estimate the corresponding singular value σ_1 by measuring $|A\mathbf{v}_1|$ on the plot;
- (iii) set the second singular vector \mathbf{v}_2 to be orthogonal to \mathbf{v}_1 by swapping components, and making one negative;
- (iv) estimate the corresponding singular value σ_2 by measuring $|A\mathbf{v}_2|$ on the plot;
- (v) compute the matrix-vector products $A\mathbf{v}_1$ and $A\mathbf{v}_2$, and confirm they are orthogonal (approximately).



Exercise 3.3.16. Use properties of the dot product to prove that when \mathbf{v}_1 and \mathbf{v} are orthogonal unit vectors the vector $\mathbf{v}_1 \cos t + \mathbf{v} \sin t$ is also a unit vector for all t (used in the proof of the SVD in section 3.3.3).

3.4 Subspaces, basis and dimension

Section Contents

3.4.1	Subspaces are lines, planes, and so on	240
3.4.2	Orthonormal bases form a foundation	249
3.4.3	Is it a line? a plane? The dimension answers	259
3.4.4	Exercises	268

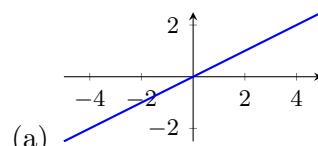
[Nature] is written in that great book which ever lies before our eyes—I mean the universe—but we cannot understand it if we do not first learn the language and grasp the symbols in which it is written. The book is written in the mathematical language, and the symbols are triangles, circles, and other geometric figures, without whose help it is impossible to comprehend a single word of it; without which one wanders in vain through a dark labyrinth.

Galileo Galilei, 1610

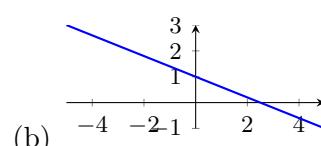
One of the most fundamental geometric structures in linear algebra are the lines or planes through the origin, and higher dimensional analogues. For example, a general solution of linear equations often involve linear combinations such as $(-2, 1, 0, 0)s + (-\frac{15}{7}, 0, \frac{9}{7}, 1)t$ (Example 2.2.24d) and $y_3\mathbf{v}_3 + y_4\mathbf{v}_4$ (Example 3.3.12): such combinations for all values of the free variables forms a plane through the origin. The aim of this section is to connect such geometric structures to the information in a singular value decomposition. Each such geometric structure is called a subspace.

3.4.1 Subspaces are lines, planes, and so on

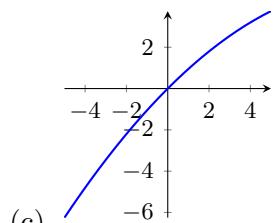
Example 3.4.1. The following graphs illustrate the concept of subspaces with structures (and imagined to infinitely extend as appropriate).



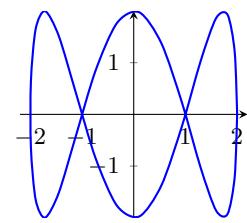
(a)
is a subspace as it is a line
through the origin.



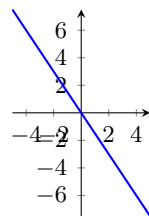
(b)
is not a subspace as it does not
include the origin.



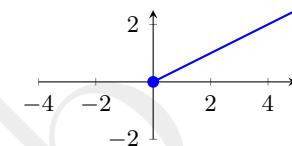
(c) is *not* a subspace as it curves.



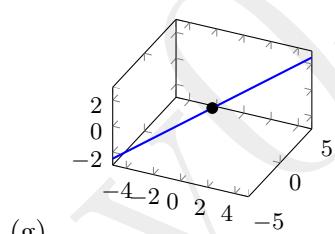
(d) is *not* a subspace as it not only curves, but does not include the origin.



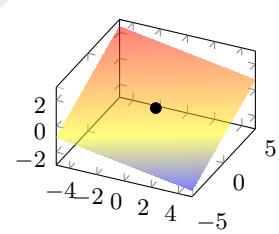
(e) is a subspace.



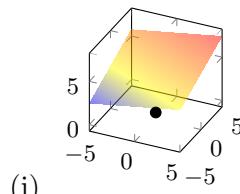
(f) where the disc indicates an end to the line, is *not* a subspace as it does not extend infinitely in both directions.



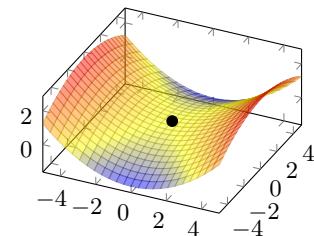
(g) is a subspace as it is a line through the origin (marked in these 3D plots).



(h) is a subspace as it is a plane through the origin.



(i) is *not* a subspace as it does not go through the origin.



(j) is *not* a subspace as it curves.

The following definition expresses precisely in algebra the concept of a subspace (this book uses the ‘blackboard bold’ font, such as \mathbb{W} and \mathbb{R} , for names of spaces and subspaces).

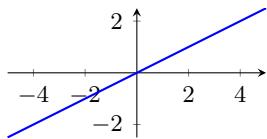
Recall that the mathematical symbol “ \in ” means “in” or “in the

set" or "is an element of the set". For examples: " $c \in \mathbb{R}$ " means " c is in the set of real numbers"; whereas " $\mathbf{v} \in \mathbb{R}^n$ " means " \mathbf{v} is a vector in \mathbb{R}^n ". Hereafter, this book uses " \in " extensively.

Definition 3.4.2. A *subspace* \mathbb{W} of \mathbb{R}^n , is a set of vectors such that $\mathbf{0} \in \mathbb{W}$ and \mathbb{W} is closed under addition and scalar multiplication: that is, for all $c \in \mathbb{R}$ and $\mathbf{u}, \mathbf{v} \in \mathbb{W}$, then both $\mathbf{u} + \mathbf{v} \in \mathbb{W}$ and $c\mathbf{u} \in \mathbb{W}$.

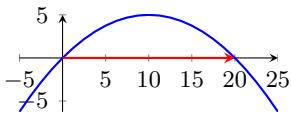
Example 3.4.3. Show why each of the following are subspaces, or not.

- (a) All vectors in the line $y = x/2$ (Example 3.4.1a).



Solution: The origin $\mathbf{0}$ is in the line $y = x/2$. The line $y = x/2$ is composed of vectors in the form $\mathbf{u} = (1, \frac{1}{2})t$ for some parameter t . Then for any $c \in \mathbb{R}$, $c\mathbf{u} = c(1, \frac{1}{2})t = (1, \frac{1}{2})(ct) = (1, \frac{1}{2})t'$ for parameter $t' = ct$; hence $c\mathbf{u}$ is in the line. Let $\mathbf{v} = (1, \frac{1}{2})s$ be another vector in the line for some parameter s , then $\mathbf{u} + \mathbf{v} = (1, \frac{1}{2})t + (1, \frac{1}{2})s = (1, \frac{1}{2})(t+s) = (1, \frac{1}{2})t'$ for parameter $t' = t+s$; hence $\mathbf{u} + \mathbf{v}$ is in the line. The three requirements of Definition 3.4.2 are met, and so this line is a subspace.

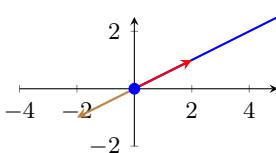
- (b) All vectors (x, y) such that $y = x - x^2/20$ (Example 3.4.1c).



Solution: To show something is not a subspace, we only need to give one instance when one of the properties fail. One instance is that the vector $(20, 0)$ is in the curve as $20 - 20^2/20 = 0$, but the scalar multiple of half of this vector $\frac{1}{2}(20, 0) = (10, 0)$ is not as $10 - 10^2/20 = 5 \neq 0$. That is, the curve is not closed under scalar multiplication and hence is not a subspace.

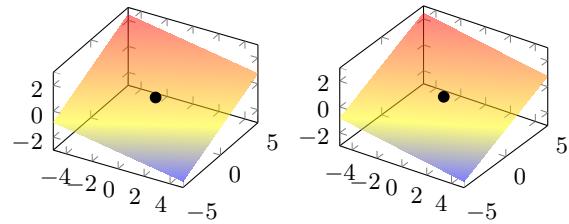
- (c) All vectors (x, y) in the line $y = x/2$ for $x, y \geq 0$ (Example 3.4.1f).

Solution: Although vectors (x, y) in the line $y = x/2$ for $x, y \geq 0$ includes the origin and is closed under addition, it fails the scalar multiplication test. For example, $\mathbf{u} = (2, 1)$ is in the line, but the scalar multiple $(-1)\mathbf{u} = (-2, -1)$ is not. Hence it is not a subspace.



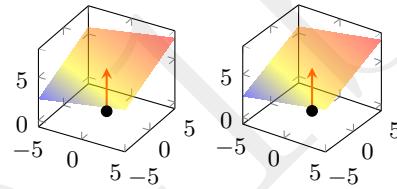
- (d) All vectors (x, y, z) in the plane $z = -x/6 + y/3$ (Example 3.4.1h).

Solution: The origin $\mathbf{0}$ is in the plane $z = -x/6 + y/3$. A vector $\mathbf{u} = (u_1, u_2, u_3)$ is in the plane provided $-u_1 + 2u_2 - 6u_3 = 0$. Consider $c\mathbf{u} = (cu_1, cu_2, cu_3)$ for which $-(cu_1) + 2(cu_2) - 6(cu_3) = c(-u_1 + 2u_2 - 6u_3) = c \times 0 = 0$ and hence must also be in the plane. Also let vector $\mathbf{v} = (v_1, v_2, v_3)$ be in the plane and consider $\mathbf{u} + \mathbf{v} = (u_1 + v_1, u_2 + v_2, u_3 + v_3)$ for which $-(u_1 + v_1) + 2(u_2 + v_2) - 6(u_3 + v_3) = -u_1 - v_1 + 2u_2 + 2v_2 - 6u_3 - 6v_3 = (-u_1 + 2u_2 - 6u_3) + (-v_1 + 2v_2 - 6v_3) = 0 + 0 = 0$ and hence must also be in the plane. The three requirements of Definition 3.4.2 are met, and so this plane is a subspace.



- (e) All vectors (x, y, z) in the plane $z = 5 + x/6 + y/3$ (Example 3.4.1i).

Solution: In this case, consider the vector $\mathbf{u} = (0, 0, 5)$ (shown in the margin): any scalar multiple, say $2\mathbf{u} = (0, 0, 10)$, is not in the plane. That is, vectors in the plane are not closed under scalar multiplication, and hence the plane is not a subspace.



- (f) $\{\mathbf{0}\}$.

Solution: The zero vector forms a trivial subspace, $\mathbb{W} = \{\mathbf{0}\}$: firstly, $\mathbf{0} \in \mathbb{W}$; secondly, the only vector in \mathbb{W} is $\mathbf{u} = \mathbf{0}$ for which every scalar multiple $c\mathbf{u} = c\mathbf{0} = \mathbf{0} \in \mathbb{W}$; and thirdly, a second vector \mathbf{v} in \mathbb{W} can only be $\mathbf{v} = \mathbf{0}$ so $\mathbf{u} + \mathbf{v} = \mathbf{0} + \mathbf{0} = \mathbf{0} \in \mathbb{W}$. The three requirements of Definition 3.4.2 are met, and so $\{\mathbf{0}\}$ is always a subspace.

- (g) \mathbb{R}^n .

Solution: Lastly, \mathbb{R}^n also is a subspace: firstly, $\mathbf{0} = (0, 0, \dots, 0) \in \mathbb{R}^n$; secondly, for $\mathbf{u} = (u_1, u_2, \dots, u_n) \in \mathbb{R}^n$, the scalar multiplication $c\mathbf{u} = c(u_1, u_2, \dots, u_n) = (cu_1, cu_2, \dots, cu_n) \in \mathbb{R}^n$; and thirdly, for $\mathbf{v} = (v_1, v_2, \dots, v_n) \in \mathbb{R}^n$, the vector addition $\mathbf{u} + \mathbf{v} = (u_1, u_2, \dots, u_n) + (v_1, v_2, \dots, v_n) = (u_1 + v_1, u_2 + v_2, \dots, u_n + v_n) \in \mathbb{R}^n$. The three requirements of Definition 3.4.2 are met, and so \mathbb{R}^n is always a subspace.

■

In summary:

- in two dimensions, \mathbb{R}^2 , subspaces are the origin $\mathbf{0}$, a line, through $\mathbf{0}$, or the entire plane \mathbb{R}^2 ;
- in three dimensions, \mathbb{R}^3 , subspaces are the origin $\mathbf{0}$, a line through $\mathbf{0}$, a plane through $\mathbf{0}$, or the entire space \mathbb{R}^3 ;

- and analogously for higher dimensions, \mathbb{R}^n .

Recall that the set of all linear combinations of a set of vectors, such as $(-2, 1, 0, 0)s + (-\frac{15}{7}, 0, \frac{9}{7}, 1)t$ (Example 2.2.24d), is called the span of that set (Definition 2.3.7).

Theorem 3.4.4. *Let $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ be vectors in \mathbb{R}^n , then $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is a subspace.*

Proof. Denote $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ by \mathbb{W} ; we aim to prove it is a subspace (Definition 3.4.2). First, $\mathbf{0} = 0\mathbf{v}_1 + 0\mathbf{v}_2 + \dots + 0\mathbf{v}_n$ and so the zero vector $\mathbf{0} \in \mathbb{W}$. Now let $\mathbf{u}, \mathbf{v} \in \mathbb{W}$ then by Definition 2.3.7 there are coefficients a_1, a_2, \dots, a_n and b_1, b_2, \dots, b_n such that

$$\begin{aligned}\mathbf{u} &= a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \dots + a_k\mathbf{v}_k, \\ \mathbf{v} &= b_1\mathbf{v}_1 + b_2\mathbf{v}_2 + \dots + b_k\mathbf{v}_k.\end{aligned}$$

Secondly, consequently

$$\begin{aligned}\mathbf{u} + \mathbf{v} &= a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \dots + a_k\mathbf{v}_k \\ &\quad + b_1\mathbf{v}_1 + b_2\mathbf{v}_2 + \dots + b_k\mathbf{v}_k \\ &= (a_1 + b_1)\mathbf{v}_1 + (a_2 + b_2)\mathbf{v}_2 + \dots + (a_k + b_k)\mathbf{v}_k \\ &\in \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\} = \mathbb{W}.\end{aligned}$$

Thirdly, for any scalar c ,

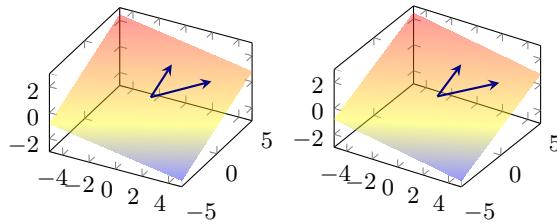
$$\begin{aligned}c\mathbf{u} &= c(a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \dots + a_k\mathbf{v}_k) \\ &= ca_1\mathbf{v}_1 + ca_2\mathbf{v}_2 + \dots + ca_k\mathbf{v}_k \\ &\in \text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\} = \mathbb{W}.\end{aligned}$$

Hence $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is a subspace. \square

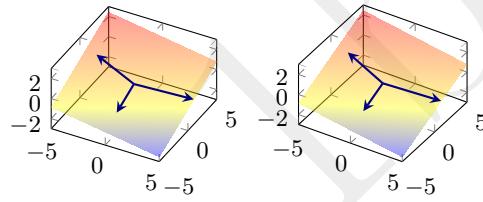
Example 3.4.5. $\text{span}\{(1, \frac{1}{2})\}$ is the subspace $y = x/2$. The reason is that a vector $\mathbf{u} \in \text{span}\{(1, \frac{1}{2})\}$ only if there is some constant a_1 such that $\mathbf{u} = a_1(1, \frac{1}{2}) = (a_1, a_1/2)$. That is, the y -component is half the x -component and hence it lies on the line $y = x/2$.

$\text{span}\{(1, \frac{1}{2}), (-2, -1)\}$ is also the subspace $y = x/2$ since every linear combination $a_1(1, \frac{1}{2}) + a_2(-2, -1) = (a_1 - 2a_2, a_1/2 - a_2)$ satisfies that the y -component is half the x -component and hence the linear combination lies on the line $y = x/2$. \blacksquare

Example 3.4.6. The plane $z = -x/6 + y/3$ may be written as $\text{span}\{(3, 3, 1/2), (0, 3, 1)\}$, as illustrated in stereo below, since every linear combination of these two vectors fills out the plane: $a_1(3, 3, 1/2) + a_2(0, 3, 1) = (3a_1, 3a_1 + 3a_2, a_1/2 + a_2)$ and so lies in the plane as $-x/6 + y/3 - z = -\frac{1}{6}3a_1 + \frac{1}{3}(3a_1 + 3a_2) - (a_1/2 + a_2) = -\frac{1}{2}a_1 + a_1 + a_2 - \frac{1}{2}a_1 - a_2 = 0$ for all a_1 and a_2 (although such arguments do not establish that the linear combinations cover the whole plane—we need Theorem 3.4.11).

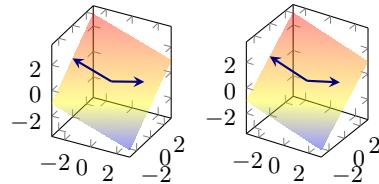


Also, $\text{span}\{(5, 1, -1/2), (0, -3, -1), (-4, 1, 1)\}$ is the plane $z = -x/6 + y/3$, as illustrated in the margin. The reason is that every linear combination of these three vectors fills out the plane: $a_1(5, 1, -1/2) + a_2(0, -3, -1) + a_3(-4, 1, 1) = (5a_1 - 4a_3, a_1 - 3a_2 + a_3, -a_1/2 - a_2 + a_3)$ and so lies in the plane as $-x/6 + y/3 - z = -\frac{1}{6}(5a_1 - 4a_3) + \frac{1}{3}(a_1 - 3a_2 + a_3) - (-a_1/2 - a_2 + a_3) = -\frac{5}{6}a_1 + \frac{2}{3}a_3 + \frac{1}{3}a_1 - a_2 + \frac{1}{3}a_3 + \frac{1}{2}a_1 + a_2 - a_3 = 0$ for all a_1, a_2 and a_3 .



Example 3.4.7. Find a set of two vectors that spans the plane $x - 2y + 3z = 0$.

Solution: Write the equation for this plane as $x = 2y - 3z$, say, then vectors in the plane are all of the form $\mathbf{u} = (x, y, z) = (2y - 3z, y, z) = (2, 1, 0)y + (-3, 0, 1)z$. That is, all vectors in the plane may be written as a linear combination of the two vectors $(2, 1, 0)$ and $(-3, 0, 1)$, hence the plane is $\text{span}\{(2, 1, 0), (-3, 0, 1)\}$ as illustrated in stereo below.



Such subspaces connect with matrices by considering a matrix whose rows or columns are the vectors appearing within the span. The columns are more often invoked than the rows.

Definition 3.4.8. (a) The **column space** of any $m \times n$ matrix A is the subspace of \mathbb{R}^m spanned by the n column vectors of A .¹⁶

¹⁶ Some of you will know that the column space is also called the range, but for the moment we just use column space.

- (b) **row space** of any $m \times n$ matrix A is the subspace of \mathbb{R}^n spanned by the m row vectors of A .

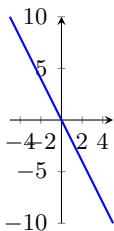
Example 3.4.9. Examples 3.4.5–3.4.7 provide some cases.

- From Example 3.4.5, the column space of $A = \begin{bmatrix} 1 & -2 \\ \frac{1}{2} & -1 \end{bmatrix}$ is the line $y = x/2$.

The row space of this matrix A is $\text{span}\{(1, -2), (\frac{1}{2}, -1)\}$ which is the set of all vectors of the form $(1, -2)s + (\frac{1}{2}, -1)t = (s + t/2, -2s - t) = (1, -2)(s + t/2) = (1, -2)t'$ is the line $y = -2x$ as illustrated in the margin. That the row space and the column space are both lines, albeit different lines, is not a coincidence (Theorem 3.4.26).

- Example 3.4.6 shows that the column space of matrix

$$B = \begin{bmatrix} 3 & 0 \\ 3 & 3 \\ \frac{1}{2} & 1 \end{bmatrix}$$



is the plane $z = -x/6 + y/3$ in \mathbb{R}^3 .

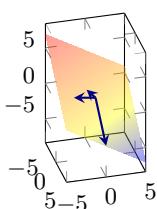
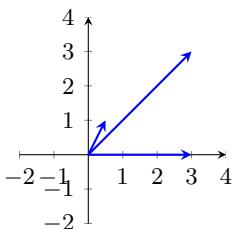
The row space of matrix B is $\text{span}\{(3, 0), (3, 3), (\frac{1}{2}, 1)\}$ which is a subspace of \mathbb{R}^2 , whereas the column space is a subspace of \mathbb{R}^3 . Here the span is all of \mathbb{R}^2 as for each $(x, y) \in \mathbb{R}^2$ choose the linear combination $\frac{x-y}{3}(3, 0) + \frac{y}{3}(3, 3) + 0(\frac{1}{2}, 1) = (x - y + y, 0 + y) = (x, y)$ so all of the \mathbb{R}^2 plane is the span. That the column space and the row space are both planes is no coincidence (Theorem 3.4.26).

- Example 3.4.6 also shows that the column space of matrix

$$C = \begin{bmatrix} 5 & 0 & -4 \\ 1 & -3 & 1 \\ -\frac{1}{2} & -1 & 1 \end{bmatrix}$$

is also the plane $z = -x/6 + y/3$ in \mathbb{R}^3 .

Now, $\text{span}\{(5, 0, -4), (1, -3, 1), (-\frac{1}{2}, -1, 1)\}$ is the row space of matrix C . It is not readily apparent but we can check that this space is the plane $4x + 3y + 5z = 0$ as illustrated in the margin. To see this, consider all linear combinations $a_1(5, 0, -4) + a_2(1, -3, 1) + a_3(-\frac{1}{2}, -1, 1) = (5a_1 + a_2 - a_3/2, -3a_2 - a_3, -4a_1 + a_2 + a_3)$ satisfy $4x + 3y + 5z = 4(5a_1 + a_2 - a_3/2) + 3(-3a_2 - a_3) + 5(-4a_1 + a_2 + a_3) = 20a_1 + 4a_2 - 2a_3 - 9a_2 - 3a_3 - 20a_1 + 5a_2 + 5a_3 = 0$. Again, it is no coincidence that the row and column spaces of C are both planes (Theorem 3.4.26). ■



Example 3.4.10. Is vector $\mathbf{b} = (-0.6, 0, -2.1, 1.9, 1.2)$ in the column space of matrix

$$A = \begin{bmatrix} 2.8 & -3.1 & 3.4 \\ 4.0 & 1.7 & 0.8 \\ -0.4 & -0.1 & 4.4 \\ 1.0 & -0.4 & -4.7 \\ -0.3 & 1.9 & 0.7 \end{bmatrix} ?$$

What about vector $\mathbf{c} = (15.2, 5.4, 3.8, -1.9, -3.7)$?

Solution: The question is: can we find a linear combination of the columns of A which equals vector \mathbf{b} ? That is, can we find some vector \mathbf{x} such that $A\mathbf{x} = \mathbf{b}$? Answer using our knowledge of linear equations.

Let's use Procedure 3.3.13 in Matlab/Octave.

- (a) Compute an SVD of this 5×3 matrix with

```
A=[2.8 -3.1 3.4
   4.0  1.7  0.8
  -0.4 -0.1  4.4
   1.0 -0.4 -4.7
  -0.3  1.9  0.7]
[U,S,V]=svd(A)
```

to find (2 d.p.)

```
U =
-0.58  0.49  0.53 -0.07  0.37
-0.17  0.69 -0.65 -0.04 -0.25
-0.56 -0.28 -0.10  0.74 -0.22
 0.57  0.43  0.21  0.66  0.14
-0.04 -0.15 -0.49  0.10  0.85
S =
 7.52      0      0
    0  4.91      0
    0      0  3.86
    0      0      0
    0      0      0
V = ...
```

- (b) Then solve $Uz = \mathbf{b}$ with $z=U'*[-0.6;0;-2.1;1.9;1.2]$ to find (2 d.p.) $\mathbf{z} = (2.55, 0.92, -0.29, -0.15, 1.54)$.
- (c) Now the diagonal matrix S has three non-zero singular values, and the last two rows are zero. So to be able to solve $Sy = \mathbf{b}$ we need the last two components of \mathbf{b} to be zero. They are not, so the system is not solvable. Hence there is no linear combination of the columns of A that gives us vector \mathbf{b} . Consequently, vector \mathbf{b} is not in the column space of A .
- (d) For the vector \mathbf{c} solve $Uz = \mathbf{c}$ with

```
z=U'*[15.2;5.4;3.8;-1.9;-3.7]
```



to find $\mathbf{z} = (-12.800, 9.876, 5.533, 0.000, 0.000)$. Since the last two entries in vector \mathbf{z} are zero, corresponding to the zero rows of S , a solution exists to $S\mathbf{y} = \mathbf{z}$. Hence a solution exists to $A\mathbf{x} = \mathbf{c}$. Consequently, vector \mathbf{c} is in the column space of A .

(Incidentally, you may check that $\mathbf{c} = 2\mathbf{a}_1 - 2\mathbf{a}_2 + \mathbf{a}_3$.)

■

Another subspace associated with matrices is the set of possible solutions to a homogeneous system of linear equations.

Theorem 3.4.11. *For any $m \times n$ matrix A , define the set $\text{null}(A)$ to be all the solutions \mathbf{x} of the homogeneous system $A\mathbf{x} = \mathbf{0}$. The set $\text{null}(A)$ is a subspace of \mathbb{R}^n called the **nullspace** of A .*

Proof. First, $A\mathbf{0} = \mathbf{0}$ so $\mathbf{0} \in \text{null } A$. Let $\mathbf{u}, \mathbf{v} \in \text{null } A$. Then by the distributivity of matrix-vector multiplication (Theorem 3.1.18), $A(\mathbf{u} + \mathbf{v}) = A\mathbf{u} + A\mathbf{v} = \mathbf{0} + \mathbf{0} = \mathbf{0}$ and so $\mathbf{u} + \mathbf{v} \in \text{null } A$. Lastly, by the associativity and commutativity of scalar multiplication (Theorem 3.1.16), for any $c \in \mathbb{R}$, $A(c\mathbf{u}) = Ac\mathbf{u} = cA\mathbf{u} = c(\mathbf{0}) = \mathbf{0}$ and so $c\mathbf{u} \in \text{null } A$. Hence $\text{null } A$ is a subspace (Definition 3.4.2). □

Example 3.4.12. • Example 2.2.24a showed that the only solution of the homogeneous system $\begin{cases} 3x_1 - 3x_2 = 0 \\ -x_1 - 7x_2 = 0 \end{cases}$ is $\mathbf{x} = \mathbf{0}$. Thus its set of solutions is $\{\mathbf{0}\}$ which is a subspace (Example 3.4.3f). Thus $\{\mathbf{0}\}$ is the nullspace of matrix $\begin{bmatrix} 3 & -3 \\ -1 & -7 \end{bmatrix}$.

- Recall the homogeneous system of linear equations from Example 2.2.24d has solutions $\mathbf{x} = (-2s - \frac{15}{7}t, s, \frac{9}{7}t, t) = (-2, 1, 0, 0)s + (-\frac{15}{7}, 0, \frac{9}{7}, 1)t$ for arbitrary s and t . That is, the set of solutions is $\text{span}\{(-2, 1, 0, 0), (-\frac{15}{7}, 0, \frac{9}{7}, 1)\}$. Since the set is a span (Theorem 3.4.4), the set of solutions is a subspace of \mathbb{R}^4 . Thus this set of solutions is the nullspace of the matrix $\begin{bmatrix} 1 & 2 & 4 & -3 \\ 1 & 2 & -3 & 6 \end{bmatrix}$.
- In contrast, Example 2.2.21 shows that the set of solutions of the non-homogeneous system $\begin{cases} -2v + 3w = -1, \\ 2u + v + w = -1. \end{cases}$ is $(u, v, w) = (-\frac{3}{4} - \frac{1}{4}t, \frac{1}{2} + \frac{3}{2}t, t) = (-\frac{3}{4}, \frac{1}{2}, 0) + (-\frac{1}{4}, \frac{3}{2}, 1)t$ over all values of parameter t . But there is no value of parameter t giving $\mathbf{0}$ as a solution: for the last component to be zero requires $t = 0$, but when $t = 0$ neither of the other components are zero, so they cannot all be zero. Since the origin $\mathbf{0}$ is not in the set of solutions, the set does not form a subspace, as is appropriate for a non-homogeneous system.

Example 3.4.13. Given the matrix

$$A = \begin{bmatrix} 3 & 1 & 0 \\ -5 & -1 & -4 \end{bmatrix},$$

is vector $\mathbf{v} = (-2, 6, 1)$ in the null space of A ? What about vector $\mathbf{w} = (1, -3, 2)$?

Solution: To test a given vector, just multiply by the matrix and see if the result is zero.

- $A\mathbf{v} = \begin{bmatrix} 3 \cdot (-2) + 1 \cdot 6 + 0 \cdot 1 \\ -5 \cdot (-2) - 1 \cdot 6 - 4 \cdot 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \mathbf{0}$, so $\mathbf{v} \in \text{null } A$.
- $A\mathbf{w} = \begin{bmatrix} 3 \cdot 1 + 1 \cdot (-3) + 0 \cdot 2 \\ -5 \cdot 1 - 1 \cdot (-3) - 4 \cdot 2 \end{bmatrix} = \begin{bmatrix} 0 \\ -10 \end{bmatrix} \neq \mathbf{0}$, so \mathbf{w} is not in the nullspace.

Summary Three ways that subspaces arise from a matrix are as the column space, row space, and nullspace.

3.4.2 Orthonormal bases form a foundation

The importance of orthogonal basis functions in interpolation and approximation cannot be overstated.

(Cuyt 2015, §5.3)

Given that subspaces arise frequently in linear algebra, and that there are many ways of representing the same subspace (as seen in some examples), is there an ‘efficient and canonical’ way of representing subspaces so we can reasonably distinguish whether two subspaces are the same or different? The next definition and theorems largely answer this question.

The aim is to mainly use an orthonormal set of vectors to span a subspace—the virtue is that orthonormal sets have many practically useful properties (for example, they underpin JPEG images, modes of vibration, and weather forecasting). Recall that an orthonormal set is composed of vectors that are both at right-angles to each other (their dot products are zero) and all of unit length (Definition 3.2.31).

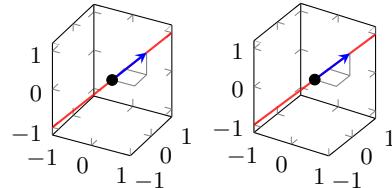
Definition 3.4.14. An *orthonormal basis* for a subspace \mathbb{W} of \mathbb{R}^n is an orthonormal set of vectors that span \mathbb{W} .

Example 3.4.15. Recall that \mathbb{R}^n is itself a subspace of \mathbb{R}^n (Example 3.4.3g).

- (a) The n standard unit vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ in \mathbb{R}^n form a set of n orthonormal vectors. They span the subspace \mathbb{R}^n as every vector in \mathbb{R}^n can be written as a linear combination $\mathbf{x} = (x_1, x_2, \dots, x_n) = x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + \dots + x_n\mathbf{e}_n$. Hence the set of standard unit vectors in \mathbb{R}^n are an orthonormal basis for the subspace \mathbb{R}^n .
- (b) The n columns $\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_n$ of an $n \times n$ orthogonal matrix Q also form an orthonormal basis for the subspace \mathbb{R}^n . The reasons are: first, Theorem 3.2.39b establishes the column vectors of Q are orthonormal; and second they span the subspace \mathbb{R}^n as for every vector $\mathbf{x} \in \mathbb{R}^n$ there exists a linear combination $\mathbf{x} = c_1\mathbf{q}_1 + c_2\mathbf{q}_2 + \dots + c_n\mathbf{q}_n$ obtained by solving $Q\mathbf{c} = \mathbf{x}$ through calculating $\mathbf{c} = Q^T\mathbf{x}$ since Q^T is the inverse of an orthogonal matrix Q (Theorem 3.2.39c). ■

Example 3.4.16. Find an orthonormal basis for the line $x = y = z$ in \mathbb{R}^3 .

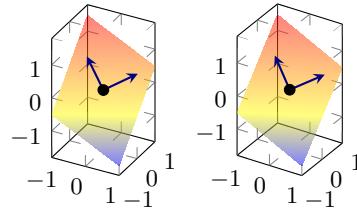
Solution: This line is a subspace as it passes through $\mathbf{0}$. A parametric description of the line is $\mathbf{x} = (x, y, z) = (t, t, t) = (1, 1, 1)t$ for every t . So the subspace is spanned by $\{(1, 1, 1)\}$. But this is not an orthonormal basis as it is not of unit length, so divide by its length $|(1, 1, 1)| = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}$. That is, $\{(1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})\}$ is an orthonormal basis for the subspace, as illustrated in stereo below.



Another orthonormal basis is the unit vector in the opposite direction, $\{(-1/\sqrt{3}, -1/\sqrt{3}, -1/\sqrt{3})\}$. ■

For subspaces that are planes in \mathbb{R}^n , orthonormal bases have more details to confirm as in the next example. The SVD then empowers us to find such bases as in the next Theorem 3.4.18.

Example 3.4.17. Confirm that the plane $-x + 2y - 2z = 0$ has an orthonormal basis $\{\mathbf{u}_1, \mathbf{u}_2\}$ where $\mathbf{u}_1 = (-\frac{2}{3}, \frac{1}{3}, \frac{2}{3})$, and $\mathbf{u}_2 = (\frac{2}{3}, \frac{2}{3}, \frac{1}{3})\}$ as illustrated in stereo below.



Solution: First, the given set is of unit vectors as the lengths are $|\mathbf{u}_1| = \sqrt{\frac{4}{9} + \frac{1}{9} + \frac{4}{9}} = 1$ and $|\mathbf{u}_2| = \sqrt{\frac{4}{9} + \frac{4}{9} + \frac{1}{9}} = 1$. Second, the set is orthonormal as their dot product is zero: $\mathbf{u}_1 \cdot \mathbf{u}_2 = -\frac{4}{9} + \frac{2}{9} + \frac{2}{9} = 0$. Third, they both lie in the plane as we check by substituting their components in the equation: for \mathbf{u}_1 , $-x+2y-2z = \frac{2}{3} + 2(\frac{1}{3}) - 2(\frac{2}{3}) = \frac{2}{3} + \frac{2}{3} - \frac{4}{3} = 0$; and for \mathbf{u}_2 , $-x+2y-2z = -\frac{2}{3} + 2(\frac{2}{3}) - 2(\frac{2}{3}) = -\frac{2}{3} + \frac{4}{3} - \frac{2}{3} = 0$. Lastly, from the parametric form of an equation for a plane (section 1.3.4) we know that all linear combinations of \mathbf{u}_1 and \mathbf{u}_2 will span the plane.

■

Theorem 3.4.18 (orthonormal basis for a span). *Let $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ be a set of n vectors in \mathbb{R}^m , then the following procedure finds an orthonormal basis for the subspace $\text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$.*

- (a) Form matrix $A := [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$.
- (b) Factorise A into its SVD, $A = USV^T$, let \mathbf{u}_j denote the columns of U (singular vectors), and let $r = \text{rank } A$ be the number of nonzero singular values (Definition 3.3.16).
- (c) Then $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ is an orthonormal basis for the subspace $\text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$.

Proof. The argument corresponds to Procedure 3.3.13. Consider any point $\mathbf{b} \in \text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$. Because \mathbf{b} is in the span, there exist coefficients x_1, x_2, \dots, x_n such that

$$\begin{aligned}
 \mathbf{b} &= \mathbf{a}_1 x_1 + \mathbf{a}_2 x_2 + \cdots + \mathbf{a}_n x_n \\
 &= A\mathbf{x} \quad (\text{by matrix-vector product §3.1.2}) \\
 &= USV^T\mathbf{x} \quad (\text{by the SVD of } A) \\
 &= US\mathbf{y} \quad (\text{for } \mathbf{y} = V^T\mathbf{x}) \\
 &= U\mathbf{z} \quad (\text{for } \mathbf{z} = (z_1, z_2, \dots, z_r, 0, \dots, 0) = S\mathbf{y}) \\
 &= \mathbf{u}_1 z_1 + \mathbf{u}_2 z_2 + \cdots + \mathbf{u}_r z_r \quad (\text{by matrix-vector product}) \\
 &\in \text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}.
 \end{aligned}$$

These equalities also hold in reverse due to the invertibility of U and V , and with $y_i = z_i/\sigma_i$ for $i = 1, 2, \dots, r$. Hence a point is in $\text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ if and only if it is in $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$. Lastly, U is an orthogonal matrix, hence the set of columns $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ is an orthonormal set and so forms an orthonormal basis for $\text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$. □

Example 3.4.19. Compute an orthonormal basis for $\text{span}\{(1, \frac{1}{2}), (-2, -1)\}$.

Solution: Form the matrix whose columns are the given vectors

$$A = \begin{bmatrix} 1 & -2 \\ \frac{1}{2} & -1 \end{bmatrix},$$

then ask Matlab/Octave for the SVD and interpret.

```
A=[1 -2; 1/2 -1]
[U,S,V]=svd(A)
```

The computed SVD is (V is immaterial here)

```
U =
-0.8944   -0.4472
-0.4472    0.8944
S =
2.5000      0
0      0.0000
V = ...
```

There is one non-zero singular value—the matrix has rank one—so an orthonormal basis for the span is the first column of matrix U , namely the set $\{(-0.89, -0.45)\}$ (2 d.p.). That is, every vector in $\text{span}\{(1, \frac{1}{2}), (-2, -1)\}$ can be written as $(-0.89, -0.45)t$ for some t : hence the span is a line.

■

Example 3.4.20. Recall that Example 3.4.6 found the plane $z = -x/6 + y/3$ could be written as $\text{span}\{(3, 3, 1/2), (0, 3, 1)\}$ or as $\text{span}\{(5, 1, -1/2), (0, -3, -1), (-4, 1, 1)\}$. Use each of these spans to find two different orthonormal bases for the plane.

Solution: • Form the matrix whose columns are the given vectors

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

then ask Matlab/Octave for the SVD and interpret. It is often easier to form the matrix in Matlab/Octave by entering the vectors as rows and then transposing:

```
A=[3 3 1/2;0 3 1] '
[U,S,V]=svd(A)
```

The computed SVD is (2 d.p.)

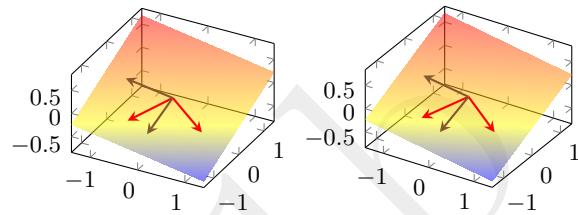
```
U =
-0.51   0.85   0.16
-0.84  -0.44  -0.31
-0.20  -0.29   0.94
```



```
S =
 4.95      0
    0  1.94
    0      0
```

V = ...

There are two non-zero singular values—the matrix has rank two—so an orthonormal basis for the plane is the set of the first two columns of matrix U , namely $(-0.51, -0.84, -0.20)$ and $(0.85, -0.44, -0.29)$. These basis vectors are illustrated as the red vectors in stereo below.



- Similarly, form the matrix

$$B = \begin{bmatrix} 5 & 0 & -4 \\ 1 & -3 & 1 \\ -\frac{1}{2} & -1 & 1 \end{bmatrix}$$

then ask Matlab/Octave for the SVD and interpret. Form the matrix in Matlab/Octave by entering the vectors as rows and then transposing:

```
B=[5 1 -1/2; 0 -3 -1; -4 1 1]'
```

```
[U,S,V]=svd(B)
```

The computed SVD is (2 d.p.)

```
U =
 -0.99  -0.04   0.16
 -0.01  -0.95  -0.31
  0.16  -0.31   0.94
S =
 6.49      0      0
    0  3.49      0
    0      0  0.00
V = ...
```

There are two non-zero singular values—the matrix has rank two—so an orthonormal basis for the plane spanned by the three vectors is the set of the first two columns of matrix U , namely the vectors $(-0.99, -0.01, 0.16)$ and $(-0.04, -0.95, -0.31)$. These are the brown pair of vectors in the above stereo illustration.



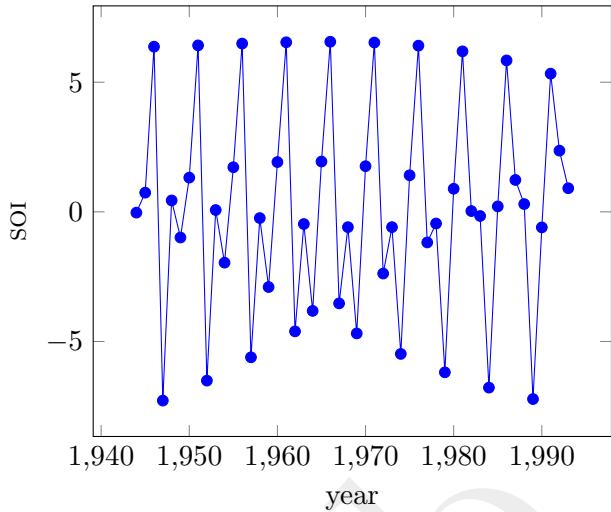


Figure 3.1: yearly average SOI over fifty years ('smoothed' somewhat for the purposes of the example). The nearly regular behaviour suggests it should be predictable.

Example 3.4.21 (data reduction). Every four or five years the phenomenon of El Nino makes a large impact on the world's weather: from drought in Australia to floods in South America. We would like to predict El Nino in advance to save lives and economies. El Nino is correlated significantly with the difference in atmospheric pressure between Darwin and Tahiti—the so-called Southern Oscillation Index (soi). This example seeks patterns in the soi in order to be able to predict the soi and hence predict El Nino.

Figure 3.1 plots the yearly average SOI each year for fifty years up to 1993. A strong regular structure is apparent, but there are significant variations and complexities in the year-to-year signal. The challenge of this example is to explore the full details of this signal.

Let's use a general technique called a Singular Spectrum Analysis. Consider a window of ten years of the SOI, and let the window 'slide' across the data to give us many 'local' pictures of the evolution in time. For example, Figure 3.2 plots six windows (each displaced vertically for clarity) each of length ten years. As the 'window' slides across the fifty year data of Figure 3.1 there are 41 such local views of the data of length ten years. Let's invoke the concept of subspaces to detect regularity in the data via these windows.

The fundamental property is that if the data has regularities, then it should lie in some subspace. We detect such subspaces using the SVD of a matrix.

- First form the 41 data windows of length ten into a matrix of size 10×41 . The numerical values of the SOI data of Figure 3.1 are the following:

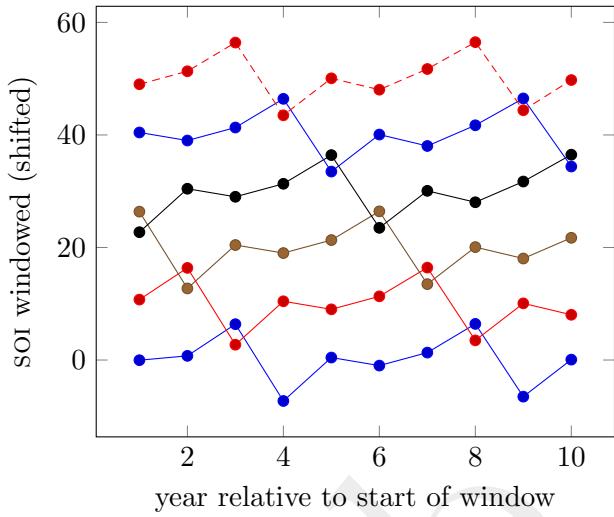


Figure 3.2: the first six windows of the SOI data of Figure 3.1—displaced vertically for clarity. Each window is of length ten years: lowest, the first window is data 1944–1953; second lowest, the second is 1945–1954; third lowest, covers 1946–1955; and so on to the 41st window is data 1984–1993, not shown.

```
year=(1944:1993)'
soi=[-0.03; 0.74; 6.37; -7.28; 0.44; -0.99; 1.32
    6.42; -6.51; 0.07; -1.96; 1.72; 6.49; -5.61
    -0.24; -2.90; 1.92; 6.54; -4.61; -0.47; -3.82
    1.94; 6.56; -3.53; -0.59; -4.69; 1.76; 6.53
    -2.38; -0.59; -5.48; 1.41; 6.41; -1.18; -0.45
    -6.19; 0.89; 6.19; 0.03; -0.16; -6.78; 0.21; 5.84
    1.23; 0.30; -7.22; -0.60; 5.33; 2.36; 0.91 ]
```

- Second form the 10×41 matrix of the windows of the data: the first seven columns being

```
A =
Columns 1 through 7
-0.03    0.74    6.37   -7.28    0.44   -0.99    1.32
    0.74    6.37   -7.28    0.44   -0.99    1.32    6.42
    6.37   -7.28    0.44   -0.99    1.32    6.42   -6.51
   -7.28    0.44   -0.99    1.32    6.42   -6.51    0.07
    0.44   -0.99    1.32    6.42   -6.51    0.07   -1.96
   -0.99    1.32    6.42   -6.51    0.07   -1.96    1.72
    1.32    6.42   -6.51    0.07   -1.96    1.72    6.49
    6.42   -6.51    0.07   -1.96    1.72    6.49   -5.61
   -6.51    0.07   -1.96    1.72    6.49   -5.61   -0.24
    0.07   -1.96    1.72    6.49   -5.61   -0.24   -2.90
```

Figure 3.2 plots the first six of these columns. The simplest way to form this matrix in Matlab/Octave—useful for all such shifting windows of data—is to invoke the `hankel()` function:

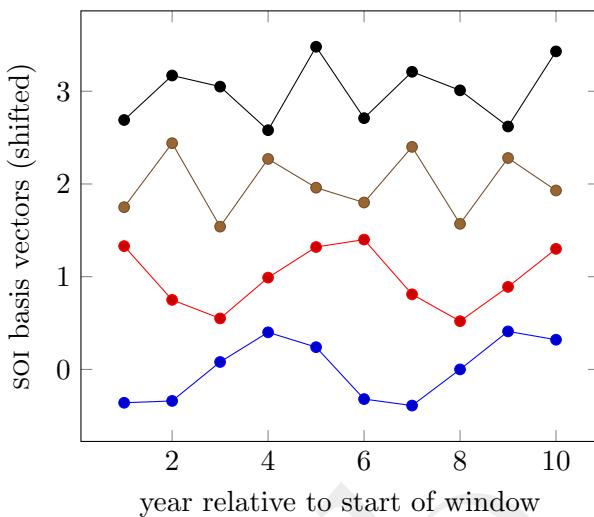


Figure 3.3: first four singular vectors of the SOI data—displaced vertically for clarity. The bottom two form a pair to show a five year cycle. The top two are a pair that show a two–three year cycle. The combination of these two cycles leads to the structure of the SOI in Figure 3.1.

```
A=hankel(soi(1:10),soi(10:50))
```

In Matlab/Octave the command `hankel(s(1:w),s(w:n))` forms the $w \times (n - w + 1)$ so-called Hankel matrix

$$\begin{bmatrix} s_1 & s_2 & s_3 & \cdots & s_{n-w} & s_{n-w+1} \\ s_2 & s_3 & \vdots & & s_{n-w+1} & \vdots \\ s_3 & \vdots & s_w & & \vdots & \vdots \\ \vdots & s_w & s_{w+1} & & \vdots & s_{n-1} \\ s_w & s_{w+1} & s_{w+2} & \cdots & s_{n-1} & s_n \end{bmatrix}$$

- Lastly, compute the SVD of the matrix of these windows:

```
[U,S,V]=svd(A);
singValues=diag(S)
plot(U(:,1:4))
```

The computed singular values are 44.63, 43.01, 39.37, 36.69, 0.03, 0.03, 0.02, 0.02, 0.02, 0.01. In practice, treat the six small singular values as zero. Since there are four non-zero singular values, the windows of data lie in a subspace spanned by the first four columns of U .

That is, all the structure seen in the fifty year SOI data of Figure 3.1 can be expressed in terms of the orthonormal basis of four ten-year vectors plotted in Figure 3.3. This analysis implies the SOI data is composed of two cycles of different frequencies.¹⁷ ■

¹⁷ However, I ‘smoothed’ the SOI data for the purposes of this example. The real



Example 3.4.20 obtained two different orthonormal bases for the one plane. Although the bases are different, they both had the same number of vectors. The next theorem establishes that this same number always occurs.

Theorem 3.4.22. *Any two orthonormal bases for a given subspace have the same number of vectors.*

Proof. Let $\mathcal{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ and $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s\}$ be any two orthonormal bases for a subspace in \mathbb{R}^n . Prove the number of vectors $r = s$ by contradiction. First assume $r < s$ (\mathcal{U} has less vectors than \mathcal{V}). Since \mathcal{U} is an orthonormal basis for the subspace every vector in \mathcal{V} can be written as a linear combination of vectors in \mathcal{U} with some coefficients a_{ij} :

$$\begin{aligned}\mathbf{v}_1 &= \mathbf{u}_1 a_{11} + \mathbf{u}_2 a_{21} + \cdots + \mathbf{u}_r a_{r1}, \\ \mathbf{v}_2 &= \mathbf{u}_1 a_{12} + \mathbf{u}_2 a_{22} + \cdots + \mathbf{u}_r a_{r2}, \\ &\vdots \\ \mathbf{v}_s &= \mathbf{u}_1 a_{1s} + \mathbf{u}_2 a_{2s} + \cdots + \mathbf{u}_r a_{rs}.\end{aligned}$$

Write each of these, such as the first one, in the form

$$\mathbf{v}_1 = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_r] \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{r1} \end{bmatrix} = U\mathbf{a}_1,$$

where matrix $U = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_r]$. Then the $n \times s$ matrix

$$\begin{aligned}V &= [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_s] \\ &= [U\mathbf{a}_1 \ U\mathbf{a}_2 \ \cdots \ U\mathbf{a}_s] \\ &= U[\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_s] = UA\end{aligned}$$

for the $r \times s$ matrix $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_s]$. By assumption $r < s$ and so Theorem 2.2.25 assures that the homogeneous system $A\mathbf{x} = \mathbf{0}$ has infinitely many solutions, choose any non-trivial solution $\mathbf{x} \neq \mathbf{0}$. Consider

$$\begin{aligned}V\mathbf{x} &= UA\mathbf{x} \quad (\text{from above}) \\ &= U\mathbf{0} \quad (\text{since } A\mathbf{x} = \mathbf{0}) \\ &= \mathbf{0}.\end{aligned}$$

Since \mathcal{V} is an orthonormal set, the matrix V is orthogonal (Theorem 3.2.39). Then premultiplying $V\mathbf{x} = \mathbf{0}$ by V^T gives $V^TV\mathbf{x} =$

so I data is considerably noisier. Also we would use 600 monthly averages not 50 yearly averages: so a ten year window would be a window of 120 months, and the matrix would be considerably larger 120×481 . Nonetheless, the conclusions with the real data, and justified by Chapter 5, are much the same.

$V^T \mathbf{0}$ which simplifies to $I_s \mathbf{x} = \mathbf{0}$; that is, $\mathbf{x} = \mathbf{0}$. But this is a contradiction, so we cannot have $r < s$.

Second, a corresponding argument establishes we cannot have $s < r$. Hence all orthonormal bases of a given subspace must have the same number of vectors. \square

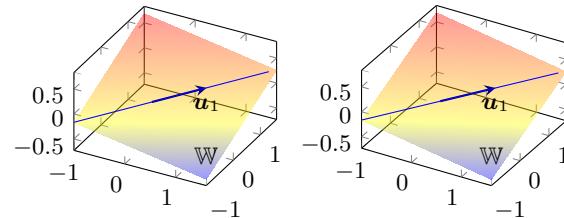
The following optional theorem and proof settles an existential issue.

An existential issue How do we know that every subspace has an orthonormal basis? We know many subspaces, such as row and column spaces, have an orthonormal basis because they are the span of rows and columns of a matrix, and then Theorem 3.4.18 assures us they have an orthonormal basis. But do all subspaces have an orthonormal basis? The following theorem certifies that they do.

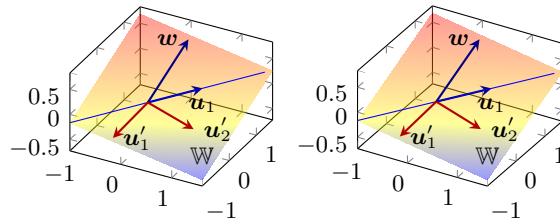
Theorem 3.4.23 (existence of basis). *Let \mathbb{W} be a subspace of \mathbb{R}^n , then there exists an orthonormal basis for \mathbb{W} .*

Proof. If the subspace is $\mathbb{W} = \{\mathbf{0}\}$, then $\mathbb{W} = \text{span}(\emptyset)$ gives a basis and this trivial case is done.

Then for other subspaces $\mathbb{W} \neq \{\mathbf{0}\}$ there exists a non-zero vector $\mathbf{w} \in \mathbb{W}$. Normalising the vector to $\mathbf{u}_1 = \mathbf{w}/|\mathbf{w}|$ all scalar multiples $c\mathbf{u}_1 = cw/|\mathbf{w}| = (c/|\mathbf{w}|)\mathbf{w} \in \mathbb{W}$ by closure of \mathbb{W} under scalar multiplication (Definition 3.4.2). Hence $\text{span}\{\mathbf{u}_1\} \subseteq \mathbb{W}$ (as illustrated below for an example). Consequently, either $\text{span}\{\mathbf{u}_1\} = \mathbb{W}$ and we are done, or we repeat the following step until the space \mathbb{W} is spanned.



Given orthonormal vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ such that the set $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\} \subset \mathbb{W}$, so $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\} \neq \mathbb{W}$. Then there must exist a vector $\mathbf{w} \in \mathbb{W}$ which is not in the set $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\}$. By the closure of subspace \mathbb{W} under addition and scalar multiplication (Definition 3.4.2), the set $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k, \mathbf{w}\} \subseteq \mathbb{W}$. Theorem 3.4.18, on the orthonormal basis for a span, then assures us that an SVD gives an orthonormal basis $\{\mathbf{u}'_1, \mathbf{u}'_2, \dots, \mathbf{u}'_{k+1}\}$ for the set $\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k, \mathbf{w}\} \subseteq \mathbb{W}$ (as illustrated for an example). Consequently, either $\text{span}\{\mathbf{u}'_1, \mathbf{u}'_2, \dots, \mathbf{u}'_{k+1}\} = \mathbb{W}$ and we are done, or we repeat the process of this paragraph with k bigger by one.



The process must terminate because $\mathbb{W} \subseteq \mathbb{R}^n$. If the process ever repeats until $k = n$, then we know $\mathbb{W} = \mathbb{R}^n$ and we are done as \mathbb{R}^n is spanned by the n standard unit vectors. \square

Ensemble simulation makes better weather forecasts Near the end of the twentieth century weather forecasts were becoming amazingly good at predicting the chaotic weather days in advance. However, there were notable failures: occasionally the weather forecast would give no hint of storms that developed (such as the severe 1999 storm in Sydney¹⁸) Why?

Occasionally the weather is both near a ‘tipping point’ where small changes may cause a storm, and where the errors in measuring the current weather are of the size of the necessary changes. Then the storm would be within the possibilities, but it would not be forecast if the measurements were, by chance error, the ‘other side’ of the tipping point. Meteorologists now mostly overcome this problem by executing on their computers an ensemble of simulations of, say, a hundred different forecast simulations (Roulstone & Norbury 2013, pp.274–80, e.g.). Such a set of 100 simulations essentially lie in a subspace spanned by 100 vectors in the vastly larger space, say $\mathbb{R}^{1000,000,000}$, of the maybe billion variables in the weather model. But what happens in the computational simulations is that the ensemble of simulations degenerate in time: so the meteorologists continuously ‘renormalise’ the ensemble of simulations by rewriting the ensemble in terms of an *orthonormal basis* of 100 vectors. Such an orthonormal basis for the ensemble reasonably ensures unusual storms are retained in the range of possibilities explored by the ensemble forecast.

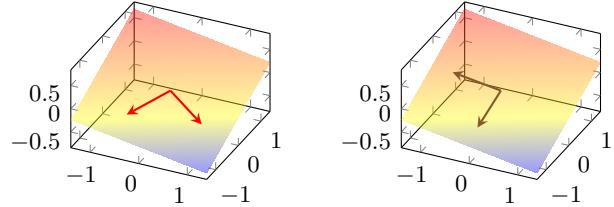
3.4.3 Is it a line? a plane? The dimension answers

physical dimension. It is an intuitive notion that appears to go back to an archaic state before Greek geometry, yet deserves to be taken up again. *Mandelbrot (1982)*

One of the beauties of an orthonormal basis is that, being orthonormal, they look just like a rotated version of the standard unit vectors. That is, the two orthonormal basis of a plane could form the two ‘standard unit vectors’ of a coordinate system in that plane.

¹⁸ http://en.wikipedia.org/wiki/1999_Sydney_hailstorm [April 2015]

Example 3.4.20 found the plane $z = -x/6 + y/3$ could have the following two orthonormal bases: either of these orthonormal bases, or indeed any other pair of orthonormal vectors, could act as a pair of standard unit vectors of the given planar subspace.



Similarly in other dimensions for other subspaces. Just as \mathbb{R}^n is called n -dimensional and has n standard unit vectors, so we analogously define the dimension of any subspace.

Definition 3.4.24. Let \mathbb{W} be a subspace of \mathbb{R}^n . The number of vectors in an orthonormal basis for \mathbb{W} is called the **dimension** of \mathbb{W} , denoted $\dim \mathbb{W}$. By convention, $\dim \{\mathbf{0}\} = 0$.

Example 3.4.25.

- Example 3.4.16 finds that the linear subspace $x = y = z$ is spanned by the orthonormal basis $\{(1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})\}$. With one vector in the basis, the line is one dimensional.
- Example 3.4.17 finds that the planar subspace $-x+2y-2z = 0$ is spanned by the orthonormal basis $\{\mathbf{u}_1, \mathbf{u}_2\}$ where $\mathbf{u}_1 = (-\frac{2}{3}, \frac{1}{3}, \frac{2}{3})$, and $\mathbf{u}_2 = (\frac{2}{3}, \frac{2}{3}, \frac{1}{3})\}$. With two vectors in the basis, the plane is two dimensional.
- Subspace $\mathbb{W} = \text{span}\{(5, 1, -1/2), (0, -3, -1), (-4, 1, 1)\}$ of Example 3.4.20 is found to have an orthonormal basis of the vectors $(-0.99, -0.01, 0.16)$ and $(-0.04, -0.95, -0.31)$. With two vectors in the basis, the subspace is two dimensional; that is, $\dim \mathbb{W} = 2$.
- Since the subspace \mathbb{R}^n (Example 3.4.3g) has an orthonormal basis of the n standard unit vectors, $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$, then its dimension $\dim \mathbb{R}^n = n$.
- The El Nino windowed data of Example 3.4.21 is effectively spanned by four orthonormal vectors. Despite the apparent complexity of the signal, the data effectively lies in a subspace of dimension four (that of two oscillators).

■

Theorem 3.4.26. The row space and column space of a matrix A have the same dimension. Further, given an SVD of the matrix, say $A = USV^T$, an orthonormal basis for the column space is the first rank A columns of U , and that for the row space is the first rank A columns of V .

Proof. From Definition 3.1.12 of the transpose, the rows of A are the same as the columns of A^T , and so the row space of A is the same as the column space of A^T . Hence,

$$\begin{aligned} & \text{dimension of the row space of } A \\ &= \text{dimension of the column space of } A^T \\ &= \text{rank}(A^T) \quad (\text{by Thm. 3.4.18}) \\ &= \text{rank } A \quad (\text{by Thm. 3.3.19}) \\ &= \text{dimension of the column space of } A \quad (\text{by Thm. 3.4.18}). \end{aligned}$$

Let $m \times n$ matrix A have an SVD $A = USV^T$ and $r = \text{rank } A$. Then Theorem 3.4.18 establishes that an orthonormal basis for the column space of A is the first r columns of U . Recall that $A^T = (USV^T)^T = VS^T U^T$ is an SVD for A^T (Theorem 3.3.19) and so an orthonormal basis for the column space of A^T is the first r columns of V (Theorem 3.4.18). Since the row space of A is the column space of A^T , an orthonormal basis for the row space of A is the first r columns of V . \square

Example 3.4.27. Use the SVD of the matrix A in Example 3.4.19 to compare the column space and the row space of matrix A .

Solution: Example 3.4.19 finds the matrix

$$A = \begin{bmatrix} 1 & -2 \\ \frac{1}{2} & -1 \end{bmatrix},$$

has an SVD of

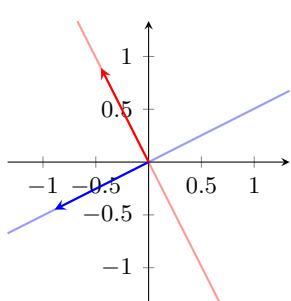
$$\mathbf{U} = \begin{bmatrix} -0.8944 & -0.4472 \\ -0.4472 & 0.8944 \end{bmatrix}$$

$$\mathbf{S} = \begin{bmatrix} 2.5000 & 0 \\ 0 & 0.0000 \end{bmatrix}$$

$$\mathbf{V} = \begin{bmatrix} -0.4472 & 0.8944 \\ 0.8944 & 0.4472 \end{bmatrix}$$

There is one non-zero singular value—the matrix has rank one—so an orthonormal basis for the column space is the first column of matrix U , namely $(-0.89, -0.45)$ (2 d.p.).

Complementing this, as there is one non-zero singular value—the matrix has rank one—so an orthonormal basis for the row space is the first column of matrix V , namely $(-0.45, 0.89)$. As illustrated in the margin, the two subspaces, the row space (red) and the column space (blue), are different but of the same dimension. (Here the orthogonality of the row and column spaces is not significant as generally the row and column spaces are not orthogonal.)

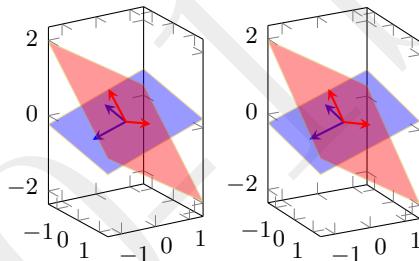


■

Example 3.4.28. Use the SVD of the matrix B in Example 3.4.20 to compare the column space and the row space of matrix B .

Solution: Recall that there are two non-zero singular values—the matrix has rank two—so an orthonormal basis for the column space is the first two columns of matrix U , namely the vectors $(-0.99, -0.01, 0.16)$ and $(-0.04, -0.95, -0.31)$.

Complementing this, as there are two non-zero singular values—the matrix has rank two—so an orthonormal basis for the row space is the set of the first two columns of matrix V , namely the vectors $(-0.78, -0.02, 0.63)$ and $(-0.28, 0.91, -0.32)$. As illustrated below in stereo, the two subspaces, the row space (red) and the column space (blue), are different but of the same dimension.



Definition 3.4.29. The **nullity** of a matrix A is the dimension of its nullspace (defined in Theorem 3.4.11), and is denoted by $\text{nullity}(A)$.

Example 3.4.30. Example 3.4.12 finds the nullspace of the two matrices

$$\begin{bmatrix} 3 & -3 \\ -1 & -7 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 1 & 2 & 4 & -3 \\ 1 & 2 & -3 & 6 \end{bmatrix}.$$

- The first matrix has nullspace $\{\mathbf{0}\}$ which has dimension zero and hence the nullity of the matrix is zero.
- The second matrix, 2×4 , has nullspace written as $\text{span}\{(-2, 1, 0, 0), (-\frac{15}{7}, 0, \frac{9}{7}, 1)\}$. Being spanned by two vectors not proportional to each other, we expect the dimension of the nullspace, the nullity, to be two. To check, compute the singular values of the matrix whose columns are these vectors: calling the matrix N for nullspace,

```
N=[-2 1 0 0; -15/7 0 9/7 1];
svd(N)
```

which computes the singular values

```
3.2485
1.3008
```



Since there are two non-zero singular values, there are two orthonormal vectors spanning the subspace, the nullspace, hence its dimension, the nullity, is two.

■

Example 3.4.31. For the matrix

$$C = \begin{bmatrix} -1 & -2 & 2 & 1 \\ -3 & 3 & 1 & 0 \\ 2 & -5 & 1 & 1 \end{bmatrix},$$

find an orthonormal basis for its nullspace and hence determine its nullity.

Solution: To find the nullspace construct a general solution to the homogeneous system $C\mathbf{x} = \mathbf{0}$ with Procedure 3.3.13.



- (a) Enter into Matlab/Octave the matrix C and compute an SVD via $[U, S, V] = \text{svd}(C)$ to find (2 d.p.)

$$\begin{aligned} U &= \\ &\begin{array}{ccc} 0.24 & 0.78 & -0.58 \\ -0.55 & 0.60 & 0.58 \\ 0.80 & 0.18 & 0.58 \end{array} \\ S &= \\ &\begin{array}{cccc} 6.95 & 0 & 0 & 0 \\ 0 & 3.43 & 0 & 0 \\ 0 & 0 & 0.00 & 0 \end{array} \\ V &= \\ &\begin{array}{cccc} 0.43 & -0.65 & 0.63 & -0.02 \\ -0.88 & -0.19 & 0.42 & 0.10 \\ 0.11 & 0.68 & 0.62 & -0.37 \\ 0.15 & 0.28 & 0.21 & 0.92 \end{array} \end{aligned}$$

- (b) Since the right-hand side is zero the solution to $U\mathbf{z} = \mathbf{0}$ is $\mathbf{z} = \mathbf{0}$.
- (c) Then, because the rank of the matrix is two, the solution to $S\mathbf{y} = \mathbf{0}$ is $\mathbf{y} = (0, 0, y_3, y_4)$ for free variables y_3 and y_4 .
- (d) The solution to $V^T\mathbf{x} = \mathbf{y}$ is $\mathbf{x} = V\mathbf{y} = \mathbf{v}_3y_3 + \mathbf{v}_4y_4 = y_3(0.63, 0.42, 0.62, 0.21) + y_4(-0.02, 0.10, -0.37, 0.92)$.

Hence $\text{span}\{(0.63, 0.42, 0.62, 0.21), (-0.02, 0.10, -0.37, 0.92)\}$ is the nullspace of matrix C . Because the columns of V are orthonormal, the two vectors appearing in this span are orthonormal and so form an orthonormal basis for the nullspace. Hence nullity $C = 2$.

■

This Example 3.4.31 indicates that the nullity is determined by the number of zero columns in the diagonal matrix S of an SVD. Conversely, the rank of a matrix is determined by the number of non-zero columns in the diagonal matrix S of an SVD. Put these two facts together in general and we get the following theorem that helps characterise solutions of linear equations.

Theorem 3.4.32 (rank theorem). *Let A be an $m \times n$ matrix, then $\text{rank } A + \text{nullity } A = n$, the number of columns of A .*

Proof. Set $r = \text{rank } A$. By Procedure 3.3.13 a general solution to the homogeneous system $Ax = \mathbf{0}$ involves $n - r$ free variables y_{r+1}, \dots, y_n in the linear combination form $\mathbf{v}_{r+1}y_{r+1} + \dots + \mathbf{v}_ny_n$. Hence the nullspace is $\text{span}\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$. Because matrix V is orthogonal, the vectors $\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$ are orthonormal; that is, they form an orthonormal basis for the nullspace, and so the nullspace is of dimension $n - r$. Consequently, $\text{rank } A + \text{nullity } A = r + (n - r) = n$. \square

Example 3.4.33. Use SVDS to determine the rank and nullity of each of the given matrices.

(a)
$$\begin{bmatrix} 1 & -1 & 2 \\ 2 & -2 & 4 \end{bmatrix}$$

Solution: Enter the matrix into Matlab/Octave and compute the singular values:

```
A=[1 -1 2
  2 -2 4]
svd(A)
```

The resultant singular values are

```
5.4772
0.0000
```

The one non-zero singular value indicates $\text{rank } A = 1$. Since the matrix has three columns, the nullity—the dimension of the nullspace—is $3 - 1 = 2$.

(b)
$$\begin{bmatrix} 1 & -1 & -1 \\ 1 & 0 & -1 \\ -1 & 3 & 1 \end{bmatrix}$$

Solution: Enter the matrix into Matlab/Octave and compute the singular values:

```
B=[1 -1 -1
  1 0 -1
 -1 3 1]
svd(B)
```

The resultant singular values are



```
3.7417
1.4142
0.0000
```

The two non-zero singular values indicate $\text{rank } B = 2$. Since the matrix has three columns, the nullity—the dimension of the nullspace—is $3 - 2 = 1$.

$$(c) \begin{bmatrix} 0 & 0 & -1 & -3 & 2 \\ -2 & -2 & 1 & 0 & 1 \\ 1 & -1 & 2 & 8 & -2 \\ -1 & 1 & 0 & -2 & -2 \\ -3 & -1 & 0 & -5 & 1 \end{bmatrix}$$

Solution: Enter the matrix into Matlab/Octave and compute the singular values:

```
C=[0 0 -1 -3 2
-2 -2 1 -0 1
1 -1 2 8 -2
-1 1 -0 -2 -2
-3 -1 -0 -5 1]
svd(C)
```

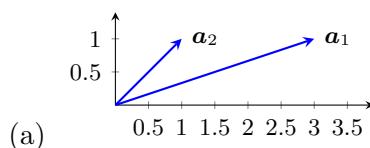
The resultant singular values are

```
10.8422
4.0625
3.1532
0.0000
0.0000
```

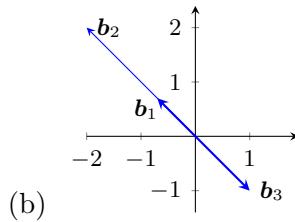
Three non-zero singular values indicate $\text{rank } C = 3$. Since the matrix has five columns, the nullity—the dimension of the nullspace—is $5 - 3 = 2$.



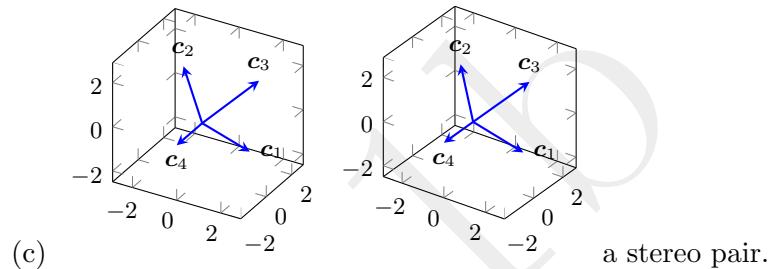
Example 3.4.34. Each of the following graphs plot all the column vectors of a matrix. What is the nullity of each of the matrices? Give reasons.



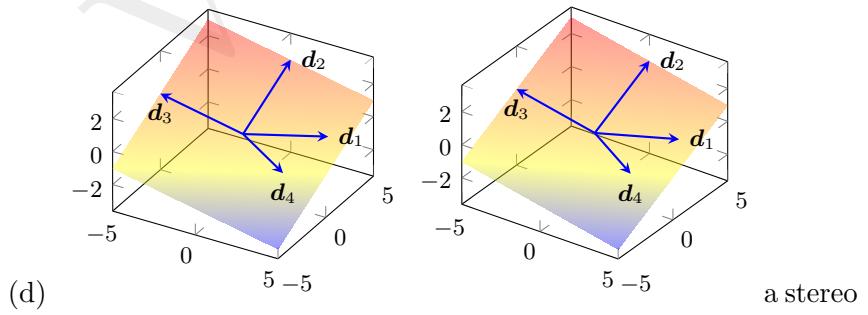
Solution: Zero. The two column vectors in the plane must come from a 2×2 matrix A . Since the two columns are at a non-trivial angle, every point in the plane may be written as a linear combination of a_1 and a_2 , hence the column space of A is \mathbb{R}^2 . Consequently, $\text{rank } A = 2$. From the rank theorem: nullity $A = n - \text{rank } A = 2 - 2 = 0$.



Solution: Two. The three column vectors in the plane must come from a 2×3 matrix B . The three vectors are all in a line, so the column space of matrix B is a line. Consequently, $\text{rank } B = 1$. From the rank theorem: nullity $B = n - \text{rank } B = 3 - 1 = 2$.



Solution: One. The four column vectors in 3D space must come from a 3×4 matrix C . Since the four columns do not all lie in a line or plane, every point in space may be written as a linear combination of c_1, c_2, \dots, c_4 , hence the column space of C is \mathbb{R}^3 . Consequently, $\text{rank } C = 3$. From the rank theorem: nullity $C = n - \text{rank } C = 4 - 3 = 1$.



Solution: Two. The four column vectors in 3D space must come from a 3×4 matrix D . Since the four columns all lie in a plane (as suggested by the drawn plane), and linear combinations can give every point in the plane, hence the column space of D has dimension two. Consequently, $\text{rank } D = 2$. The rank theorem gives nullity $D = n - \text{rank } D = 4 - 2 = 2$.

■

The recognition of these new concepts associated with matrices and linear equations, then empowers us to extend the list of exact

properties that ensure a system of linear equations has a unique solution.

Theorem 3.4.35 (Unique Solutions: version 2). *Let A be an $n \times n$ square matrix. Extending Theorem 3.3.21, the following statements are equivalent:*

- (a) A is invertible;
- (b) $A\mathbf{x} = \mathbf{b}$ has a unique solution for every $\mathbf{b} \in \mathbb{R}^n$;
- (c) $A\mathbf{x} = \mathbf{0}$ has only the zero solution;
- (d) all n singular values of A are nonzero;
- (e) $\text{rank } A = n$;
- (f) $\text{nullity } A = 0$;
- (g) the column vectors of A span \mathbb{R}^n ;
- (h) the row vectors of A span \mathbb{R}^n .

Proof. Theorem 3.3.21 establishes the equivalence of the statements 3.4.35a–3.4.35e. We prove the equivalence of these with the statements 3.4.35f–3.4.35h.

3.4.35e \iff 3.4.35f : The Rank Theorem 3.4.32 assures us that $\text{nullity } A = 0$ iff $\text{rank } A = n$.

3.4.35b \implies 3.4.35g : By 3.4.35b every $\mathbf{b} \in \mathbb{R}^n$ can be written as $\mathbf{b} = A\mathbf{x}$. But $A\mathbf{x}$ is a linear combination of the columns of A and so \mathbf{b} is in the span of the columns. Hence the column vectors of A span \mathbb{R}^n .

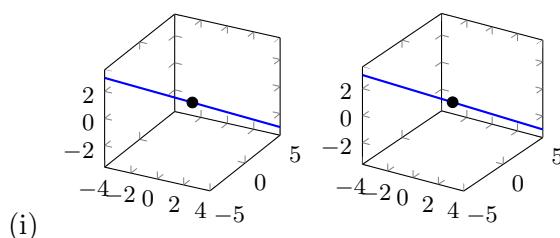
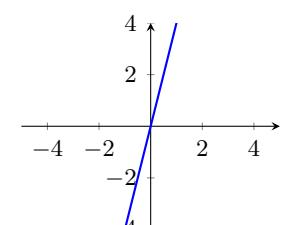
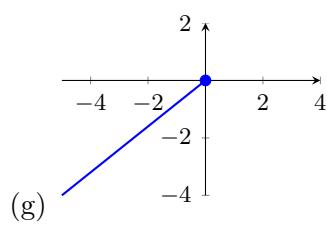
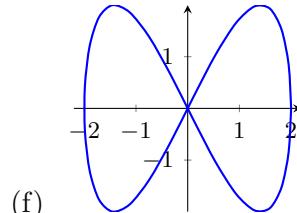
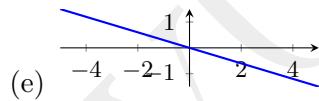
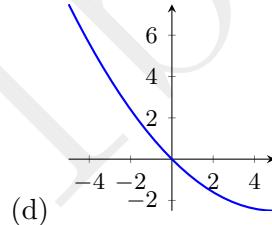
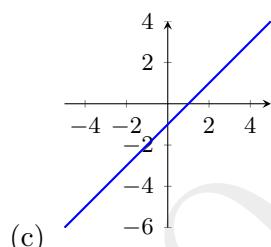
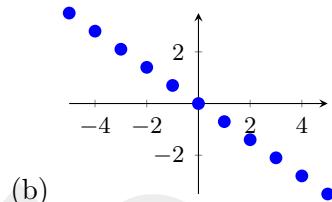
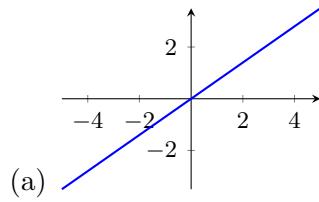
3.4.35g \implies 3.4.35e : Suppose $\text{rank } A = r$ reflecting r non-zero singular values in an SVD $A = USV^T$. Theorem 3.4.18 assures us the column space of A has orthonormal basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$. But the column space is \mathbb{R}^n (statement 3.4.35g) which also has the orthonormal basis of the n standard unit vectors. Theorem 3.4.22 assures us that the number of basis vectors must be the same; that is, $\text{rank } A = r = n$.

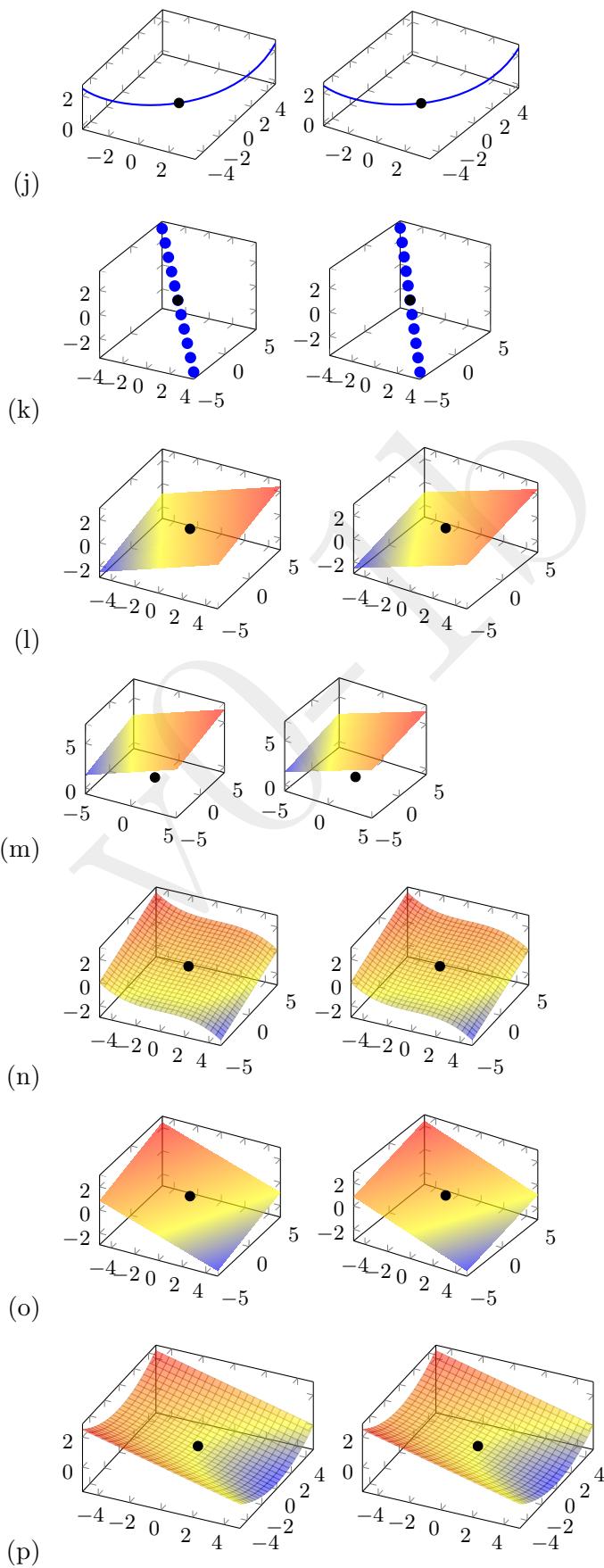
3.4.35e \iff 3.4.35h : Theorem 3.3.19 asserts $\text{rank}(A^T) = \text{rank } A$, so the statement 3.4.35e implies $\text{rank}(A^T) = n$, and so statement 3.4.35g asserts the columns of A^T span \mathbb{R}^n . But the columns of A^T are the rows of A so the rows of A span \mathbb{R}^n . Conversely, if the rows of A span \mathbb{R}^n , then so do the columns of A^T , hence $\text{rank}(A^T) = n$ which by Theorem 3.3.19 implies $\text{rank } A = n$.

□

3.4.4 Exercises

Exercise 3.4.1. Use your intuitive notion of a subspace to decide whether each of the following drawn sets (3D in stereo pair) is a subspace, or not.





Exercise 3.4.2. Use Definition 3.4.2 to decide whether each of the following is a subspace, or not. Give reasons.

- (a) All vectors in the line $y = 2x$.
- (b) All vectors in the line $3.2y = 0.8x$.
- (c) All vectors $(x, y) = (t, 2 + t)$ for all real t .
- (d) All vectors $(1.3n, -3.4n)$ for all integer n .
- (e) All vectors $(x, y) = (-3.3 - 0.3t, 2.4 - 1.8t)$ for all real t .
- (f) $\text{span}\{(6, -1), (1, 2)\}$
- (g) All vectors $(x, y) = (6 - 3t, t - 2)$ for all real t .
- (h) The vectors $(2, 1, -3)t + (5, -\frac{1}{2}, 2)s$ for all real s, t .
- (i) The vectors $(0.9, 2.4, 1)t - (0.2, 0.6, 0.3)s$ for all real s, t .
- (j) All vectors (x, y) such that $y = x^3$.
- (k) All vectors (x, y, z) such that $x = 2t$, $y = t^2$ and $z = t/2$ for all t .
- (l) The vectors $(t, n, 2t + 3n)$ for real t and integer n .
- (m) $\text{span}\{(0, -1, 1), (-1, 0, 2)\}$
- (n) The vectors $(2.7, 2.6, -0.8, 2.1)s + (0.5, 0.1, -1, 3.3)t$ for all real s, t .
- (o) The vectors $(1.4, 2.3, 1.5, 4) + (1.2, -0.8, -1.2, 2)t$ for all real t .
- (p) The vectors $(t^3, 2t^3)$ for all real t (tricky).
- (q) The vectors $(t^2, 3t^2)$ for all real t (tricky).

Exercise 3.4.3. Let \mathbb{W}_1 and \mathbb{W}_2 be any two subspaces of \mathbb{R}^n (Definition 3.4.2).

- (a) Use the definition to prove that the intersection of \mathbb{W}_1 and \mathbb{W}_2 is also a subspace of \mathbb{R}^n .
- (b) Give an example to prove that the union of \mathbb{W}_1 and \mathbb{W}_2 is not necessarily a subspace of \mathbb{R}^n .

Exercise 3.4.4. For each of the following matrices, partially solve linear equations to determine whether the given vector \mathbf{b}_j is in the column space, and to determine if the given vector \mathbf{r}_j is in the row space of the matrix. Work small problems by hand, and address larger problems with Matlab/Octave. Record your working or Matlab/Octave commands and output.

$$(a) A = \begin{bmatrix} 2 & 1 \\ 5 & 4 \end{bmatrix}, \mathbf{b}_1 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}, \mathbf{r}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$(b) B = \begin{bmatrix} -2 & 1 \\ 4 & -2 \end{bmatrix}, \mathbf{b}_2 = \begin{bmatrix} 1 \\ -3 \end{bmatrix}, \mathbf{r}_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(c) \ C = \begin{bmatrix} 1 & -1 \\ -3 & 4 \\ -3 & 5 \end{bmatrix}, \ b_3 = \begin{bmatrix} 5 \\ 0 \\ -1 \end{bmatrix}, \ r_3 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$(d) \ D = \begin{bmatrix} -2 & -4 & -5 \\ -6 & -2 & 1 \end{bmatrix}, \ b_4 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \ r_4 = \begin{bmatrix} 2 \\ -6 \\ -11 \end{bmatrix}$$

$$(e) \ E = \begin{bmatrix} 3 & 2 & 4 \\ 1 & 6 & 0 \\ 1 & -2 & 2 \end{bmatrix}, \ b_5 = \begin{bmatrix} 10 \\ 2 \\ 4 \end{bmatrix}, \ r_5 = \begin{bmatrix} 0 \\ -1 \\ 2 \end{bmatrix}$$

$$(f) \ F = \begin{bmatrix} 0 & -1 & -4 \\ -2 & 0 & 4 \\ 7 & -1 & -3 \\ 1 & 1 & 3 \end{bmatrix}, \ b_6 = \begin{bmatrix} 2 \\ 3 \\ 1 \\ 3 \end{bmatrix}, \ r_6 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$(g) \ G = \begin{bmatrix} -2 & -1 & 5 & 4 \\ 0 & -3 & 1 & -1 \\ 3 & 3 & 4 & 3 \end{bmatrix}, \ b_7 = \begin{bmatrix} 1 \\ 2 \\ -4 \end{bmatrix}, \ r_7 = \begin{bmatrix} 0 \\ -3 \\ 1 \\ -1 \end{bmatrix}$$

$$(h) \ H = \begin{bmatrix} -2 & 1 & 1 & -1 \\ 2 & -2 & -1 & 0 \\ -2 & 1 & 1 & -2 \\ 2 & 5 & -1 & -1 \end{bmatrix}, \ b_8 = \begin{bmatrix} 2 \\ -1 \\ 3 \\ 0 \end{bmatrix}, \ r_8 = \begin{bmatrix} -1 \\ -2 \\ 0 \\ -1 \end{bmatrix}$$

$$(i) \ I = \begin{bmatrix} 1.0 & 0.8 & 2.1 & 1.4 \\ 0.6 & -0.1 & 2.1 & 1.8 \\ 0.1 & -0.1 & 2.1 & 1.2 \\ 1.7 & -1.1 & -1.9 & 2.9 \end{bmatrix}, \ b_9 = \begin{bmatrix} 4.3 \\ 1.2 \\ 0.5 \\ 0.3 \end{bmatrix}, \ r_9 = \begin{bmatrix} 0.0 \\ 1.5 \\ -0.5 \\ 1.0 \end{bmatrix}$$

Exercise 3.4.5. In each of the following, is the given vector in the nullspace of the given matrix?

$$(a) \ A = \begin{bmatrix} -11 & -2 & 5 \\ -1 & 1 & 1 \end{bmatrix}, \ p = \begin{bmatrix} -7 \\ 6 \\ -13 \end{bmatrix}$$

$$(b) \ B = \begin{bmatrix} 3 & -3 & 2 \\ 1 & 1 & -3 \end{bmatrix}, \ q = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$$

$$(c) \ C = \begin{bmatrix} -5 & 0 & -2 \\ -6 & 2 & -2 \\ 0 & -5 & -1 \end{bmatrix}, \ r = \begin{bmatrix} -2 \\ -1 \\ 5 \end{bmatrix}$$

$$(d) \ D = \begin{bmatrix} -3 & -2 & 0 & 2 \\ 5 & 0 & 1 & -2 \\ 4 & -4 & 4 & 2 \end{bmatrix}, \ s = \begin{bmatrix} 6 \\ 1 \\ -10 \\ 10 \end{bmatrix}$$

$$(e) E = \begin{bmatrix} -3 & 2 & 3 & 1 \\ -3 & -2 & -1 & 4 \\ 6 & 1 & -1 & -1 \end{bmatrix}, t = \begin{bmatrix} 2 \\ -4 \\ 1 \\ -2 \end{bmatrix}$$

$$(f) F = \begin{bmatrix} -4 & -2 & 2 & -2 \\ 2 & -1 & -2 & 1 \\ 0 & 2 & 1 & 0 \\ 0 & 0 & -8 & -2 \end{bmatrix}, u = \begin{bmatrix} 11 \\ -2 \\ 4 \\ -16 \end{bmatrix}$$

Exercise 3.4.6. Given the SVDs of Exercises 3.3.2 and 3.3.3, write down an orthonormal basis for the span of the following sets of vectors.

- (a) $(-\frac{9}{5}, -4), (\frac{12}{5}, -3)$
- (b) $(-0.96, -0.72), (1.28, 0.96)$
- (c) $(7, -22, -4)/39, (-34, -38, -53)/78$
- (d) $(4, 4, 2)/33, (4, 4, 2)/11, (6, 6, 3)/11$
- (e) $(-\frac{2}{5}, \frac{11}{9}, \frac{31}{90}, \frac{4}{9}), (-\frac{2}{5}, \frac{11}{9}, \frac{31}{90}, \frac{4}{9}), (-\frac{26}{45}, -\frac{1}{3}, \frac{17}{90}, -\frac{2}{9}), (\frac{26}{45}, \frac{1}{3}, -\frac{17}{90}, \frac{2}{9})$

Exercise 3.4.7. Given any $m \times n$ matrix A .

- (a) Explain how Theorem 3.4.18 uses an SVD to find an orthonormal basis for the column space of A .
- (b) How does the same SVD give the orthonormal basis $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$ for the row space of A ? Justify your answer.
- (c) Why does the same SVD also give the orthonormal basis $\{\mathbf{v}_{r+1}, \dots, \mathbf{v}_n\}$ for the nullspace of A ? Justify.

Exercise 3.4.8. For each of the following matrices, compute an SVD with Matlab/Octave, and then use the properties of Exercise 3.4.7 to write down an orthonormal basis for the column space, the row space, and the nullspace of the matrix. (The bases, especially for the nullspace, may differ in detail depending upon your version of Matlab/Octave.)

$$(a) \begin{bmatrix} 19 & -36 & -18 \\ -3 & 12 & 6 \\ -17 & 48 & 24 \end{bmatrix}$$

$$(b) \begin{bmatrix} -12 & 0 & -4 \\ -30 & -6 & 4 \\ 34 & 22 & 8 \\ -50 & -10 & 12 \end{bmatrix}$$

$$(c) \begin{bmatrix} 0 & 0 & 0 & 0 \\ 4 & 10 & 1 & -3 \\ 2 & 6 & 0 & -2 \\ -2 & -4 & -1 & 1 \end{bmatrix}$$





(d)
$$\begin{bmatrix} -13 & 9 & 10 & -4 & -6 \\ -7 & 27 & -2 & 4 & -10 \\ -4 & 0 & 4 & 4 & -4 \\ -4 & -18 & 10 & -8 & 5 \end{bmatrix}$$



(e)
$$\begin{bmatrix} 1 & -2 & 3 & 9 \\ -1 & 5 & 0 & 0 \\ 0 & 3 & 3 & 9 \\ 2 & -9 & 1 & 3 \\ 1 & -7 & -2 & -6 \end{bmatrix}$$



(f)
$$\begin{bmatrix} 9 & 3 & 0 & -9 \\ 15 & 1 & 24 & -15 \\ -12 & -4 & 0 & 12 \\ 9 & 3 & 0 & -9 \\ -3 & -1 & 0 & 3 \\ 11 & 5 & -8 & -11 \end{bmatrix}$$



(g)
$$\begin{bmatrix} -8 & 17 & 7 & -51 & 20 \\ 5 & -2 & -1 & 15 & -2 \\ 15 & -30 & -15 & 75 & -30 \\ -2 & -1 & -5 & -33 & 8 \end{bmatrix}$$



(h)
$$\begin{bmatrix} -128 & 6 & 55 & -28 & -1 \\ 20 & 12 & -31 & 18 & -3 \\ -12 & -30 & 39 & -24 & 7 \\ -1 & 6 & -1 & 7 & -3 \end{bmatrix}$$

Exercise 3.4.9. For each of the matrices in Exercise 3.4.8, from your computed bases write down the dimension of the column space, the row space and the nullspace. Comment on how these confirm the rank theorem 3.4.32.

Exercise 3.4.10. What are the possible values for $\text{nullity}(A)$ in the following cases?

- (a) A is a 2×5 matrix.
- (b) A is a 3×3 matrix.
- (c) A is a 3×2 matrix.
- (d) A is a 4×6 matrix.
- (e) A is a 4×4 matrix.
- (f) A is a 6×5 matrix.

Exercise 3.4.11 (Cowen (1997)). Alice and Bob are taking linear algebra. One of the problems in their homework assignment is to find the nullspace of a 4×5 matrix A . In each of the following cases: are their answers consistent with each other? Give reasons.

- (a) Alice's answer is that the nullspace is spanned by $(-2, -2, 0, 2, -6)$, $(1, 5, 4, -3, 11)$, $(3, 5, 2, -4, 13)$, and $(0, -2, -2, 1, -4)$. Bob's answer is that the nullspace is spanned by $(1, 1, 0, -1, 3)$, $(-2, 0, 2, 1, -2)$, and $(-1, 3, 4, 1, 5)$.

- (b) Alice's answer is that the nullspace is spanned by $(2, -3, 1, -2, -5), (2, -7, 2, -1, -6), (1, -2, 1, 1, 0), (3, -6, 3, 3, 0)$. Bob's answer is that the nullspace is spanned by $(1, -2, 1, 1, 0), (0, 4, -1, -1, 1), (1, -1, 0, -3, -5)$.
- (c) Alice's answer is that the nullspace is spanned by $(-2, 0, -2, 4, -5), (0, 2, -2, 2, -2), (0, -2, 2, -2, 2), (-4, -12, 8, -4, 2)$. Bob's answer is that the nullspace is spanned by $(0, 2, -2, 2, -2), (-2, -4, 2, 0, -1), (1, 0, -1, 1, -3)$.
- (d) Alice's answer is that the nullspace is spanned by $(-1, 0, 0, 0, 0), (5, 3, -2, 5, 1), (-5, 1, 0, -6, -2), (4, -2, 0, 1, 8)$. Bob's answer is that the nullspace is spanned by $(1, -2, 0, -3, 4), (2, -1, 1, 3, 2), (3, 0, -1, 2, 3)$.

Exercise 3.4.12. Prove that if the columns of a matrix A are orthonormal, then they must form an orthonormal basis for the column space of A .

Exercise 3.4.13. Let A be any $m \times n$ matrix. Use an SVD to prove that every vector in the row space of A is orthogonal to every vector in the nullspace.

Exercise 3.4.14. Bachlin et al. [*IEEE Transactions on Information Technology in Biomedicine*, 14(2), 2010] explored the walking gait of people with Parkinson's Disease. Among many measurements, they measured the vertical ankle acceleration of the people when they walked. Figure 3.4 shows ten seconds of just one example: use the so-called Singular Spectrum Analysis to find the regular structures in this complex data.

Following Example 3.4.21:



(a) enter the data into Matlab/Octave;

```
time=(0.0625:0.125:9.85)'
acc=[5.34; 0.85; -1.90; -1.39; -0.99; 5.64;
-1.76; -5.90; 4.74; 1.85; -2.09; -1.16; -1.58;
5.19; -0.27; -6.54; 3.94; 2.70; -2.36; -0.94;
-1.68; 4.61; 0.79; -6.90; 3.23; 3.59; -2.65;
-0.99; -1.65; 4.09; 1.78; -7.11; 2.26; 4.27;
-2.53; -0.84; -1.84; 3.34; 2.62; -6.74; 1.54;
4.16; -2.29; -0.50; -1.97; 2.80; 2.92; -6.37;
1.09; 4.17; -2.05; -0.44; -2.03; 2.08; 3.91;
-5.84; -0.78; 4.98; -1.28; -0.94; -1.86; 0.50;
5.40; -4.19; -3.88; 5.45; 0.44; -1.71; -1.59;
-0.90; 5.86; -1.95; -5.95; 4.75; 1.90; -2.06;
-1.21; -1.61; 5.16]
```

- (b) use the `hankel()` function to form a matrix whose 66 columns are 66 'windows' of accelerations, each of length fourteen data points (of length 1.75 s);

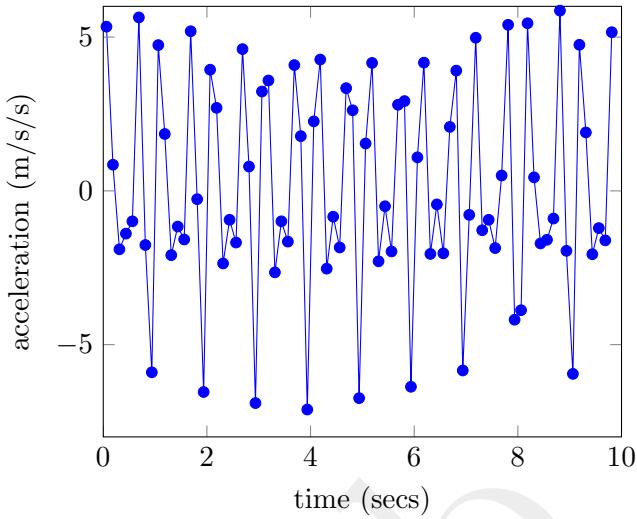


Figure 3.4: vertical acceleration of the ankle of a person walking normally, a person who has Parkinson's Disease. The data is recorded 0.125 s apart (here subsampled and smoothed for the purposes of the exercise).

- (c) compute an SVD of the matrix, and explain why the windows of measured accelerations are close to lying in a four dimensional subspace;
- (d) plot orthonormal basis vectors for the four dimensional subspace of the windowed accelerations.

Exercise 3.4.15. Consider $m \times n$ matrices bordered by zeros in the block form

$$E = \begin{bmatrix} F & O_{k \times n-\ell} \\ O_{m-k \times \ell} & O_{m-k \times n-\ell} \end{bmatrix}$$

where F is some $k \times \ell$ matrix. Given matrix F has an SVD, find an SVD of matrix E , and hence prove that $\text{rank } E = \text{rank } F$.

Exercise 3.4.16. For compatibly sized matrices A and B , use their SVDs, and the result of the previous Exercise 3.4.15 (applied to the matrix $S_A V_A^T U_B S_B$), to prove that $\text{rank}(AB) \leq \text{rank } A$ and that $\text{rank}(AB) \leq \text{rank } B$.

3.5 Project to solve inconsistent equations

Section Contents

3.5.1	Make a minimal change to the problem . . .	276
3.5.2	Compute the smallest appropriate solution .	290
3.5.3	Orthogonal projection resolves vector components	296
	Project onto a direction	297
	Project onto a subspace	299
	Orthogonal decomposition separates	310
3.5.4	Exercises	323

Agreement with experiment is the sole criterion of truth
for a physical theory.

Pierre Duhem, 1906

As well as being fundamental to engineering, scientific and computational inference, approximately solving inconsistent equations also introduces the linear transformation of projection.

The scientific method is to infer general laws from data and then validate the laws. This section addresses the inference of general laws from data. A big challenge is that data is typically corrupted by noise and errors. So this section shows how the singular value decomposition (SVD) leads to understanding the so-called least square methods.

3.5.1 Make a minimal change to the problem

Example 3.5.1 (rationalise contradictory weights). I weighed myself the other day. I weighed myself four times, each time separated by a few minutes: the scales reported my weight in kg as 84.8, 84.1, 84.7 and 84.4. What sense can we make of this apparently contradictory data? Traditionally we just average and say my weight is $x \approx (84.8 + 84.1 + 84.7 + 84.4)/4 = 84.5$ kg. Let's see this same answer from a new linear algebra view.

In the linear algebra view my weight x is an unknown and the four experimental measurements give four equations for this one unknown:

$$x = 84.8, \quad x = 84.1, \quad x = 84.7, \quad x = 84.4.$$

Despite being manifestly impossible to satisfy all four equations, let's see what linear algebra can do for us. Linear algebra writes these four equations as the matrix-vector system

$$Ax = \mathbf{b}, \quad \text{namely} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} x = \begin{bmatrix} 84.8 \\ 84.1 \\ 84.7 \\ 84.4 \end{bmatrix}.$$

The linear algebra Procedure 3.3.13 is to ‘solve’ this system, despite its contradictions, via an SVD and some intermediaries:

$$Ax = U \underbrace{S V^T x}_{=z} = b.$$

- (a) You are given that this matrix A of ones has an SVD of (perhaps check the columns of U are orthonormal)

$$A = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix}^T = USV^T.$$

- (b) Solve $Uz = b$ by computing

$$z = U^T b = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 84.8 \\ 84.1 \\ 84.7 \\ 84.4 \end{bmatrix} = \begin{bmatrix} 169 \\ -0.1 \\ 0.2 \\ 0.5 \end{bmatrix}.$$

- (c) Now try to solve $Sy = z$, that is,

$$\begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} y = \begin{bmatrix} 169 \\ -0.1 \\ 0.2 \\ 0.5 \end{bmatrix}.$$

But we cannot because the last three components in the equation are impossible: we cannot satisfy any of

$$0y = -0.1, \quad 0y = 0.2, \quad 0y = 0.5.$$

Instead of seeking an *exact* solution, ask what is the *smallest change* we can make to $z = (169, -0.1, 0.2, 0.5)$ so that we can report a solution to a slightly different problem? Answer: we *have to* adjust the last three components to zero. Also, any adjustment to the first component is unnecessary, would make the change to z bigger than necessary, and so we do not adjust the first component. Hence we find a solution to a slightly different problem by solving

$$\begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} y = \begin{bmatrix} 169 \\ 0 \\ 0 \\ 0 \end{bmatrix},$$

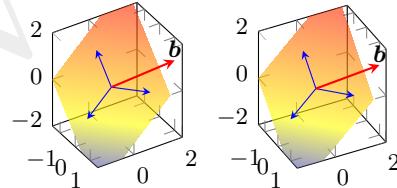
with solution $y = 84.5$. We treat this solution to a slightly different problem, as an *approximate* solution to the original problem.

- (d) Lastly, solve $V^T x = y$ by computing $x = Vy = 1y = y = 84.5 \text{ kg}$ (upon including the physical units).

This linear algebra procedure recovers the traditional answer of averaging measurements. \blacksquare

The answer to the previous Example 3.5.1 illustrates how traditional averaging emerges from trying to make sense of apparently inconsistent information. Importantly, the principle of making the smallest possible change to the intermediary \mathbf{z} is equivalent to making the smallest possible change to the original data vector \mathbf{b} . The reason is that $\mathbf{b} = U\mathbf{z}$ for an orthogonal matrix U : since U is an orthogonal matrix, multiplication by U preserves distances and angles (Theorem 3.2.39) and so the smallest possible change to \mathbf{b} is the same as the smallest possible change to \mathbf{z} . Scientists and engineers implicitly use this same ‘smallest change’ approach to approximately solve many sorts of inconsistent linear equations.

Example 3.5.2. Recall the table tennis player rating Example 3.3.11. There we found that we could not solve the equations to find some ratings because the equations were inconsistent. In our new terminology of the previous Section 3.4, the right-hand side vector \mathbf{b} is not in the column space of the matrix A (Definition 3.4.8): the stereo picture below illustrates the 2D column space spanned by the three columns of A and that the vector \mathbf{b} lies outside the column space.



Now reconsider step 3 in Example 3.3.11.

- (a) We need to interpret and ‘solve’ $S\mathbf{y} = \mathbf{z}$ which here is

$$\begin{bmatrix} 1.7321 & 0 & 0 \\ 0 & 1.7321 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{y} = \begin{bmatrix} -2.0412 \\ -2.1213 \\ 0.5774 \end{bmatrix}.$$

The third line of this system says $0y_3 = 0.5774$ which is impossible for any y_3 : we cannot have zero on the left-hand side equalling 0.5774 on the right-hand side. Instead of seeking an *exact* solution, ask what is the *smallest change* we can make to $\mathbf{z} = (-2.0412, -2.1213, 0.5774)$ so that we can report a solution, albeit to a slightly different problem? Answer: we must adjust the last component to zero. But any adjustment to the first two components is unnecessary, would make the change bigger than necessary, and so we do not adjust the

first two components. Hence find an approximate solution to the player ratings via solving

$$\begin{bmatrix} 1.7321 & 0 & 0 \\ 0 & 1.7321 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{y} = \begin{bmatrix} -2.0412 \\ -2.1213 \\ 0 \end{bmatrix}.$$

Here a general soution is $\mathbf{y} = (-1.1785, -1.2247, y_3)$ from $\mathbf{y}=\mathbf{z}(1:2) ./ \text{diag}(\mathbf{S}(1:2, 1:2))$. Varying the free variable y_3 gives equally good solutions as approximations.

- (b) Lastly, solve $V^T \mathbf{x} = \mathbf{y}$, via computing $\mathbf{x}=V(:, 1:2)*\mathbf{y}$, to determines

$$\begin{aligned} \mathbf{x} = V\mathbf{y} &= \begin{bmatrix} 0.0000 & -0.8165 & 0.5774 \\ -0.7071 & 0.4082 & 0.5774 \\ 0.7071 & 0.4082 & 0.5774 \end{bmatrix} \begin{bmatrix} -1.1785 \\ -1.2247 \\ y_3 \end{bmatrix} \\ &= \begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} + \frac{y_3}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}. \end{aligned}$$

As before, it is only the relative ratings that are important so we choose any particular (approximate) solution by setting y_3 to anything we like, such as zero. ■

The reliability and likely error of such approximate solutions are the province of Statistics courses. We focus on the geometry and linear algebra of obtaining the ‘best’ approximate solution.

Procedure 3.5.3 (approximate solution). *Obtain the so-called ‘least square’ approximate solution(s) of inconsistent equations $A\mathbf{x} = \mathbf{b}$ using an SVD and via intermediate unknowns:*

1. factorise $A = USV^T$ and set $r = \text{rank } A$ (remembering that relatively small singular values are effectively zero);
2. solve $U\mathbf{z} = \mathbf{b}$ by $\mathbf{z} = U^T\mathbf{b}$;
3. disregard the equations for $i = r + 1, \dots, m$ as errors, set $y_i = z_i/\sigma_i$ for $i = 1, \dots, r$ (as these $\sigma_i > 0$), and otherwise y_i is free for $i = r + 1, \dots, n$;
4. solve $V^T \mathbf{x} = \mathbf{y}$ to obtain a general approximate solution as $\mathbf{x} = V\mathbf{y}$.

Example 3.5.4 (round robin tournament). Consider four players (or teams) that play in a round robin sporting event: Anne, Bob, Chris and Dee. Table 3.4 summarises the results of the six games played. From these results estimate the relative player ratings of the four players. As in many real-life situations, the information appears contradictory such as Anne beats Bob who beats Dee who in turn

Table 3.4: the results of six games played in a round robin: the scores are games/goals/points scored by each when playing the others. For example, Dee beat Anne 3 to 1.

	Anne	Bob	Chris	Dee
Anne	-	3	3	1
Bob	2	-	2	4
Chris	0	1	-	2
Dee	3	0	3	-

beats Anne. Assume that the rating x_i of player i is to reflect, as best we can, the difference in scores upon playing player j : that is, pose the difference in ratings, $x_i - x_j$, should equal the difference in the scores when they play.

Solution: The first stage is to model the results by idealised mathematical equations. From Table 3.4 six games were played with the following scores. Each game then generates the shown ideal equation for the difference between two ratings.

- Anne beats Bob 3-2, so $x_1 - x_2 = 3 - 2 = 1$.
- Anne beats Chris 3-0, so $x_1 - x_3 = 3 - 0 = 3$.
- Bob beats Chris 2-1, so $x_2 - x_3 = 2 - 1 = 1$.
- Anne is beaten by Dee 1-3, so $x_1 - x_4 = 1 - 3 = -2$.
- Bob beats Dee 4-0, so $x_2 - x_4 = 4 - 0 = 4$.
- Chris is beaten by Dee 2-3, so $x_3 - x_4 = 2 - 3 = -1$.

These six equations form the linear system $A\mathbf{x} = \mathbf{b}$ where

$$A = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 3 \\ 1 \\ -2 \\ 4 \\ -1 \end{bmatrix}.$$

We cannot satisfy all these equations exactly, so we have to accept an approximate solution that estimates the ratings as best we can. The second stage is to proceed using an SVD obtained computationally.

- (a) Enter the matrix A and vector \mathbf{b} into Matlab/Octave with

```
A=[1 -1 0 0
   1 0 -1 0
   0 1 -1 0
   1 0 0 -1
   0 1 0 -1
   0 0 1 -1 ]
b=[1;3;1;-2;4;-1]
```



Then factorise matrix $A = USV^T$ with $[U, S, V] = \text{svd}(A)$ (2 d.p.):

$$\begin{aligned} U &= \\ &\begin{matrix} 0.31 & -0.26 & -0.58 & -0.26 & 0.64 & -0.15 \\ 0.07 & 0.40 & -0.58 & 0.06 & -0.49 & -0.51 \\ -0.24 & 0.67 & 0.00 & -0.64 & 0.19 & 0.24 \\ -0.38 & -0.14 & -0.58 & 0.21 & -0.15 & 0.66 \\ -0.70 & 0.13 & 0.00 & 0.37 & 0.45 & -0.40 \\ -0.46 & -0.54 & -0.00 & -0.58 & -0.30 & -0.26 \end{matrix} \\ S &= \\ &\begin{matrix} 2.00 & 0 & 0 & 0 \\ 0 & 2.00 & 0 & 0 \\ 0 & 0 & 2.00 & 0 \\ 0 & 0 & 0 & 0.00 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix} \\ V &= \\ &\begin{matrix} 0.00 & 0.00 & -0.87 & -0.50 \\ -0.62 & 0.53 & 0.29 & -0.50 \\ -0.14 & -0.80 & 0.29 & -0.50 \\ 0.77 & 0.28 & 0.29 & -0.50 \end{matrix} \end{aligned}$$

Although the first three columns of U and V may be different for you (because the first three singular values are all the same), the eventual solution is the same. The system of equations $A\mathbf{x} = \mathbf{b}$ for the ratings becomes

$$U \underbrace{S V^T}_{=z} \mathbf{x} = \mathbf{b}.$$

- (b) Solve $Uz = \mathbf{b}$ by $z = U^T \mathbf{b}$ via computing $\mathbf{z} = \mathbf{U}' * \mathbf{b}$ to get the \mathbb{R}^6 vector

$$\begin{aligned} z &= \\ &\begin{matrix} -1.27 \\ 2.92 \\ -1.15 \\ 0.93 \\ 1.76 \\ -4.07 \end{matrix} \end{aligned}$$

- (c) Now solve $Sy = z$. But the last three rows of the diagonal matrix S are zero, whereas the last three components of \mathbf{z} are non-zero: hence there is no exact solution. Instead we approximate by setting the last three components of \mathbf{z} to zero. This approximation is the *smallest change* we can make to the data of the game results that will make the results consistent.

That is, since $\text{rank } A = 3$ from the three non-zero singular values, so we approximately solve the system in Matlab/

Octave by `y=z(1:3)./diag(S(1:3,1:3)):`

```
y =
-0.63
1.46
-0.58
```

The fourth component y_4 is arbitrary.

- (d) Lastly, solve $V^T \mathbf{x} = \mathbf{y}$ as $\mathbf{x} = V\mathbf{y}$. Obtain a particular solution in Matlab/Octave by computing `x=V(:,1:3)*y:`

```
x =
0.50
1.00
-1.25
-0.25
```

Add an arbitrary multiple of the fourth column of V to get a general solution

$$\mathbf{x} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ 1 \\ -\frac{5}{4} \\ -\frac{1}{4} \end{bmatrix} + y_4 \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix}.$$

The final stage is to interpret the solution for the application. In this application the absolute ratings are not important, so we ignore y_4 (consider it zero). From the game results of Table 3.4 this analysis indicates the players' rankings are, in decreasing order, Bob, Anne, Dee, and Chris.

■

When rating players or teams based upon results, be clear the purpose. For example, is the purpose to summarise past performance? or to predict future contests? If the latter, then my limited experience suggests that one should fit the win-loss record instead of the scores.

Theorem 3.5.5 (smallest change). All approximations obtained by Procedure 3.5.3 solve the linear system $A\mathbf{x} = \tilde{\mathbf{b}}$ for the consistent right-hand side $\tilde{\mathbf{b}}$ that minimises the distance $|\tilde{\mathbf{b}} - \mathbf{b}|$.

Proof. From an SVD $A = USV^T$, Procedure 3.5.3 computes $\mathbf{z} = U^T \mathbf{b} \in \mathbb{R}^m$, that is, $\mathbf{b} = U\mathbf{z}$ as U is orthogonal. For any $\tilde{\mathbf{b}} \in \mathbb{R}^m$ let $\tilde{\mathbf{z}} = U^T \tilde{\mathbf{b}} \in \mathbb{R}^m$, that is, $\tilde{\mathbf{b}} = U\tilde{\mathbf{z}}$. Then $|\tilde{\mathbf{b}} - \mathbf{b}| = |U\tilde{\mathbf{z}} - U\mathbf{z}| = |U(\tilde{\mathbf{z}} - \mathbf{z})| = |\tilde{\mathbf{z}} - \mathbf{z}|$ as multiplication by orthogonal U preserves distances (Theorem 3.2.39). Thus minimising $|\tilde{\mathbf{b}} - \mathbf{b}|$ is equivalent

Be aware of Kenneth Arrow's Impossibility Theorem (one of the great theorems of the 20th century): *all 1D ranking systems are flawed!* Wikipedia [2014] describes the theorem this way (in the context of voting systems): that among

three or more distinct alternatives (options), no rank order voting system can convert the ranked preferences of individuals into a community-wide (complete and transitive) ranking while also meeting [four sensible] criteria ... called unrestricted domain, non-dictatorship, Pareto efficiency, and independence of irrelevant alternatives.

In rating sport players/teams:

- the “distinct alternatives” are the players/teams;
- the “ranked preferences of individuals” are the individual results of each game played; and
- the “community-wide ranking” is the assumption that we can rate each player/team by a one-dimensional numerical rating.

Arrow's theorem assures us that every such scheme must violate at least one of four sensible criteria. Every ranking scheme is thus open to criticism. But every alternative scheme will also be open to criticism by also violating the criteria.

to minimising $|\tilde{\mathbf{z}} - \mathbf{z}|$. Procedure 3.5.3 seeks to solve the diagonal system $S\mathbf{y} = \mathbf{z}$ for $\mathbf{y} \in \mathbb{R}^n$. That is, for a matrix of rank $A = r$

$$\begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_r \\ O_{(m-r) \times r} & O_{(m-r) \times (n-r)} \end{bmatrix} \mathbf{y} = \begin{bmatrix} z_1 \\ \vdots \\ z_r \\ z_{r+1} \\ \vdots \\ z_m \end{bmatrix}.$$

Procedure 3.5.3 approximately solves this inconsistent system by adjusting the right-hand side to $\tilde{\mathbf{z}} = (z_1, \dots, z_r, 0, \dots, 0) \in \mathbb{R}^m$. This change makes $|\mathbf{z} - \tilde{\mathbf{z}}|$ as small as possible because we must zero the last $(m - r)$ components of \mathbf{z} in order to obtain a consistent set of equations, and because any adjustment to the first r components of \mathbf{z} would only increase $|\mathbf{z} - \tilde{\mathbf{z}}|$. Hence the solution computed by Procedure 3.5.3 solves the consistent system $A\mathbf{x} = \tilde{\mathbf{b}}$ (with $\tilde{\mathbf{b}} = U\tilde{\mathbf{z}}$) such that $|\mathbf{b} - \tilde{\mathbf{b}}|$ is minimised. \square

Example 3.5.6 (life expectancy). Table 3.5 lists life expectancies of people born in a given year; Figure 3.5 plots the data points. Over the decades the life expectancies have increased. Let's quantify the overall trend to be able to draw, as in Figure 3.5, the best straight line to the female life expectancy. Solve the approximation problem with an SVD and confirm it gives the same solution as `A\b` in Matlab/Octave.

Table 3.5: life expectancy in years of (white) females and males born in the given years [<http://www.infoplease.com/ipa/A0005140.html>, 2014]. Used by Example 3.5.6.

year	1951	1961	1971	1981	1991	2001	2011
female	72.0	74.2	75.5	78.2	79.6	80.2	81.1
male	66.3	67.5	67.9	70.8	72.9	75.0	76.3

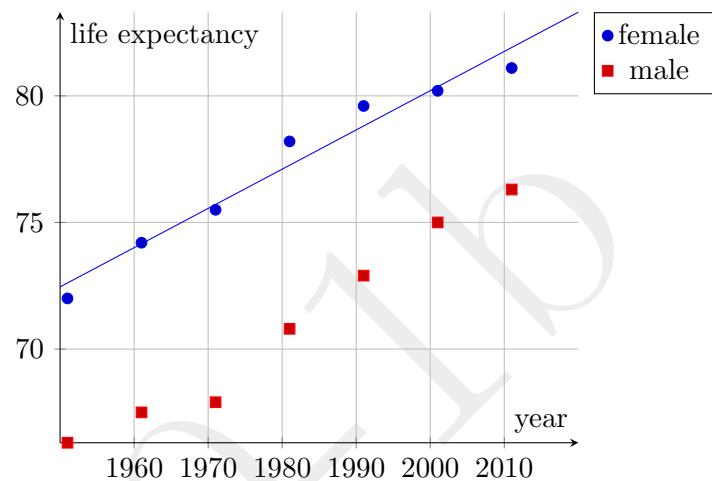


Figure 3.5: the life expectancies in years of females and males born in the given years (Table 3.5). Also plotted is the best straight line fit to the female data obtained by Example 3.5.6.

Solution: Start by posing a mathematical model: let's suppose that the life expectancy ℓ is a straight line function of year of birth: $\ell = x_1 + x_2 t$ where we need to find the coefficients x_1 and x_2 , and where t counts the number of decades since 1951, the start of the data. Table 3.5 then gives seven ideal equations to solve for x_1 and x_2 :

$$\begin{aligned}
 (1951) \quad & x_1 + 0x_2 = 72.0, \\
 (1961) \quad & x_1 + 1x_2 = 74.2, \\
 (1971) \quad & x_1 + 2x_2 = 75.5, \\
 (1981) \quad & x_1 + 3x_2 = 78.2, \\
 (1991) \quad & x_1 + 4x_2 = 79.6, \\
 (2001) \quad & x_1 + 5x_2 = 80.2, \\
 (2011) \quad & x_1 + 6x_2 = 81.1.
 \end{aligned}$$

Form these into the matrix-vector system $A\mathbf{x} = \mathbf{b}$ where

$$A = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \\ 1 & 5 \\ 1 & 6 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 72.0 \\ 74.2 \\ 75.5 \\ 78.2 \\ 79.6 \\ 80.2 \\ 81.1 \end{bmatrix}.$$

Procedure 3.5.3 then determines a best approximate solution.

- (a) Enter the matrix A and vector \mathbf{b} into Matlab/Octave, and compute an SVD of $A = USV^T$ via `[U,S,V]=svd(A)` (2 d.p.):



```

U =
 0.02  0.68 -0.38 -0.35 -0.32 -0.30 -0.27
 0.12  0.52 -0.14  0.06  0.26  0.45  0.65
 0.22  0.36  0.89 -0.09 -0.08 -0.07 -0.05
 0.32  0.20 -0.10  0.88 -0.13 -0.15 -0.16
 0.42  0.04 -0.10 -0.14  0.81 -0.23 -0.28
 0.52 -0.12 -0.09 -0.16 -0.24  0.69 -0.39
 0.62 -0.28 -0.09 -0.19 -0.29 -0.40  0.50
S =
 9.80    0
 0   1.43
 0    0
 0    0
 0    0
 0    0
 0    0
V =
 0.23  0.97
 0.97 -0.23

```

- (b) Solve $U\mathbf{z} = \mathbf{b}$ to give this first intermediary $\mathbf{z} = U^T\mathbf{b}$ via the command `z=U'*b`:

```

z =
 178.19
 100.48
 -0.05
 1.14
 1.02
 0.10
 -0.52

```

- (c) Now solve approximately $S\mathbf{y} = \mathbf{z}$. From the two non-zero singular values in S the matrix A has rank 2. So the approximation is to discard/zero (as ‘errors’) all but the first

Table 3.6: orbital periods for the eight planets of the solar system: the periods are in (Earth) days; the distance is the length of the semi-major axis of the orbits [Wikipedia, 2014].

planet	distance (Gigametres)	period (days)
Mercury	57.91	87.97
Venus	108.21	224.70
Earth	149.60	365.26
Mars	227.94	686.97
Jupiter	778.55	4332.59
Saturn	1433.45	10759.22
Uranus	2870.67	30687.15
Neptune	4498.54	60190.03

two elements of \mathbf{z} and find the best approximate \mathbf{y} via
 $\mathbf{y} = \mathbf{z}(1:2) ./ \text{diag}(\mathbf{S}(1:2, 1:2))$:

```
y =
18.19
70.31
```

(d) Solve $\mathbf{V}^T \mathbf{x} = \mathbf{y}$ by $\mathbf{x} = \mathbf{V} \mathbf{y}$ via $\mathbf{x} = \mathbf{V} * \mathbf{y}$:

```
x =
72.61
1.55
```

We soon discuss why computing $\mathbf{x} = \mathbf{A} \setminus \mathbf{b}$ gives exactly the same ‘best’ approximate solution.

Lastly, interpret the answer. The approximation gives $x_1 = 72.61$ and $x_2 = 1.55$ (2 d.p.). Since the ideal model was life expectancy $\ell = x_1 + x_2 t$ we determine a ‘best’ approximate model is $\ell \approx 72.61 + 1.55 t$ years where t is the number of decades since 1951: this is the straight line drawn in Figure 3.5. That is, females tend to live an extra 1.55 years for every decade born after 1951. For example, for females born in 2021, some seven decades after 1951, this model predicts a life expectancy of $\ell \approx 72.61 + 1.55 \times 7 = 83.46$ years.

■

Example 3.5.7 (planetary orbital periods). Table 3.6 lists each orbital period of the planets of the solar system; Figure 3.6 plots the data points as a function of the distance of the planets from the sun. Let’s infer Kepler’s law that the period grows as the distance to the power 3/2: shown by the straight line fit in Figure 3.6. Use the data for Mercury to Uranus to infer the law with an SVD, confirm it gives the same solution as $\mathbf{A} \setminus \mathbf{b}$ in Matlab/Octave, and use the fit to predict Neptune’s period from its distance.

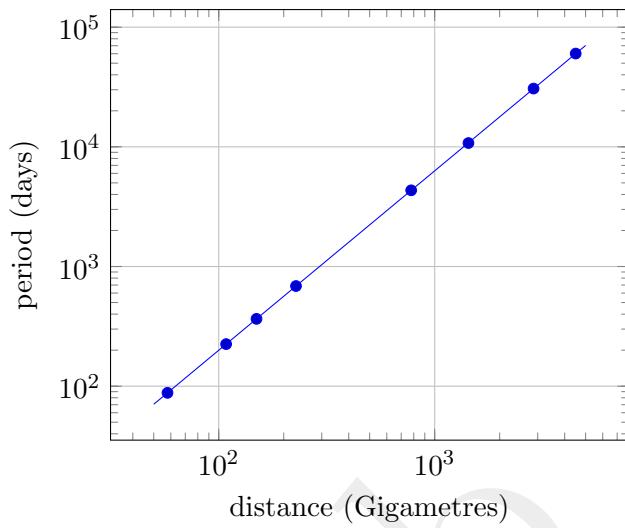


Figure 3.6: the planetary periods as a function of the distance from the data of Table 3.6: the graph is a log-log plot to show the excellent power law. Also plotted is the power law fit computed by Example 3.5.7.

Solution: Start by posing a mathematical model: Kepler's law is a power law that the i th period $p_i = c_1 d_i^{c_2}$ for some unknown coefficient c_1 and exponent c_2 . Take logarithms (to any base so let's use base 10) and seek that $\log_{10} p_i = \log_{10} c_1 + c_2 \log_{10} d_i$; that is, seek unknowns x_1 and x_2 such that $\log_{10} p_i = x_1 + x_2 \log_{10} d_i$. The first seven rows of Table 3.6 then gives seven ideal equations to solve for x_1 and x_2 :

$$\begin{aligned}x_1 + \log_{10} 57.91 x_2 &= \log_{10} 87.97, \\x_1 + \log_{10} 108.21 x_2 &= \log_{10} 224.70, \\x_1 + \log_{10} 149.60 x_2 &= \log_{10} 365.26, \\x_1 + \log_{10} 227.94 x_2 &= \log_{10} 686.97, \\x_1 + \log_{10} 778.55 x_2 &= \log_{10} 4332.59, \\x_1 + \log_{10} 1433.45 x_2 &= \log_{10} 10759.22, \\x_1 + \log_{10} 2870.67 x_2 &= \log_{10} 30687.15.\end{aligned}$$

Form these into the matrix-vector system $A\mathbf{x} = \mathbf{b}$: for simplicity recorded here to two decimal places,

$$A = \begin{bmatrix} 1 & 1.76 \\ 1 & 2.03 \\ 1 & 2.17 \\ 1 & 2.36 \\ 1 & 2.89 \\ 1 & 3.16 \\ 1 & 3.46 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1.94 \\ 2.35 \\ 2.56 \\ 2.84 \\ 3.64 \\ 4.03 \\ 4.49 \end{bmatrix}.$$

Procedure 3.5.3 then determines a best approximate solution.

- (a) Enter these matrices in Matlab/Octave by the commands, for example,



```
d=[ 57.91
    108.21
    149.60
    227.94
    778.55
    1433.45
    2870.67];
p=[ 87.97
    224.70
    365.26
    686.97
    4332.59
    10759.22
    30687.15];
A=[ones(7,1) log10(d)]
b=log10(p)
```

since the Matlab/Octave function `log10()` computes the logarithm to base 10 of each component in its argument (Table 3.2). Then compute an SVD of $A = USV^T$ via `[U,S,V]=svd(A)` (2 d.p.):

```
U =
-0.27 -0.57 -0.39 -0.38 -0.34 -0.32 -0.30
-0.31 -0.40 -0.21 -0.09  0.27  0.45  0.65
-0.32 -0.31  0.88 -0.10 -0.06 -0.04 -0.02
-0.35 -0.19 -0.11  0.90 -0.10 -0.09 -0.09
-0.41  0.14 -0.08 -0.11  0.80 -0.25 -0.30
-0.45  0.31 -0.06 -0.11 -0.26  0.67 -0.41
-0.49  0.51 -0.04 -0.11 -0.31 -0.41  0.47

S =
 7.38      0
      0  0.55
      0      0
      0      0
      0      0
      0      0
      0      0

V =
-0.35 -0.94
-0.94  0.35
```

- (b) Solve $Uz = b$ to give this first intermediary $z = U^T b$ via the command `z=U'*b`:

```
z =
-8.5507
 0.6514
```

```

0.0002
0.0004
0.0005
-0.0018
0.0012

```

- (c) Now solve approximately $S\mathbf{y} = \mathbf{z}$. From the two non-zero singular values in S the matrix A has rank two. So the approximation is to discard/zero all but the first two elements of \mathbf{z} (as an error, here all small in value). Then find the best approximate \mathbf{y} via $\mathbf{y}=\mathbf{z}(1:2) ./ \text{diag}(S(1:2,1:2))$:

```

y =
-1.1581
1.1803

```

- (d) Solve $V^T\mathbf{x} = \mathbf{y}$ by $\mathbf{x} = V\mathbf{y}$ via $\mathbf{x}=V*\mathbf{y}$:

```

x =
-0.6980
1.4991

```

Also check that computing $\mathbf{x}=A\backslash\mathbf{b}$ gives exactly the same ‘best’ approximate solution.

Lastly, interpret the answer. The approximation gives $x_1 = -0.6980$ and $x_2 = 1.4991$. Since the ideal model was the log of the period $\log_{10} p = x_1 + x_2 \log_{10} d$ we determine a ‘best’ approximate model is $\log_{10} p \approx -0.6980 + 1.4991 \log_{10} d$. Raising ten to the power of both sides gives the power law that the period $p \approx 0.2005 d^{1.4991}$ days: this is the straight line drawn in Figure 3.6. The exponent 1.4991 is within 0.1% of the exponent 3/2 that is Kepler’s law.

For example, for Neptune with a semi-major axis distance of 4498.542 Gm, extrapolating the fit predicts Neptune’s period

$$10^{-0.6980+1.4991 \log_{10} 4498.542} = 60019 \text{ days.}$$

This prediction is pleasingly close to the observed period of 60190 days. ■

Compute in Matlab/Octave. There are two separate important computational issues.

- Many books approximate solutions of $A\mathbf{x} = \mathbf{b}$ by solving the associated normal equation $(A^T A)\mathbf{x} = (A^T \mathbf{b})$. For theoretical purposes this normal equation is very useful. However, in practical computation avoid the normal equation because forming $A^T A$, and then manipulating it, is both expensive and error enhancing. For example, $\text{cond}(A^T A) = (\text{cond } A)^2$ (Exercise 3.3.14) so matrix $A^T A$ typically has a much worse condition number than matrix A (Procedure 2.2.4).

- The last two examples observe that $\mathbf{A}\backslash\mathbf{b}$ gives an answer that was identical to what the SVD procedure gives. Thus $\mathbf{A}\backslash\mathbf{b}$ can serve as a very useful short-cut to finding a best approximate solution. For non-square matrices, $\mathbf{A}\backslash\mathbf{b}$ generally does this (without comment as Matlab/Octave assumes you know what you are doing).

3.5.2 Compute the smallest appropriate solution

I'm thinking of two numbers. Their average is three.
 What are the numbers? *Cleve Moler*, The world's
 simplest impossible problem
 (1990)

The Matlab/Octave operation $\mathbf{A}\backslash\mathbf{b}$ Recall that the last few examples of the previous section observed that $\mathbf{A}\backslash\mathbf{b}$ gave an answer that was identical to the best approximate solution that the SVD procedure gave. But there are just as many circumstances when $\mathbf{A}\backslash\mathbf{b}$ is not ‘the approximate answer’ that you want. Beware.

Example 3.5.8. Use $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$ to ‘solve’ the problems of Examples 3.5.1–3.5.4.

- With Octave, observe the answer returned is the *particular* solution determined by the SVD Procedure 3.5.3 (whether approximate or exact): respectively 84.5 kg; ratings $(1, \frac{1}{3}, -\frac{4}{3})$; and ratings $(\frac{1}{2}, 1, -\frac{5}{4}, -\frac{1}{4})$.
- With Matlab, the computed answers are often different: respectively 84.5 kg (the same); ratings $(\text{NaN}, \text{Inf}, \text{Inf})$ with a warning; and ratings $(0.75, 1.25, -1, 0)$ with a warning.

How do we make sense of such differences in computed answers? ■

Recall that systems of linear equations may not have unique solutions (as in the rating examples): what does $\mathbf{A}\backslash\mathbf{b}$ compute when there are an infinite number of solutions?

- For systems of equations with the number of equations not equal to the number of variables, $m \neq n$, the Octave operation $\mathbf{A}\backslash\mathbf{b}$ computes for you the smallest solution of all valid solutions (Theorem 3.5.9): often ‘exact’ when $m < n$, or approximate when $m > n$ (Theorem 3.5.5). Using $\mathbf{A}\backslash\mathbf{b}$ is the most efficient computationally, but using the SVD helps us understand what it does.
- Matlab (R2013b) does something different with $\mathbf{A}\backslash\mathbf{b}$ in the case of fewer equations than variables, $m < n$. Matlab’s different ‘answer’ does reinforce that a choice of one solution among many is a subjective decision. But Octave’s choice of the smallest valid solution is often more appealing.

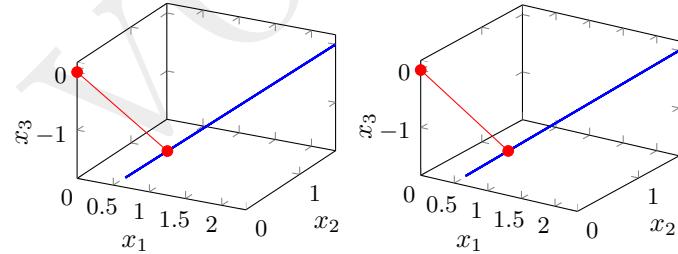
Theorem 3.5.9 (smallest solution). Obtain the smallest solution, whether exact or as an approximation, to a system of linear equations by setting to zero the free variables, $y_{r+1} = \dots = y_n = 0$, in Procedures 3.3.13 and 3.5.3.

Proof. We obtain all possible solutions, whether exact (Procedure 3.3.13) or approximate (Procedure 3.5.3), from solving $\mathbf{x} = V\mathbf{y}$. Since multiplication by orthogonal V preserves lengths (Theorem 3.2.39), the lengths of \mathbf{x} and \mathbf{y} are the same: consequently, $|\mathbf{x}|^2 = |\mathbf{y}|^2 = y_1^2 + \dots + y_r^2 + y_{r+1}^2 + \dots + y_n^2$. Now variables y_1, y_2, \dots, y_r are fixed by Procedures 3.3.13 and 3.5.3, but y_{r+1}, \dots, y_n are free to vary. Hence the smallest $|\mathbf{y}|^2$ is obtained by setting $y_{r+1} = \dots = y_n = 0$. Then this gives the particular solution $\mathbf{x} = V\mathbf{y}$ of smallest $|\mathbf{x}|$. \square

Example 3.5.10. In the table tennis ratings of Example 3.5.2 the procedure found the ratings were any of

$$\mathbf{x} = \begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} + \frac{y_3}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix},$$

as illustrated in stereo below (blue). Verify $|\mathbf{x}|$ is a minimum only when the free variable $y_3 = 0$ (a disc in the plot).



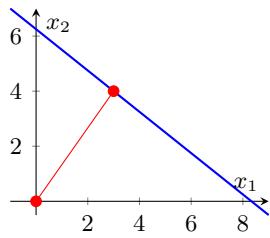
Solution:

$$\begin{aligned} |\mathbf{x}|^2 &= \mathbf{x} \cdot \mathbf{x} \\ &= \left(\begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} + \frac{y_3}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right) \cdot \left(\begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} + \frac{y_3}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right) \\ &= \begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} \cdot \begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} + \frac{2y_3}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix} + \frac{y_3^2}{3} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \\ &= \frac{26}{9} + 0y_3 + y_3^2 \end{aligned}$$

This quadratic is minimised for $y_3 = 0$. Hence the length $|\mathbf{x}|$ is minimised by the free variable $y_3 = 0$.

■

Example 3.5.11 (closest point to the origin). What is the point on the line $3x_1 + 4x_2 = 25$ that is closest to the origin? I am sure you could think of several methods, perhaps inspired by the marginal graph, but here use an SVD and Theorem 3.5.9. Confirm the Octave computation $\mathbf{A}\backslash\mathbf{b}$ gives this same closest point, but Matlab gives a different answer.



Solution: The point on the line $3x_1 + 4x_2 = 25$ closest to the origin, is the smallest solution of $3x_1 + 4x_2 = 25$. Rephrase as the matrix vector system $A\mathbf{x} = \mathbf{b}$ for matrix $A = [3 \ 4]$ and $\mathbf{b} = 25$, and apply Procedure 3.3.13.

- (a) Factorise $A = USV^T$ in Matlab/Octave via the command
 $[U,S,V]=svd([3 \ 4]):$

$$\begin{aligned} U &= \begin{bmatrix} 0.6000 & -0.8000 \\ 0.8000 & 0.6000 \end{bmatrix} \\ S &= \begin{bmatrix} 5 & 0 \end{bmatrix} \\ V &= \begin{bmatrix} 1 & \\ & 1 \end{bmatrix} \end{aligned}$$

- (b) Solve $Uz = b = 25$ which here gives $z = 25$.
(c) Solve $Sy = z = 25$ with general solution here of $y = (5, y_2)$. Obtain the smallest solution with free variable $y_2 = 0$.
(d) Solve $V^T\mathbf{x} = \mathbf{y}$ by $\mathbf{x} = Vy = V(5, 0) = (3, 4)$.

This is the smallest solution and hence the point on the line closest to the origin.

Computing $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$, which here is simply $\mathbf{x}=[3 \ 4]\backslash 25$, gives answer $\mathbf{x} = (3, 4)$ in Octave; as determined by the SVD, this point is the closest on the line to the origin. In Matlab, $\mathbf{x}=[3 \ 4]\backslash 25$ gives $\mathbf{x} = (0, 6.25)$ which the marginal graph shows is a valid solution, but not the smallest solution.

■

Example 3.5.12 (computed tomography).

A CT-scan, also called X-ray computed tomography (X-ray CT) or computerized axial tomography scan (CAT scan), makes use of computer-processed combinations of many X-ray images taken from different angles to produce cross-sectional (tomographic) images (virtual 'slices') of specific areas of a scanned object, allowing the user to see inside the object without cutting.

Wikipedia, 2015

Table 3.7: As well as the Matlab/Octave commands and operations listed in Tables 1.2, 2.3, 3.1, 3.2, and 3.3 we may invoke these functions for drawing images—functions which are otherwise not needed.

- `reshape(A,p,q)` for a $m \times n$ matrix/vector A , provided $mn = pq$, generates a $p \times q$ matrix with entries taken column-wise from A . Either p or q can be [] in which case Matlab/Octave uses $p = mn/q$ or $q = mn/p$ respectively.
- `colormap(gray)` Matlab/Octave usually draws graphs with colour, but for many images we need grayscale; this command changes the current figure to 64 shades of gray.
(`colormap(jet)` is the default, `colormap(hot)` gives colours that when also reproduced in grayscale are suitably gray, `colormap('list')` lists the available colormaps you can try)
- `imagesc(A)` where A is a $m \times n$ matrix of values draws an $m \times n$ image in the current figure window using the values of A (scaled to fit) to determine the colour from the current colormap (e.g., grayscale).
- `log(x)` where x is a matrix, vector or scalar computes the natural logarithm to the base e of each element, and returns the result(s) as a correspondingly sized matrix, vector or scalar.
- `exp(x)` where x is a matrix, vector or scalar computes the exponential of each element, and returns the result(s) as a correspondingly sized matrix, vector or scalar.

Importantly for medical diagnosis and industrial purposes, the computed answer must not have artificial features: if there is any ambiguity about the answer, then the answer shown should be the ‘greyest’—the ‘greyest’ corresponds to the mathematical smallest solution.

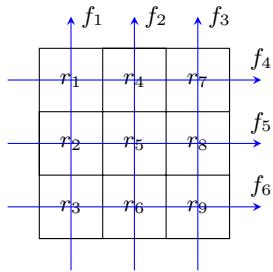
Let’s analyse a toy example.¹⁹ Suppose we divide a cross-section of a body into nine squares (large pixels) in a 3×3 grid. Inside each square the body’s material has some unknown density represented by transmission factors, r_1, r_2, \dots, r_9 as shown in the margin, that the CT-scan is to find: the fraction r_j of the incident X-ray emerges after passing through the j th square.

As indicated next in the margin, six X-ray measurements are made through the body where f_1, f_2, \dots, f_6 denote the fraction of energy in the measurements relative to the power of the X-ray beam. Thus we need to solve six equations for the nine unknown transmission factors:

r_1	r_4	r_7
r_2	r_5	r_8
r_3	r_6	r_9

$$\begin{aligned} r_1 r_2 r_3 &= f_1, & r_4 r_5 r_6 &= f_2, & r_7 r_8 r_9 &= f_3, \\ r_1 r_4 r_7 &= f_4, & r_2 r_5 r_8 &= f_5, & r_3 r_6 r_9 &= f_6. \end{aligned}$$

¹⁹ For those interested in reading further, Kress (2015) [§8] introduces the advanced, highly mathematical, approach to computerized tomography.



Computers almost always use “log” to denote the natural logarithm, so we do too. Herein unadorned “log” means the same as “ln”.

Turn such nonlinear equations into linear equations that we can handle by taking the logarithm (to any base, but here say the natural logarithm to base e) of both sides of all equations:

$$r_i r_j r_k = f_l \iff (\log r_i) + (\log r_j) + (\log r_k) = (\log f_l).$$

That is, letting new unknowns $x_i = \log r_i$ and new right-hand sides $b_i = \log f_i$, we solve six linear equations for nine unknowns:

This forms the matrix-vector system $A\mathbf{x} = \mathbf{b}$ for 6×9 matrix

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

For example, let’s find an answer for the densities when the measurements give vector $\mathbf{b} = (-0.91, -1.04, -1.54, -1.52, -1.43, -0.53)$ (all negative as they are the logarithms of fractions f_i less than one)



```
A=[1 1 1 0 0 0 0 0 0
    0 0 0 1 1 1 0 0 0
    0 0 0 0 0 0 1 1 1
    1 0 0 1 0 0 1 0 0
    0 1 0 0 1 0 0 1 0
    0 0 1 0 0 1 0 0 1]
b=[-0.91 -1.04 -1.54 -1.52 -1.43 -0.53] '
x=A\b
r=reshape(exp(x),3,3)
colormap(gray),imagesc(r)
```

The answer from Octave is (2 d.p.)

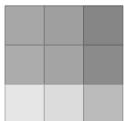
$$\mathbf{x} = (-.42, -.39, -.09, -.47, -.44, -.14, -.63, -.60, -.30).$$

These are logarithms so to get the corresponding physical transmission factors compute the exponential of each component, denoted as $\exp(\mathbf{x})$,

$$\mathbf{r} = \exp(\mathbf{x}) = (.66, .68, .91, .63, .65, .87, .53, .55, .74),$$

although it is perhaps more appealing to put these factors into the shape of the 3×3 array of pixels as in (and as illustrated in the margin)

$$\begin{bmatrix} r_1 & r_4 & r_7 \\ r_2 & r_5 & r_8 \\ r_3 & r_6 & r_9 \end{bmatrix} = \begin{bmatrix} 0.66 & 0.63 & 0.53 \\ 0.68 & 0.65 & 0.55 \\ 0.91 & 0.87 & 0.74 \end{bmatrix}.$$



Octave's answer predicts that there is less transmitting, more absorbing, denser, material to the top-right; and more transmitting, less absorbing, less dense, material to the bottom-left.

However, the answer from Matlab's $\mathbf{A}\backslash\mathbf{b}$ is (2 d.p.)

$$\mathbf{x} = (-0.91, 0, 0, -0.61, -1.43, 1.01, 0, 0, -1.54),$$

as illustrated below—the leftmost picture—which is quite different!
²⁰



Furthermore, Matlab could give other ‘answers’ as illustrated in the other pictures above. Reordering the rows in the matrix A and right-hand side \mathbf{b} does not change the system of equations. But after such reordering the answer from Matlab’s $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$ variously predicts each of the above four pictures.

The reason for such multiplicity of mathematically valid answers is that the problem is underdetermined. There are nine unknowns but only six equations, so in linear algebra there are typically an infinity of valid answers (as in Theorem 2.2.25): just five of these are illustrated above. *In this application to CT-scans* we add the additional information that we desire the answer that is the ‘greyest’, the most ‘washed out’, the answer with fewest features. Finding the answer \mathbf{x} that minimises $|\mathbf{x}|$ is a reasonable way to quantify this desire.²¹

The SVD procedure guarantees that we find such a smallest answer. Procedure 3.5.3 in Matlab/Octave gives the following process to satisfy the experimental measurements expressed in $A\mathbf{x} = \mathbf{b}$.

- (a) First, find an SVD, $A = USV^T$, via `[U,S,V]=svd(A)` and get (2 d.p.)

```
U =
-0.41 -0.00  0.82 -0.00  0.00  0.41
-0.41 -0.00 -0.41 -0.57 -0.42  0.41
-0.41 -0.00 -0.41  0.57  0.42  0.41
-0.41  0.81 -0.00  0.07 -0.09 -0.41
-0.41 -0.31 -0.00 -0.45  0.61 -0.41
-0.41 -0.50  0.00  0.38 -0.52 -0.41
```

²⁰ Matlab does give a warning in this instance (`Warning: Rank deficient, ...`), but it does not always. For example, it does not warn of issues when you ask it to solve $\frac{1}{2}(x_1 + x_2) = 3$ via `[0.5 0.5]\3`: it simply computes the ‘answer’ $\mathbf{x} = (6, 0)$.

²¹ Another possibility is to increase the number of measurements in order to increase the number of equations to match the number of unknown pixels. However, measurements are often prohibitively expensive. Further, increasing the number of measurements may tempt us to increase the resolution by having more smaller pixels: in which case we again have to deal with the same issue of more variables than known equations.



```

S =
 2.45   0   0   0   0   0   0   0   0
    0  1.73   0   0   0   0   0   0   0
    0   0  1.73   0   0   0   0   0   0
    0   0   0  1.73   0   0   0   0   0
    0   0   0   0  1.73   0   0   0   0
    0   0   0   0   0  0.00   0   0   0
V =
 -0.33  0.47  0.47  0.04 -0.05  0.03 -0.58 -0.21 -0.25
 -0.33 -0.18  0.47 -0.26  0.35 -0.36  0.49 -0.27 -0.07
 -0.33 -0.29  0.47  0.22 -0.30  0.33  0.09  0.47  0.33
 -0.33  0.47 -0.24 -0.29 -0.29 -0.48  0.11  0.37  0.26
 -0.33 -0.18 -0.24 -0.59  0.11  0.41 -0.24 -0.27  0.38
 -0.33 -0.29 -0.24 -0.11 -0.54  0.07  0.13 -0.10 -0.64
 -0.33  0.47 -0.24  0.37  0.19  0.45  0.47 -0.16 -0.00
 -0.33 -0.18 -0.24  0.07  0.59 -0.05 -0.25  0.53 -0.31
 -0.33 -0.29 -0.24  0.55 -0.06 -0.40 -0.22 -0.37  0.32

```

(b) Solve $Uz = b$ by $z=U'*b$ to find

$$z = (2.85, -0.52, 0.31, 0.05, -0.67, -0.00).$$

(c) Because the sixth singular value is zero, ignore the sixth equation: because $z_6 = 0.00$ this is only a small inconsistency error. Now set $y_i = z_i/\sigma_i$ for $i = 1, \dots, 5$ and for the smallest magnitude answer set the free variables $y_6 = y_7 = y_8 = y_9 = 0$ (Theorem 3.5.9). Obtain the non-zero values via $y=z(1:5)./diag(S(1:5,1:5))$ to find

$$y = (1.16, -0.30, 0.18, 0.03, -0.39, 0, 0, 0, 0)$$

(d) Then solve $V^T x = y$ to determine the smallest solution via $x=V(:,1:5)*y$ is $x = (-0.42, -0.39, -0.09, -0.47, -0.44, -0.14, -0.63, -0.60, -0.30)$. This is the same answer as computed by Octave's $A\backslash b$ to give the pixel image shown that has minimal artifices.



In practice, each slice of a real CT-scan would involve finding the absorption of tens of millions of pixels. That is, a CT-scan needs to best solve many systems of tens of millions of equations in tens of millions of unknowns! ■

3.5.3 Orthogonal projection resolves vector components

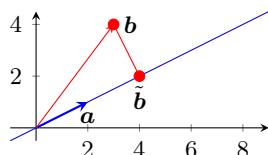
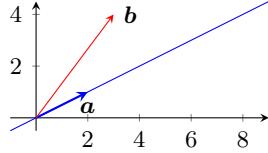
This section is optional, but does usefully support least square approximation, and provides another class of transformations for the next section 3.6 on abstract linear transformations. Orthogonal projections are also extensively used in applications.

Reconsider the task of making a minimal change to the right-hand side of a system of linear equations, and let's connect it to the so-called orthogonal projection. This important connection occurs because of the geometry that the closest point on a line or plane to another given point is the one which forms a right-angle; that is, is forms an orthogonal vector.

Project onto a direction

Example 3.5.13. Consider ‘solving’ the inconsistent system

$$\mathbf{a}x = \begin{bmatrix} 2 \\ 1 \end{bmatrix} x = \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \mathbf{b}.$$

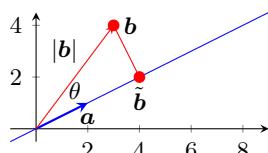


As illustrated in the margin, the impossible task is to find some multiple of the vector $\mathbf{a} = (2, 1)$ (all multiples plotted) that equals $\mathbf{b} = (3, 4)$. It cannot be done. Question: how may we change the right-hand side vector \mathbf{b} so that the task is possible? A partial answer is to replace \mathbf{b} by some vector $\tilde{\mathbf{b}}$ which is in the column space of matrix $A = [\mathbf{a}]$. But we could choose any $\tilde{\mathbf{b}}$ in the column space, so any answer is possible! Surely any answer is not acceptable. Instead, the preferred answer is to find the vector $\tilde{\mathbf{b}}$ in the column space of matrix $A = [\mathbf{a}]$ and which is closest to \mathbf{b} , as illustrated in the margin here.

The SVD approach of Procedure 3.5.3 to find $\tilde{\mathbf{b}}$ and x is the following.

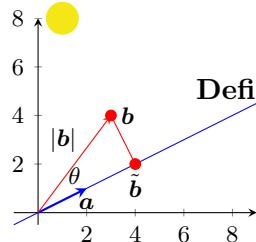
- Use $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}([2; 1])$ to find the SVD factorisation $A = \mathbf{U}\mathbf{S}\mathbf{V}^T = \begin{bmatrix} 0.89 & -0.45 \\ 0.45 & 0.89 \end{bmatrix} \begin{bmatrix} 2.24 \\ 0 \end{bmatrix} [1]^T$ (2 d.p.).
- Then $\mathbf{z} = \mathbf{U}^T \mathbf{b} = (4.47, 2.24)$.
- Treat the second component of $S\mathbf{y} = \mathbf{z}$ as an error—it is of the magnitude $|\mathbf{b} - \tilde{\mathbf{b}}|$ —to deduce $y = 4.47/2.24 = 2.00$ (2 d.p.) from the first component.
- Then $x = V^T y = 1y = 2$ solves the changed problem.

From this solution, vector $\tilde{\mathbf{b}} = \mathbf{a}x = (2, 1)2 = (4, 2)$, as is recognisable in the graphs. ■



Now let’s derive the same result but with two differences: firstly, use more elementary arguments, not the SVD; and secondly, derive the result for general vectors \mathbf{a} and \mathbf{b} (although continuing to use the same illustration). Start with the crucial observation that the closest point/vector $\tilde{\mathbf{b}}$ in the column space of $A = [\mathbf{a}]$ is such that $\mathbf{b} - \tilde{\mathbf{b}}$ is at right-angles, orthogonal, to \mathbf{a} . If $\mathbf{b} - \tilde{\mathbf{b}}$ were not orthogonal, then we would be able to slide $\tilde{\mathbf{b}}$ along the line spanned by \mathbf{a} to reduce the length of $\mathbf{b} - \tilde{\mathbf{b}}$. Thus we form a right-angle triangle with hypotenuse of length $|\mathbf{b}|$ and angle θ as shown in the margin. Trigonometry then gives the adjacent length $|\tilde{\mathbf{b}}| = |\mathbf{b}| \cos \theta$. But the angle θ is that between the given vectors \mathbf{a} and \mathbf{b} , so the dot product gives the cosine as $\cos \theta = \mathbf{a} \cdot \mathbf{b} / (|\mathbf{a}| |\mathbf{b}|)$ (Theorem 1.3.4). Hence the adjacent length $|\tilde{\mathbf{b}}| = |\mathbf{b}| \mathbf{a} \cdot \mathbf{b} / (|\mathbf{a}| |\mathbf{b}|) = \mathbf{a} \cdot \mathbf{b} / |\mathbf{a}|$. To approximately solve $\mathbf{a}x = \mathbf{b}$, replace the inconsistent $\mathbf{a}x = \mathbf{b}$ by the consistent $\mathbf{a}x = \tilde{\mathbf{b}}$: then as it is a scalar $x = |\tilde{\mathbf{b}}| / |\mathbf{a}| = \mathbf{a} \cdot \mathbf{b} / |\mathbf{a}|^2$. For Example 3.5.13, this gives ‘solution’ $x = (2, 1) \cdot (3, 4) / (2^2 + 1^2) = 10/5 = 2$ as before.

A crucial part of such solutions is the general formula for $\tilde{\mathbf{b}} = \mathbf{a}x = \mathbf{a}(\mathbf{a} \cdot \mathbf{b})/|\mathbf{a}|^2$. Geometrically the formula gives the ‘shadow’ $\tilde{\mathbf{b}}$ of vector \mathbf{b} when projected by a ‘sun’ high above the line of the vector \mathbf{a} , as illustrated schematically in the margin. As such, the formula is called an orthogonal projection.



Definition 3.5.14 (orthogonal projection onto 1D). Let $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and vector $\mathbf{u} \neq \mathbf{0}$, then the **orthogonal projection** of \mathbf{v} onto \mathbf{u} is

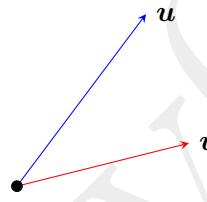
$$\text{proj}_{\mathbf{u}}(\mathbf{v}) := \mathbf{u} \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}|^2}. \quad (3.5a)$$

In the special but common case when \mathbf{u} is a unit vector,

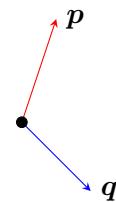
$$\text{proj}_{\mathbf{u}}(\mathbf{v}) := \mathbf{u}(\mathbf{u} \cdot \mathbf{v}). \quad (3.5b)$$

Example 3.5.15. For the following pairs of vectors: draw the named orthogonal projection; and for the given inconsistent system, determine whether the ‘best’ approximate solution is in the range $x < -1$, $-1 < x < 0$, $0 < x < 1$, or $1 < x$.

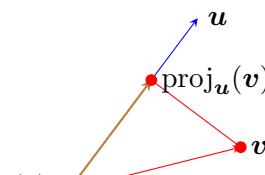
(a) $\text{proj}_{\mathbf{u}}(\mathbf{v})$ and $\mathbf{u}x = \mathbf{v}$



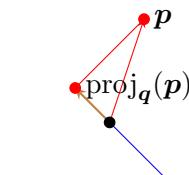
(b) $\text{proj}_{\mathbf{q}}(\mathbf{p})$ and $\mathbf{q}x = \mathbf{p}$



Solution:



(a) Draw a line perpendicular to \mathbf{u} that passes through the tip of \mathbf{v} . Then $\text{proj}_{\mathbf{u}}(\mathbf{v})$ is as shown. To ‘best solve’ $\mathbf{u}x = \mathbf{v}$, approximate the equation $\mathbf{u}x = \mathbf{v}$ by $\mathbf{u}x = \text{proj}_{\mathbf{u}}(\mathbf{v})$. Since $\text{proj}_{\mathbf{u}}(\mathbf{v})$ is smaller than \mathbf{u} and the same direction, $0 < x < 1$.



(b) Vector \mathbf{q} in $\text{proj}_{\mathbf{q}}(\mathbf{p})$ gives the direction of a line, so we can and do project onto the negative direction of \mathbf{q} . To ‘best solve’ $\mathbf{q}x = \mathbf{p}$, approximate the equation $\mathbf{q}x = \mathbf{p}$ by $\mathbf{q}x = \text{proj}_{\mathbf{q}}(\mathbf{p})$. Since $\text{proj}_{\mathbf{q}}(\mathbf{p})$ is smaller than \mathbf{q} and in the opposite direction, $-1 < x < 0$.

Example 3.5.16. For the following pairs of vectors: compute the given orthogonal projection; and hence find the ‘best’ approximate solution to the given inconsistent system.

- (a) Find $\text{proj}_{\mathbf{u}}(\mathbf{v})$ for vectors $\mathbf{u} = (3, 4)$ and $\mathbf{v} = (4, 1)$, and hence best solve $\mathbf{u}x = \mathbf{v}$.

Solution:

$$\text{proj}_{\mathbf{u}}(\mathbf{v}) = (3, 4) \frac{(3, 4) \cdot (4, 1)}{|(3, 4)|^2} = (3, 4) \frac{16}{25} = \left(\frac{48}{25}, \frac{64}{25}\right).$$

Approximate equation $\mathbf{u}x = \mathbf{v}$ by $\mathbf{u}x = \text{proj}_{\mathbf{u}}(\mathbf{v})$, that is, $(3, 4)x = \left(\frac{48}{25}, \frac{64}{25}\right)$ with solution $x = \frac{16}{25}$ (from either component).

- (b) Find $\text{proj}_{\mathbf{s}}(\mathbf{r})$ for vectors $\mathbf{r} = (1, 3)$ and $\mathbf{s} = (2, -2)$, and hence best solve $\mathbf{s}x = \mathbf{r}$.

Solution:

$$\text{proj}_{\mathbf{s}}(\mathbf{r}) = (2, -2) \frac{(2, -2) \cdot (1, 3)}{|(2, -2)|^2} = (2, -2) \frac{-4}{8} = (-1, 1).$$

Approximate equation $\mathbf{s}x = \mathbf{r}$ by $\mathbf{s}x = \text{proj}_{\mathbf{s}}(\mathbf{r})$, that is, $(2, -2)x = (-1, 1)$ with solution $x = -1/2$ (from either component).

- (c) Find $\text{proj}_{\mathbf{p}}(\mathbf{q})$ for vectors $\mathbf{p} = \left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right)$ and $\mathbf{q} = (3, 2, 1)$, and best solve $\mathbf{p}x = \mathbf{q}$.

Solution: Vector \mathbf{r} is a unit vector, so we use the simpler formula that

$$\begin{aligned} \text{proj}_{\mathbf{r}}(\mathbf{q}) &= \left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right) \left[\left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right) \cdot (3, 2, 1) \right] \\ &= \left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right) \left[1 + \frac{4}{3} + \frac{2}{3} \right] \\ &= \left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right) 3 = (1, 2, 2). \end{aligned}$$

Then ‘best solve’ equation $\mathbf{p}x = \mathbf{q}$ by the approximation $\mathbf{p}x = \text{proj}_{\mathbf{p}}(\mathbf{q})$, that is, $\left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right)x = (1, 2, 2)$ with solution $x = 3$ (from any component). ■

Project onto a subspace

The previous subsection develops a geometric view of the ‘best’ solution to the inconsistent system $\mathbf{a}x = \mathbf{b}$. The discussion introduced that the conventional ‘best’ solution—that determined by Procedure 3.5.3—is to replace \mathbf{b} by its projection $\text{proj}_{\mathbf{a}}(\mathbf{b})$, namely to solve $\mathbf{a}x = \text{proj}_{\mathbf{a}}(\mathbf{b})$. The rationale is that this is the *smallest* change to the right-hand side that enables the equation to be solved. This subsection introduces that solving inconsistent equations in more variables involves the analogous projection onto a subspace.

Definition 3.5.17 (project onto a subspace). *Let \mathbb{W} be a k -dimensional subspace of \mathbb{R}^n with an orthonormal basis $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k\}$. For any vector $\mathbf{v} \in \mathbb{R}^n$, the **orthogonal projection** of vector \mathbf{v} onto subspace \mathbb{W} is*

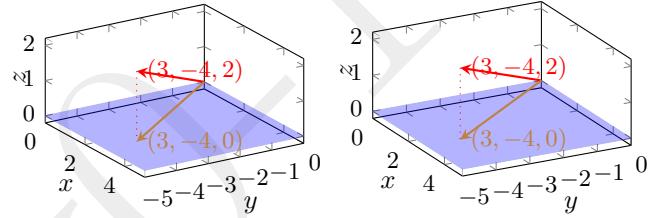
$$\text{proj}_{\mathbb{W}}(\mathbf{v}) = \mathbf{w}_1(\mathbf{w}_1 \cdot \mathbf{v}) + \mathbf{w}_2(\mathbf{w}_2 \cdot \mathbf{v}) + \cdots + \mathbf{w}_k(\mathbf{w}_k \cdot \mathbf{v}). \quad (3.6)$$

Example 3.5.18. (a) Let \mathbb{X} be the xy -plane in xyz -space, find $\text{proj}_{\mathbb{X}}(3, -4, 2)$.

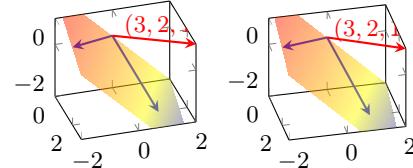
Solution: An orthogonal basis for the xy -plane (blue plane in the stereo picture below) are the two unit vectors $\mathbf{i} = (1, 0, 0)$ and $\mathbf{j} = (0, 1, 0)$. Hence

$$\begin{aligned} \text{proj}_{\mathbb{X}}(3, -4, 2) &= \mathbf{i}(\mathbf{i} \cdot (3, -4, 2)) + \mathbf{j}(\mathbf{j} \cdot (3, -4, 2)) \\ &= \mathbf{i}(3 + 0 + 0) + \mathbf{j}(0 - 4 + 0) \\ &= (3, -4, 0) \quad (\text{shown in brown}). \end{aligned}$$

That is, just set the third component of $(3, -4, 2)$ to zero.

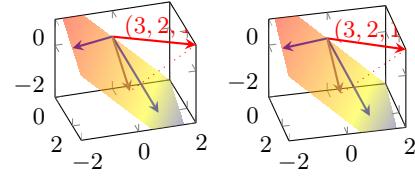


(b) For the subspace $\mathbb{W} = \text{span}\{(2, -2, 1), (2, 1, -2)\}$, determine $\text{proj}_{\mathbb{W}}(3, 2, 1)$ (these vectors and subspace are illustrated below).



Solution: Although the two vectors in the span are orthogonal (blue in the stereo picture above), they are not unit vectors. Normalise the vectors by dividing by their length $\sqrt{2^2 + (-2)^2 + 1^2} = \sqrt{2^2 + 1^2 + (-2)^2} = 3$ to find the vectors $\mathbf{w}_1 = (\frac{2}{3}, -\frac{2}{3}, \frac{1}{3})$ and $\mathbf{w}_2 = (\frac{2}{3}, \frac{1}{3}, -\frac{2}{3})$ are an orthonormal basis for \mathbb{W} (plane). Hence

$$\begin{aligned} \text{proj}_{\mathbb{W}}(3, 2, 1) &= \mathbf{w}_1(\mathbf{w}_1 \cdot (3, 2, 1)) + \mathbf{w}_2(\mathbf{w}_2 \cdot (3, 2, 1)) \\ &= \mathbf{w}_1(2 - \frac{4}{3} + \frac{1}{3}) + \mathbf{w}_2(2 + \frac{2}{3} - \frac{2}{3}) \\ &= \mathbf{w}_1 + 2\mathbf{w}_2 \\ &= (\frac{2}{3}, -\frac{2}{3}, \frac{1}{3}) + 2(\frac{2}{3}, \frac{1}{3}, -\frac{2}{3}) \\ &= (2, 0, -1) \quad (\text{shown in brown below}). \end{aligned}$$

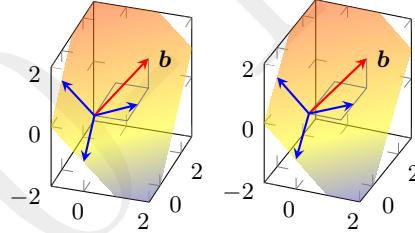


- (c) Recall the table tennis ranking Examples 3.3.11 and 3.5.2. To rank the players we seek to solve the matrix-vector system, $A\mathbf{x} = \mathbf{b}$,

$$\begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \\ 0 & 1 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}.$$

Letting \mathbb{A} denote the column space of matrix A , determine $\text{proj}_{\mathbb{A}}(\mathbf{b})$.

Solution: We need to find an orthonormal basis for the column space (the illustrated plane spanned by the three shown column vectors)—an SVD gives it to us.



Example 3.3.10 found an SVD $A = USV^T$, in Matlab/Octave via `[U,S,V]=svd(A)`, to be

```

U =
    0.4082   -0.7071    0.5774
   -0.4082   -0.7071   -0.5774
   -0.8165   -0.0000    0.5774

S =
    1.7321         0         0
        0    1.7321         0
        0         0    0.0000

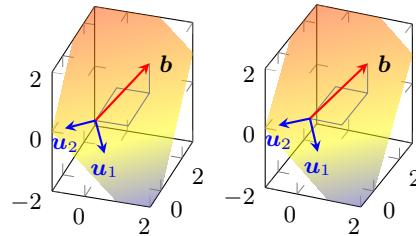
V =
    0.0000   -0.8165    0.5774
   -0.7071    0.4082    0.5774
    0.7071    0.4082    0.5774

```

Since there are only two non-zero singular values, the column space \mathbb{A} is 2D and spanned by the first two orthonormal columns of matrix U : that is, an orthonormal basis for \mathbb{A} is the two vectors (as illustrated below)

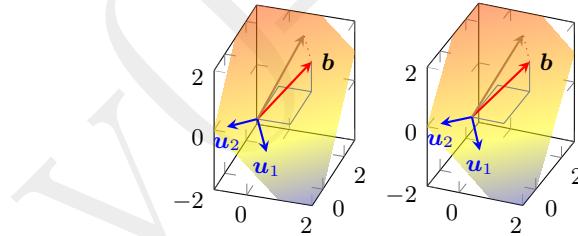
$$\mathbf{u}_1 = \begin{bmatrix} 0.4082 \\ -0.4082 \\ -0.8165 \end{bmatrix} = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 \\ -1 \\ -2 \end{bmatrix},$$

$$\mathbf{u}_2 = \begin{bmatrix} -0.7071 \\ -0.7071 \\ -0.0000 \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ -1 \\ 0 \end{bmatrix}.$$



Hence

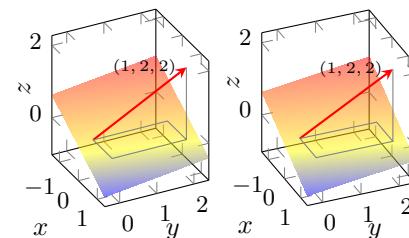
$$\begin{aligned} & \text{proj}_{\mathbb{A}}(1, 2, 2) \\ &= \mathbf{u}_1(\mathbf{u}_1 \cdot (1, 2, 2)) + \mathbf{u}_2(\mathbf{u}_2 \cdot (1, 2, 2)) \\ &= \mathbf{u}_1(1 - 2 - 4)/\sqrt{6} + \mathbf{u}_2(-1 - 2 + 0)/\sqrt{2} \\ &= -\frac{5}{\sqrt{6}}\mathbf{u}_1 - \frac{3}{\sqrt{2}}\mathbf{u}_2 \\ &= \frac{1}{6}(-5, 5, 10) + \frac{1}{2}(3, 3, 0) \\ &= \frac{1}{3}(2, 7, 5) \quad (\text{shown in brown below}). \end{aligned}$$



- (d) Find the projection of the vector $(1, 2, 2)$ onto the plane $2x - \frac{1}{2}y + 4z = 6$.

Solution: This plane is not a subspace as it does not pass through the origin. Definition 3.5.17 only defines projection onto a subspace so we cannot answer this problem (as yet).

- (e) Use an SVD to find the projection of the vector $(1, 2, 2)$ onto the plane $2x - \frac{1}{2}y + 4z = 0$ (illustrated below).



Solution: This plane does pass through the origin so it forms a subspace, call it \mathbb{P} (illustrated above). To project we need two orthonormal basis vectors. Recall that a normal to the plane is its vectors of coefficients, here $(2, -\frac{1}{2}, 4)$, so we

need to find two orthonormal vectors which are orthogonal to $(2, -\frac{1}{2}, 4)$. Further, recall that the columns of an orthogonal matrix are orthonormal (Theorem 3.2.39b), so use an SVD to find orthonormal vectors to $(2, -\frac{1}{2}, 4)$. In Matlab/Octave, set column matrix $\mathbf{n}=[2; -1/2; 4]$ then compute an SVD with $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{n})$ to find

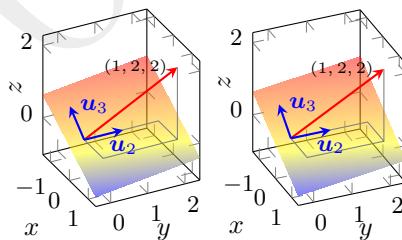


```

U =
-0.4444   0.1111  -0.8889
 0.1111   0.9914   0.0684
 -0.8889   0.0684   0.4530
S =
 4.5000
 0
 0
V = -1

```

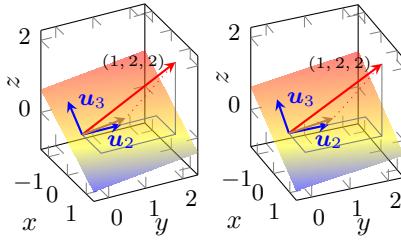
The first column $\mathbf{u}_1 = (-4, 1, -8)/9$ of orthogonal matrix U is in the direction of a normal to the plane as it must since it must be in the span of $(2, -\frac{1}{2}, 4)$. Since matrix U is orthogonal, the last two columns (say \mathbf{u}_2 and \mathbf{u}_3 , drawn in blue below) are not only orthonormal, but also orthogonal to \mathbf{u}_1 and hence an orthonormal basis for the plane \mathbb{P} .



Hence

$$\begin{aligned}
& \text{proj}_{\mathbb{P}}(1, 2, 2) \\
&= \mathbf{u}_2(\mathbf{u}_2 \cdot (1, 2, 2)) + \mathbf{u}_3(\mathbf{u}_3 \cdot (1, 2, 2)) \\
&= 2.2308 \mathbf{u}_2 + 0.1539 \mathbf{u}_3 \\
&= 2.2308(0.1111, 0.9914, 0.0684) \\
&\quad + 0.1539(-0.8889, 0.0684, 0.4530) \\
&= (0.1111, 2.2222, 0.2222) \\
&= \frac{1}{9}(1, 10, 2) \quad (\text{shown in brown below}).
\end{aligned}$$

This answer may be computed in Matlab/Octave via the two dot products $\mathbf{cs}=\mathbf{U}(:, 2:3)'*[1; 2; 2]$, giving the two coefficients 2.2308 and 0.1539, and then the linear combination $\text{proj}=\mathbf{U}(:, 2:3)*\mathbf{cs}$.



■

Example 3.5.18c determines the orthogonal projection of the given table tennis results $\mathbf{b} = (1, 2, 2)$ onto the column space of matrix A is the vector $\tilde{\mathbf{b}} = \frac{1}{3}(2, 7, 5)$. Recall that in Example 3.5.2, Procedure 3.5.3 gives the ‘approximate’ solution of the impossible $A\mathbf{x} = \mathbf{b}$ to be $\mathbf{x} = (1, \frac{1}{3}, -\frac{4}{3})$. Now see that $A\mathbf{x} = (1 - \frac{1}{3}, 1 - (-\frac{4}{3}), \frac{1}{3} - (-\frac{4}{3})) = (\frac{2}{3}, \frac{7}{3}, \frac{5}{3}) = \tilde{\mathbf{b}}$. That is, the approximate solution method of Procedure 3.5.3 solved the problem $A\mathbf{x} = \text{proj}_{\mathbb{A}}(\mathbf{b})$. The following theorem confirms this is no accident: orthogonally projecting the right-hand side onto the column space of the matrix in a system of linear equations is equivalent to solving the system with a smallest change to the right-hand side that makes it consistent.

Theorem 3.5.19. *The ‘least square’ solution(s) of the system $A\mathbf{x} = \mathbf{b}$ determined by Procedure 3.5.3 is(are) the solution(s) of $A\mathbf{x} = \text{proj}_{\mathbb{A}}(\mathbf{b})$ where \mathbb{A} denotes the column space of A .*

Proof. For any $m \times n$ matrix A , Procedure 3.5.3 first finds an SVD $A = USV^T$ and sets $r = \text{rank } A$. Second, it computes $\mathbf{z} = U^T\mathbf{b}$ but disregards z_i for $i = r + 1, \dots, m$ as errors. That is, instead of using $\mathbf{z} = U^T\mathbf{b}$ Procedure 3.5.3 solves the equations with $\tilde{\mathbf{z}} = (z_1, z_2, \dots, z_r, 0, \dots, 0)$. This vector $\tilde{\mathbf{z}}$ corresponds to a modified right-hand side $\tilde{\mathbf{b}}$ satisfying $\tilde{\mathbf{z}} = U^T\tilde{\mathbf{b}}$; that is, $\tilde{\mathbf{b}} = U\tilde{\mathbf{z}}$ as matrix U is orthogonal. Recalling \mathbf{u}_i denotes the i th column of U and that components $z_i = \mathbf{u}_i \cdot \mathbf{b}$ from $\mathbf{z} = U^T\mathbf{b}$, the matrix-vector product $\tilde{\mathbf{b}} = U\tilde{\mathbf{z}}$ is the linear combination (Example 2.3.4)

$$\begin{aligned}\tilde{\mathbf{b}} &= \mathbf{u}_1\tilde{z}_1 + \mathbf{u}_2\tilde{z}_2 + \cdots + \mathbf{u}_r\tilde{z}_r + \mathbf{u}_{r+1}0 + \cdots + \mathbf{u}_m0 \\ &= \mathbf{u}_1(\mathbf{u}_1 \cdot \mathbf{b}) + \mathbf{u}_2(\mathbf{u}_2 \cdot \mathbf{b}) + \cdots + \mathbf{u}_r(\mathbf{u}_r \cdot \mathbf{b}) \\ &= \text{proj}_{\text{span}\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}}(\mathbf{b}),\end{aligned}$$

by Definition 3.5.17 since the columns \mathbf{u}_i of U are orthonormal (Theorem 3.2.39). Theorem 3.4.26 establishes that this span is the column space \mathbb{A} of matrix A . Hence, $\tilde{\mathbf{b}} = \text{proj}_{\mathbb{A}}(\mathbf{b})$ and so Procedure 3.5.3 solves the system $A\mathbf{x} = \text{proj}_{\mathbb{A}}(\mathbf{b})$. □

Example 3.5.20. Recall Example 3.5.1 rationalises four apparently contradictory weighings: in kg the weighings are 84.8, 84.1, 84.7

and 84.4. Denoting the ‘uncertain’ weight by x , we write these weighings as the inconsistent matrix-vector system

$$Ax = \mathbf{b}, \quad \text{namely } \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} x = \begin{bmatrix} 84.8 \\ 84.1 \\ 84.7 \\ 84.4 \end{bmatrix}.$$

Let’s see that the orthogonal projection of the right-hand side onto the column space of A is the same as the minimal change of Example 3.5.1, which in turn is the well known average.

To find the orthogonal projection, observe matrix A has one column $\mathbf{a}_1 = (1, 1, 1, 1)$ so by Definition 3.5.14 the orthogonal projection

$$\begin{aligned} & \text{proj}_{\text{span}\{\mathbf{a}_1\}}(84.8, 84.1, 84.7, 84.4) \\ &= \mathbf{a}_1 \frac{\mathbf{a}_1 \cdot (84.8, 84.1, 84.7, 84.4)}{|\mathbf{a}_1|^2} \\ &= \mathbf{a}_1 \frac{84.8 + 84.1 + 84.7 + 84.4}{1 + 1 + 1 + 1} \\ &= \mathbf{a}_1 \times 84.5 \\ &= (84.5, 84.5, 84.5, 84.5). \end{aligned}$$

The projected system $Ax = (84.5, 84.5, 84.5, 84.5)$ is now consistent, with solution $x = 84.5$ kg. As in Example 3.5.1, this solution is the well-known averaging of the four weights. ■

Example 3.5.21. Recall the round robin tournament amongst four players of Example 3.5.4. To estimate the player ratings of the four players from the results of six matches we want to solve the inconsistent system $Ax = \mathbf{b}$ where

$$A = \begin{bmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 3 \\ 1 \\ -2 \\ 4 \\ -1 \end{bmatrix}.$$

Let’s see that the orthogonal projection of \mathbf{b} onto the column space of A is the same as the minimal change of Example 3.5.4.

An SVD finds an orthonormal basis for the column space \mathbb{A} of matrix A : Example 3.5.4 uses the SVD (2 d.p.)

$\mathbf{U} =$

$$\begin{array}{ccccccc} 0.31 & -0.26 & -0.58 & -0.26 & 0.64 & -0.15 \\ 0.07 & 0.40 & -0.58 & 0.06 & -0.49 & -0.51 \\ -0.24 & 0.67 & 0.00 & -0.64 & 0.19 & 0.24 \\ -0.38 & -0.14 & -0.58 & 0.21 & -0.15 & 0.66 \\ -0.70 & 0.13 & 0.00 & 0.37 & 0.45 & -0.40 \end{array}$$



```

-0.46 -0.54 -0.00 -0.58 -0.30 -0.26
S =
2.00      0      0      0
0   2.00      0      0
0      0  2.00      0
0      0      0  0.00
0      0      0      0
0      0      0      0
V = ...

```

As there are three non-zero singular values in S , the first three columns of U are an orthonormal basis for the column space \mathbb{A} . Letting u_j denote the columns of U , Definition 3.5.17 gives the orthogonal projection (2 d.p.)

$$\begin{aligned}
\text{proj}_{\mathbb{A}}(\mathbf{b}) &= u_1(u_1 \cdot \mathbf{b}) + u_2(u_2 \cdot \mathbf{b}) + u_3(u_3 \cdot \mathbf{b}) \\
&= -1.27 u_1 + 2.92 u_2 - 1.15 u_3 \\
&= (-0.50, 1.75, 2.25, 0.75, 1.25, -1.00).
\end{aligned}$$

Compute these three dot products in Matlab/Octave with $\mathbf{cs}=U(:,1:3) * \mathbf{b}$, and then compute the linear combination with $\text{proj}_{\mathbb{A}}=\mathbf{U}(:,1:3)*\mathbf{cs}$. To confirm that Procedure 3.5.3 solves $A\mathbf{x} = \text{proj}_{\mathbb{A}}(\mathbf{b})$ we check that the ratings found by Example 3.5.4, $\mathbf{x} = (\frac{1}{2}, 1, -\frac{5}{4}, -\frac{1}{4})$, satisfy $A\mathbf{x} = \text{proj}_{\mathbb{A}}(\mathbf{b})$: in Matlab/Octave compute $A*[0.50; 1.00; -1.25; -0.25]$ and see the product is $\text{proj}_{\mathbb{A}}(\mathbf{b})$. ■

Section 3.6 uses orthogonal projection as an example of a linear transformation. The section shows that a linear transformation always correspond to multiplying by a matrix, which for orthogonal projection is here WW^T .

There is an useful feature of Examples 3.5.18e and 3.5.21. In both we use Matlab/Octave to compute the projection in two steps: letting matrix W denote the matrix of appropriate columns of orthogonal U (respectively $W = U(:,2:3)$ and $W = U(:,1:3)$), first the examples compute $\mathbf{cs}=W^*\mathbf{b}$, that is, the vector $\mathbf{c} = W^T\mathbf{b}$; and second the examples compute $\text{proj}=W*\mathbf{cs}$, that is, $\text{proj}(\mathbf{b}) = W\mathbf{c}$. Combining these two steps into one (using associativity) gives

$$\text{proj}_{\mathbb{W}}(\mathbf{b}) = W\mathbf{c} = W(W^T)\mathbf{b} = (WW^T)\mathbf{b}.$$

The useful feature is that the orthogonal projection formula of Definition 3.5.17 is equivalent to the multiplication by matrix (WW^T) for an appropriate matrix W .

Theorem 3.5.22 (orthogonal projection matrix). *Let \mathbb{W} be a k -dimensional subspace of \mathbb{R}^n with an orthonormal basis $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k\}$, then for any vector $\mathbf{v} \in \mathbb{R}^n$, the orthogonal projection*

$$\text{proj}_{\mathbb{W}}(\mathbf{v}) = (WW^T)\mathbf{v} \tag{3.7}$$

for the $n \times k$ matrix $W = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_k]$.

Proof. Directly from Definition 3.5.17,

$$\text{proj}_{\mathbb{W}}(\mathbf{v}) = \mathbf{w}_1(\mathbf{w}_1 \cdot \mathbf{v}) + \mathbf{w}_2(\mathbf{w}_2 \cdot \mathbf{v}) + \dots + \mathbf{w}_k(\mathbf{w}_k \cdot \mathbf{v})$$

$$\begin{aligned}
 & \text{(using that } \mathbf{w} \cdot \mathbf{v} = \mathbf{w}^T \mathbf{v}, \text{ Example 3.1.13)} \\
 &= \mathbf{w}_1 \mathbf{w}_1^T \mathbf{v} + \mathbf{w}_2 \mathbf{w}_2^T \mathbf{v} + \cdots + \mathbf{w}_k \mathbf{w}_k^T \mathbf{v} \\
 &= (\mathbf{w}_1 \mathbf{w}_1^T + \mathbf{w}_2 \mathbf{w}_2^T + \cdots + \mathbf{w}_k \mathbf{w}_k^T) \mathbf{v}.
 \end{aligned}$$

Let the components of the vector $\mathbf{w}_j = (w_{1j}, w_{2j}, \dots, w_{nj})$, then from the matrix product Definition 3.1.8, the k products in the sum

$$\begin{aligned}
 & \mathbf{w}_1 \mathbf{w}_1^T + \mathbf{w}_2 \mathbf{w}_2^T + \cdots + \mathbf{w}_k \mathbf{w}_k^T \\
 &= \begin{bmatrix} w_{11}w_{11} & w_{11}w_{21} & \cdots & w_{11}w_{n1} \\ w_{21}w_{11} & w_{21}w_{21} & \cdots & w_{21}w_{n1} \\ \vdots & \vdots & & \vdots \\ w_{n1}w_{11} & w_{n1}w_{21} & \cdots & w_{n1}w_{n1} \end{bmatrix} \\
 &+ \begin{bmatrix} w_{12}w_{12} & w_{12}w_{22} & \cdots & w_{12}w_{n2} \\ w_{22}w_{12} & w_{22}w_{22} & \cdots & w_{22}w_{n2} \\ \vdots & \vdots & & \vdots \\ w_{n2}w_{12} & w_{n2}w_{22} & \cdots & w_{n2}w_{n2} \end{bmatrix} \\
 &+ \cdots \\
 &+ \begin{bmatrix} w_{1k}w_{1k} & w_{1k}w_{2k} & \cdots & w_{1k}w_{nk} \\ w_{2k}w_{1k} & w_{2k}w_{2k} & \cdots & w_{2k}w_{nk} \\ \vdots & \vdots & & \vdots \\ w_{nk}w_{1k} & w_{nk}w_{2k} & \cdots & w_{nk}w_{nk} \end{bmatrix}.
 \end{aligned}$$

So the (i, j) th entry of this sum is

$$\begin{aligned}
 & w_{i1}w_{j1} + w_{i2}w_{j2} + \cdots + w_{ik}w_{jk} \\
 &= w_{i1}(W^T)_{1j} + w_{i2}(W^T)_{2j} + \cdots + w_{ik}(W^T)_{kj},
 \end{aligned}$$

which, from Definition 3.1.8 again, is the (i, j) th entry of the product WW^T . Hence $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (WW^T)\mathbf{v}$. \square

Example 3.5.23. Find the matrices of the following orthogonal projections (from Example 3.5.16), and use the matrix to find the given projection.

- (a) $\text{proj}_{\mathbf{u}}(\mathbf{v})$ for vector $\mathbf{u} = (3, 4)$ and $\mathbf{v} = (4, 1)$.

Solution: First, normalise \mathbf{u} to the unit vector $\mathbf{w} = \mathbf{u}/|\mathbf{u}| = (3, 4)/5$. Second, the matrix is

$$WW^T = \mathbf{w}\mathbf{w}^T = \begin{bmatrix} \frac{3}{5} \\ \frac{4}{5} \end{bmatrix} \begin{bmatrix} \frac{3}{5} & \frac{4}{5} \end{bmatrix} = \begin{bmatrix} \frac{9}{25} & \frac{12}{25} \\ \frac{12}{25} & \frac{16}{25} \end{bmatrix}.$$

Then the projection

$$\text{proj}_{\mathbf{u}}(\mathbf{v}) = (WW^T)\mathbf{v} = \begin{bmatrix} \frac{9}{25} & \frac{12}{25} \\ \frac{12}{25} & \frac{16}{25} \end{bmatrix} \begin{bmatrix} 4 \\ 1 \end{bmatrix} = \begin{bmatrix} 48/25 \\ 64/25 \end{bmatrix}$$

- (b) $\text{proj}_s(r)$ for vector $s = (2, -2)$ and $r = (1, 1)$.

Solution: Normalise s to the unit vector $w = s/|s| = (2, -2)/(2\sqrt{2}) = (1, -1)/\sqrt{2}$, then the matrix is

$$WW^T = ww^T = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix}.$$

Consequently the projection

$$\text{proj}_s(r) = (WW^T)r = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \mathbf{0}.$$

- (c) $\text{proj}_p(q)$ for vector $p = (\frac{1}{3}, \frac{2}{3}, \frac{2}{3})$ and $q = (3, 3, 0)$.

Solution: Vector p is already a unit vector so the matrix is

$$WW^T = pp^T = \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \\ \frac{2}{3} \end{bmatrix} \begin{bmatrix} \frac{1}{3} & \frac{2}{3} & \frac{2}{3} \end{bmatrix} = \begin{bmatrix} \frac{1}{9} & \frac{2}{9} & \frac{2}{9} \\ \frac{2}{9} & \frac{4}{9} & \frac{4}{9} \\ \frac{2}{9} & \frac{4}{9} & \frac{4}{9} \end{bmatrix}.$$

Then the projection

$$\text{proj}_p(q) = (WW^T)q = \begin{bmatrix} \frac{1}{9} & \frac{2}{9} & \frac{2}{9} \\ \frac{2}{9} & \frac{4}{9} & \frac{4}{9} \\ \frac{2}{9} & \frac{4}{9} & \frac{4}{9} \end{bmatrix} \begin{bmatrix} 3 \\ 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}.$$

■

Example 3.5.24. Find the matrices of the following orthogonal projections (from Example 3.5.16).

- (a) $\text{proj}_{\mathbb{X}}(v)$ where \mathbb{X} is the xy -plane in xyz -space.

Solution: The two unit vectors $i = (1, 0, 0)$ and $j = (0, 1, 0)$ form an orthogonal basis, so matrix

$$W = [i \ j] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix},$$

hence the matrix of the projection is

$$WW^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

- (b) $\text{proj}_{\mathbb{W}}(v)$ for the subspace $\mathbb{W} = \text{span}\{(2, -2, 1), (2, 1, -2)\}$.

Solution: Now $\mathbf{w}_1 = (\frac{2}{3}, -\frac{2}{3}, \frac{1}{3})$ and $\mathbf{w}_2 = (\frac{2}{3}, \frac{1}{3}, -\frac{2}{3})$ form an orthonormal basis for \mathbb{W} , so matrix

$$W = [\mathbf{w}_1 \ \mathbf{w}_2] = \frac{1}{3} \begin{bmatrix} 2 & 2 \\ -2 & 1 \\ 1 & -2 \end{bmatrix},$$

hence the matrix of the projection is

$$\begin{aligned} WW^T &= \frac{1}{3} \begin{bmatrix} 2 & 2 \\ -2 & 1 \\ 1 & -2 \end{bmatrix} \frac{1}{3} \begin{bmatrix} 2 & -2 & 1 \\ 2 & 1 & -2 \end{bmatrix} \\ &= \frac{1}{9} \begin{bmatrix} 8 & -2 & -2 \\ -2 & 5 & -4 \\ -2 & -4 & 5 \end{bmatrix}. \end{aligned}$$

(c) The orthogonal projection onto the column space of matrix

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

Solution: The SVD of Example 3.5.18c determines an orthonormal basis is $\mathbf{u}_1 = (1, -1, -2)/\sqrt{6}$ and $\mathbf{u}_2 = (-1, -1, 0)/\sqrt{2}$. Hence the matrix of the projection is

$$\begin{aligned} WW^T &= \begin{bmatrix} \frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} \\ -\frac{2}{\sqrt{6}} & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{6}} \\ -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{2}{3} & \frac{1}{3} & -\frac{1}{3} \\ \frac{1}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{1}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix}. \end{aligned}$$

Alternatively, recall the SVD of matrix A from Example 3.3.10, and recall that the first two columns of \mathbf{U} are the orthonormal basis vectors. Hence matrix $W = \mathbf{U}(:, 1:2)$ and so Matlab/Octave computes the matrix of the projection, WW^T , via `WWT=U(:,1:2)*U(:,1:2)'` to give the answer

```
WWT =
0.6667    0.3333   -0.3333
0.3333    0.6667    0.3333
-0.3333    0.3333    0.6667
```

(d) The orthogonal projection onto the plane $2x - \frac{1}{2}y + 4z = 0$.

Solution: The SVD of Example 3.5.18e determines an orthonormal basis is the last two columns of



```
U =
-0.4444  0.1111 -0.8889
 0.1111  0.9914  0.0684
-0.8889  0.0684  0.4530
```



Hence Matlab/Octave computes the matrix of the projection with $\text{WWT}=\text{U}(:,2:3)*\text{U}(:,2:3)'$, giving the answer

```
WWT =
 0.8025  0.0494 -0.3951
 0.0494  0.9877  0.0988
-0.3951  0.0988  0.2099
```

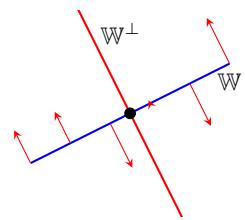
■

Orthogonal decomposition separates

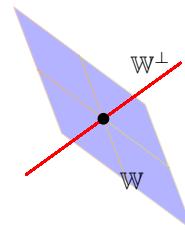
Because orthogonal projection has such a close connection to the geometry underlying important tasks such as ‘least square’ approximation (Theorem 3.5.19), this section develops further some orthogonal properties.

For any subspace \mathbb{W} of interest, it is often useful to be able to discuss the set of vectors orthogonal to all those in \mathbb{W} , called the orthogonal complement. Such a set forms a subspace, called \mathbb{W}^\perp (read as “ \mathbb{W} perp”), as illustrated below and defined by Definition 3.5.25.

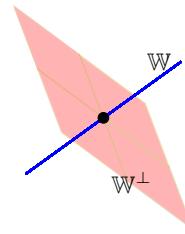
- Given the blue subspace \mathbb{W} in \mathbb{R}^2 (the origin is a black dot), consider the set of all vectors at right-angles to \mathbb{W} (drawn arrows). Move the base of these vectors to the origin, and then they all lie in the red subspace \mathbb{W}^\perp .
1. Given the blue subspace \mathbb{W} in \mathbb{R}^2 (the origin is a black dot), consider the set of all vectors at right-angles to \mathbb{W} (drawn arrows). Move the base of these vectors to the origin, and then they all lie in the red subspace \mathbb{W}^\perp .



2. Given the blue plane subspace \mathbb{W} in \mathbb{R}^3 (the origin is a black dot), the red line subspace \mathbb{W}^\perp contains all vectors orthogonal to \mathbb{W} (when drawn with their base at the origin).



3. Conversely, given the blue line subspace \mathbb{W} in \mathbb{R}^3 (the origin is a black dot), the red plane subspace \mathbb{W}^\perp contains all vectors orthogonal to \mathbb{W} (when drawn with their base at the origin).



Definition 3.5.25 (orthogonal complement). Let \mathbb{W} be a k -dimensional subspace of \mathbb{R}^n . The set of all vectors $\mathbf{u} \in \mathbb{R}^n$ (together with $\mathbf{0}$) that are each orthogonal to all vectors in \mathbb{W} is called the **orthogonal complement** \mathbb{W}^\perp (“ W -perp”); that is,

$$\mathbb{W}^\perp = \{\mathbf{u} \in \mathbb{R}^n : \mathbf{u} \cdot \mathbf{w} = 0 \text{ for all } \mathbf{w} \in \mathbb{W}\}.$$

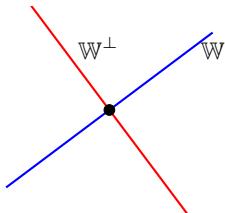
Example 3.5.26 (orthogonal complement).

- (a) Given the subspace $\mathbb{W} = \text{span}\{(3, 4)\}$, find its orthogonal complement \mathbb{W}^\perp .

Solution: Every vector in \mathbb{W} is of the form $\mathbf{w} = (3c, 4c)$. For any vector $\mathbf{v} = (u, v) \in \mathbb{R}^2$ the dot product

$$\mathbf{w} \cdot \mathbf{v} = (3c, 4c) \cdot (u, v) = c(3u + 4v),$$

is zero for all c if and only if $3u + 4v = 0$. That is, when $u = -4v/3$. Hence $\mathbf{v} = (-\frac{4}{3}v, v) = (-\frac{4}{3}, 1)v$, and so $\mathbb{W}^\perp = \text{span}\{(-\frac{4}{3}, 1)\}$.

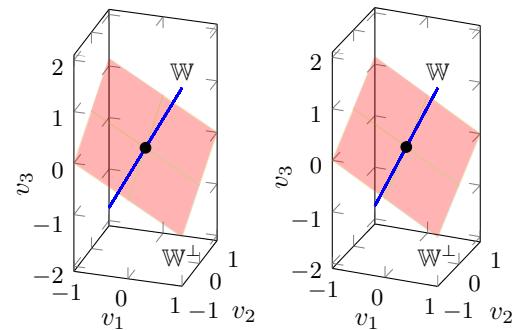


- (b) Describe the orthogonal complement \mathbb{X}^\perp to the subspace $\mathbb{X} = \text{span}\{(4, -4, 7)\}$.

Solution: Every vector in \mathbb{W} is of the form $\mathbf{w} = (4, -4, 7)c$. Seek all vectors \mathbf{v} such that $\mathbf{w} \cdot \mathbf{v} = 0$: for vectors $\mathbf{v} = (v_1, v_2, v_3)$ the inner product

$$\mathbf{w} \cdot \mathbf{v} = c(4, -4, 7) \cdot (v_1, v_2, v_3) = c(4v_1 - 4v_2 + 7v_3)$$

is zero for all c if and only if $4v_1 - 4v_2 + 7v_3 = 0$; that is, the orthogonal complement is all vectors \mathbf{v} in the plane $4v_1 - 4v_2 + 7v_3 = 0$ (illustrated in stereo below).



- (c) Describe the orthogonal complement of the set $\mathbb{W} = \{(t, t^2) : t \in \mathbb{R}\}$.

Solution: It does not exist as an orthogonal complement is only defined for a subspace, and the parabola (t, t^2) is not a subspace.

- (d) Determine the orthogonal complement of the subspace $\mathbb{W} = \text{span}\{(2, -2, 1), (2, 1, -2)\}$.

Solution: Let $\mathbf{w}_1 = (2, -2, 1)$ and $\mathbf{w}_2 = (2, 1, -2)$ then all vectors $\mathbf{w} \in \mathbb{W}$ are of the form $\mathbf{w} = c_1\mathbf{w}_1 + c_2\mathbf{w}_2$ for all c_1 and c_2 . Every vector $\mathbf{v} \in \mathbb{W}^\perp$ must satisfy, for all c_1 and c_2 ,

$$\mathbf{w} \cdot \mathbf{v} = (c_1\mathbf{w}_1 + c_2\mathbf{w}_2) \cdot \mathbf{v} = c_1\mathbf{w}_1 \cdot \mathbf{v} + c_2\mathbf{w}_2 \cdot \mathbf{v} = 0.$$

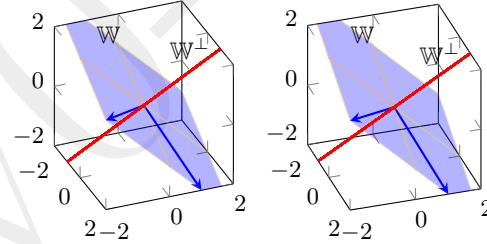
The only way to be zero for all c_1 and c_2 is for both $\mathbf{w}_1 \cdot \mathbf{v} = 0$ and $\mathbf{w}_2 \cdot \mathbf{v} = 0$. For vectors $\mathbf{v} = (v_1, v_2, v_3)$ these two equations become the pair

$$2v_1 - 2v_2 + v_3 = 0 \quad \text{and} \quad 2v_1 + v_2 - 2v_3 = 0.$$

Adding twice the second to the first, and subtracting the first from the second give the equivalent pair

$$6v_1 - 3v_3 = 0 \quad \text{and} \quad 3v_2 - 3v_3 = 0.$$

Both are satisfied for all $v_3 = t$ with $v_1 = t/2$ and $v_2 = t$. Therefore all possible \mathbf{v} in the complement \mathbb{W}^\perp are those in the form of the line $\mathbf{v} = (\frac{1}{2}t, t, t)$. That is, $\mathbb{W}^\perp = \text{span}\{(\frac{1}{2}, 1, 1)\}$ (as illustrated below in stereo).



Example 3.5.27. Prove $\{\mathbf{0}\}^\perp = \mathbb{R}^n$ and $(\mathbb{R}^n)^\perp = \{\mathbf{0}\}$.

Solution: • The only vector in $\{\mathbf{0}\}$ is $\mathbf{w} = \mathbf{0}$. Since all vectors $\mathbf{v} \in \mathbb{R}^n$ satisfy $\mathbf{w} \cdot \mathbf{v} = \mathbf{0} \cdot \mathbf{v} = 0$, by Definition 3.5.25 $\{\mathbf{0}\}^\perp = \mathbb{R}^n$.

- Certainly, $\mathbf{0} \in (\mathbb{R}^n)^\perp$ as $\mathbf{w} \cdot \mathbf{0} = 0$ for all vectors $\mathbf{w} \in \mathbb{R}^n$. Establish there are no others by contradiction. Assume a non-zero vector $\mathbf{v} \in (\mathbb{R}^n)^\perp$. Now set $\mathbf{w} = \mathbf{v} \in \mathbb{R}^n$, then $\mathbf{w} \cdot \mathbf{v} = \mathbf{v} \cdot \mathbf{v} = |\mathbf{v}|^2 \neq 0$ as \mathbf{v} is non-zero. Consequently, a non-zero \mathbf{v} cannot be in the complement. Thus $(\mathbb{R}^n)^\perp = \{\mathbf{0}\}$.

These examples find orthogonal complements that are lines, planes, or the entire space. These suggest that an orthogonal complement is generally a subspace as proved next.

Theorem 3.5.28 (orthogonal complement is subspace). *For any subspace \mathbb{W} of \mathbb{R}^n , the orthogonal complement \mathbb{W}^\perp is a subspace of \mathbb{R}^n . Further, the intersection $\mathbb{W} \cap \mathbb{W}^\perp = \{\mathbf{0}\}$; that is, the zero vector is the only vector in both \mathbb{W} and \mathbb{W}^\perp .*

Proof. Recall the Definition 3.4.2 of a subspace: we need to establish \mathbb{W}^\perp has the zero vector, and is closed under addition and scalar multiplication.

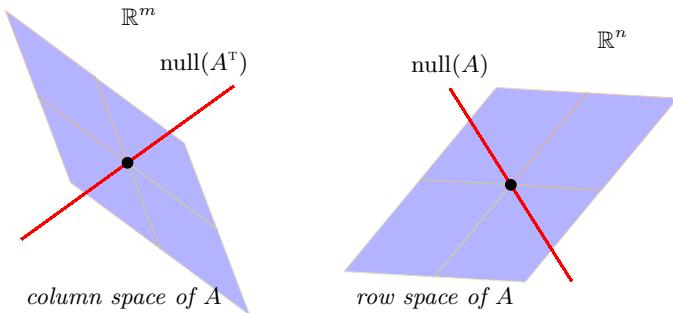
- For all $\mathbf{w} \in \mathbb{W}$, $\mathbf{0} \cdot \mathbf{w} = 0$ and so $\mathbf{0} \in \mathbb{W}^\perp$.
- Let $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{W}^\perp$, then for all $\mathbf{w} \in \mathbb{W}$ the dot product $(\mathbf{v}_1 + \mathbf{v}_2) \cdot \mathbf{w} = \mathbf{v}_1 \cdot \mathbf{w} + \mathbf{v}_2 \cdot \mathbf{w} = 0 + 0 = 0$ and so $\mathbf{v}_1 + \mathbf{v}_2 \in \mathbb{W}^\perp$.
- Let scalar $c \in \mathbb{R}$ and $\mathbf{v} \in \mathbb{W}^\perp$, then for all $\mathbf{w} \in \mathbb{W}$ the dot product $(c\mathbf{v}) \cdot \mathbf{w} = c(\mathbf{v} \cdot \mathbf{w}) = c0 = 0$ and so $c\mathbf{v} \in \mathbb{W}^\perp$.

Hence, by Definition 3.4.2, \mathbb{W}^\perp is a subspace.

Further, as they are both subspaces, the zero vector is in both \mathbb{W} and \mathbb{W}^\perp . Let vector \mathbf{u} be any vector in both \mathbb{W} and \mathbb{W}^\perp . As $\mathbf{u} \in \mathbb{W}^\perp$, by Definition 3.5.25 $\mathbf{u} \cdot \mathbf{w} = 0$ for all $\mathbf{w} \in \mathbb{W}$. But $\mathbf{u} \in \mathbb{W}$ also, so using this for \mathbf{w} in the previous equation gives $\mathbf{u} \cdot \mathbf{u} = 0$; that is, $|\mathbf{u}|^2 = 0$. Hence vector \mathbf{u} has to be the zero vector (Theorem 1.1.11). That is, $\mathbb{W} \cap \mathbb{W}^\perp = \{\mathbf{0}\}$. \square

More specifically, orthogonal complements often arise and are usefully written as the nullspace of a matrix.

Theorem 3.5.29 (nullspace complementarity). *For any $m \times n$ matrix A , the column space of A has $\text{null}(A^T)$ as its orthogonal complement in \mathbb{R}^m . That is, identifying the columns of matrix $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$, and denoting the column space by $\mathbb{A} = \text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$, then the orthogonal complement $\mathbb{A}^\perp = \text{null}(A^T)$. Further, $\text{null}(A)$ in \mathbb{R}^n is the orthogonal complement of the row space of A .*



Proof. First, by Definition 3.5.25, any vector $\mathbf{v} \in \mathbb{A}^\perp$ is orthogonal to all vectors in the column space of A , in particular it is orthogonal to the columns of A :

$$\mathbf{a}_1 \cdot \mathbf{v} = 0, \ \mathbf{a}_2 \cdot \mathbf{v} = 0, \ \dots, \ \mathbf{a}_k \cdot \mathbf{v} = 0$$

$$\begin{aligned}
 &\iff \mathbf{a}_1^T \mathbf{v} = 0, \mathbf{a}_2^T \mathbf{v} = 0, \dots, \mathbf{a}_k^T \mathbf{v} = 0 \\
 &\iff \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_k^T \end{bmatrix} \mathbf{v} = \mathbf{0} \\
 &\iff A^T \mathbf{v} = \mathbf{0} \\
 &\iff \mathbf{v} \in \text{null}(A^T).
 \end{aligned}$$

That is, $\mathbb{A}^\perp \subseteq \text{null}(A^T)$. Second, for any $\mathbf{v} \in \text{null}(A^T)$, recall that by Definition 3.4.8 for any vector \mathbf{w} in the column space of A , there exists a linear combination $\mathbf{w} = c_1 \mathbf{a}_1 + c_2 \mathbf{a}_2 + \dots + c_n \mathbf{a}_n$. Then

$$\begin{aligned}
 \mathbf{w} \cdot \mathbf{v} &= (c_1 \mathbf{a}_1 + c_2 \mathbf{a}_2 + \dots + c_n \mathbf{a}_n) \cdot \mathbf{v} \\
 &= c_1(\mathbf{a}_1 \cdot \mathbf{v}) + c_2(\mathbf{a}_2 \cdot \mathbf{v}) + \dots + c_n(\mathbf{a}_n \cdot \mathbf{v}) \\
 &= c_1 0 + c_2 0 + \dots + c_n 0 \quad (\text{from above } \iff) \\
 &= 0,
 \end{aligned}$$

and so by Definition 3.5.25 vector $\mathbf{v} \in \mathbb{A}^\perp$; that is, $\text{null}(A^T) \subseteq \mathbb{A}^\perp$. Putting these two together, $\text{null}(A^T) = \mathbb{A}^\perp$.

Lastly, that the $\text{null}(A)$ in \mathbb{R}^n is the orthogonal complement of the row space of A follows from applying the above result to the matrix A^T . \square

Example 3.5.30.

- (a) The subspace $\mathbb{W} = \text{span}\{(2, -1)\}$, find its orthogonal complement \mathbb{W}^\perp .

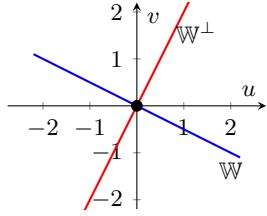
Solution: Here the subspace \mathbb{W} is the column space of matrix

$$W = \begin{bmatrix} 2 \\ -1 \end{bmatrix}.$$

To find $\mathbb{W}^\perp = \text{null}(W^T)$, solve $W^T \mathbf{v} = \mathbf{0}$, that is, for vectors $\mathbf{v} = (u, v)$

$$\begin{bmatrix} 2 & -1 \end{bmatrix} \mathbf{v} = 2u - v = 0.$$

All solutions are $v = 2u$ (as illustrated). Hence $\mathbf{v} = (u, 2u) = (1, 2)u$, and so $\mathbb{W}^\perp = \text{span}\{(1, 2)\}$.

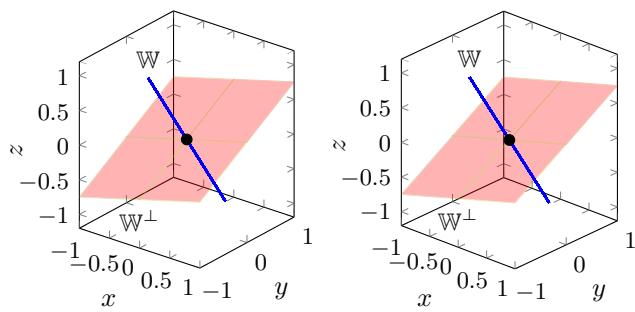


- (b) Describe the subspace of \mathbb{R}^3 whose orthogonal complement is the plane $-\frac{1}{2}x - y + 2z = 0$.

Solution: The equation of the plane in \mathbb{R}^3 may be written

$$\begin{bmatrix} -\frac{1}{2} & -1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = 0, \quad \text{that is } W^T \mathbf{v} = 0$$

for matrix $W = [\mathbf{w}_1]$ and vectors $\mathbf{w}_1 = (-\frac{1}{2}, -1, 2)$ and $\mathbf{v} = (x, y, z)$. Since the plane is the nullspace of matrix W^T , the plane must be the orthogonal complement of the line $\mathbb{W} = \text{span}\{\mathbf{w}_1\}$ (as illustrated).



- (c) Find the orthogonal complement to the column space of matrix

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

Solution: Recall from Section 2.1 that for such small problems we find all solutions of $A^T \mathbf{v} = \mathbf{0}$ by algebraic elimination; that is,

$$\begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix} \mathbf{v} = \mathbf{0} \iff \begin{cases} v_1 + v_2 = 0, \\ -v_1 + v_3 = 0, \\ -v_2 - v_3 = 0, \end{cases} \iff \begin{cases} v_2 = -v_1, \\ v_3 = v_1, \\ -v_2 - v_3 = v_1 - v_1 = 0. \end{cases}$$

Therefore all solutions of $A^T \mathbf{v} = \mathbf{0}$ are of the form $v_1 = t$, $v_2 = -v_1 = -t$ and $v_3 = v_1 = t$; that is, $\mathbf{v} = (1, -1, 1)t$. Hence the orthogonal complement is $\text{span}\{(1, -1, 1)\}$.

- (d) Describe the orthogonal complement of the subspace spanned by the four vectors $(1, 1, 0, 1, 0, 0)$, $(-1, 0, 1, 0, 1, 0)$, $(0, -1, -1, 0, 0, 1)$ and $(0, 0, 0, -1, -1, -1)$.

Solution: Arrange these vectors as the four columns of a matrix, say

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

then seek $\text{null}(A^T)$, the solutions of $A^T \mathbf{x} = \mathbf{0}$. Adapt Procedure 3.3.13 to solve $A^T \mathbf{x} = \mathbf{0}$:

- Example 3.5.4 computed an SVD $A = USV^T$ for this matrix A , which gives the SVD $A^T = VS^TU^T$ for the transpose where (2 d.p.)



```

U =
  0.31 -0.26 -0.58 -0.26  0.64 -0.15
  0.07  0.40 -0.58  0.06 -0.49 -0.51
 -0.24  0.67  0.00 -0.64  0.19  0.24
 -0.38 -0.14 -0.58  0.21 -0.15  0.66
 -0.70  0.13  0.00  0.37  0.45 -0.40
 -0.46 -0.54 -0.00 -0.58 -0.30 -0.26

S =
  2.00    0    0    0
    0  2.00    0    0
    0    0  2.00    0
    0    0    0  0.00
    0    0    0    0
    0    0    0    0

V = ...

```

- ii. $Vz = \mathbf{0}$ determines $z = \mathbf{0}$.
- iii. $S^T y = z = \mathbf{0}$ determines $y_1 = y_2 = y_3 = 0$ as there are three non-zero singular values, and y_4 , y_5 and y_6 are free variables; that is, $\mathbf{y} = (0, 0, 0, y_4, y_5, y_6)$.
- iv. Denoting the columns of U by $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_6$, the solutions of $U^T x = y$ are $x = Uy = \mathbf{u}_4 y_4 + \mathbf{u}_5 y_5 + \mathbf{u}_6 y_6$.

That is, the orthogonal complement is the three dimensional subspace $\text{span}\{\mathbf{u}_4, \mathbf{u}_5, \mathbf{u}_6\}$ in \mathbb{R}^6 , where (2 d.p.)

$$\begin{aligned}\mathbf{u}_4 &= (-0.26, 0.06, -0.64, 0.21, 0.37, -0.58), \\ \mathbf{u}_5 &= (0.64, -0.49, 0.19, -0.15, 0.45, -0.30), \\ \mathbf{u}_6 &= (-0.15, -0.51, 0.24, 0.66, -0.40, -0.26).\end{aligned}$$

■

In the previous Example 3.5.30d there are three non-zero singular values in the first three rows of S . These three nonzero singular values determine that the first three columns of U form a basis for the column space of A . The example argues that the remaining three columns of U form a basis for the orthogonal complement of the column space. That is, all six of the columns of the orthogonal U are used in either the column space or its complement. This is generally true.

Example 3.5.31. Recall the cases of Example 3.5.30.

- 3.5.30a : $\dim \mathbb{W} + \dim \mathbb{W}^\perp = 1 + 1 = 2 = \dim \mathbb{R}^2$.
- 3.5.30b : $\dim \mathbb{W} + \dim \mathbb{W}^\perp = 1 + 2 = 3 = \dim \mathbb{R}^3$.
- 3.5.30c : $\dim \mathbb{W} + \dim \mathbb{W}^\perp = 2 + 1 = 3 = \dim \mathbb{R}^3$.
- 3.5.30d : $\dim \mathbb{W} + \dim \mathbb{W}^\perp = 3 + 3 = 6 = \dim \mathbb{R}^6$.

■

Recall the Rank Theorem 3.4.32 connects the dimension of a space with the dimensions of a nullspace and column space of a matrix. Since a subspace is closely connected to matrices, and its orthogonal complement is connected to nullspaces, then the Rank Theorem should say something general here.

Theorem 3.5.32. *Let \mathbb{W} be a subspace of \mathbb{R}^n , then $\dim \mathbb{W} + \dim \mathbb{W}^\perp = n$; equivalently, $\dim \mathbb{W}^\perp = n - \dim \mathbb{W}$.*

Proof. Let the columns of a matrix W form an orthonormal basis for the subspace \mathbb{W} (Theorem 3.4.23 asserts a basis exists). Theorem 3.5.29 establishes that $\mathbb{W}^\perp = \text{null}(W^T)$. Equating dimensions of both sides,

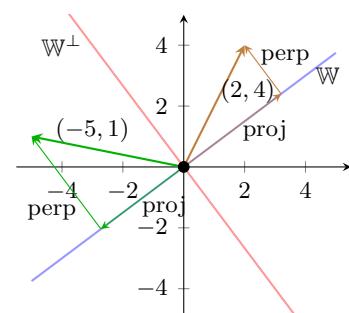
$$\begin{aligned}\dim \mathbb{W}^\perp &= \text{nullity}(W^T) \quad (\text{from Defn. 3.4.29}) \\ &= n - \text{rank}(W^T) \quad (\text{from Rank Thm. 3.4.32}) \\ &= n - \text{rank}(W) \quad (\text{from Thm. 3.3.19}) \\ &= n - \dim \mathbb{W} \quad (\text{from Thm. 3.4.18}),\end{aligned}$$

as required. □

Since the dimension of the space is the sum of the dimension of a subspace plus the dimension of its orthogonal complement, surely we must be able to separate vectors into two corresponding components.

Example 3.5.33. Recall from Example 3.5.26a that subspace $\mathbb{W} = \text{span}\{(3, 4)\}$ has orthogonal complement $\mathbb{W}^\perp = \text{span}\{(-4, 3)\}$, as illustrated below.

As shown, for example, write the brown vector $(2, 4) = (3, 2, 2, 4) + (-1.2, 1.6) = \text{proj}_{\mathbb{W}}(2, 4) + \text{perp}$, where here the vector $\text{perp} = (-1.2, 1.6) \in \mathbb{W}^\perp$. Indeed, any vector can be written as a component in subspace \mathbb{W} and a component in the orthogonal complement \mathbb{W}^\perp (Theorem 3.5.38).



For example, write the green vector $(-5, 1) = (-2.72, -2.04) + (-2.28, 3.04) = \text{proj}_{\mathbb{W}}(-5, 1) + \text{perp}$, where in this case the vector $\text{perp} = (-2.28, 3.04) \in \mathbb{W}^\perp$. ■

Further, such a separation can be done for any pair of complementary subspaces \mathbb{W} and \mathbb{W}^\perp within any space \mathbb{R}^n . To proceed, let's define what is meant by "perp" in such a context.

Definition 3.5.34 (perpendicular component). *Let \mathbb{W} be a subspace of \mathbb{R}^n .*

*For any vector $\mathbf{v} \in \mathbb{R}^n$, the **perpendicular component** of \mathbf{v} to \mathbb{W} is the vector $\text{perp}_{\mathbb{W}}(\mathbf{v}) := \mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})$.*

Example 3.5.35. (a) Let the subspace \mathbb{W} be the span of $(-2, -3, 6)$. Find the perpendicular component to \mathbb{W} of the vector $(4, 1, 3)$. Verify the perpendicular component lies in the plane $-2x - 3y + 6z = 0$.

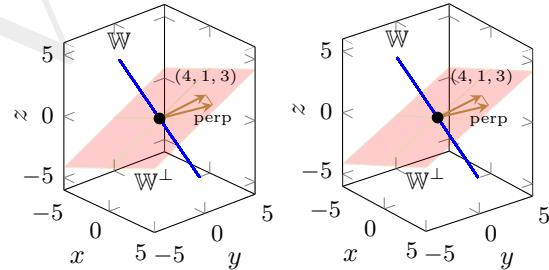
Solution: Projection is easiest with a unit vector. Obtain a unit vector to span \mathbb{W} by normalising the basis vector to $\mathbf{w}_1 = (-2, -3, 6)/\sqrt{2^2 + 3^2 + 6^2} = (-2, -3, 6)/7$. Then

$$\begin{aligned}\text{perp}_{\mathbb{W}}(4, 1, 3) &= (4, 1, 3) - \mathbf{w}_1(\mathbf{w}_1 \cdot (4, 1, 3)) \\ &= (4, 1, 3) - \mathbf{w}_1(-8 - 3 + 18)/7 \\ &= (4, 1, 3) - \mathbf{w}_1 = (30, 10, 15)/7.\end{aligned}$$

For $(x, y, z) = (30, 10, 15)/7$ we find

$$-2x - 3y + 6z = \frac{1}{7}(-60 - 30 + 90) = \frac{1}{7}0 = 0.$$

Hence $\text{perp}_{\mathbb{W}}(4, 1, 3)$ lies in the plane $-2x - 3y + 6z = 0$ (which is the orthogonal complement \mathbb{W}^\perp , as illustrated in stereo below).

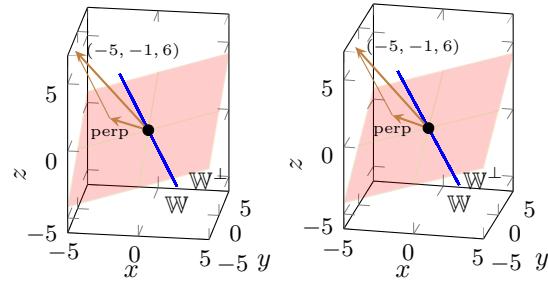


(b) For the vector $(-5, -1, 6)$ find its perpendicular component to the subspace \mathbb{W} spanned by $(-2, -3, 6)$. Verify the perpendicular component lies in the plane $-2x - 3y + 6z = 0$.

Solution: As in the previous case, use the basis vector $\mathbf{w}_1 = (-2, -3, 6)/7$. Then

$$\begin{aligned}\text{perp}_{\mathbb{W}}(-5, -1, 6) &= (-5, -1, 6) - \mathbf{w}_1(\mathbf{w}_1 \cdot (-5, -1, 6)) \\ &= (-5, -1, 6) - \mathbf{w}_1(10 + 3 + 36)/7 \\ &= (-5, -1, 6) - \mathbf{w}_1 7 = (-3, 2, 0).\end{aligned}$$

For $(x, y, z) = (-3, 2, 0)$ we find $-2x - 3y + 6z = 6 - 6 + 0 = 0$. Hence $\text{perp}_{\mathbb{W}}(-5, -1, 6)$ lies in the plane $-2x - 3y + 6z = 0$ (which is the orthogonal complement \mathbb{W}^\perp , as illustrated below in stereo).



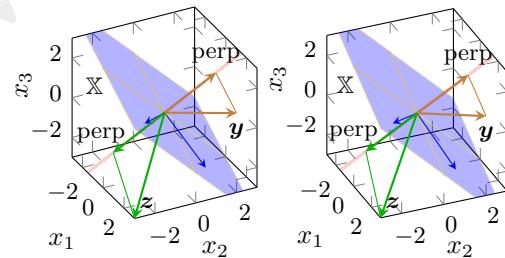
- (c) Let the subspace $\mathbb{X} = \text{span}\{(2, -2, 1), (2, 1, -2)\}$. Determine the perpendicular component of each of the two vectors $\mathbf{y} = (3, 2, 1)$ and $\mathbf{z} = (3, -3, -3)$.

Solution: Computing $\text{proj}_{\mathbb{X}}$ needs an orthonormal basis for \mathbb{X} (Definition 3.5.17). The two vectors in the span are orthogonal, so normalise them to $\mathbf{w}_1 = (2, -2, 1)/3$ and $\mathbf{w}_2 = (2, 1, -2)/3$.

- Then for the first vector $\mathbf{y} = (3, 2, 1)$,

$$\begin{aligned}\text{perp}_{\mathbb{X}}(\mathbf{y}) &= \mathbf{y} - \text{proj}_{\mathbb{X}}(\mathbf{y}) \\ &= \mathbf{y} - \mathbf{w}_1(\mathbf{w}_1 \cdot \mathbf{y}) - \mathbf{w}_2(\mathbf{w}_2 \cdot \mathbf{y}) \\ &= \mathbf{y} - \mathbf{w}_1(6 - 4 + 1)/3 - \mathbf{w}_2(6 + 2 - 2)/3 \\ &= \mathbf{y} - \mathbf{w}_1 - 2\mathbf{w}_2 \\ &= (3, 2, 1) - (2, -2, 1)/3 - (4, 2, -4)/3 \\ &= (1, 2, 2)\end{aligned}$$

(as illustrated below in brown).



- For the second vector $\mathbf{z} = (3, -3, -3)$ (in green in the picture above),

$$\begin{aligned}\text{perp}_{\mathbb{X}}(\mathbf{z}) &= \mathbf{z} - \text{proj}_{\mathbb{X}}(\mathbf{z}) \\ &= \mathbf{z} - \mathbf{w}_1(\mathbf{w}_1 \cdot \mathbf{z}) - \mathbf{w}_2(\mathbf{w}_2 \cdot \mathbf{z}) \\ &= \mathbf{z} - \mathbf{w}_1(6 + 6 - 3)/3 - \mathbf{w}_2(6 - 3 + 6)/3 \\ &= \mathbf{z} - 3\mathbf{w}_1 - 3\mathbf{w}_2 \\ &= (2, -2, -2) - (2, -2, 1) - (2, 1, -2) \\ &= (-1, -2, -2).\end{aligned}$$

■

As seen in all these examples, the perpendicular component of a vector always lies in the orthogonal complement to the subspace (as suggested by the naming).

Theorem 3.5.36 (perpendicular component is orthogonal). *Let \mathbb{W} be a subspace of \mathbb{R}^n and let \mathbf{v} be any vector in \mathbb{R}^n , then the perpendicular component $\text{perp}_{\mathbb{W}}(\mathbf{v}) \in \mathbb{W}^\perp$.*

Proof. Let vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k$ form an orthonormal basis for the subspace \mathbb{W} (the basis exists by Theorem 3.4.23). Let the $n \times k$ matrix $W = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_k]$ so subspace \mathbb{W} is the column space of matrix W , then Theorem 3.5.29 asserts we just need to check that $W^T \text{perp}_{\mathbb{W}}(\mathbf{v}) = \mathbf{0}$. Consider

$$\begin{aligned} W^T \text{perp}_{\mathbb{W}}(\mathbf{v}) &= W^T [\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})] && (\text{from Defn. 3.5.34}) \\ &= W^T [\mathbf{v} - (WW^T)\mathbf{v}] && (\text{from Thm. 3.5.22}) \\ &= W^T\mathbf{v} - W^T(WW^T)\mathbf{v} && (\text{by distributivity}) \\ &= W^T\mathbf{v} - (W^TW)W^T\mathbf{v} && (\text{by associativity}) \\ &= W^T\mathbf{v} - I_k W^T\mathbf{v} && (\text{only if } W^TW = I_k) \\ &= W^T\mathbf{v} - W^T\mathbf{v} = \mathbf{0}. \end{aligned}$$

Hence $\text{perp}_{\mathbb{W}}(\mathbf{v}) \in \text{null}(W^T)$ and so is in \mathbb{W}^\perp (by Theorem 3.5.29).

But this proof only holds if $W^TW = I_k$. To establish this identity, use the same argument as in the proof of Theorem 3.2.39a \iff 3.2.39b:

$$\begin{aligned} W^TW &= \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_k^T \end{bmatrix} [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_k] \\ &= \begin{bmatrix} \mathbf{w}_1^T \mathbf{w}_1 & \mathbf{w}_1^T \mathbf{w}_2 & \dots & \mathbf{w}_1^T \mathbf{w}_k \\ \mathbf{w}_2^T \mathbf{w}_1 & \mathbf{w}_2^T \mathbf{w}_2 & \dots & \mathbf{w}_2^T \mathbf{w}_k \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{w}_k^T \mathbf{w}_1 & \mathbf{w}_k^T \mathbf{w}_2 & \dots & \mathbf{w}_k^T \mathbf{w}_k \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{w}_1 \cdot \mathbf{w}_1 & \mathbf{w}_1 \cdot \mathbf{w}_2 & \dots & \mathbf{w}_1 \cdot \mathbf{w}_k \\ \mathbf{w}_2 \cdot \mathbf{w}_1 & \mathbf{w}_2 \cdot \mathbf{w}_2 & \dots & \mathbf{w}_2 \cdot \mathbf{w}_k \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{w}_k \cdot \mathbf{w}_1 & \mathbf{w}_k \cdot \mathbf{w}_2 & \dots & \mathbf{w}_k \cdot \mathbf{w}_k \end{bmatrix} \\ &= I_k \end{aligned}$$

as vectors $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k$ are an orthonormal set (from Definition 3.2.31, $\mathbf{w}_i \cdot \mathbf{w}_j = 0$ for $i \neq j$ and $|\mathbf{w}_i|^2 = \mathbf{w}_i \cdot \mathbf{w}_i = 1$). \square

Example 3.5.37. The previous examples' calculation of the perpendicular component confirm that $\mathbf{v} = \text{proj}_{\mathbb{W}}(\mathbf{v}) + \text{perp}_{\mathbb{W}}(\mathbf{v})$, where we now know that $\text{perp}_{\mathbb{W}}$ is orthogonal to \mathbb{W} :

3.5.33 : $(2, 4) = (3.2, 2.4) + (-1.2, 1.6)$ and
 $(-5, 1) = (-2.72, -2.04) + (-2.28, 3.04);$

3.5.35b : $(-5, -1, 6) = (-2, -3, 6) + (-3, 2, 0);$

3.5.35c : $(3, 2, 1) = (2, 0, -1) + (1, 2, 2)$ and
 $(3, -3, -3) = (4, -1, -1) + (-1, -2, -2).$

■

Given any subspace \mathbb{W} , this theorem indicates that every vector can be written as a sum of two vectors: one in the subspace \mathbb{W} ; and one in its orthogonal complement \mathbb{W}^\perp .

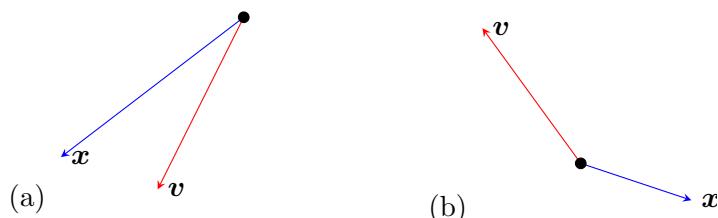
Theorem 3.5.38 (orthogonal decomposition). *Let \mathbb{W} be a subspace of \mathbb{R}^n and vector $\mathbf{v} \in \mathbb{R}^n$, then there exist unique vectors $\mathbf{w} \in \mathbb{W}$ and $\mathbf{n} \in \mathbb{W}^\perp$ such that vector $\mathbf{v} = \mathbf{w} + \mathbf{n}$; this particular sum is called an **orthogonal decomposition** of \mathbf{v} .*

Proof. • First establish existence. By Definition 3.5.34, $\text{perp}_{\mathbb{W}}(\mathbf{v}) = \mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})$, so it follows that $\mathbf{v} = \text{proj}_{\mathbb{W}}(\mathbf{v}) + \text{perp}_{\mathbb{W}}(\mathbf{v}) = \mathbf{w} + \mathbf{n}$ when we set $\mathbf{w} = \text{proj}_{\mathbb{W}}(\mathbf{v}) \in \mathbb{W}$ and $\mathbf{n} = \text{perp}_{\mathbb{W}}(\mathbf{v}) \in \mathbb{W}^\perp$.

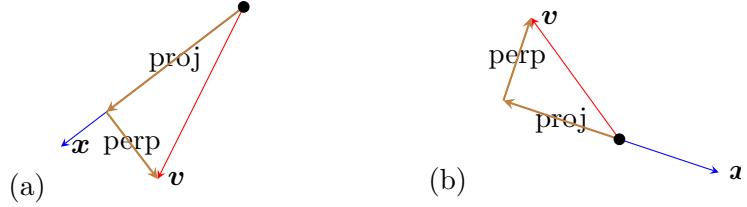
- Second establish uniqueness by contradiction. Suppose there is another decomposition $\mathbf{v} = \mathbf{w}' + \mathbf{n}'$ where $\mathbf{w}' \in \mathbb{W}$ and $\mathbf{n}' \in \mathbb{W}^\perp$. Then $\mathbf{w} + \mathbf{n} = \mathbf{v} = \mathbf{w}' + \mathbf{n}'$. Rearranging gives $\mathbf{w} - \mathbf{w}' = \mathbf{n}' - \mathbf{n}$. By closure of a subspace under vector addition (Definition 3.4.2), the left-hand side is in \mathbb{W} and the right-hand side is in \mathbb{W}^\perp , so the two sides must be both in \mathbb{W} and \mathbb{W}^\perp . The zero vector is the only common vector to the two subspaces (Theorem 3.5.28), so $\mathbf{w} - \mathbf{w}' = \mathbf{n}' - \mathbf{n} = \mathbf{0}$, and hence both $\mathbf{w} = \mathbf{w}'$ and $\mathbf{n} = \mathbf{n}'$. That is, the decomposition must be unique.

□

Example 3.5.39. For each pair of the shown subspaces $\mathbb{X} = \text{span}\{\mathbf{x}\}$ and vectors \mathbf{v} , draw the decomposition of vector \mathbf{v} into the sum of vectors in \mathbb{X} and \mathbb{X}^\perp .



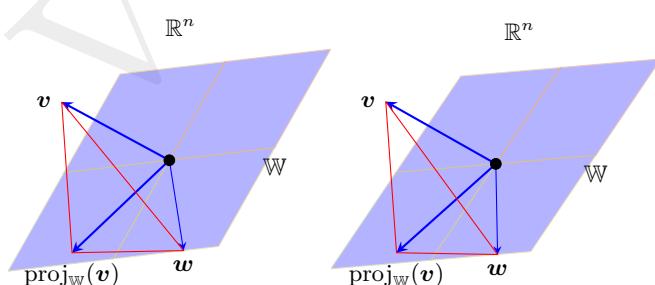
Solution: In each case, the two brown vectors shown are the decomposition, with $\text{proj} \in \mathbb{X}$ and $\text{perp} \in \mathbb{X}^\perp$.



■

In two or even three dimensions, that a decomposition has such a nice physical picture is appealing. What is powerful is that the same decomposition works in any number of dimensions: it works no matter how complicated the scenario, no matter how much data. In particular, the next theorem gives a geometric view of the ‘least square’ solution of Procedure 3.5.3: in that procedure the minimal change of the right-hand side \mathbf{b} to make the linear equation $A\mathbf{x} = \mathbf{b}$ consistent (Theorem 3.5.5) is also to be viewed as the projection of the right-hand side \mathbf{b} to the *closest* point in the columns space of the matrix. That is, the ‘least square’ procedure solves $A\mathbf{x} = \text{proj}_{\mathbb{A}}(\mathbf{b})$.

Theorem 3.5.40 (best approximation). *Given any vector \mathbf{v} in \mathbb{R}^n , and any subspace \mathbb{W} in \mathbb{R}^n , then $\text{proj}_{\mathbb{W}}(\mathbf{v})$ is the closest vector in \mathbb{W} to \mathbf{v} ; that is, $|\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})| \leq |\mathbf{v} - \mathbf{w}|$ for all $\mathbf{w} \in \mathbb{W}$.*



Proof. For any vector $\mathbf{w} \in \mathbb{W}$, consider the triangle formed by the three vectors $\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})$, $\mathbf{v} - \mathbf{w}$ and $\mathbf{w} - \text{proj}_{\mathbb{W}}(\mathbf{v})$ (the stereo illustration above schematically plots this triangle in red). This is a right-angle triangle as $\mathbf{w} - \text{proj}_{\mathbb{W}}(\mathbf{v}) \in \mathbb{W}$ by closure of the subspace \mathbb{W} , and as $\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v}) = \text{perp}_{\mathbb{W}}(\mathbf{v}) \in \mathbb{W}^\perp$. Then Pythagoras tells us

$$\begin{aligned} |\mathbf{v} - \mathbf{w}|^2 &= |\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})|^2 + |\mathbf{w} - \text{proj}_{\mathbb{W}}(\mathbf{v})|^2 \\ &\geq |\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})|^2. \end{aligned}$$

Hence $|\mathbf{v} - \mathbf{w}| \geq |\mathbf{v} - \text{proj}_{\mathbb{W}}(\mathbf{v})|$ for all $\mathbf{w} \in \mathbb{W}$. □

3.5.4 Exercises

Exercise 3.5.1. During an experiment on the strength of beams, you and your partner measure the length of a crack in the beam. With vernier callipers you measure the crack as 17.8 mm long, whereas your partner measures it as 18.4 mm long.

- Write this information as a simple matrix-vector equation for the as yet to be decided length x , and involving the matrix $A = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$.
- Confirm that an SVD of the matrix is

$$A = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} \sqrt{2} \\ 0 \end{bmatrix} [1]^T.$$

- Use the SVD to ‘best’ solve the inconsistent equations and estimate the length of the crack is $x \approx 18.1$ mm—the average of the two measurements.

Exercise 3.5.2. In measuring the amount of butter to use in cooking a recipe you weigh a container to have 207 g (grams), then a bit later weigh it at 211 g. Wanting to be more accurate you weigh the butter container a third time and find 206 g.

- Write this information as a simple matrix-vector equation for the as yet to be decided weight x , and involving the matrix $B = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}$.
- Confirm that an SVD of the matrix is

$$B = \begin{bmatrix} \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & 0 & -\frac{2}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \end{bmatrix} \begin{bmatrix} \sqrt{3} \\ 0 \\ 0 \end{bmatrix} [1]^T.$$

- Use the SVD to ‘best’ solve the inconsistent equations and estimate the butter container weighs $x \approx 208$ g—the average of the three measurements.

Exercise 3.5.3. Consider three sporting teams that play each other in a round robin event: Newark, Yonkers, and Edison: Yonkers beat Newark, 2 to 0; Edison beat Newark 5 to 2; and Edison beat Yonkers 3 to 2. Assuming the teams can be rated, and based upon the scores, write three equations that ideally relate the team ratings. Use Procedure 3.5.3 to estimate the ratings.

Table 3.8: the results of six matches played in a round robin: the scores are games/goals/points scored by each when playing the others. For example, Clapham beat Acton 4 to 2. Exercise 3.5.5 rates these teams.

	Acton	Barbican	Clapham	Dalston
Acton	-	2	2	6
Barbican	2	-	2	6
Clapham	4	4	-	5
Dalston	3	1	0	-

Table 3.9: the results of ten matches played in a round robin: the scores are games/goals/points scored by each when playing the others. For example, Atlanta beat Concord 3 to 2. Exercise 3.5.6 rates these teams.

	Atlanta	Boston	Concord	Denver	Frankfort
Atlanta	-	3	3	2	5
Boston	2	-	2	3	8
Concord	2	7	-	6	1
Denver	2	2	1	-	5
Frankfort	2	3	6	7	-

Exercise 3.5.4. Consider three sporting teams that play each other in a round robin event: Adelaide, Brisbane, and Canberra: Adelaide beat Brisbane, 5 to 1; Canberra beat Adelaide 5 to 0; and Brisbane beat Canberra 2 to 1. Assuming the teams can be rated, and based upon the scores, write three equations that ideally relate the team ratings. Use Procedure 3.5.3 to estimate the ratings.

Exercise 3.5.5. Consider four sporting teams that play each other in a round robin event: Acton, Barbican, Clapham, and Dalston. Table 3.8 summarises the results of the six matches played. Assuming the teams can be rated, and based upon the scores, write six equations that ideally relate the team ratings. Use Procedure 3.5.3 to estimate the ratings.

Exercise 3.5.6. Consider five sporting teams that play each other in a round robin event: Atlanta, Boston, Concord, Denver, and Frankfort. Table 3.9 summarises the results of the ten matches played. Assuming the teams can be rated, and based upon the scores, write ten equations that ideally relate the team ratings. Use Procedure 3.5.3 to estimate the ratings.

Exercise 3.5.7. Consider six sporting teams in a weekly competition: Algeria, Botswana, Chad, Djibouti, Ethiopia, and Gabon. In the first week of competition Algeria beat Botswana 3 to 0, Chad and Djibouti drew 3 all, and Ethiopia beat Gabon 4 to 2. In the second week of competition Chad beat Algeria 4 to 2, Botswana beat

Table 3.10: the body weight and heat production of various mammals (Kleiber 1947). Recall that numbers written as xEn denote the number $x \cdot 10^n$.

animal	body weight (kg)	heat prod. (kcal/day)
mouse	1.95E-2	3.06E+0
rat	2.70E-1	2.61E+1
cat	3.62E+0	1.56E+2
dog	1.28E+1	4.35E+2
goat	2.58E+1	7.50E+2
sheep	5.20E+1	1.14E+3
cow	5.34E+2	7.74E+3
elephant	3.56E+3	4.79E+4

Ethiopia 4 to 2, Djibouti beat Gabon 4 to 3. In the third week of competition Algeria beat Ethiopia 4 to 1, Botswana beat Djibouti 3 to 1, Chad drew with Gabon 2 all. Assuming the teams can be rated, and based upon the scores after the first three weeks, write nine equations that ideally relate the ratings of the six teams. Use Procedure 3.5.3 to estimate the ratings.

Discover power laws Exercises 3.5.8–3.5.11 use log-log plots as examples of the scientific inference of some surprising patterns in nature. These are simple examples of what, in modern parlance, might be termed ‘data mining’, ‘knowledge discovery’ or ‘artificial intelligence’.

Exercise 3.5.8. Table 3.10 lists data on the body weight and heat production of various mammals. As in Example 3.5.7, use this data to discover Kleiber’s power law that $(\text{heat}) \propto (\text{weight})^{3/4}$. Graph the data on a log-log plot, fit a straight line, check the correspondence between neglected parts of the right-hand side and the quality of the graphical fit, describe the power law.



Exercise 3.5.9. Table 3.11 lists data on river lengths and basin areas of some Russian rivers. As in Example 3.5.7, use this data to discover Hack’s exponent in the power law that $(\text{length}) \propto (\text{area})^{0.58}$. Graph the data on a log-log plot, fit a straight line, check the correspondence between neglected parts of the right-hand side and the quality of the graphical fit, describe the power law.



Exercise 3.5.10. Find for another country some river length and basin area data akin to that of Exercise 3.5.9. Confirm, or otherwise, Hack’s exponent for your data. Write a short report.



Exercise 3.5.11. The area-length relationship of a river is expected to be $(\text{length}) \propto (\text{area})^{1/2}$, so it is a puzzle as to why one consistently

Table 3.11: river length and basin area for some Russian rivers (Arnold 2014, p.154).

river	basin area (km ²)	length (km)
Moscow	17640	502
Protva	4640	275
Vorya	1160	99
Dubna	5474	165
Istra	2120	112
Nara	2170	156
Pakhra	2720	129
Skhodnya	259	47
Volgusha	265	40
Pekhorka	513	42
Setun	187	38
Yauza	452	41

Table 3.12: given a measuring stick of some length, compute the length of the west coast of Britain (Mandelbrot 1982, Plate 33).

stick length (km)	coast length (km)
10.4	2845
30.2	2008
99.6	1463
202.	1138
532.	929
933.	914

finds Hack's exponent (e.g., Exercise 3.5.9). The puzzle may be answered by the surprising notion that rivers do not have a well defined length! L. F. Richardson first established this remarkable notion for coastlines.

Table 3.12 lists data on the length of the west coast of Britain computed by using measuring sticks of various lengths: as one uses a smaller and smaller measuring stick, more and more bays and inlets are resolved and measured which increases the computed coast length. As in Example 3.5.7, use this data to discover the power law that the coast length $\propto (\text{stick})^{-1/4}$. Hence as the measuring stick length goes to 'zero', the coast length goes to 'infinity'! Graph the data on a log-log plot, fit a straight line, check the correspondence between neglected parts of the right-hand side and the quality of the graphical fit, describe the power law.



Exercise 3.5.12.

Table 3.13 lists nine of the US universities ranked by

Table 3.13: a selection of nine of the US universities ranked in 2013 by *The Center for Measuring University Performance* [http://mup.asu.edu/research_data.html]. Among others, these particular nine universities are listed by the Center in the following order. The other three columns give just three of the attributes used to create their ranked list.

Institution	Research fund(M\$)	Faculty awards	Median SAT U/G
Stanford University	868	45	1455
Yale University	654	45	1500
University of California, San Diego	1004	35	1270
University of Pittsburgh, Pittsburgh	880	22	1270
Vanderbilt University	535	19	1440
Pennsylvania State University, University Park	677	20	1195
Purdue University, West Lafayette	520	22	1170
University of Utah	410	12	1110
University of California, Santa Barbara	218	11	1205

I do not condone nor endorse such naive one dimensional ranking of complex multi-faceted institutions. This exercise simply illustrates a technique that deconstructs such a credulous endeavour.

an organisation in 2013, in the order they list. The table also lists three of the attributes used to generate the ranked list. Find a formula that approximately reproduces the listed ranking from the three given attributes.

- (a) Pose the rank of the i th institution is a linear function of the attributes and a constant, say the rank $i = x_1 f_i + x_2 a_i + x_3 s_i + x_4$ where f_i denotes the funding, a_i denotes the awards, and s_i denotes the SAT.
- (b) Form a system of nine equations that we would ideally solve to find the coefficients $\mathbf{x} = (x_1, x_2, x_3, x_4)$.
- (c) Enter the data into Matlab/Octave and find a best approximate solution (you should find the formula is roughly that $\text{rank} \approx 97 - 0.01f_i - 0.07a_i - 0.01s_i$).
- (d) Discuss briefly how well the approximation reproduces the ranking of the list.



Exercise 3.5.13. For each of the following lines and planes, use an SVD to find the point closest to the origin in the line or plane. For the lines in 2D, draw a graph to show the answer is correct.

(a) $5x_1 - 12x_2 = 169$

(b) $x_1 - 2x_2 = 5$

(c) $-x + y = -1$

(d) $-2p - 3q = 5$

(e) $2x_1 - 3x_2 + 6x_3 = 7$

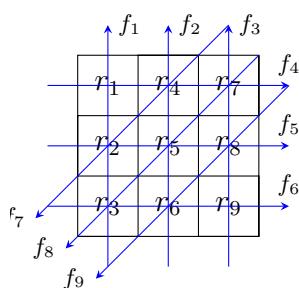
(f) $x_1 + 4x_2 - 8x_3 = 27$

(g) $2u_2 - 5u_2 - 3u_3 = -2$

(h) $q_1 + q_2 - 5q_3 = 2$

Exercise 3.5.14. Following the computed tomography Example 3.5.12, predict the densities in the body if the fraction of X-ray energy measured in the six paths is $\mathbf{f} = (0.9, 0.2, 0.8, 0.9, 0.8, 0.2)$ respectively. Draw an image of your predictions. Which region is the most absorbing (least transmitting)?

Exercise 3.5.15. In an effort to remove the need for requiring the ‘smallest’, most washed out, CT-scan, you make three more measurements, as illustrated in the margin, so that you obtain nine equations for the nine unknowns.



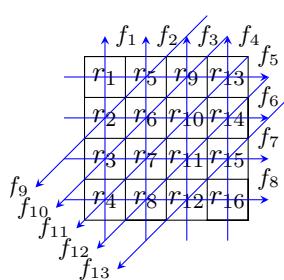
- (a) Write down the nine equations for the transmission factors in terms of the fraction of X-ray energy measured after passing through the body. Take logarithms to form a system of linear equations.

- (b) Encode the matrix A of the system and check `rcond(A)`: curses, `rcond` is terrible, so we must still use an SVD.

- (c) Suppose the measured fractions of X-ray energy are $\mathbf{f} = (0.05, 0.35, 0.33, 0.31, 0.05, 0.36, 0.07, 0.32, 0.51)$. Use an SVD to find the ‘grayest’ transmission factors consistent with the measurements.

- (d) Which part of the body is predicted to be the most absorbing?

Exercise 3.5.16. Use a little higher resolution in computed tomography: suppose the two dimensional ‘body’ is notionally divided into sixteen regions as illustrated in the margin. Suppose a CT-scan takes thirteen measurements of the intensity of an X-ray after passing through the shown paths, and that the fraction of the X-ray energy that is measured is $\mathbf{f} = (0.29, 0.33, 0.07, 0.35, 0.36, 0.07, 0.31, 0.32, 0.62, 0.40, 0.06, 0.47, 0.58)$.



- (a) Write down the thirteen equations for the sixteen transmission factors in terms of the fraction of X-ray energy measured after passing through the body. Take logarithms to form a system of linear equations.

- (b) Encode the matrix A of the system and find it has rank twelve.

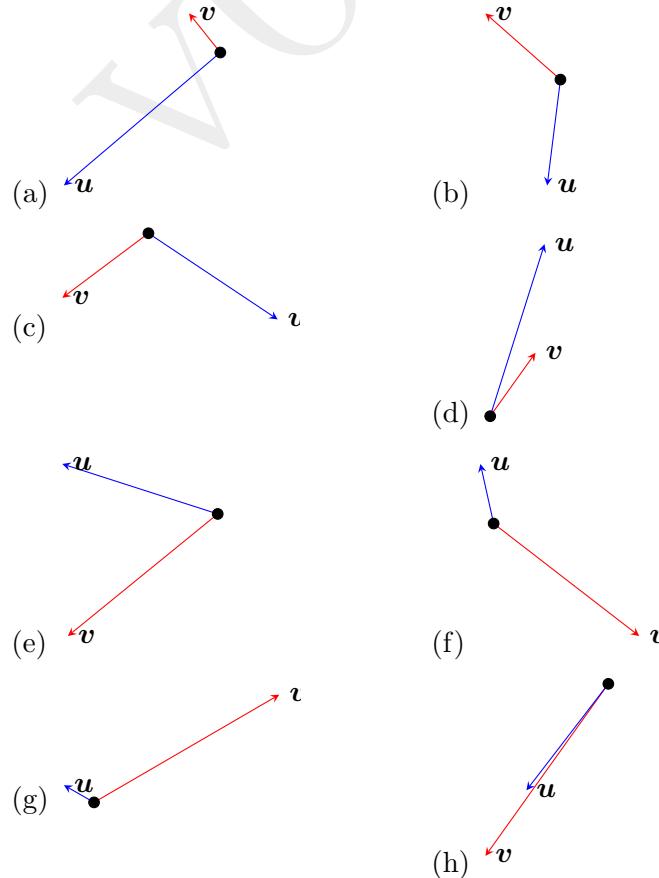
- (c) Use an SVD to find the ‘grayest’ transmission factors consistent with the measurements.

- (d) In which square pixel is the ‘lump’ of dense material?

Exercise 3.5.17. This exercise is for those who, in Calculus courses, have studied constrained optimisation with Lagrange multipliers. The aim is to derive how to use the SVD to find the vector \mathbf{x} that minimises $|A\mathbf{x} - \mathbf{b}|$ such that the magnitude $|\mathbf{x}| \leq \alpha$ for some given magnitude α .

- Given vector $\mathbf{z} \in \mathbb{R}^n$ and $n \times n$ diagonal matrix $S = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$, with $\sigma_1, \sigma_2, \dots, \sigma_n > 0$. In one of two cases, use a Lagrange multiplier λ to find the vector \mathbf{y} (as a function of λ and \mathbf{z}) that minimises $|S\mathbf{y} - \mathbf{z}|^2$ such that $|\mathbf{y}|^2 \leq \alpha^2$ for some given magnitude α : show that the multiplier λ satisfies a polynomial equation of degree n .
- What can be further deduced if one or more $\sigma_j = 0$?
- Use an SVD of $n \times n$ matrix A to find the vector $\mathbf{x} \in \mathbb{R}^n$ that minimises $|A\mathbf{x} - \mathbf{b}|$ such that the magnitude $|\mathbf{x}| \leq \alpha$ for some given magnitude α . Use that multiplication by orthogonal matrices preserves lengths.

Exercise 3.5.18. For each pair of vectors, draw the orthogonal projection $\text{proj}_{\mathbf{u}}(\mathbf{v})$.



Exercise 3.5.19. For the following pairs of vectors: compute the orthogonal projection $\text{proj}_{\mathbf{u}}(\mathbf{v})$; and hence find the ‘best’ approximate solution to the inconsistent system $\mathbf{u}x = \mathbf{v}$.

$$(a) \quad \mathbf{u} = (2, 1), \mathbf{v} = (2, 0) \qquad (b) \quad \mathbf{u} = (4, -1), \mathbf{v} = (-1, 1)$$

$$(c) \quad \mathbf{u} = (6, 0), \mathbf{v} = (-1, -1) \qquad (d) \quad \mathbf{u} = (2, -2), \mathbf{v} = (-1, 2)$$

$$(e) \quad \mathbf{u} = (4, 5, -1), \\ \mathbf{v} = (-1, 2, -1) \qquad (f) \quad \mathbf{u} = (-3, 2, 2), \\ \mathbf{v} = (0, 1, -1)$$

$$(g) \quad \mathbf{u} = (0, 2, 0), \\ \mathbf{v} = (-2, 1, 1) \qquad (h) \quad \mathbf{u} = (-1, -7, 5), \\ \mathbf{v} = (1, 1, -1)$$

$$(i) \quad \mathbf{u} = (2, 4, 0, -1), \\ \mathbf{v} = (0, 2, -1, 0) \qquad (j) \quad \mathbf{u} = (3, -6, -3, -2), \\ \mathbf{v} = (-1, 1, 0, 1)$$

$$(k) \quad \mathbf{u} = (1, 2, 1, -1, -4), \\ \mathbf{v} = (1, -1, 2, -2, 1) \qquad (l) \quad \mathbf{u} = (-2, 2, -1, 3, 2), \\ \mathbf{v} = (-1, 2, 2, 2, 0)$$

Exercise 3.5.20. For each of the following subspaces \mathbb{W} (given as the span of orthogonal vectors), and the given vectors \mathbf{v} , find the orthogonal projection $\text{proj}_{\mathbb{W}}(\mathbf{v})$.

$$(a) \quad \mathbb{W} = \text{span}\{(-6, -6, 7), (2, -9, -6)\}, \mathbf{v} = (0, 1, -2) \qquad (b) \quad \mathbb{W} = \text{span}\{(4, -7, -4), (1, -4, 8)\}, \mathbf{v} = (0, -4, -1)$$

$$(c) \quad \mathbb{W} = \text{span}\{(-6, -3, -2), (-2, 6, -3)\}, \mathbf{v} = (3, -2, -3) \qquad (d) \quad \mathbb{W} = \text{span}\{(1, 8, -4), (-8, -1, -4)\}, \mathbf{v} = (-2, 2, 0)$$

$$(e) \quad \mathbb{W} = \text{span}\{(-1, 2, -2), (-2, 1, 2), (2, 2, 1)\}, \\ \mathbf{v} = (3, -1, 1) \qquad (f) \quad \mathbb{W} = \text{span}\{(-2, 4, -2, 5), (-5, -2, -4, -2)\}, \\ \mathbf{v} = (1, -2, -1, -3)$$

$$(g) \quad \mathbb{W} = \text{span}\{(6, 2, -4, 5), (-5, 2, -4, 2)\}, \\ \mathbf{v} = (3, 3, 2, 7) \qquad (h) \quad \mathbb{W} = \text{span}\{(-1, 3, 1, 5), (-3, -1, -5, 1)\}, \\ \mathbf{v} = (-3, 2, 3, -2)$$

$$(i) \quad \mathbb{W} = \text{span}\{(-1, 5, 3, -1), (-1, -1, 1, -1)\}, \\ \mathbf{v} = (0, -2, -5, -5) \qquad (j) \quad \mathbb{W} = \text{span}\{(-1, 1, -1, 1), (-1, 1, 1, -1), (1, 1, 1, 1)\}, \\ \mathbf{v} = (0, 1, 1, 2)$$

$$(k) \quad \begin{aligned} \mathbb{W} &= \text{span}\{(1, 4, -2, -2), \\ &\quad (2, -2, -4, -5), (-4, 1, 4, -4), (-2, 4, 2, 5)\}, \\ &\quad (-4, 4, 1, -4), (2, -1, 4, -2)\}, \quad \mathbf{v} = (-2, -4, 3, -1) \\ \mathbf{v} &= (2, -3, 1, 0) \end{aligned}$$

Exercise 3.5.21. For each of the following matrices, compute an SVD in Matlab/Octave to find an orthonormal basis for the column space of the matrix, and then compute the matrix of the orthogonal projection onto the column space.

$$(a) A = \begin{bmatrix} 0 & -2 & 4 \\ 4 & -1 & -14 \\ 1 & -1 & -2 \end{bmatrix}$$

$$(b) B = \begin{bmatrix} -3 & 4 \\ -1 & 5 \\ -3 & -1 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} -3 & 11 & 6 \\ 12 & 19 & 3 \\ -30 & 5 & 15 \end{bmatrix}$$

$$(d) D = \begin{bmatrix} -8 & 4 & -2 \\ -24 & 12 & -6 \\ -16 & 8 & -4 \end{bmatrix}$$

$$(e) E = \begin{bmatrix} -3 & 0 & -5 \\ -1 & -4 & 1 \end{bmatrix}$$

$$(f) F = \begin{bmatrix} -5 & 5 & 5 \\ 4 & -4 & -4 \\ -1 & 1 & 1 \\ 5 & -5 & -5 \end{bmatrix}$$

$$(g) G = \begin{bmatrix} 12 & 0 & 10 & 5 \\ -26 & -5 & 5 & 0 \\ -1 & -2 & -16 & 1 \\ -29 & -9 & 29 & 8 \end{bmatrix}$$

$$(h) H = \begin{bmatrix} -12 & 4 & 8 & 16 & 8 \\ 15 & -5 & -10 & -20 & -10 \end{bmatrix}$$

$$(i) \mathcal{I} = \begin{bmatrix} 1 & 26 & -13 & 10 \\ -13 & 2 & 9 & 10 \\ -4 & -2 & 4 & 2 \\ -21 & 32 & 1 & 28 \\ -1 & -9 & 5 & -3 \end{bmatrix}$$

$$(j) J = \begin{bmatrix} 51 & -15 & -19 & -35 & 11 \\ -7 & 2 & 5 & 6 & -5 \\ 14 & -17 & -2 & -8 & -4 \\ 10 & -12 & -2 & -6 & -2 \\ -40 & 30 & 14 & 27 & -4 \end{bmatrix}$$

Exercise 3.5.22. Generally, each of the following systems of equations are inconsistent. Use your answers to the previous Exercise 3.5.21 to find the right-hand side vector \mathbf{b}' that is the closest vector to the given right-hand side among all the vectors in the column space of the matrix. What is the magnitude of the difference between \mathbf{b}' and the given right-hand side? Hence write down a system of *consistent* equations that best approximates the original system.

$$(a) \begin{bmatrix} 0 & -2 & 4 \\ 4 & -1 & -14 \\ 1 & -1 & -2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 6 \\ -19 \\ -3 \end{bmatrix} \quad (b) \begin{bmatrix} 0 & -2 & 4 \\ 4 & -1 & -14 \\ 1 & -1 & -2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 2 \\ -8 \\ -1 \end{bmatrix}$$

$$(c) \begin{bmatrix} -3 & 4 \\ -1 & 5 \\ -3 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 9 \\ 11 \\ -1 \end{bmatrix} \quad (d) \begin{bmatrix} -3 & 4 \\ -1 & 5 \\ -3 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -1 \\ 2 \\ -3 \end{bmatrix}$$

$$(e) \begin{bmatrix} -3 & 11 & 6 \\ 12 & 19 & 3 \\ -30 & 5 & 15 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 3 \\ 5 \\ -3 \end{bmatrix} \quad (f) \begin{bmatrix} -3 & 11 & 6 \\ 12 & 19 & 3 \\ -30 & 5 & 15 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 5 \\ 27 \\ -14 \end{bmatrix}$$

$$(g) \begin{bmatrix} -3 & 0 & -5 \\ -1 & -4 & 1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -9 \\ 10 \end{bmatrix} \quad (h) \begin{bmatrix} -3 & 0 & -5 \\ -1 & -4 & 1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 6 \\ 3 \end{bmatrix}$$

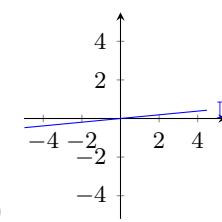
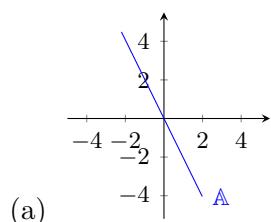
$$(i) \begin{bmatrix} -5 & 5 & 5 \\ 4 & -4 & -4 \\ -1 & 1 & 1 \\ 5 & -5 & -5 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 5 \\ -6 \\ 1 \\ -6 \end{bmatrix} \quad (j) \begin{bmatrix} -5 & 5 & 5 \\ 4 & -4 & -4 \\ -1 & 1 & 1 \\ 5 & -5 & -5 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -6 \\ 6 \\ -2 \\ 5 \end{bmatrix}$$

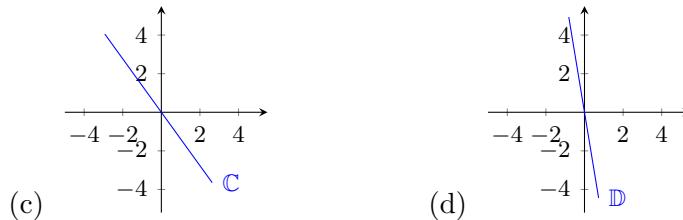
$$(k) \begin{bmatrix} 12 & 0 & 10 & 5 \\ -26 & -5 & 5 & 0 \\ -1 & -2 & -16 & 1 \\ -29 & -9 & 29 & 8 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 4 \\ -45 \\ 27 \\ -98 \end{bmatrix} \quad (l) \begin{bmatrix} 12 & 0 & 10 & 5 \\ -26 & -5 & 5 & 0 \\ -1 & -2 & -16 & 1 \\ -29 & -9 & 29 & 8 \end{bmatrix} \mathbf{x} = \begin{bmatrix} -11 \\ -4 \\ 18 \\ -37 \end{bmatrix}$$

Exercise 3.5.23. Theorems 3.5.5 and 3.5.19, examples and Exercise 3.5.22 solve an inconsistent system of equations by some specific ‘best approximation’ that forms a consistent system of equations to solve. Describe briefly the key idea of this ‘best approximation’. Discuss other possibilities for a ‘best approximation’ that might be developed.

Exercise 3.5.24. For any matrix A , suppose you know an orthonormal basis for the column space of A . Form the matrix W from all the vectors of the orthonormal basis. What is the result of the product $(WW^T)A$? Explain why.

Exercise 3.5.25. For each of the following subspaces, draw its orthogonal complement on the plot.





Exercise 3.5.26. Describe the orthogonal complement of each of the sets given below, if the set has one.

- (a) $\mathbb{A} = \text{span}\{(-1, 2)\}$
- (b) $\mathbb{B} = \text{span}\{(5, -1)\}$
- (c) $\mathbb{C} = \text{span}\{(1, 9, -9)\}$
- (d) \mathbb{D} is the plane $-4x_1 + 4x_2 + 5x_3 = 0$
- (e) \mathbb{E} is the plane $5x + 2y + 3z = 3$
- (f) $\mathbb{F} = \text{span}\{(-5, 5, -3), (-2, 1, 1)\}$
- (g) $\mathbb{G} = \text{span}\{(-2, 2, 8), (5, 3, 5)\}$
- (h) $\mathbb{H} = \text{span}\{(6, 5, 1, -3)\}$

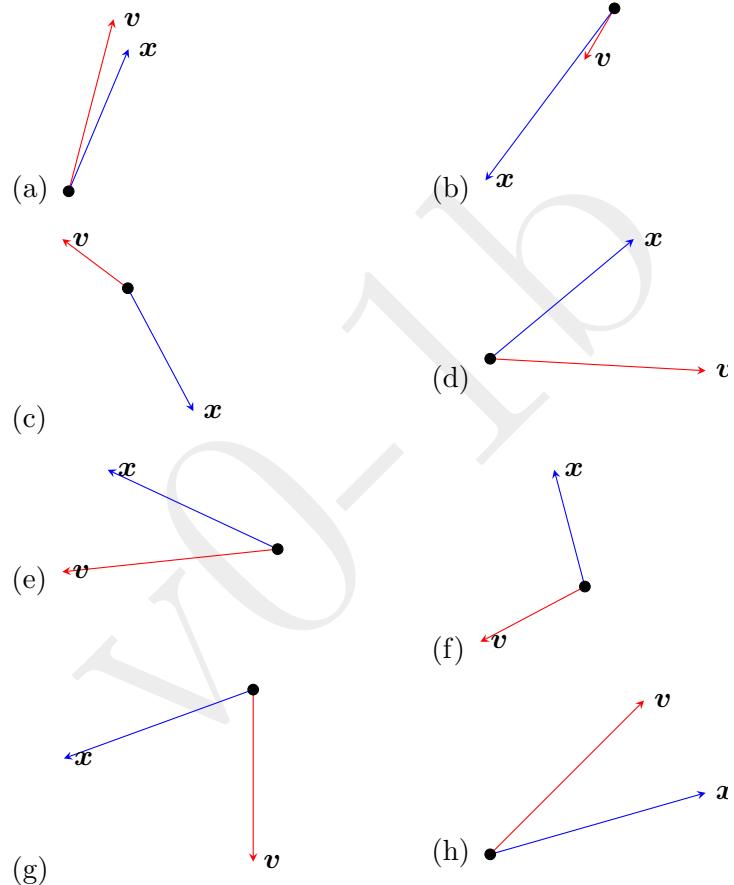
Exercise 3.5.27. Compute, using Matlab/Octave when necessary, an orthonormal basis for the orthogonal complement, if it exists, to each of the following sets. Use that the an orthogonal complement is the nullspace of the transpose of a matrix of column vectors (Theorem 3.5.29).

- (a) The \mathbb{R}^3 vectors in the plane $-6x + 2y - 3z = 0$
- (b) The \mathbb{R}^3 vectors in the plane $x + 4y + 8z = 0$
- (c) The \mathbb{R}^3 vectors in the plane $3x + 3y + 2z = 9$
- (d) The span of vectors $(-3, 11, -25), (24, 32, -40), (-8, -8, 8)$.
- (e) The span of vectors $(3, -2, 1), (-3, 2, -1), (-9, 6, -3), (-6, 4, -2)$
- (f) The span of vectors $(26, -2, -4, 20), (23, -3, 2, 6), (2, -2, 8, -16), (21, -5, 12, -16)$
- (g) The span of vectors $(7, -5, 1, -6, -4), (6, -4, -2, -8, -4), (-5, 5, -15, -10, 0), (8, -6, 4, -4, -4)$

- (h) The column space of matrix $\begin{bmatrix} 2 & -1 & 2 & 6 \\ -9 & 11 & -12 & -22 \\ -7 & -6 & -15 & -46 \\ 7 & -23 & 2 & -14 \\ 0 & -2 & 2 & 0 \end{bmatrix}$

- (i) The intersection in \mathbb{R}^4 of the two hyper-planes $4x_1 + x_2 - 2x_3 + 5x_4 = 0$ and $-4x_1 - x_2 - 7x_3 + 2x_4 = 0$.
- (j) The intersection in \mathbb{R}^4 of the two hyper-planes $-3x_1 + x_2 + 4x_3 - 7x_4 = 0$ and $-6x_2 - x_3 - 2x_4 = 0$.

Exercise 3.5.28. For the subspace $\mathbb{X} = \text{span}\{\mathbf{x}\}$ and the vector \mathbf{v} , draw the decomposition of \mathbf{v} into the sum of vectors in \mathbb{X} and \mathbb{X}^\perp .



Exercise 3.5.29. For each of the following vectors, find the perpendicular component to the subspace $\mathbb{W} = \text{span}\{(4, -4, 7)\}$. Verify that the perpendicular component lies in the plane $4x - 4y + 7z = 0$.

- | | |
|-------------------|-------------------|
| (a) $(4, 2, 4)$ | (b) $(0, 1, -2)$ |
| (c) $(0, -2, -2)$ | (d) $(-2, -1, 1)$ |
| (e) $(5, 1, 5)$ | (f) (p, q, r) |

Exercise 3.5.30. For each of the following vectors, find the perpendicular component to the subspace $\mathbb{W} = \text{span}\{(1, 5, 5, 7), (-5, 1, -7, 5)\}$.

- | | |
|----------------------|---------------------|
| (a) $(1, 2, -1, -1)$ | (b) $(-2, 4, 5, 0)$ |
|----------------------|---------------------|

$$(c) (2, -6, 1, -3) \quad (d) (p, q, r, s)$$

Exercise 3.5.31. Let \mathbb{W} be a subspace of \mathbb{R}^n and let \mathbf{v} be any vector in \mathbb{R}^n .

Prove that $\text{perp}_{\mathbb{W}}(\mathbf{v}) = (I_n - WW^T)\mathbf{v}$ where the columns of the matrix W are an orthonormal basis for \mathbb{W} .

Exercise 3.5.32. For each of the following vectors in \mathbb{R}^2 , write the vector as the orthogonal decomposition with respect to the subspace $\mathbb{W} = \text{span}\{(3, 4)\}$.

$$(a) (-2, 4) \quad (b) (-3, 3)$$

$$(c) (0, 0) \quad (d) (3, 1)$$

Exercise 3.5.33. For each of the following vectors in \mathbb{R}^3 , write the vector as the orthogonal decomposition with respect to the subspace $\mathbb{W} = \text{span}\{(3, -6, 2)\}$.

$$(a) (-5, 4, -5) \quad (b) (0, 5, -1)$$

$$(c) (1, -1, -2) \quad (d) (-3, 1, -1)$$

Exercise 3.5.34. For each of the following vectors in \mathbb{R}^4 , write the vector as the orthogonal decomposition with respect to the subspace $\mathbb{W} = \text{span}\{(3, -1, 9, 3), (-9, 3, 3, 1)\}$.

$$(a) (5, -5, 1, -3) \quad (b) (-4, -2, 5, 5)$$

$$(c) (2, -1, -4, -3) \quad (d) (5, 4, 0, 3)$$

Exercise 3.5.35. The vector $(-3, 4)$ has an orthogonal decomposition $(1, 2) + (-4, 2)$. Draw in \mathbb{R}^2 the possibilities for the subspace \mathbb{W} and its orthogonal complement.

Exercise 3.5.36. The vector $(2, 0, -3)$ in \mathbb{R}^3 has an orthogonal decomposition $(2, 0, 0) + (0, 0, -3)$. Describe the possibilities for the subspace \mathbb{W} and its orthogonal complement.

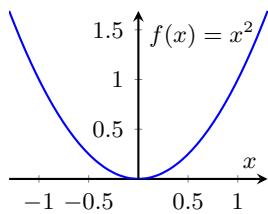
Exercise 3.5.37. The vector $(0, -2, 5, 0)$ in \mathbb{R}^4 has an orthogonal decomposition $(0, -2, 0, 0) + (0, 0, 5, 0)$. Describe the possibilities for the subspace \mathbb{W} and its orthogonal complement.

3.6 Introducing linear transformations

Section Contents

3.6.1	Matrices characterise linear transforms	341
3.6.2	The pseudo-inverse of a matrix	345
3.6.3	Function composition connects to matrix inverse	353
3.6.4	Exercises	359

This optional section unifies the transformation examples seen so far, and forms a foundation for more abstract algebra.



Recall the function notation such as $f(x) = x^2$ means that for each $x \in \mathbb{R}$, the function $f(x)$ gives a result in \mathbb{R} , namely the value x^2 , as plotted in the margin. We often write $f : \mathbb{R} \rightarrow \mathbb{R}$ to denote this functionality: that is, $f : \mathbb{R} \rightarrow \mathbb{R}$ means function f transforms any given real number into another real number by some rule.

There is analogous functionality in multiple dimensions with vectors: given any vector

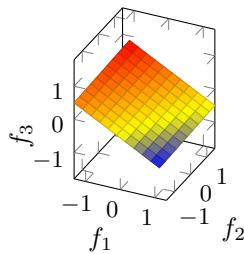
- multiplication by a diagonal matrix stretches and shrinks the vector (section 3.2.2);
- multiplication by an orthogonal matrix rotates the vector (section 3.2.3); and
- projection finds the vector's components in a subspace (section 3.5.3).

Correspondingly, we use the notation $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ to mean function f transforms a given vector with n components (in \mathbb{R}^n) into another vector with m components (in \mathbb{R}^m) according to some rule. For example, suppose the function $f(\mathbf{x})$ is to denote multiplication by the matrix

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

Then the function

$$f(\mathbf{x}) = \begin{bmatrix} 1 & -\frac{1}{3} \\ \frac{1}{2} & -1 \\ -1 & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 - x_2/3 \\ x_1/2 - x_2 \\ -x_1 - x_2/2 \end{bmatrix}.$$



That is, here $f : \mathbb{R}^2 \rightarrow \mathbb{R}^3$. Given any vector in the 2D-plane, the function f , also called a transformation, returns a vector in 3D-space. The marginal plot illustrates the subspace formed by $f(\mathbf{x})$ for all 2D vectors \mathbf{x} .

There is a major difference between ‘curvaceous’ functions like the parabola above, and matrix multiplication functions such as rotation

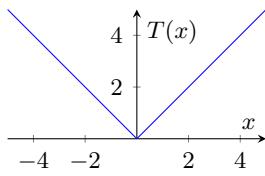
and projection. The difference is that linear algebra empowers many practical results in the latter case.

Definition 3.6.1. A transformation/function $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called a **linear transformation** if

- (a) $T(\mathbf{u} + \mathbf{v}) = T(\mathbf{u}) + T(\mathbf{v})$ for all $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, and
- (b) $T(c\mathbf{v}) = cT(\mathbf{v})$ for all $\mathbf{v} \in \mathbb{R}^n$ and all scalars c .

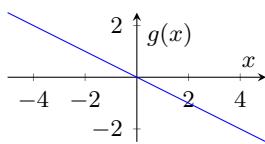
Example 3.6.2 (1D cases). (a) Show that function $f : \mathbb{R} \rightarrow \mathbb{R}$ where $f(x) = x^2$ is not a linear transformation.

Solution: To test Property 3.6.1a, consider $f(x+y) = (x+y)^2 = x^2 + 2xy + y^2 = f(x) + 2xy + f(y) \neq f(x) + f(y)$ for all x and y (it is equal if either are zero, but the test requires equality to hold for all x and y). Alternatively one could test Property 3.6.1b and consider $f(cx) = (cx)^2 = c^2x^2 = c^2f(x) \neq cf(x)$ for all c . Either of these prove that f is not a linear transform.



- (b) Is the function $T(x) = |x|$ ($T : \mathbb{R} \rightarrow \mathbb{R}$) a linear transform?

Solution: To prove not it is sufficient to find just one instance when Definition 3.6.1 fails. Let $u = -1$ and $v = 2$, then $T(u+v) = |-1+2| = |1| = 1$ whereas $T(u) + T(v) = |-1| + |2| = 1 + 2 = 3 \neq T(u+v)$ so the function T fails the additivity and so is not a linear transform.



- (c) Is the function $g : \mathbb{R} \rightarrow \mathbb{R}$ such that $g(x) = -x/2$ a linear transform?

Solution: Because the graph of g is a straight line (as in the marginal picture) we suspect it is a linear transform. Thus check the properties in full generality:

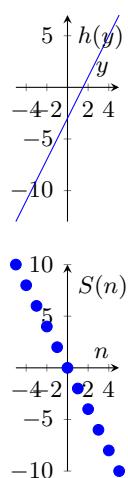
3.6.1a : for all $u, v \in \mathbb{R}$, $g(u+v) = -(u+v)/2 = -u/2 - v/2 = (-u/2) + (-v/2) = g(u) + g(v)$;

3.6.1b : for all $u, c \in \mathbb{R}$, $g(cu) = -(cu)/2 = c(-u/2) = cg(u)$.

Hence g is a linear transformation.

- (d) Show that the function $h(y) = 2y - 3$, $h : \mathbb{R} \rightarrow \mathbb{R}$, is not a linear transform.

Solution: Because the graph of $h(y)$ is a straight line we suspect it may be a linear transform (as shown in the margin). To prove not it is enough to find one instance when Definition 3.6.1 fails. Let $u = 0$ and $c = 2$, then $h(cu) = h(2 \cdot 0) = h(0) = -3$ whereas $ch(u) = 2h(0) = 2 \cdot (-3) = -6 \neq h(cu)$ so the function g fails the multiplication rule and hence is not a linear transform. (This function fails because linear transforms have to pass through the origin.)



- (e) Is the function $S : \mathbb{N} \rightarrow \mathbb{N}$ given by $S(n) = -2n$ a linear transform?

Solution: No, because the function S is here only defined for integers \mathbb{N} (as plotted in the margin) whereas Definition 3.6.1 requires the function to be defined for all reals. ²²

■

Example 3.6.3 (higher-D cases). (a) Let function $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be $T(x, y, z) = (y, z, x)$. Is T a linear transformation?

Solution: Yes, because:

3.6.1a : for all $\mathbf{u} = (x, y, z)$ and $\mathbf{v} = (x', y', z')$ in \mathbb{R}^3 consider $T(\mathbf{u} + \mathbf{v}) = T(x + x', y + y', z + z') = (y + y', z + z', x + x') = (y, z, x) + (y', z', x') = T(x, y, z) + T(x', y', z') = T(\mathbf{u}) + T(\mathbf{v})$;

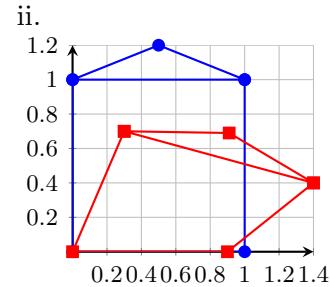
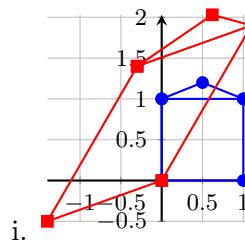
3.6.1b : for all $\mathbf{u} = (x, y, z)$ and scalars c consider $T(c\mathbf{u}) = T(cx, cy, cz) = (cy, cz, cx) = c(y, z, x) = cT(x, y, z) = cT(\mathbf{u})$.

Hence, T is a linear transformation.

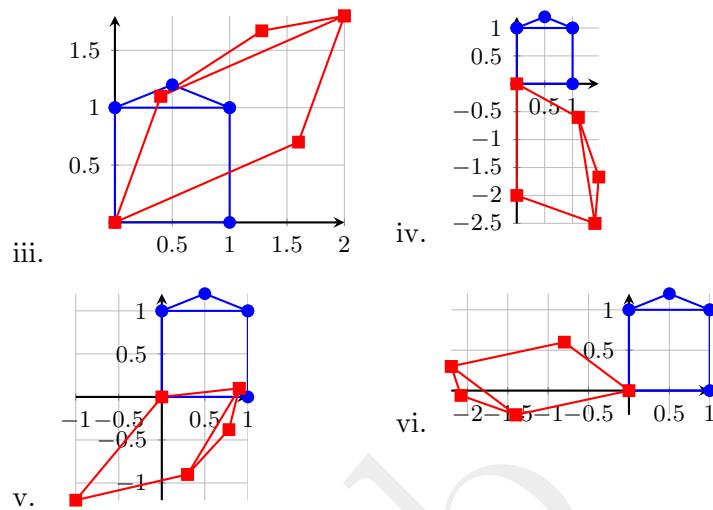
- (b) Consider the function $f(x, y, z) = x + y + 1$, $f : \mathbb{R}^3 \rightarrow \mathbb{R}$: is f a linear transformation?

Solution: No. For example, choose $\mathbf{u} = \mathbf{0}$ and scalar $c = 2$ then $f(c\mathbf{u}) = f(2 \cdot \mathbf{0}) = f(\mathbf{0}) = 1$ whereas $cf(\mathbf{u}) = 2f(\mathbf{0}) = 2 \cdot 1 = 2$. Hence f fails the scalar multiplication property 3.6.1b.

- (c) Which of the following illustrated transformations of the plane *cannot* be that of a linear transformation? In each illustration of a transformation T , the four corners of the blue unit square ((0, 0), (1, 0), (1, 1) and (0, 1)), are transformed to the four corners of the red figure ($T(0, 0)$, $T(1, 0)$, $T(1, 1)$ and $T(0, 1)$ —the ‘roof’ of the unit square clarifies which side goes where).



²² More advanced linear algebra generalises the definition of a linear transformation to non-reals, but not here.

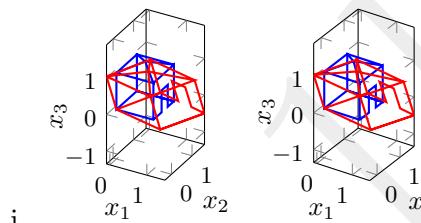


Solution: To test we check the addition property 3.6.1a. First, with $\mathbf{u} = \mathbf{v} = \mathbf{0}$ Definition 3.6.1a requires $T(\mathbf{0} + \mathbf{0}) = T(\mathbf{0}) + T(\mathbf{0})$, but the left-hand side is just $T(\mathbf{0})$ which cancels with one on the right-hand side to leave that a linear transform has to satisfy $T(\mathbf{0}) = \mathbf{0}$: all the shown transforms satisfy $T(\mathbf{0}) = \mathbf{0}$ as the (blue) origin point is transformed to the (red) origin point. Second, with $\mathbf{u} = (1, 0)$, $\mathbf{v} = (1, 0)$ and $\mathbf{u} + \mathbf{v} = (1, 1)$ Definition 3.6.1a requires $T(1, 1) = T(1, 0) + T(0, 1)$: let's see which do not pass this test.

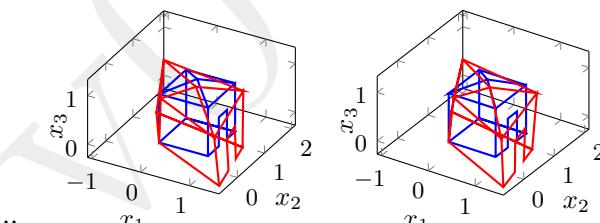
- i. Here $T(1, 1) \approx (-0.3, 1.4)$, whereas $T(1, 0) + T(0, 1) \approx (-1.4, -0.5) + (1.1, 1.9) = (-0.3, 1.4) \approx T(1, 1)$ so this may be a linear transformation.
- ii. Here $T(1, 1) \approx (1.4, 0.4)$, whereas $T(1, 0) + T(0, 1) \approx (0.9, 0) + (0.3, 0.7) = (1.2, 0.7) \not\approx T(1, 1)$ so this *cannot* be a linear transformation.
- iii. Here $T(1, 1) \approx (2.0, 1.8)$, whereas $T(1, 0) + T(0, 1) \approx (1.6, 0.7) + (0.4, 1.1) = (2.0, 1.8) \approx T(1, 1)$ so this may be a linear transformation.
- iv. Here $T(1, 1) \approx (1.4, -2.5)$, whereas $T(1, 0) + T(0, 1) \approx (0, -2) + (1.1, -0.6) = (1.1, -2.6) \not\approx T(1, 1)$ so this *cannot* be a linear transformation.
- v. Here $T(1, 1) \approx (0.3, -0.9)$, whereas $T(1, 0) + T(0, 1) \approx (-1, -1.2) + (0.9, 0.1) = (-0.1, -1.1) \not\approx T(1, 1)$ so this *cannot* be a linear transformation.
- vi. Here $T(1, 1) \approx (-2.2, 0.3)$, whereas $T(1, 0) + T(0, 1) \approx (-0.8, 0.6) + (-1.4, -0.3) = (-2.2, 0.3) \approx T(1, 1)$ so this may be a linear transformation.

(The ones that pass this test may fail other tests: all we are sure of is that those that fail such tests *cannot* be linear transforms.)

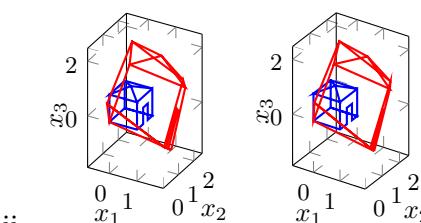
- (d) The previous Example 3.6.3c illustrated that a linear transform of the square seems to transform the unit square to a rhombus: if a function transforms the unit square to something that is not a rhombus, then the function cannot be a linear transform. Analogously in higher dimensions: for example, if a function transforms the unit cube to something which is not a parallelepiped, then the function is not a linear transform. Using this information, which of the following illustrated functions, $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$, *cannot* be a linear transform? Each of these stereo illustrations plot the unit cube in blue (with a ‘roof’ and ‘door’ to help orientate), and the transform of the unit cube in red (with its transformed ‘roof’ and ‘door’).



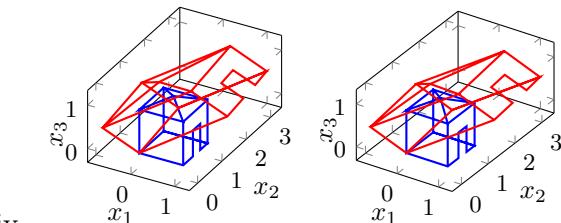
i. This *may* be a linear transform as the transform of the unit cube looks like a parallelepiped.



ii. This *cannot* be a linear transform as the unit cube transforms to something not a parallelepiped.



iii. This *cannot* be a linear transform as the unit cube transforms to something not a parallelepiped.



iv. This *may* be a linear transform as the transform of the unit cube looks like a parallelepiped.

Example 3.6.4. For any given nonzero vector $\mathbf{w} \in \mathbb{R}^n$, prove that the projection $P : \mathbb{R}^n \rightarrow \mathbb{R}^n$ by $P(\mathbf{u}) = \text{proj}_{\mathbf{w}}(\mathbf{u})$ is a linear transformation (as a function of \mathbf{u}). But, for any given nonzero vector $\mathbf{u} \in \mathbb{R}^n$, prove that the projection $Q : \mathbb{R}^n \rightarrow \mathbb{R}^n$ by $Q(\mathbf{w}) = \text{proj}_{\mathbf{w}}(\mathbf{u})$ is not a linear transformation (as a function of \mathbf{w}).

Solution: • Consider the two properties of Definition 3.6.1 for the function P .

$$\begin{aligned} 3.6.1a : \quad & \text{for all } \mathbf{u}, \mathbf{v} \in \mathbb{R}^n, \text{ from Defintion 3.5.14 for a projection,} \\ & P(\mathbf{u} + \mathbf{v}) = \text{proj}_{\mathbf{w}}(\mathbf{u} + \mathbf{v}) = \mathbf{w}[\mathbf{w} \cdot (\mathbf{u} + \mathbf{v})]/|\mathbf{w}|^2 = \mathbf{w}[(\mathbf{w} \cdot \mathbf{u}) + (\mathbf{w} \cdot \mathbf{v})]/|\mathbf{w}|^2 = \mathbf{w}(\mathbf{w} \cdot \mathbf{u})/|\mathbf{w}|^2 + \mathbf{w}(\mathbf{w} \cdot \mathbf{v})/|\mathbf{w}|^2 = \\ & \text{proj}_{\mathbf{w}}(\mathbf{u}) + \text{proj}_{\mathbf{w}}(\mathbf{v}) = P(\mathbf{u}) + P(\mathbf{v}); \end{aligned}$$

$$\begin{aligned} 3.6.1b : \quad & \text{for all } \mathbf{u} \in \mathbb{R}^n \text{ and scalars } c, P(c\mathbf{u}) = \text{proj}_{\mathbf{w}}(c\mathbf{u}) = \mathbf{w}[\mathbf{w} \cdot (c\mathbf{u})]/|\mathbf{w}|^2 = \mathbf{w}[c(\mathbf{w} \cdot \mathbf{u})]/|\mathbf{w}|^2 = c[\mathbf{w}(\mathbf{w} \cdot \mathbf{u})/|\mathbf{w}|^2] = \\ & c \text{proj}_{\mathbf{w}}(\mathbf{u}) = cP(\mathbf{u}). \end{aligned}$$

Hence, the projection P is a linear transformation.

- Now consider $Q(\mathbf{w}) = \text{proj}_{\mathbf{w}}(\mathbf{u})$. For any $\mathbf{u}, \mathbf{w} \in \mathbb{R}^n$ let's check $Q(2\mathbf{w}) = \text{proj}_{(2\mathbf{w})}(\mathbf{u}) = (2\mathbf{w})[(2\mathbf{w}) \cdot \mathbf{u}]/|2\mathbf{w}|^2 = 4\mathbf{w}(\mathbf{w} \cdot \mathbf{u})/(4|\mathbf{u}|^2) = \mathbf{w}(\mathbf{w} \cdot \mathbf{u})/|\mathbf{u}|^2 = \text{proj}_{\mathbf{w}}(\mathbf{u}) = Q(\mathbf{w}) \neq 2Q(\mathbf{w})$ and so the projection is not a linear transformation when considered as a function of the direction of the transformation \mathbf{w} for some given \mathbf{u} .

3.6.1 Matrices characterise linear transforms

One important class of linear transformations are the transforms that can be written as matrix multiplications. The reason for the importance is that Theorem 3.6.8 establishes all linear transforms may be written as matrix multiplications! This in turn justifies why we define matrix multiplication to be as it is (section 3.1.2): *matrix multiplication is defined just so that all linear transformations are encompassed by them.*

Example 3.6.5. But first, the following Theorem 3.6.6 proves, among many other possibilities, that the following are linear transformations:

- stretching/shrinking along coordinate axes as these are multiplication by a diagonal matrix (section 3.2.2);
- rotations and/or reflections as they arise as multiplications by an orthogonal matrix (section 3.2.3);

- orthogonal projection onto a subspace as all such projections may be expressed as multiplication by a matrix (the matrix WW^T in Theorem 3.5.22).

■

Theorem 3.6.6. *Let A be a given $m \times n$ matrix and define the transformation $T_A : \mathbb{R}^n \rightarrow \mathbb{R}^m$ by the matrix multiplication $T_A(\mathbf{x}) := A\mathbf{x}$ for all $\mathbf{x} \in \mathbb{R}^n$. Then T_A is a linear transformation.*

Proof. Let vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and scalar $c \in \mathbb{R}$ and consider the two properties of Definition 3.6.1.

3.6.1a. By the distributivity of matrix-vector multiplication (Theorem 3.1.16), $T_A(\mathbf{u} + \mathbf{v}) = A(\mathbf{u} + \mathbf{v}) = A\mathbf{u} + A\mathbf{v} = T_A(\mathbf{u}) + T_A(\mathbf{v})$.

3.6.1b. By commutativity of scalar multiplication (Theorem 3.1.18), $T_A(c\mathbf{u}) = A(c\mathbf{u}) = c(A\mathbf{u}) = cT_A(\mathbf{u})$.

Hence T_A is a linear transformation. □

Example 3.6.7. Prove that a matrix multiplication with a nonzero shift \mathbf{b} , $S : \mathbb{R}^n \rightarrow \mathbb{R}^m$ where $S(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ for vector $\mathbf{b} \neq \mathbf{0}$, is not a linear transform.

Solution: Just consider the addition property 3.6.1a for the zero vectors $\mathbf{u} = \mathbf{v} = \mathbf{0}$: on the one hand, $S(\mathbf{u} + \mathbf{v}) = S(\mathbf{0} + \mathbf{0}) = S(\mathbf{0}) = A\mathbf{0} + \mathbf{b} = \mathbf{b}$; on the other hand $S(\mathbf{u}) + S(\mathbf{v}) = S(\mathbf{0}) + S(\mathbf{0}) = A\mathbf{0} + \mathbf{b} + A\mathbf{0} + \mathbf{b} = 2\mathbf{b}$. Hence when the shift \mathbf{b} is nonzero, there are vectors for which $S(\mathbf{u} + \mathbf{v}) \neq S(\mathbf{u}) + S(\mathbf{v})$ and so S is not a linear transform.

■

Now let's establish the important converse to Theorem 3.6.6: every linear transformation can be written as a matrix multiplication.

Theorem 3.6.8. *Let $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a linear transformation. Then T is the transformation corresponding to the $m \times n$ matrix*

$$A = [T(\mathbf{e}_1) \ T(\mathbf{e}_2) \ \cdots \ T(\mathbf{e}_n)]$$

where \mathbf{e}_j are the standard unit vectors in \mathbb{R}^n . This matrix A , often denoted $[T]$, is called the **standard matrix** of the linear transformation T .

Proof. Let \mathbf{x} be any vector in \mathbb{R}^n : then $\mathbf{x} = (x_1, x_2, \dots, x_n) = x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + \cdots + x_n\mathbf{e}_n$ for standard unit vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$. Then

$$T(\mathbf{x}) = T(x_1\mathbf{e}_1 + x_2\mathbf{e}_2 + \cdots + x_n\mathbf{e}_n)$$

$$\begin{aligned}
 & \text{(using the identity of Exercise 3.6.6)} \\
 & = x_1 T(\mathbf{e}_1) + x_2 T(\mathbf{e}_2) + \cdots + x_n T(\mathbf{e}_n) \\
 & = [T(\mathbf{e}_1) \ T(\mathbf{e}_2) \ \cdots \ T(\mathbf{e}_n)] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\
 & = A\mathbf{x}
 \end{aligned}$$

for matrix A of the theorem. Since $T(\mathbf{e}_1), T(\mathbf{e}_2), \dots, T(\mathbf{e}_n)$ are n (column) vectors in \mathbb{R}^m , the matrix A is $m \times n$. \square

Example 3.6.9. (a) Find the standard matrix of the linear transformation $T : \mathbb{R}^3 \rightarrow \mathbb{R}^4$ where $T(x, y, z) = (y, z, x, 3x - 2y + z)$.

Solution: We need to find the transform of the three standard unit vectors in \mathbb{R}^3 :

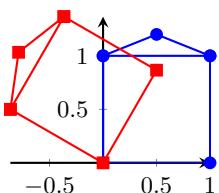
$$\begin{aligned}
 T(\mathbf{e}_1) &= T(1, 0, 0) = (0, 0, 1, 3); \\
 T(\mathbf{e}_2) &= T(0, 1, 0) = (1, 0, 0, -2); \\
 T(\mathbf{e}_3) &= T(0, 0, 1) = (0, 1, 0, 1).
 \end{aligned}$$

Form the standard matrix with these as its three columns, in order,

$$[T] = [T(\mathbf{e}_1) \ T(\mathbf{e}_2) \ T(\mathbf{e}_3)] = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 3 & -2 & 1 \end{bmatrix}.$$

(b) Find the standard matrix of the rotation of the plane by 60° about the origin.

Solution: Denote the rotation of the plane by the function $R : \mathbb{R}^2 \rightarrow \mathbb{R}^2$. Since $60^\circ = \frac{\pi}{3}$ then, as illustrated in the margin,



$$\begin{aligned}
 R(\mathbf{e}_1) &= (\cos \frac{\pi}{3}, \sin \frac{\pi}{3}) = (\frac{1}{2}, \frac{\sqrt{3}}{2}), \\
 R(\mathbf{e}_2) &= (-\sin \frac{\pi}{3}, \cos \frac{\pi}{3}) = (-\frac{\sqrt{3}}{2}, \frac{1}{2}).
 \end{aligned}$$

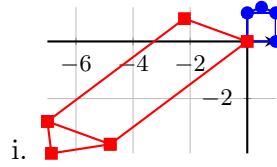
Form the standard matrix with these as its columns, in order,

$$[R] = [R(\mathbf{e}_1) \ R(\mathbf{e}_2)] = \begin{bmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix}.$$

(c) Find the standard matrix of the rotation about the point $(1, 0)$ of the plane by 45° .

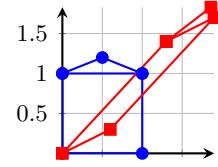
Solution: Since the origin $(0, 0)$ is transformed by the rotation to $(1, -1)$ which is nonzero, this transform cannot be of the form $A\mathbf{x}$, so cannot have a standard matrix, and hence is not a linear transform.

- (d) Estimate the standard matrix for each of the illustrated transforms given they transform the unit square as shown.



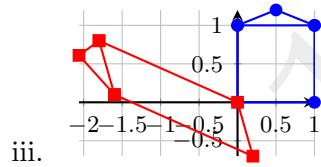
Solution: Here

$T(1, 0) \approx (-2.2, 0.8)$ and
 $T(0, 1) \approx (-4.8, -3.6)$ so
the approximate standard
matrix is $\begin{bmatrix} -2.2 & -4.8 \\ 0.8 & -3.6 \end{bmatrix}$.



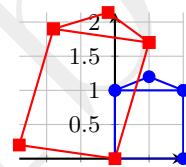
Solution: Here

$T(1, 0) \approx (0.6, 0.3)$ and
 $T(0, 1) \approx (1.3, 1.4)$ so the
approximate standard
matrix is $\begin{bmatrix} 0.6 & 1.3 \\ 0.3 & 1.4 \end{bmatrix}$.



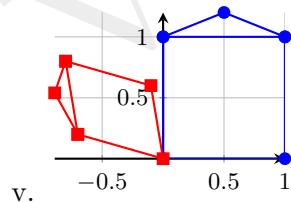
Solution: Here

$T(1, 0) \approx (0.2, -0.7)$ and
 $T(0, 1) \approx (-1.8, 0.8)$ so the
approximate standard
matrix is $\begin{bmatrix} 0.2 & -1.8 \\ -0.7 & 0.8 \end{bmatrix}$.



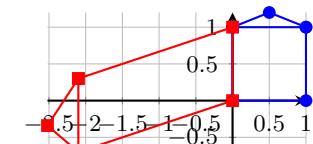
Solution: Here

$T(1, 0) \approx (-1.4, 0.2)$ and
 $T(0, 1) \approx (0.5, 1.7)$ so the
approximate standard
matrix is $\begin{bmatrix} -1.4 & 0.5 \\ 0.2 & 1.7 \end{bmatrix}$.



Solution: Here

$T(1, 0) \approx (-0.1, 0.6)$ and
 $T(0, 1) \approx (-0.7, 0.2)$ so the
approximate standard
matrix is $\begin{bmatrix} -0.1 & -0.7 \\ 0.6 & 0.2 \end{bmatrix}$.



Solution: Here

$T(1, 0) \approx (0, 1.0)$ and
 $T(0, 1) \approx (-2.1, -0.7)$ so the
approximate standard
matrix is $\begin{bmatrix} 0 & -2.1 \\ 1.0 & -0.7 \end{bmatrix}$.

Example 3.6.10. For a fixed scalar a , let the function $H : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be $H(\mathbf{u}) = a\mathbf{u}$. Show that H is a linear transformation, and then find its standard matrix.

Solution: Let $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and c be any scalar. Function H is a linear transformation because

- $H(\mathbf{u} + \mathbf{v}) = a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v} = H(\mathbf{u}) + H(\mathbf{v})$, and

- $H(c\mathbf{u}) = a(c\mathbf{u}) = (ac)\mathbf{u} = (ca)\mathbf{u} = c(a\mathbf{u}) = cH(\mathbf{u})$.

To find the standard matrix consider

$$\begin{aligned} H(\mathbf{e}_1) &= a\mathbf{e}_1 = (a, 0, 0, \dots, 0), \\ H(\mathbf{e}_2) &= a\mathbf{e}_2 = (0, a, 0, \dots, 0), \\ &\vdots \\ H(\mathbf{e}_n) &= a\mathbf{e}_n = (0, 0, \dots, 0, a). \end{aligned}$$

Hence the standard matrix $[H] = \text{diag}(a, a, \dots, a) = aI_n$.

$$\begin{aligned} [H] &= [H(\mathbf{e}_1) \ H(\mathbf{e}_2) \ \cdots \ H(\mathbf{e}_n)] \\ &= \begin{bmatrix} a & 0 & \cdots & 0 \\ 0 & a & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & a \end{bmatrix} = aI_n. \end{aligned}$$

■

Consider this last Example 3.6.10 in the case $a = 1$: then $H(\mathbf{u}) = \mathbf{u}$ is the identity and so the example shows that the standard matrix of the identity transformation is I_n .

3.6.2 The pseudo-inverse of a matrix

In solving inconsistent linear equations, $A\mathbf{x} = \mathbf{b}$ for some given A , This subsection is an optional extension.

Procedure 3.5.3 finds a solution \mathbf{x} that depends upon the right-hand side \mathbf{b} . That is, any given \mathbf{b} is transformed by the procedure to some result \mathbf{x} . This section establishes the resulting solution given by the procedure is a linear transformation of \mathbf{b} , and hence there must be a matrix, say A^+ , corresponding to the procedure and hence giving a solution $\mathbf{x} = A^+\mathbf{b}$. We call the matrix A^+ the pseudo-inverse of A .

Example 3.6.11. Find the pseudo-inverse of the matrix $A = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$.

Solution: Apply Procedure 3.5.3 to solve $Ax = \mathbf{b}$ for any right-hand side \mathbf{b} .

(a) This matrix has an SVD

$$A = \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 5 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \end{bmatrix}^T = USV^T.$$

(b) Hence $\mathbf{z} = U^T\mathbf{b} = \begin{bmatrix} \frac{3}{5}b_1 + \frac{4}{5}b_2 \\ -\frac{4}{5}b_1 + \frac{3}{5}b_2 \end{bmatrix}$.

- (c) To solve the diagonal system $Sy = z$ consider $\begin{bmatrix} 5 \\ 0 \end{bmatrix} y = \begin{bmatrix} \frac{3}{5}b_1 + \frac{4}{5}b_2 \\ -\frac{4}{5}b_1 + \frac{3}{5}b_2 \end{bmatrix}$: approximate by neglecting the second component in the equations and just set $y = \frac{3}{25}b_1 + \frac{4}{25}b_2$.
- (d) Then the procedure's solution is $x = Vy = 1(\frac{3}{25}b_1 + \frac{4}{25}b_2) = \frac{3}{25}b_1 + \frac{4}{25}b_2$.

That is, for all right-hand side vectors b , this least square solution

$$x = A^+b \quad \text{for pseudo-inverse } A^+ = \left[\begin{array}{cc} \frac{3}{25} & \frac{4}{25} \end{array} \right].$$

■

A pseudo-inverse A^+ of a non-invertible matrix A is only an ‘inverse’ because the pseudo-inverse builds in extra information that we *choose* to be desirable. This extra information rationalises all the contradictions encountered in trying to construct an inverse of a non-invertible matrix. Namely, we *choose* to desire that the pseudo-inverse solves the *nearest* consistent system to the one specified, and we *choose* the smallest of all possibilities then allowed. However, although there are many situations where these choices are useful, beware that there are also many situations where such choices are not appropriate.

Theorem 3.6.12 (pseudo-inverse). *For a given $m \times n$ matrix A and in the context of wanting to solve $Ax = b$, Procedure 3.5.3 forms a linear transformation $T : \mathbb{R}^m \rightarrow \mathbb{R}^n$ to find the smallest solution x (Theorem 3.5.9) to the closest consistent system $Ax = \tilde{b}$ (Theorem 3.5.5). This linear transformation has an $n \times m$ standard matrix A^+ called the **pseudo-inverse**, or **Moore–Penrose inverse**, of matrix A .*

Proof. First, for each right-hand side vector b in \mathbb{R}^m , the procedure gives a result x in \mathbb{R}^n and so is some function $T : \mathbb{R}^m \rightarrow \mathbb{R}^n$. We proceed to confirm that Procedure 3.5.3 satisfies the two defining properties of a linear transformation (Definition 3.6.1). For each of any two right-hand side vectors $b', b'' \in \mathbb{R}^m$ let Procedure 3.5.3 generate similarly dashed intermediaries (z', z'', y', y'') through to two corresponding least square solutions $x' = T(b') \in \mathbb{R}^n$ and $x'' = T(b'') \in \mathbb{R}^n$. Throughout let the matrix A have SVD $A = USV^T$ and set $r = \text{rank } A$.

3.6.1a : To check $T(b' + b'') = T(b') + T(b'')$, apply the procedure when the right-hand side is $b = b' + b''$:

2. solve $Uz = b$ by $z = U^T b = U^T(b' + b'') = U^T b' + U^T b'' = z' + z''$;
3. • for $i = 1, \dots, r$ set $y_i = z_i/\sigma_i = (z'_i + z''_i)/\sigma_i = z'_i/\sigma_i + z''_i/\sigma_i = y'_i + y''_i$, and

- for $i = r+1, \dots, n$ set the free variables $y_i = 0 = 0 + 0 = y'_i + y''_i$ to obtain the smallest solution (Theorem 3.5.9);;

and hence $\mathbf{y} = \mathbf{y}' + \mathbf{y}''$;

4. solve $\mathbf{x} = V^T \mathbf{y}$ with $\mathbf{x} = V\mathbf{y} = V(\mathbf{y}' + \mathbf{y}'') = V\mathbf{y}' + V\mathbf{y}'' = \mathbf{x}' + \mathbf{x}''$.

Since result $\mathbf{x} = \mathbf{x}' + \mathbf{x}''$, thus $T(\mathbf{b}' + \mathbf{b}'') = T(\mathbf{b}') + T(\mathbf{b}'')$.

3.6.1b : To check $T(c\mathbf{b}') = cT(\mathbf{b}')$ for any scalar c , apply the procedure when the right-hand side is $\mathbf{b} = c\mathbf{b}'$:

2. solve $U\mathbf{z} = \mathbf{b}$ by $\mathbf{z} = U^T \mathbf{b} = U^T(c\mathbf{b}') = cU^T \mathbf{b}' = cz'$;
 3. • for $i = 1, \dots, r$ set $y_i = z_i/\sigma_i = (cz'_i)/\sigma_i = c(z'_i/\sigma_i) = cy'_i$, and
 - for $i = r+1, \dots, n$ set the free variables $y_i = 0 = c0 = cy'_i$ to obtain the smallest solution (Theorem 3.5.9),
- and hence $\mathbf{y} = cy'$;

4. solve $\mathbf{x} = V^T \mathbf{y}$ with $\mathbf{x} = V\mathbf{y} = V(c\mathbf{y}') = cV\mathbf{y}' = cx'$.

Since the result $\mathbf{x} = cx'$, consequently $T(c\mathbf{b}') = cT(\mathbf{b}')$.

Since Procedure 3.5.3, denoted by T , is a linear transform $T : \mathbb{R}^m \rightarrow \mathbb{R}^n$, Theorem 3.6.8 assures us T has a corresponding $n \times m$ standard matrix which we denote A^+ and call the pseudo-inverse. \square

Example 3.6.13. Find the pseudo-inverse of the matrix $A = [5 \ 12]$.

Solution: Apply Procedure 3.5.3 to solve $Ax = b$ for any right-hand side b .

(a) This matrix has an SVD

$$A = [5 \ 12] = [1] [13 \ 0] \begin{bmatrix} \frac{5}{13} & -\frac{12}{13} \\ \frac{12}{13} & \frac{5}{13} \end{bmatrix}^T = USV^T.$$

(b) Hence $z = U^T b = 1b = b$.

(c) The diagonal system $S\mathbf{y} = z$ becomes $[13 \ 0] \mathbf{y} = b$ with general solution $\mathbf{y} = (b/13, y_2)$. The smallest of these solutions is $\mathbf{y} = (b/13, 0)$.

(d) Then the procedure's result is

$$\mathbf{x} = V\mathbf{y} = \begin{bmatrix} \frac{5}{13} & -\frac{12}{13} \\ \frac{12}{13} & \frac{5}{13} \end{bmatrix} \begin{bmatrix} b/13 \\ 0 \end{bmatrix} = \begin{bmatrix} \frac{5}{169}b \\ -\frac{12}{169}b \end{bmatrix}.$$

That is, for all right-hand sides b , this procedure's result is

$$\mathbf{x} = A^+b \quad \text{for pseudo-inverse } A^+ = \begin{bmatrix} \frac{5}{169} \\ -\frac{12}{169} \end{bmatrix}.$$

■

Example 3.6.14. Recall that Example 3.5.1 explored how to determine weight from four apparently contradictory measurements. The exploration showed that Procedure 3.5.3 agrees with the traditional method of simple averaging. Let's see that the pseudo-inverse implements the simple average of the four weights.

Recall that Example 3.5.1 sought to solve the inconsistent system

$$Ax = \mathbf{b}, \quad \text{namely} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} x = \begin{bmatrix} 84.8 \\ 84.1 \\ 84.7 \\ 84.4 \end{bmatrix}.$$

To find the pseudo-inverse, do this for arbitrary right-hand side \mathbf{b} .

(a) Recall this matrix A of ones has an SVD of

$$A = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} [1]^T = USV^T.$$

(b) Solve $Uz = \mathbf{b}$ by computing

$$\begin{aligned} z &= U^T \mathbf{b} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \mathbf{b} \\ &= \begin{bmatrix} \frac{1}{2}b_1 + \frac{1}{2}b_2 + \frac{1}{2}b_3 + \frac{1}{2}b_4 \\ \frac{1}{2}b_1 + \frac{1}{2}b_2 - \frac{1}{2}b_3 - \frac{1}{2}b_4 \\ \frac{1}{2}b_1 - \frac{1}{2}b_2 - \frac{1}{2}b_3 + \frac{1}{2}b_4 \\ \frac{1}{2}b_1 - \frac{1}{2}b_2 + \frac{1}{2}b_3 - \frac{1}{2}b_4 \end{bmatrix}. \end{aligned}$$

(c) Now try to solve $Sy = z$, that is,

$$\begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} y = \begin{bmatrix} \frac{1}{2}b_1 + \frac{1}{2}b_2 + \frac{1}{2}b_3 + \frac{1}{2}b_4 \\ \frac{1}{2}b_1 + \frac{1}{2}b_2 - \frac{1}{2}b_3 - \frac{1}{2}b_4 \\ \frac{1}{2}b_1 - \frac{1}{2}b_2 - \frac{1}{2}b_3 + \frac{1}{2}b_4 \\ \frac{1}{2}b_1 - \frac{1}{2}b_2 + \frac{1}{2}b_3 - \frac{1}{2}b_4 \end{bmatrix}.$$

Instead of seeking an *exact* solution, we *have to* adjust the last three components to zero. Hence we find a solution to a slightly different problem by solving

$$\begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} y = \begin{bmatrix} \frac{1}{2}b_1 + \frac{1}{2}b_2 + \frac{1}{2}b_3 + \frac{1}{2}b_4 \\ 0 \\ 0 \\ 0 \end{bmatrix},$$

with solution $y = \frac{1}{4}b_1 + \frac{1}{4}b_2 + \frac{1}{4}b_3 + \frac{1}{4}b_4$.

(d) Lastly, solve $V^T x = y$ by computing

$$x = V y = 1y = \frac{1}{4}b_1 + \frac{1}{4}b_2 + \frac{1}{4}b_3 + \frac{1}{4}b_4 = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix} \mathbf{b}.$$

Hence the pseudo-inverse of matrix A is $A^+ = \begin{bmatrix} \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$. This pseudo-inverse implements the traditional answer of averaging measurements. ■

Example 3.6.15. Recall that Example 3.5.2 rates three table tennis players, Anne, Bob and Chris. The rating involved solving the inconsistent system $Ax = \mathbf{b}$ for the particular matrix and vector

$$\begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \\ 0 & 1 & -1 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix}.$$

Find the pseudo-inverse of this matrix A . Use the pseudo-inverse to rate the players in the cases of Examples 3.3.10 and 3.5.2.

Solution: To find the pseudo-inverse, follow Procedure 3.5.3 with a general right-hand side vector \mathbf{b} .

(a) Compute an SVD $A = USV^T$ in Matlab/Octave with `[U,S,V]=svd(A)`:

$$\begin{aligned} \mathbf{U} &= \\ &\begin{array}{ccc} 0.4082 & -0.7071 & 0.5774 \\ -0.4082 & -0.7071 & -0.5774 \\ -0.8165 & -0.0000 & 0.5774 \end{array} \\ \mathbf{S} &= \\ &\begin{array}{ccc} 1.7321 & 0 & 0 \\ 0 & 1.7321 & 0 \\ 0 & 0 & 0.0000 \end{array} \\ \mathbf{V} &= \\ &\begin{array}{ccc} 0.0000 & -0.8165 & 0.5774 \\ -0.7071 & 0.4082 & 0.5774 \\ 0.7071 & 0.4082 & 0.5774 \end{array} \end{aligned}$$

Upon recognising various square-roots, these matrices are

$$\begin{aligned} U &= \begin{bmatrix} \frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{3}} \\ -\frac{2}{\sqrt{6}} & 0 & \frac{1}{\sqrt{3}} \end{bmatrix}, \\ S &= \begin{bmatrix} \sqrt{3} & 0 & 0 \\ 0 & \sqrt{3} & 0 \\ 0 & 0 & 0 \end{bmatrix}, \end{aligned}$$

$$V = \begin{bmatrix} 0 & -\frac{2}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \end{bmatrix}.$$

The system of equations for the ratings becomes

$$\underbrace{Ax}_{=z} = U \underbrace{S V^T \overbrace{x}^{=y}}_{=z} = b.$$

(b) As U is orthogonal, $Uz = b$ has unique solution

$$z = U^T b = \begin{bmatrix} \frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{6}} \\ -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{bmatrix} b.$$

(c) Now solve $Sy = z$. But S has a troublesome zero on the diagonal. So interpret the equation $Sy = z$ in detail as

$$\begin{bmatrix} \sqrt{3} & 0 & 0 \\ 0 & \sqrt{3} & 0 \\ 0 & 0 & 0 \end{bmatrix} y = \begin{bmatrix} \frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{6}} \\ -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ \frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \end{bmatrix} b :$$

i. the first line requires $y_1 = \frac{1}{\sqrt{3}} \left[\frac{1}{\sqrt{6}} \ -\frac{1}{\sqrt{6}} \ -\frac{2}{\sqrt{6}} \right] b = \frac{1}{3} \left[\frac{1}{\sqrt{2}} \ -\frac{1}{\sqrt{2}} \ -\sqrt{2} \right]$;

ii. the second line requires $y_2 = \frac{1}{\sqrt{3}} \left[-\frac{1}{\sqrt{2}} \ -\frac{1}{\sqrt{2}} \ 0 \right] b = \left[-\frac{1}{\sqrt{6}} \ -\frac{1}{\sqrt{6}} \ 0 \right] b$;

iii. the third line requires $0y_3 = \left[\frac{1}{\sqrt{3}} \ -\frac{1}{\sqrt{3}} \ \frac{1}{\sqrt{3}} \right] b$ which generally cannot be satisfied, so we set $y_3 = 0$ to get the *smallest solution of the system after projecting b onto the column space of A* .

(d) Finally, as V is orthogonal, $V^T x = y$ has the solution $x = Vy$ (unique for each valid y):

$$\begin{aligned} x &= Vy \\ &= \begin{bmatrix} 0 & -\frac{2}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \end{bmatrix} \begin{bmatrix} \frac{1}{3\sqrt{2}} & -\frac{1}{3\sqrt{2}} & -\frac{\sqrt{2}}{3} \\ -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{6}} & 0 \\ 0 & 0 & 0 \end{bmatrix} b \\ &= \frac{1}{3} \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix} b \end{aligned}$$

Hence the pseudo-inverse of A is

$$A^+ = \frac{1}{3} \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix}.$$

- In Example 3.3.10, Anne beat Bob 3-2 games; Anne beat Chris 3-1; Bob beat Chris 3-2 so the right-hand side vector is $\mathbf{b} = (1, 2, 1)$. The procedure's ratings are then, as before,

$$\mathbf{x} = A^+ \mathbf{b} = \frac{1}{3} \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}.$$

- In Example 3.5.2, Bob instead beat Chris 3-1 so the right-hand side vector is $\mathbf{b} = (1, 2, 2)$. The procedure's ratings are then, as before,

$$\mathbf{x} = A^+ \mathbf{b} = \frac{1}{3} \begin{bmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 0 & -1 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \\ \frac{1}{3} \\ -\frac{4}{3} \end{bmatrix}.$$

■

Theorem 3.6.16. *In the case of an $m \times n$ matrix A with $\text{rank } A = n$ (so $m \geq n$), the pseudo-inverse $A^+ = (A^T A)^{-1} A^T$.*

Proof. Apply Procedure 3.5.3 to the system $A\mathbf{x} = \mathbf{b}$ for an arbitrary $\mathbf{b} \in \mathbb{R}^m$.

1. Let $m \times n$ matrix A have SVD $A = U S V^T$. Since $\text{rank } A = n$ there are n nonzero singular values on the diagonal of $m \times n$ matrix S , and so $m \geq n$.
2. Solve $U\mathbf{z} = \mathbf{b}$ with $\mathbf{z} = U^T \mathbf{b} \in \mathbb{R}^m$.
3. Approximately solve $S\mathbf{y} = \mathbf{z}$ by setting $y_i = z_i/\sigma_i$ for $i = 1, \dots, n$ and neglecting the last $(m - n)$ equations. Setting the $n \times m$ matrix

$$S^+ := \begin{bmatrix} 1/\sigma_1 & 0 & \cdots & 0 & 0 & \cdots & 0 \\ 0 & 1/\sigma_2 & & 0 & 0 & \cdots & 0 \\ \vdots & & \ddots & \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1/\sigma_n & 0 & \cdots & 0 \end{bmatrix}$$

this is identical to setting $\mathbf{y} = S^+ \mathbf{z} = S^+ U^T \mathbf{b}$.

4. Solve $V^T \mathbf{x} = \mathbf{y}$ with $\mathbf{x} = V\mathbf{y} = VS^+U^T\mathbf{b}$.

Hence the pseudo-inverse is $A^+ = VS^+U^T$.

Let's find $(A^TA)^{-1}A^T$ is the same expression. First, since $A^T = (USV^T)^T = VS^TU^T$,

$$A^TA = VS^TU^TUSV^T = VS^TSV^T = V(S^TS)V^T$$

where $n \times n$ matrix $(S^TS) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$. Since $\text{rank } A = n$, the singular values $\sigma_1, \sigma_2, \dots, \sigma_n > 0$ and so matrix (S^TS) is invertible as it is square and diagonal with all nonzero elements in the diagonal, and the inverse is $(S^TS)^{-1} = \text{diag}(1/\sigma_1^2, 1/\sigma_2^2, \dots, 1/\sigma_n^2)$ (Theorem 3.2.21). Second, the $n \times n$ matrix (A^TA) is invertible as it has SVD $V(S^TS)V^T$ with n nonzero singular values (Theorem 3.3.21d), and so

$$\begin{aligned} (A^TA)^{-1}A^T &= (V(S^TS)V^T)^{-1}VS^TU^T \\ &= (V^T)^{-1}(S^TS)^{-1}V^{-1}VS^TU^T \\ &= V(S^TS)^{-1}V^TUV^T \\ &= V(S^TS)^{-1}S^TU^T \\ &= VS^+U^T, \end{aligned}$$

where the last equality follows because

$$\begin{aligned} &(S^TS)^{-1}S^T \\ &= \text{diag}(1/\sigma_1^2, 1/\sigma_2^2, \dots, 1/\sigma_n^2) \text{diag}_{n \times m}(\sigma_1, \sigma_2, \dots, \sigma_n) \\ &= \text{diag}_{n \times m}(1/\sigma_1, 1/\sigma_2, \dots, 1/\sigma_n) = S^+. \end{aligned}$$

Hence the pseudo-inverse $A^+ = VS^+U^T = (A^TA)^{-1}A^T$. \square

Theorem 3.6.17. *If A is an invertible matrix, then the pseudo-inverse $A^+ = A^{-1}$, the inverse.*

Proof. If A is invertible it must be square, say $n \times n$, and of rank $A = n$ (Theorem 3.3.21e). Further, A^T is invertible with inverse $(A^{-1})^T$ (Theorem 3.2.11d). Then the expression from Theorem 3.6.16 for the pseudo-inverse gives

$$A^+ = (A^TA)^{-1}A^T = A^{-1}(A^T)^{-1}A^T = A^{-1}I_n = A^{-1}.$$

\square

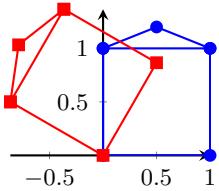
Computer considerations Except for easy cases, we (almost) never explicitly compute the pseudo-inverse of a matrix. In practical computation, forming A^TA and then manipulating it is both expensive and error enhancing: for example, $\text{cond}(A^TA) = (\text{cond } A)^2$ so matrix A^TA typically has a much worse condition number than matrix A . Computationally there are (almost) always better ways to proceed, such as Procedure 3.5.3. Like an inverse, a pseudo-inverse is a theoretical device, rarely a practical tool.

A main point of this subsection is to illustrate how a complicated procedure is conceptually expressible as a linear transformation, and so has associated matrix properties such as being equivalent to multiplication by the pseudo-inverse.

3.6.3 Function composition connects to matrix inverse

To achieve a complex goal we typically decompose the goal into a set of smaller tasks and achieve those tasks one after another. The analogy in linear algebra is that we often apply linear transforms one after another to build up or solve a complex problem. This section introduces how applying a sequence of linear transforms is equivalent to one grand linear transform.

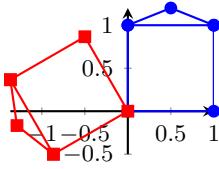
Example 3.6.18 (simple rotation). Recall Example 3.6.9b on rotation by 60° (illustrated in the margin) with its standard matrix



$$[R] = \begin{bmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix}.$$

Consider two successive rotations by 60° : show that the standard matrix of the resultant rotation by 120° is the same as the product $[R][R]$.

Solution: On the one hand, rotation by $120^\circ = 2\pi/3$, call it S , transforms the unit vectors as (illustrated in the margin)



$$\begin{aligned} S(\mathbf{e}_1) &= (\cos \frac{2\pi}{3}, \sin \frac{2\pi}{3}) = \left(-\frac{1}{2}, \frac{\sqrt{3}}{2}\right), \\ S(\mathbf{e}_2) &= (-\sin \frac{2\pi}{3}, \cos \frac{2\pi}{3}) = \left(-\frac{\sqrt{3}}{2}, -\frac{1}{2}\right). \end{aligned}$$

Form the standard matrix with these as its columns, in order,

$$[S] = [S(\mathbf{e}_1) \quad S(\mathbf{e}_2)] = \begin{bmatrix} -\frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & -\frac{1}{2} \end{bmatrix}.$$

On the other hand, the matrix multiplication

$$\begin{aligned} [R][R] &= \begin{bmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix} \begin{bmatrix} \frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{4} - \frac{3}{4} & -\frac{\sqrt{3}}{4} - \frac{\sqrt{3}}{4} \\ \frac{\sqrt{3}}{4} + \frac{\sqrt{3}}{4} & -\frac{3}{4} + \frac{1}{4} \end{bmatrix} \\ &= \begin{bmatrix} -\frac{1}{2} & -\frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} & -\frac{1}{2} \end{bmatrix} = [S]. \end{aligned}$$

That is, multiplying the two matrices is equivalent to performing the two rotations in succession: the next theorem confirms this is generally true.

■

Theorem 3.6.19. Let $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $S : \mathbb{R}^m \rightarrow \mathbb{R}^p$ be linear transformations. Recalling the **composition** of functions is $(S \circ T)(\mathbf{v}) = S(T(\mathbf{v}))$, then $S \circ T : \mathbb{R}^n \rightarrow \mathbb{R}^p$ is a linear transformation with standard matrix $[S \circ T] = [S][T]$.

Proof. Let matrix $A = [S]$ ($p \times m$) and $B = [T]$ ($m \times n$). Let \mathbf{u} be any vector in \mathbb{R}^n , then $(S \circ T)(\mathbf{u}) = S(T(\mathbf{u})) = S(B\mathbf{u}) = A(B\mathbf{u}) = (AB)\mathbf{u}$ (using associativity, Theorem 3.1.18c). Hence the effect of $S \circ T$ is identical to multiplication by the $p \times n$ matrix (AB) . It is thus a matrix transformation, which is consequently linear (Theorem 3.6.6), and its standard matrix $[S \circ T] = AB = [S][T]$. \square

Example 3.6.20. Consider the linear transform $T : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ defined by $T(x_1, x_2, x_3) := (3x_1 + x_2, -x_2 - 7x_3)$, and the linear transform $S : \mathbb{R}^2 \rightarrow \mathbb{R}^4$ defined by $S(y_1, y_2) = (-y_1, -3y_1 + 2y_2, 2y_1 - y_2, 2y_2)$. Find the standard matrix of the linear transform $S \circ T$, and also that of $T \circ S$.

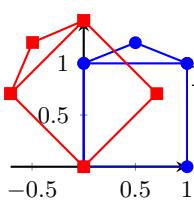
Solution: From the given formulas for the two given linear transforms we write down the standard matrices

$$[T] = \begin{bmatrix} 3 & 1 & 0 \\ 0 & -1 & -7 \end{bmatrix} \quad \text{and} \quad [S] = \begin{bmatrix} -1 & 0 \\ -3 & 2 \\ 2 & -1 \\ 0 & 2 \end{bmatrix}.$$

First, Theorem 3.6.19 assures us the standard matrix of the composition

$$\begin{aligned} [S \circ T] &= [S][T] \\ &= \begin{bmatrix} -1 & 0 \\ -3 & 2 \\ 2 & -1 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 3 & 1 & 0 \\ 0 & -1 & -7 \end{bmatrix} \\ &= \begin{bmatrix} -3 & -1 & 0 \\ -9 & -5 & -14 \\ 6 & 3 & 7 \\ 0 & -2 & -14 \end{bmatrix}. \end{aligned}$$

However, second, the standard matrix of $T \circ S$ does not exist because it would require the multiplication of a 2×3 matrix by a 4×2 matrix, and such a multiplication is not defined. The failure is rooted earlier in the question because $S : \mathbb{R}^2 \rightarrow \mathbb{R}^4$ and $T : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ so a result of S , which is in \mathbb{R}^4 , cannot be used as an argument to T , which must be in \mathbb{R}^3 : the lack of a defined multiplication is a direct reflection of this incompatibility in ' $T \circ S$ ' which means $T \circ S$ cannot exist. \blacksquare



Example 3.6.21. Find the standard matrix of the transformation of the plane that first rotates by 45° about the origin, and then reflects in the vertical axis.

Solution: Two possible solutions are the following.

- Let R denote the rotation about the origin of the plane by 45° , $R : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ (illustrated in the margin). Its standard matrix is

$$[R] = [R(1, 0) \ R(0, 1)] = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}.$$

Let F denote the reflection in the vertical axis of the plane, $F : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ (illustrated in the margin). Its standard matrix is

$$[F] = [F(1, 0) \ F(0, 1)] = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Then the standard matrix of the composition

$$[F \circ R] = [F][R] = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

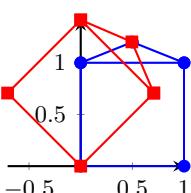
- Alternatively, just consider the action of the two component transforms on the standard unit vectors.

- $(F \circ R)(1, 0) = F(R(1, 0))$ which first rotates $(1, 0)$ to point to the top-right, then reflects in the vertical axis to point to the top-left and thus $(F \circ R)(1, 0) = (-\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$.
- $(F \circ R)(0, 1) = F(R(0, 1))$ which first rotates $(0, 1)$ to point to the top-left, then reflects in the vertical axis to point to the top-right and thus $(F \circ R)(0, 1) = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$.

Then the standard matrix of the composition (as illustrated in the margin)

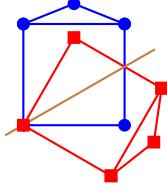
$$[F \circ R] = [(F \circ R)(1, 0) \ (F \circ R)(0, 1)] = \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

As an extension, check that although $R \circ F$ is defined, it is different to $F \circ R$: the difference corresponds to the non-commutativity of matrix multiplication (section 3.1.3).



Having introduced and characterised the composition of linear transformations, we now discuss when two transforms composed together end up ‘cancelling’ each other out.

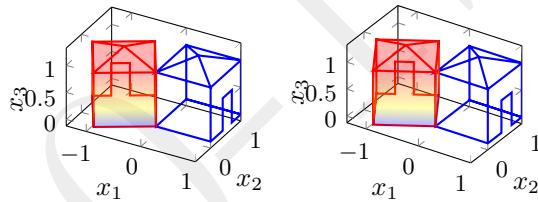
Example 3.6.22 (inverse transforms). (a) Let S be rotation of the plane by 60° , and T be rotation of the plane by -60° . Then $S \circ T$ is first rotation by -60° by T , and second rotation by 60° by S , results together in the plane being unchanged. Because $S \circ T$ is effectively the identity transformation, we call these rotations the inverse transform of each other.



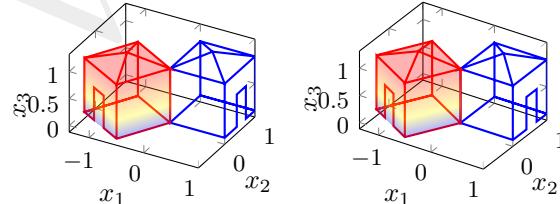
- (b) Let R be reflection of the plane in the line at 30° to the horizontal (illustrated in the margin). Then $R \circ R$ is first reflection in the line at 30° by R , and second another reflection in the line at 30° by R , results together in the plane being unchanged. Because $R \circ R$ is effectively the identity transformation, the reflection R is its own inverse. ■

Definition 3.6.23. Let S and T be linear transforms from \mathbb{R}^n to \mathbb{R}^n (the same dimension). If $S \circ T = T \circ S = I$, the identity transformation, then S and T are **inverse transformations** of each other. Further, we say S and T are **invertible**.

Example 3.6.24. Let $S : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be rotation about the vertical axis by 120° (as illustrated in stereo below),



and let $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be rotation about the vertical axis by 240° (below).



Argue that $S \circ T = T \circ S = I$ the identity and so S and T are inverse transformations of each other.

Solution: A basic argument is that rotation by 120° together with a rotation of 240° about the same axis, in either order, is the same as a rotation by 360° about the axis. But a 360° rotation leaves everything unchanged and so must be the identity.

Alternatively one could dress up the argument with some algebra as in the following. First consider $T \circ S$:

- the vertical unit vector e_3 is unchanged by both S and T so $(T \circ S)(e_3) = T(S(e_3)) = e_3$;
- the unit vector e_1 is rotated 120° by S , and then by 240° by T which is a total of 360° , that is, it is rotated back to itself so $(T \circ S)(e_1) = T(S(e_1)) = e_1$; and

- the unit vector e_2 is rotated 120° by S , and then by 240° by T which is a total of 360° , that is, it is rotated back to itself so $(T \circ S)(e_2) = T(S(e_2)) = e_2$.

Form these results into a matrix to deduce the standard matrix

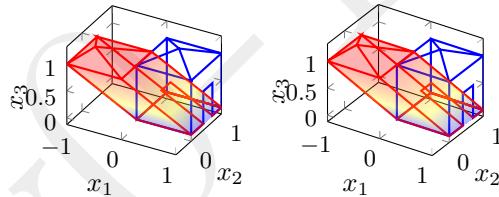
$$[T \circ S] = [e_1 \ e_2 \ e_3] = I_3$$

which is the standard matrix of the identity transform. Hence $T \circ S = I$ the identity.

Second, an exactly corresponding argument gives $S \circ T = I$. By Definition 3.6.23, S and T are inverse transforms of each other.

■

Example 3.6.25. In some violent weather a storm passes and the strong winds lean a house sideways as in the shear transformation illustrated below.



Estimate the standard matrix of the shear transformation shown. To restore the house back upright, we need to shear it an equal amount in the opposite direction: hence write down the standard matrix of the inverse shear. Confirm that the product of the two standard matrices is the standard matrix of the identity.

Solution: As shown, the unit vectors e_1 and e_2 are unchanged by the storm S . However, the vertical unit vector e_3 is sheared by the storm to $S(e_3) = (-1, -0.5, 1)$. Hence the standard matrix of the storm S is

$$[S] = \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -\frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix}.$$

To restore, R , the house upright we need to shear in the opposite direction, so the restoration shear has $R(e_3) = (1, \frac{1}{2}, 1)$; that is, the standard matrix of the inverse is

$$[R] = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & \frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix}.$$

Multiplying these matrices together gives

$$[R \circ S] = [R][S]$$

$$\begin{aligned}
 &= \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & \frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -\frac{1}{2} \\ 0 & 0 & 1 \end{bmatrix} \\
 &= \begin{bmatrix} 1+0+0 & 0+0+0 & -1+0+1 \\ 0+0+0 & 0+1+0 & 0-\frac{1}{2}+\frac{1}{2} \\ 0+0+0 & 0+0+0 & 0+0+1 \end{bmatrix} \\
 &= I_3.
 \end{aligned}$$

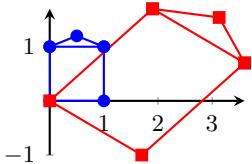
This is the standard matrix of the identity transform. ■

Because of the exact correspondence between linear transformations and matrix multiplication, the inverse of a transform exactly corresponds to the inverse of a matrix. In the last Example 3.6.25, because $[R][S] = I_3$ we know that the matrices $[R]$ and $[S]$ are inverses of each other. Correspondingly, the transforms R and S are inverses of each other.

Theorem 3.6.26. *Let $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be an invertible linear transformation. Then its standard matrix $[T]$ is invertible, and $[T^{-1}] = [T]^{-1}$.*

Proof. As T is invertible, let the symbol T^{-1} denote its inverse. Since both are linear transforms, they both have standard matrices, $[T]$ and $[T^{-1}]$. Then $[T \circ T^{-1}] = [T][T^{-1}]$; but also $[T \circ T^{-1}] = [I] = I_n$; so $[T][T^{-1}] = I_n$. Similarly, $[T^{-1} \circ T] = [T^{-1} \circ T] = [I] = I_n$. Consequently the matrices $[T]$ and $[T^{-1}]$ are the inverses of each other. □

Example 3.6.27. Estimate the standard matrix of the linear transform T illustrated in the margin. Then use Theorem 3.2.6 to determine the standard matrix of its inverse transform T^{-1} . Hence sketch how the inverse transforms the unit square and write a sentence or two about how the sketch confirms it is a reasonable inverse.

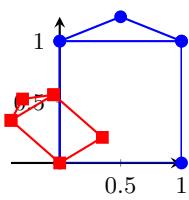


Solution: The illustrated transform shows $T(1, 0) \approx (1.7, -1)$ and $T(0, 1) \approx (1.9, 1.7)$ hence its standard matrix is

$$[T] \approx \begin{bmatrix} 1.7 & 1.9 \\ -1 & 1.7 \end{bmatrix}.$$

Using Theorem 3.2.6 the inverse of this matrix is, since its determinant = $1.7 \cdot 1.7 - (-1) \cdot 1.9 = 4.79$,

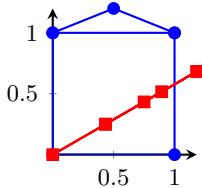
$$[T^{-1}] = [T]^{-1} = \frac{1}{4.79} \begin{bmatrix} 1.7 & -1.9 \\ 1 & 1.7 \end{bmatrix} \approx \begin{bmatrix} 0.35 & -0.40 \\ 0.21 & 0.35 \end{bmatrix}$$



This matrix determines that the inverse maps the corners of the unit square as $T^{-1}(1, 0) = (0.35, 0.21)$, $T^{-1}(0, 1) = (-0.40, 0.35)$ and $T^{-1}(1, 1) = (-0.05, 0.56)$. Hence the unit square is transformed as

shown in the margin. The original transform, roughly, rotated the unit square clockwise and stretched it: the sketch shows the inverse roughly rotates the unit square anti-clockwise and shrinks it. Thus the inverse does indeed undo the action of the original transform. ■

Example 3.6.28. Determine if the orthogonal projection of the plane onto the line at 30° to the horizontal (illustrated in the margin) is an invertible transform; if it is find its inverse.



Solution: Recall that Theorem 3.5.22 gives the matrix of an orthogonal projection as WW^T where columns of W are an orthonormal basis for the projected space. Here the projected space is the line at 30° to the horizontal (illustrated in the margin) which has orthonormal basis of the one vector $\mathbf{w} = (\cos 30^\circ, \sin 30^\circ) = (\frac{\sqrt{3}}{2}, \frac{1}{2})$. Hence the standard matrix of the projection is

$$\mathbf{w}\mathbf{w}^T = \begin{bmatrix} \frac{\sqrt{3}}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} \begin{bmatrix} \frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix} = \begin{bmatrix} \frac{3}{4} & \frac{\sqrt{3}}{4} \\ \frac{\sqrt{3}}{4} & \frac{1}{4} \end{bmatrix}.$$

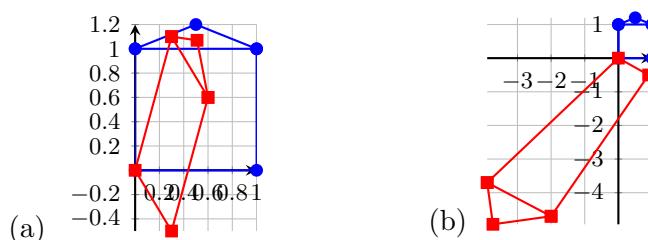
From Theorem 3.2.6 this matrix is invertible only if the determinant ($\det = ad - bc$) is nonzero, but here

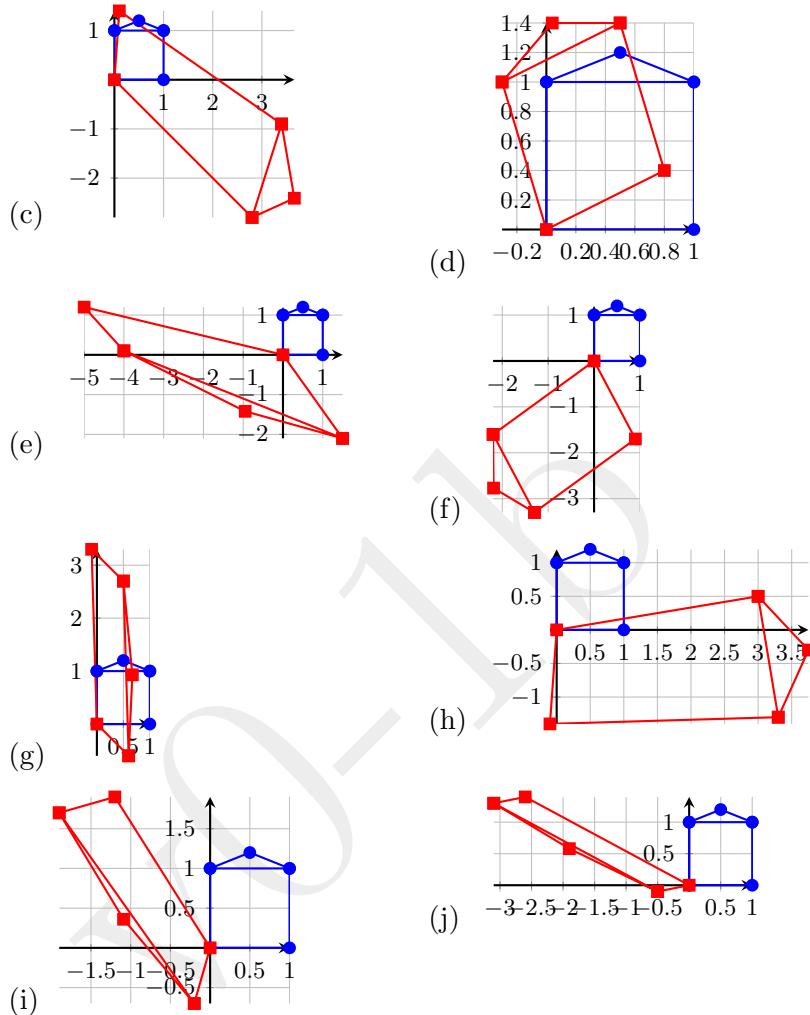
$$\det(\mathbf{w}\mathbf{w}^T) = \frac{3}{4} \cdot \frac{1}{4} - \frac{\sqrt{3}}{4} \cdot \frac{\sqrt{3}}{4} = \frac{3}{16} - \frac{3}{16} = 0.$$

Since its standard matrix is not invertible, the given orthogonal projection is also not invertible (as the illustration shows, the projection ‘squashes’ the plane onto the line, which cannot be uniquely undone, and hence is not invertible). ■

3.6.4 Exercises

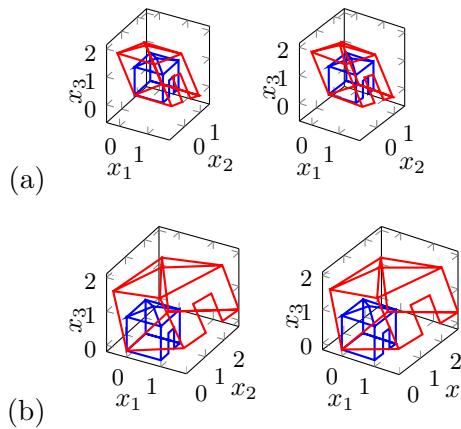
Exercise 3.6.1. Which of the following illustrated transformations of the plane *cannot* be that of a linear transformation? In each illustration of a transformation T , the four corners of the blue unit square ($(0,0)$, $(1,0)$, $(1,1)$ and $(0,1)$), are mapped to the four corners of the red figure ($T(0,0)$, $T(1,0)$, $T(1,1)$ and $T(0,1)$ —the ‘roof’ of the unit square clarifies which side goes where).

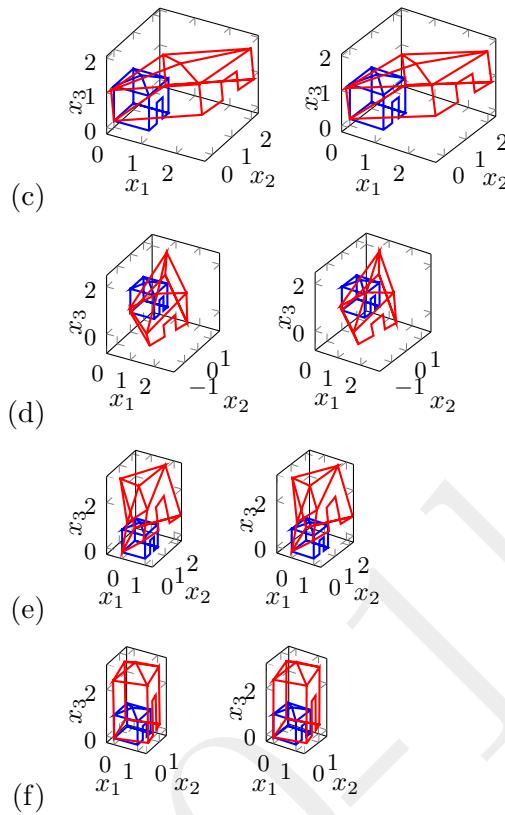




Exercise 3.6.2. Consider the transforms of Example 3.6.1: for those transforms that *may* be linear transforms, assume they are and so estimate roughly the standard matrix of each such linear transform.

Exercise 3.6.3. Consider the following illustrated transforms of \mathbb{R}^3 . Which *cannot* be that of a linear transformation?





Exercise 3.6.4. Let $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a linear transformation. Prove from the Definition 3.6.1 that $T(\mathbf{0}) = \mathbf{0}$ and $T(\mathbf{u} - \mathbf{v}) = T(\mathbf{u}) - T(\mathbf{v})$ for all $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$.

Exercise 3.6.5 (equivalent definition). Consider a function $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Prove that T is a linear transformation (Definition 3.6.1) if and only if $T(c_1\mathbf{v}_1 + c_2\mathbf{v}_2) = c_1T(\mathbf{v}_1) + c_2T(\mathbf{v}_2)$ for all $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^n$ and all scalars c_1, c_2 .

Exercise 3.6.6. Given $T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a linear transformation, use induction to prove that for all k

$$T(c_1\mathbf{u}_1 + c_2\mathbf{u}_2 + \cdots + c_k\mathbf{u}_k) = c_1T(\mathbf{u}_1) + c_2T(\mathbf{u}_2) + \cdots + c_kT(\mathbf{u}_k)$$

for all scalars c_1, c_2, \dots, c_k and all vectors $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k$ in \mathbb{R}^n .

Exercise 3.6.7. Consider each of the following vector transformations: if it is a linear transformation, then write down its standard matrix.

- (a) $A(x, y, z) = (-3x + 2y, -z)$
- (b) $B(x, y) = (0, 3x + 7y, -2y, -3y)$
- (c) $C(x_1, x_2, \dots, x_5) = (2x_1 + x_2 - 2x_3, 7x_1 + 7x_4)$
- (d) $D(x, y) = (2x + 3, -5y + 3, 2x - 4y + 3, 0, -6x)$
- (e) $E(p, q, r, s) = (-3p - 4r, -s, 0, p + r + 6s, 5p + 6q - s)$
- (f) $F(x, y) = (5x + 4y, x^2, 2y, -4x, 0)$

$$(g) \quad G(x_1, x_2, x_3, x_4) = (-x_1 + 4x_2 + e^{x_3}, 8x_4)$$

$$(h) \quad H(u_1, u_2, \dots, u_5) = (7u_1 - 9u_3 - u_5, 3u_1 - 3u_4)$$

Exercise 3.6.8. Use Procedure 3.5.3 to derive that the pseudo-inverse of the general 2×1 matrix $A = \begin{bmatrix} a \\ b \end{bmatrix}$ is the 1×2 matrix $A^+ = \begin{bmatrix} a \\ \frac{a}{a^2+b^2} & \frac{b}{a^2+b^2} \end{bmatrix}$. Further, what is the pseudo-inverse of the general 1×2 matrix $\begin{bmatrix} a & b \end{bmatrix}$?

Exercise 3.6.9. Consider the general $m \times n$ diagonal matrix of rank r ,

$$S = \begin{bmatrix} \sigma_1 & \cdots & 0 & & & \\ \vdots & \ddots & \vdots & & O_{r \times (n-r)} & \\ 0 & \cdots & \sigma_r & & & \\ & & & & & \\ O_{(m-r) \times r} & & & O_{(m-r) \times (n-r)} & & \end{bmatrix},$$

equivalently $S = \text{diag}_{m \times n}(\sigma_1, \sigma_2, \dots, \sigma_r, 0, \dots, 0)$. Derive that, based upon Procedure 3.5.3, the pseudo-inverse of S is the $n \times m$ diagonal matrix of rank r ,

$$S^+ = \begin{bmatrix} 1/\sigma_1 & \cdots & 0 & & & \\ \vdots & \ddots & \vdots & & O_{r \times (m-r)} & \\ 0 & \cdots & 1/\sigma_r & & & \\ & & & & & \\ O_{(n-r) \times r} & & & O_{(n-r) \times (m-r)} & & \end{bmatrix},$$

equivalently $S^+ = \text{diag}_{n \times m}(1/\sigma_1, 1/\sigma_2, \dots, 1/\sigma_r, 0, \dots, 0)$.

Exercise 3.6.10. Given any $m \times n$ matrix A with SVD $A = USV^T$, use the result of Exercise 3.6.9 to establish that the pseudo-inverse $A^+ = VS^+U^T$.

Exercise 3.6.11. Use Matlab/Octave and the identity of Exercise 3.6.10 to compute the pseudo-inverse of each of the following matrices.

$$(a) \quad A = \begin{bmatrix} 0.2 & 0.3 \\ 0.2 & 1.5 \\ 0.1 & -0.3 \end{bmatrix}$$

$$(b) \quad B = \begin{bmatrix} 2.5 & 0.6 & 0.3 \\ 0.5 & 0.4 & 0.2 \end{bmatrix}$$

$$(c) \quad C = \begin{bmatrix} 0.1 & -0.5 \\ 0.6 & -3.0 \\ 0.4 & -2.0 \end{bmatrix}$$

$$(d) \quad D = \begin{bmatrix} 0.3 & -0.3 & -1.2 \\ 0.1 & 0.3 & 1.4 \\ -3.1 & -0.5 & -3.8 \\ 1.5 & -0.3 & -0.6 \end{bmatrix}$$





$$(e) E = \begin{bmatrix} 4.1 & 1.8 & -0.4 & 0.0 & -0.1 & -1.4 \\ -3.3 & -3.9 & 0.6 & -2.2 & 0.5 & 0.1 \\ -0.9 & -1.9 & 0.6 & -2.2 & 0.1 & 0.5 \\ -4.3 & -3.6 & 0.8 & -2.0 & 0.3 & 1.2 \end{bmatrix}$$



$$(f) F = \begin{bmatrix} -0.6 & -1.3 & -1.2 & -1.9 & 1.6 & 1.6 \\ -0.7 & 0.6 & -0.2 & 0.9 & -0.6 & -0.7 \\ 0.0 & 0.2 & 0.7 & 1.1 & -0.4 & -0.6 \\ 0.8 & 0.1 & -1.6 & -0.9 & -0.5 & -0.8 \\ -0.5 & 0.9 & 1.3 & 1.7 & -0.3 & -0.6 \end{bmatrix}$$



$$(g) G = \begin{bmatrix} 0.0 & -1.6 & 1.2 & -0.4 & -0.4 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & -1.5 & 0.5 & 0.5 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & -1.6 & 1.2 & -0.4 & -0.4 \end{bmatrix}$$



$$(h) H = \begin{bmatrix} -1.9 & -1.8 & -1.9 & -1.8 & -1.0 & -1.9 \\ 0.4 & 2.5 & -0.5 & 1.9 & 0.8 & 2.4 \\ -0.3 & -0.3 & 1.3 & -0.5 & -0.1 & -0.4 \\ 0.7 & 0.5 & 1.1 & 0.5 & 0.4 & 0.6 \end{bmatrix}$$

Exercise 3.6.12. Prove that in the case of an $m \times n$ matrix A with $\text{rank } A = m$ (so $m \leq n$), the pseudo-inverse is the $n \times m$ matrix $A^+ = A^T(AA^T)^{-1}$.

Exercise 3.6.13. Use Theorem 3.6.16 and the identity in Exercise 3.6.12 to prove that $(A^+)^+ = A$ in the case when $m \times n$ matrix A has $\text{rank } A = n$. (Be careful as many plausible looking steps are incorrect.)

Exercise 3.6.14. Confirm that the composition of the two linear transforms in \mathbb{R}^2 has a standard matrix that is the same as multiplying the two standard matrices of the specified linear transforms.

- (a) Rotation by 30° followed by rotation by 60° . (b) Rotation by 120° followed by rotation by -60° (clockwise 60°).

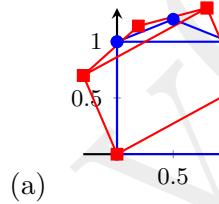
- (c) Reflection in the x -axis followed by reflection in the line $y = x$. (d) Reflection in the line $y = x$ followed by reflection in the x -axis.

- (e) Reflection in the line $y = x$ followed by rotation by 90° . (f) Reflection in the line $y = \sqrt{3}x$ followed by rotation by -30° (clockwise 30°).

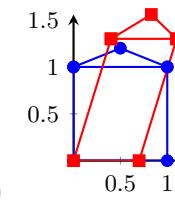
Exercise 3.6.15. For each of the following pairs of linear transformations S and T , if possible determine the standard matrices of the compositions $S \circ T$ and $T \circ S$.

- (a) $S(x) = (-5x, 2x, -x, -x)$ and $T(y_1, y_2, y_3, y_4) = -4y_1 - 3y_2 - 4y_3 + 5y_4$
- (b) $S(x_1, x_2, x_3, x_4) = (-2x_2 - 3x_3 + 6x_4, -3x_1 + 2x_2 - 4x_3 + 3x_4)$
and $T(y) = (-4y, 0, 0, -y)$
- (c) $S(x, y) = (-5x - 3y, -5x + 5y)$ and $T(z_1, z_2, z_3, z_4) = (3z_1 + 3z_2 - 2z_3 - 2z_4, 4z_1 + 7z_2 + 4z_3 + 3z_4)$
- (d) $S(x, y, z) = 5x - y + 4z$ and $T(p) = (6p, -3p, 3p, p)$
- (e) $S(u_1, u_2, u_3, u_4) = (-u_1 + 2u_2 + 2u_3 - 3u_4, -3u_1 + 3u_2 + 4u_3 + 3u_4, -u_2 - 2u_3)$ and $T(x, y, z) = (-2y - 4z, -4x + 2y + 2z, x + 3y + 4z, 2x - 2z)$
- (f) $S(p, q, r, s) = (5p - r - 2s, q - r + 2s, 7p + q + s)$ and $T(x, y) = (y, 2x + 3y, -2x + 2y, -5x - 4y)$

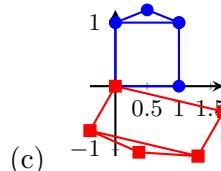
Exercise 3.6.16. For each of the illustrated transforms, estimate the standard matrix of the linear transform. Then use Theorem 3.2.6 to determine the standard matrix of its inverse transform. Hence sketch how the inverse transforms the unit square and write a sentence or two about how the sketch confirms it is a reasonable inverse.



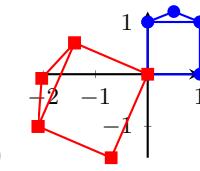
(a)



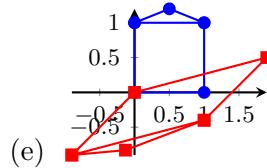
(b)



(c)



(d)



(e)

Answers to selected exercises

3.1.1b : Only B and D .

3.1.2a : A , 4×2 ; B , 1×3 ; C , 3×2 ; D , 3×4 ; E , 2×2 ; F , 2×1 .

3.1.2c : AE , AF , BC , BD , CE , CF , DA , E^2 , EF , FB .

3.1.4 : $\mathbf{b}_1 = (7.6, -1.1, 2.6, -1.5, -0.2)$, $\mathbf{b}_2 = (-1.1, -9.3, 6.9, -7.5, 5.5)$,
 $\mathbf{b}_3 = (-0.7, 0.1, 1.2, 3.7, -0.9)$, $\mathbf{b}_4 = (-4.5, 8.2, -3.6, 2.6, 2.4)$;
 $b_{13} = -0.7$, $b_{31} = 2.6$, $b_{42} = -7.5$.

3.1.6a : $A + B = \begin{bmatrix} 3 & 2 & -1 \\ 0 & -5 & -9 \\ -8 & 6 & -1 \end{bmatrix}$, $A - B = \begin{bmatrix} 1 & 0 & -1 \\ -8 & 7 & 3 \\ 4 & -2 & -1 \end{bmatrix}$

3.1.6c : $P + Q = \begin{bmatrix} -3 & 2 & 0 \\ 9 & -7 & 0 \\ 0 & 0 & -2 \end{bmatrix}$, $P - Q = \begin{bmatrix} -1 & 8 & 2 \\ -3 & 1 & 4 \\ -6 & 6 & -4 \end{bmatrix}$

3.1.7a : $-2A = \begin{bmatrix} 6 & 4 \\ -8 & 4 \\ -4 & 8 \end{bmatrix}$, $2A = \begin{bmatrix} -6 & -4 \\ 8 & -4 \\ 4 & -8 \end{bmatrix}$, $3A = \begin{bmatrix} -9 & -6 \\ 12 & -6 \\ 6 & -12 \end{bmatrix}$.

3.1.7c : $-4U = \begin{bmatrix} 15.6 & 1.2 & 11.6 \\ -12.4 & 15.6 & 4. \\ -12.4 & 26. & -3.6 \end{bmatrix}$, $2U = \begin{bmatrix} 7.8 & 0.6 & 5.8 \\ -6.2 & 7.8 & 2. \\ -6.2 & 13. & -1.8 \end{bmatrix}$,
 $4U = \begin{bmatrix} -15.6 & -1.2 & -11.6 \\ 12.4 & -15.6 & -4. \\ 12.4 & -26. & 3.6 \end{bmatrix}$.

3.1.10a : $A\mathbf{p} = \begin{bmatrix} -9 \\ -13 \end{bmatrix}$, $A\mathbf{q} = \begin{bmatrix} 4 \\ -16 \end{bmatrix}$, $A\mathbf{r} = \begin{bmatrix} -15 \\ 11 \end{bmatrix}$.

3.1.10c : $C\mathbf{u} = \begin{bmatrix} 6 \\ 3 \end{bmatrix}$, $C\mathbf{v} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$, $C\mathbf{w} = \begin{bmatrix} 24 \\ -5 \end{bmatrix}$.

3.1.11a : $A\mathbf{u} = (7, -5)$, $A\mathbf{v} = (-6, 3)$, $A\mathbf{w} = (9, -6)$.

3.1.11c : $C\mathbf{x}_1 = (-4.41, 9.66)$, $C\mathbf{x}_2 = (1.42, -1.56)$, $C\mathbf{x}_3 = (-0.47, -0.38)$.

3.1.12a : $P\mathbf{u} = (1, 1.4)$, $P\mathbf{v} = (-3.6, 1.7)$, $P\mathbf{w} = (0.1, -2.3)$. Reflection in the horizontal axis.

3.1.12c : $R\mathbf{x}_1 = (-4.4, -0.8)$, $R\mathbf{x}_2 = (5, 0)$, $R\mathbf{x}_3 = (-0.2, 3.6)$. Rotation (by 36.87°).

3.1.13a : $\begin{bmatrix} -9 & 4 \\ -13 & -16 \end{bmatrix}$, $\begin{bmatrix} -9 & -15 \\ -13 & 11 \end{bmatrix}$, $\begin{bmatrix} 4 & -15 \\ -16 & 11 \end{bmatrix}$, $\begin{bmatrix} -9 & 4 & -15 \\ -13 & -16 & 11 \end{bmatrix}$.

3.1.13c : $\begin{bmatrix} 6 & 3 \\ 3 & 4 \end{bmatrix}$, $\begin{bmatrix} 24 & 3 \\ -5 & 4 \end{bmatrix}$, $\begin{bmatrix} 24 & 6 \\ -5 & 3 \end{bmatrix}$, $\begin{bmatrix} 24 & 3 & 6 \\ -5 & 4 & 3 \end{bmatrix}$.

3.1.17 : Let x_j be the number of females of age $(j - 1)$ years. $L =$

$\begin{bmatrix} 0 & 0 & 2 & 2 \\ \frac{2}{3} & 0 & 0 & 0 \\ 0 & \frac{2}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \end{bmatrix}$. After one year $\mathbf{x}' = (72, 12, 6, 9)$, two years $\mathbf{x}'' = (30, 48, 8, 3)$, three years $\mathbf{x}''' = (22, 20, 32, 4)$. Increasing.

3.1.18b : $\begin{bmatrix} 3 & -5 \\ -4 & 2 \\ -2 & -3 \\ 2 & 3 \end{bmatrix}$

3.1.18d : $\begin{bmatrix} 3 \\ 1 \\ -2 \\ -3 \end{bmatrix}$

3.1.18f : $\begin{bmatrix} -4 & -5.1 & 0.3 \\ -5.1 & -7.4 & -3 \\ 0.3 & -3 & 2.6 \end{bmatrix}$, symmetric

3.1.18h : $\begin{bmatrix} 1.7 & 0.7 & 0.6 \\ -0.2 & -0.3 & 3 \\ -0.4 & -0.4 & -2.2 \end{bmatrix}$

3.2.1a : Inverse.

3.2.1c : Not inverse.

3.2.1e : Inverse.

3.2.1g : Not inverse.

3.2.1i : Inverse.

3.2.1k : Inverse.

3.2.2b : No inverse.

3.2.2d : $\begin{bmatrix} -1/4 & -1/4 \\ 1/8 & -3/8 \end{bmatrix}$

3.2.2f : $\begin{bmatrix} -0.8974 & 0.5769 \\ -0.5128 & -0.3846 \end{bmatrix}$

3.2.2h : $\begin{bmatrix} 10/9 & -10/9 \\ 2/3 & 4/3 \end{bmatrix}$

3.2.3b : $(p, q) = (1, 2)$

3.2.3d : $(x, y, z) = (2, 13/4, -9)$

3.2.3f : $\mathbf{x} = (-1, 1, 1, 4)$

3.2.3h : $(a, d, c, b) = (3, -7/3, 13/3, -8/3)$

3.2.4b : Not enough information

3.2.4d : Not invertible.

3.2.4f : Not enough information.

3.2.7b : $\begin{bmatrix} -1/64 & 0 \\ 0 & -1/64 \end{bmatrix}$

3.2.7d : $\begin{bmatrix} -4 & 16 \\ -80 & -4 \end{bmatrix}$

3.2.7f : $\begin{bmatrix} 10 & -6 & -6 \\ -2 & 6 & 0 \\ -2 & 2 & 4 \end{bmatrix}$

3.2.7h : $\begin{bmatrix} 25 & -96 & -46 & 32 \\ -6 & 25 & 11 & -8 \\ 0 & -40 & -15 & 16 \\ 18 & -76 & -41 & 29 \end{bmatrix}$

3.2.8b : Not diagonal.

3.2.8d : Not diagonal.

3.2.8f : Not diagonal.

3.2.8h : $\text{diag}_{2 \times 3}(0, 2)$

3.2.8j : Diagonal only when $a = d = 0$.

3.2.9b : No solution.

3.2.9d : $\mathbf{x} = (0, -1/4, 1, s, t)$ for all s, t

3.2.9f : $(w, x, y, z) = (6, -5, s, t)$ for all s, t

3.2.9h : No solution.

3.2.10b : $\text{diag}(1, 1.5)$

3.2.10d : $\text{diag}(-0.3, 0.3)$

3.2.10f : Not diagonal.

3.2.10h : Not diagonal.

3.2.10j : Not diagonal.

3.2.11b : Not diagonal.

3.2.11d : $\text{diag}(1.2, 0.3, 0.6)$

3.2.11f : $\text{diag}(0.6, 1.2, 0.6)$

3.2.13b : Orthonormal.

3.2.13d : Orthonormal.

3.2.13f : Not orthogonal.

3.2.13h : Orthogonal.

3.2.14b : Orthogonal.

3.2.14d : Orthonormal.

- 3.2.14f : Not orthogonal.
- 3.2.15b : Orthogonal set, divide each by nine.
- 3.2.15d : Not orthogonal set.
- 3.2.15f : Orthogonal set, divide each by five.
- 3.2.15h : Not orthogonal set.
- 3.2.16b : Orthogonal matrix.
- 3.2.16d : Orthogonal matrix.
- 3.2.16f : Orthogonal matrix.
- 3.2.16h : Orthogonal matrix.
- 3.2.16j : Not orthogonal matrix.
- 3.2.17a : $\theta = 21.04^\circ$
- 3.2.17c : $\theta = 33.69^\circ$
- 3.2.17e : $\theta = 60^\circ$
- 3.2.17g : $\theta = 42.79^\circ$
- 3.2.18a : $(x, y) = (1.4, 5.2)$
- 3.2.18c : $(x, y) = (5/\sqrt{2}, 1/\sqrt{2})$
- 3.2.18e : $(p, q, r) = (6, -1, -5)$
- 3.2.18g : $(a, b, c) = (8, 32, -1)/11$
- 3.2.18i : $\mathbf{y} = (0, -3.3, -0.6, 2.5)$
- 3.2.22a : No—the square is deformed.
- 3.2.22c : No—the square is squashed (and rotated/reflected).
- 3.2.22e : No—the square is stretched.
- 3.2.22g : Yes—the square is rotated and reflected.
- 3.2.23a : Yes—the cube appears rotated.
- 3.2.23c : Yes—the cube appears rotated.
- 3.2.23e : Yes—the cube appears rotated.
- 3.2.23g : No—the cube is deformed.
- 3.3.2b : $\mathbf{x} = (0, \frac{3}{2})$
- 3.3.2d : $\mathbf{x} = (-2, 2)$
- 3.3.2f : $\mathbf{x} = (\frac{5}{14}, \frac{15}{14}, \frac{15}{28}) + (-\frac{6}{7}, \frac{3}{7}, -\frac{2}{7})s + (\frac{3}{7}, \frac{2}{7}, -\frac{6}{7})t$
- 3.3.2h : No solution.
- 3.3.2j : No solution.
- 3.3.2l : $\mathbf{x} = (2, -\frac{7}{4}, -\frac{9}{2})$

3.3.3b : No solution.

3.3.3d : $\mathbf{x} = (27, -12, -24, 9) + (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})t$

3.3.3f : $\mathbf{x} = (-\frac{33}{2}, \frac{25}{2}, \frac{17}{2}, \frac{9}{2})$

3.3.4b : $\mathbf{x} = (-1, -1, 3, -3)$

3.3.4d : $\mathbf{x} = (-4, -9, 0, 7)$

3.3.4f : $\mathbf{x} = (0.6, 8.2, 9.8, -0.6) + (-0.1, 0.3, -0.3, -0.9)t$

3.3.4h : $\mathbf{x} = (0, -0.3, -3.1, 4.9)$

3.3.4j : $\mathbf{x} = (0.18, 3.35, -4.86, 0.33, 0.35) + (0.91, -0.34, -0.21, -0.09, -0.07)t$ (2 d.p.)

3.3.6 : 1. cond = 5/3, rank = 2;

2. cond = 1, rank = 2;

3. cond = ∞ , rank = 1;

4. cond = 2, rank = 2;

5. cond = 2, rank = 2;

6. cond = ∞ , rank = 1;

7. cond = 2, rank = 2;

8. cond = 2, rank = 2;

9. cond = 2, rank = 3;

10. cond = ∞ , rank = 2;

11. cond = ∞ , rank = 1;

12. cond = 9, rank = 3;

3.3.10 : The theorem applies to the square matrix systems of Exercise 3.3.2 (a)–(d), (i)–(l), and of Exercise 3.3.3 (e) and (f). The cases with no zero singular value, full rank, have a unique solution. The cases with a zero singular value, rank less than n , either have no solution or an infinite number.

3.3.15b : $\mathbf{v}_1 \approx (0.1, 1.0)$, $\sigma_1 \approx 1.7$, $\mathbf{v}_2 \approx (1.0, -0.1)$, $\sigma_1 \approx 0.3$.

3.3.15d : $\mathbf{v}_1 \approx (0.9, -0.4)$, $\sigma_1 \approx 2.3$, $\mathbf{v}_2 \approx (0.4, 0.9)$, $\sigma_1 \approx 0.3$.

3.4.1b : Not a subspace.

3.4.1d : Not a subspace.

3.4.1f : Not a subspace.

3.4.1h : Subspace.

3.4.1j : Not a subspace.

3.4.11 : Subspace.

- 3.4.1n : Not a subspace.
- 3.4.1p : Not a subspace.
- 3.4.2b : Subspace.
- 3.4.2d : Not a subspace.
- 3.4.2f : Subspace.
- 3.4.2h : Subspace.
- 3.4.2j : Not a subspace.
- 3.4.2l : Not a subspace.
- 3.4.2n : Subspace.
- 3.4.2p : Subspace.
- 3.4.4a : \mathbf{b}_5 is in column space; \mathbf{r}_5 is in row space.
- 3.4.4c : \mathbf{b}_5 is not in column space; \mathbf{r}_5 is in row space.
- 3.4.4e : \mathbf{b}_5 is in column space; \mathbf{r}_5 is not in row space.
- 3.4.4g : \mathbf{b}_5 is in column space; \mathbf{r}_5 is in row space.
- 3.4.4i : \mathbf{b}_5 is in column space; \mathbf{r}_5 is in row space.
- 3.4.5b : no
- 3.4.5d : yes
- 3.4.5f : yes
- 3.4.6b : $(\frac{4}{5}, \frac{3}{5})$
- 3.4.6d : $(2, 2, 1)/3$
- 3.4.8a : (2 d.p.) column space $\{(-0.61, 0.19, 0.77), (-0.77, -0.36, -0.52)\}$; row space $\{(-0.35, 0.84, 0.42), (-0.94, -0.31, -0.15)\}$; nullspace $\{(0, 0.45, -0.89)\}$.
- 3.4.8c : (2 d.p.) column space $\{(0.00, -0.81, -0.48, 0.33), (0.00, -0.08, 0.66, 0.74)\}$; row space $\{(-0.35, -0.89, -0.08, 0.27), (-0.48, 0.17, -0.80, -0.32)\}$; nullspace $\{(-0.78, 0.36, 0.43, 0.28), (-0.19, -0.22, 0.41, -0.86)\}$.
- 3.4.8e : (2 d.p.) column space $\{(-0.53, -0.11, -0.64, 0.01, 0.54), (0.40, -0.37, 0.03, 0.76, 0.35)\}$; row space $\{(0.01, -0.34, -0.30, -0.89), (0.21, -0.92, 0.11, 0.32)\}$; nullspace $\{(-0.06, -0.01, 0.95, -0.31), (0.98, 0.20, 0.04, -0.08)\}$.
- 3.4.8g : (2 d.p.) column space $\{(0.52, -0.14, -0.79, 0.28), (-0.02, 0.18, -0.37, -0.91), (-0.31, -0.94, -0.09, -0.14)\}$; row space $\{(-0.16, 0.29, 0.13, -0.87, 0.33), (-0.15, 0.67, 0.58, 0.40, 0.17), (-0.75, -0.12, 0.19, -0.11, -0.61)\}$; nullspace $\{(0.62, 0.05, 0.45, -0.25, -0.59), (-0.05, -0.67, 0.63, 0.02, 0.38)\}$.
- 3.4.9 : 1, 1, 1; 3, 3, 0; 2, 2, 2; 3, 3, 2; 2, 2, 2; 2, 2, 2; 3, 3, 2; 4, 4, 1.

3.4.10b : 0,1,2,3

3.4.10d : 2,3,4,5,6

3.4.10f : 0,1,2,3,4,5

3.4.11b : yes.

3.4.11d : no.

3.5.4 : To within an arbitrary constant: Adelaide, -0.33; Brisbane, -1.00; Canberra, 1.33.

3.5.6 : To within an arbitrary constant: Atlanta, 1.0; Boston, 0.0; Concord, 0.8; Denver, -1.6; Frankfort, -0.2.

3.5.13a : $\mathbf{x} = (5, -12)$

3.5.13c : $(x, y) = (\frac{1}{2}, -\frac{1}{2})$

3.5.13e : $\mathbf{x} = (\frac{2}{7}, -\frac{3}{7}, \frac{6}{7})$

3.5.13g : $\mathbf{u} = (-0.1053, 0.2632, 0.1579)$

3.5.14 : The middle bottom is most absorbing.

3.5.16 : Pixel ten is the most absorbing, $r_{10} \approx 0.25$.

3.5.19b : (2 d.p.) $\text{proj}_{\mathbf{u}}(\mathbf{v}) = (-1.18, 0.29)$, $x = -0.29$

3.5.19d : $\text{proj}_{\mathbf{u}}(\mathbf{v}) = (-1.5, 1.5)$, $x = -0.75$

3.5.19f : $\text{proj}_{\mathbf{u}}(\mathbf{v}) = \mathbf{0}$, $x = 0$

3.5.19h : (2 d.p.) $\text{proj}_{\mathbf{u}}(\mathbf{v}) = (0.17, 1.21, -0.87)$, $x = -0.17$

3.5.19j : (2 d.p.) $\text{proj}_{\mathbf{u}}(\mathbf{v}) = (-0.57, 1.14, 0.57, 0.38)$, $x = -0.19$

3.5.19l : (2 d.p.) $\text{proj}_{\mathbf{u}}(\mathbf{v}) = (-0.91, 0.91, -0.45, 1.36, 0.91)$, $x = 0.45$

3.5.20b : (2 d.p.) $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (1.68, -3.16, -0.79)$

3.5.20d : (2 d.p.) $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (-1.21, 1.21, -1.38)$

3.5.20f : (2 d.p.) $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (0.02, -2.24, 0.20, -2.71)$

3.5.20h : (2 d.p.) $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (0.78, 0.44, 1.44, 0)$

3.5.20j : $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (0.5, 1.5, 0.5, 1.5)$

3.5.20l : (2 d.p.) $\text{proj}_{\mathbb{W}}(\mathbf{v}) = (-2.06, -4.01, 2.94, -1.00)$

3.5.21b : (2 d.p.)
$$\begin{bmatrix} 0.57 & 0.40 & 0.29 \\ 0.40 & 0.63 & -0.27 \\ 0.29 & -0.27 & 0.80 \end{bmatrix}$$

3.5.21d : (2 d.p.)
$$\begin{bmatrix} 0.07 & 0.21 & 0.14 \\ 0.21 & 0.64 & 0.43 \\ 0.14 & 0.43 & 0.29 \end{bmatrix}$$

3.5.21f : (2 d.p.)
$$\begin{bmatrix} 0.37 & -0.30 & 0.07 & -0.37 \\ -0.30 & 0.24 & -0.06 & 0.30 \\ 0.07 & -0.06 & 0.01 & -0.07 \\ -0.37 & 0.30 & -0.07 & 0.37 \end{bmatrix}$$

3.5.21h : (2 d.p.)
$$\begin{bmatrix} 0.39 & -0.49 \\ -0.49 & 0.61 \end{bmatrix}$$

3.5.21j : (2 d.p.)
$$\begin{bmatrix} 0.99 & 0.06 & -0.03 & -0.06 & -0.05 \\ 0.06 & 0.54 & 0.36 & 0.13 & 0.32 \\ -0.03 & 0.36 & 0.52 & 0.29 & -0.19 \\ -0.06 & 0.13 & 0.29 & 0.18 & -0.20 \\ -0.05 & 0.32 & -0.19 & -0.20 & 0.76 \end{bmatrix}$$

3.5.22b : (2 d.p.) $\mathbf{b}' = (2.08, -7.95, -1.21)$, difference 0.23

3.5.22d : (2 d.p.) $\mathbf{b}' = (-0.65, 1.68, -3.24)$, difference 0.53

3.5.22f : (2 d.p.) $\mathbf{b}' = (8.21, 25.40, -14.96)$, difference 3.71

3.5.22h : (2 d.p.) $\mathbf{b}' = (6, 3)$, difference 0

3.5.22j : (2 d.p.) $\mathbf{b}' = (-6.38, 5.00, -1.25, 6.38)$, difference 1.90

3.5.22l : (2 d.p.) $\mathbf{b}' = (-12.77, -6.03, 18.22, -35.92)$, difference 2.91

3.5.26b : The line $y = 5x$

3.5.26d : The line $\text{span}\{(-4, 4, 5)\}$

3.5.26f : The line $\text{span}\{(8, 11, 5)\}$

3.5.26h : The hyper-plane $6x_1 + 5x_2 + x_3 - 3x_4 = 0$

3.5.27b : $\{\left(\frac{1}{9}, \frac{4}{9}, \frac{8}{9}\right)\}$ is one possibility.

3.5.27d : $\{(-0.41, 0.82, 0.41)\}$ (2 d.p.).

3.5.27f : $\{(-0.12, -0.98, -0.17, 0.02), (-0.20, -0.11, 0.87, 0.43)\}$ is one possibility (2 d.p.).

3.5.27h : $\{(0.89, 0.20, 0.02, 0.02, 0.41), (-0.23, 0.68, -0.49, 0.45, 0.17)\}$ is one possibility (2 d.p.).

3.5.27j : $\{(0.45 - 0.220.830.25), (-0.82 - 0.200.260.47)\}$ is one possibility (2 d.p.).

3.5.29b : $\frac{1}{9}(8, 1, -4)$

3.5.29d : $\frac{1}{27}(-58, -23, 20)$

3.5.29f : $\frac{1}{81}(65p + 16q - 28r, 16p + 65q + 28r, -28p + 28q + 32r)$

3.5.30b : $(-3.48, 2.06, 1.38, -1.96)$

3.5.30d : $\frac{1}{100}(74p - 40r + 18s, 74q - 18r - 40s, -40p - 18q + 26r, 18p - 40q + 26s)$

3.5.32b : $(-3, 3) = (\frac{9}{25}, \frac{12}{25}) + (-\frac{84}{25}, \frac{63}{25})$

3.5.32d : $(3, 1) = \left(\frac{39}{25}, \frac{52}{25}\right) + \left(\frac{36}{25}, -\frac{27}{25}\right)$

3.5.33b : $(-1.96, 3.92, -1.31) + (1.96, 1.08, 0.31)$ (2 d.p.)

3.5.33d : $(-1.04, 2.08, -0.69) + (-1.96, -1.08, -0.31)$ (2 d.p.)

3.5.34b : $(-3, 1, 6, 2) + (-1, -3, -1, 3)$

3.5.34d : $(3.3, -1.1, 0.9, 0.3) + (1.7, 5.1, -0.9, 2.7)$

3.5.36 : Use xyz -space. Either \mathbb{W} is x -axis, or xy -plane, and \mathbb{W}^\perp corresponding complement, or vice-versa.

3.6.1a : Maybe a LT, with standard matrix $\begin{bmatrix} 0.3 & 0.3 \\ -0.5 & 1.1 \end{bmatrix}$

3.6.1c : Not a LT.

3.6.1e : Not a LT.

3.6.1g : Maybe a LT, with standard matrix $\begin{bmatrix} -0.1 & 0.6 \\ 3.3 & -0.6 \end{bmatrix}$

3.6.1i : Not a LT.

3.6.3a : Maybe a LT.

3.6.3c : Maybe a LT.

3.6.3e : Not a LT.

3.6.7a : $[E] = \begin{bmatrix} -3 & 2 & 0 \\ 0 & 0 & -1 \end{bmatrix}$

3.6.7c : $[E] = \begin{bmatrix} 2 & 1 & -2 & 0 & 0 \\ 7 & 0 & 0 & 7 & 0 \end{bmatrix}$

3.6.7e : $[E] = \begin{bmatrix} -3 & 0 & -4 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 6 \\ 5 & 6 & 0 & -1 \end{bmatrix}$

3.6.7g : Not a LT.

3.6.11a : $E^+ = \begin{bmatrix} 3.52 & -0.08 & 3.11 \\ -0.36 & 0.63 & -0.55 \end{bmatrix}$ (2 d.p.)

3.6.11c : $E^+ = \begin{bmatrix} 0.01 & 0.04 & 0.03 \\ -0.04 & -0.22 & -0.15 \end{bmatrix}$ (2 d.p.)

3.6.11e : $E^+ = \begin{bmatrix} 0.17 & -0.03 & 0.16 & -0.04 \\ -0.12 & -0.33 & 0.16 & 0.05 \\ -0.02 & -0.15 & 0.18 & 0.06 \\ -0.10 & 0.18 & -0.43 & -0.10 \\ 0.04 & 0.15 & -0.13 & -0.04 \\ -0.21 & -0.58 & 0.49 & 0.19 \end{bmatrix}$ (2 d.p.)

$$3.6.11g : E^+ = \begin{bmatrix} 0.00 & 0.00 & 0.00 & 0.00 & 0.00 \\ -0.10 & 0.00 & 0.13 & 0.00 & -0.10 \\ 0.08 & 0.00 & -0.10 & 0.00 & 0.08 \\ -0.03 & 0.00 & 0.03 & 0.00 & -0.03 \\ -0.03 & 0.00 & 0.03 & 0.00 & -0.03 \end{bmatrix} \quad (2 \text{ d.p.})$$

3.6.14a : Equivalent to rotation by 90° with matrix $\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$.

3.6.14c : Equivalent to rotation by 90° with matrix $\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$.

3.6.14e : Equivalent to reflection in the y -axis with matrix $\begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$.

$$3.6.15a : [S \circ T] = \begin{bmatrix} 20 & 15 & 20 & -25 \\ -8 & -6 & -8 & 10 \\ 4 & 3 & 4 & -5 \\ 4 & 3 & 4 & -5 \end{bmatrix} \text{ and } [T \circ S] = [13]$$

3.6.15c : $[S \circ T] = \begin{bmatrix} -27 & -36 & -2 & 1 \\ 5 & 20 & 30 & 25 \end{bmatrix}$ and $[T \circ S] =$ does not exist.

$$3.6.15e : [S \circ T] = \begin{bmatrix} -12 & 12 & 22 \\ -2 & 24 & 28 \\ 2 & -8 & -10 \end{bmatrix}, [T \circ S] = \begin{bmatrix} 6 & -2 & 0 & -6 \\ -2 & -4 & -4 & 18 \\ -10 & 7 & 6 & 6 \\ -2 & 6 & 8 & -6 \end{bmatrix}$$

3.6.16a : Inverse $\approx \begin{bmatrix} 0.74 & 0.32 \\ -0.63 & 1.16 \end{bmatrix}$

3.6.16c : Inverse $\approx \begin{bmatrix} 0.52 & -0.30 \\ -0.30 & -1.26 \end{bmatrix}$

3.6.16e : Inverse $\approx \begin{bmatrix} 0.71 & -0.71 \\ 0.40 & -1.51 \end{bmatrix}$

4 Eigenvalues and eigenvectors of symmetric matrices

Chapter Contents

4.1	Introduction to eigenvalues and eigenvectors	377
4.1.1	Systematically find eigenvalues and eigenvectors	384
4.1.2	Exercises	396
4.2	Beautiful properties for symmetric matrices	403
4.2.1	Matrix powers maintain eigenvectors	403
4.2.2	Symmetric matrices are orthogonally diagonalisable	408
4.2.3	Change orthonormal basis to classify quadratics	418
4.2.4	Exercises	425

Recall (Section 3.1.2) that a symmetric matrix A is a square matrix such that $A^T = A$, that is, $a_{ij} = a_{ji}$. For example, of the following two matrices, the first is symmetric, but the second is not:

$$\begin{bmatrix} -2 & 4 & 0 \\ 4 & 2 & -3 \\ 0 & -3 & 1 \end{bmatrix}; \quad \begin{bmatrix} -1 & 3 & 0 \\ 1 & 1 & 0 \\ 0 & -3 & 1 \end{bmatrix}.$$

Example 4.0.1. Compute some SVDS of random symmetric matrices, observe the columns of U are always \pm the columns of V (well, almost always).

Solution: Repeat as often as you like for any size of square matrix that you like (one example is recorded here to two decimal places).

- (a) Generate in Matlab/Octave some random symmetric matrix by adding a random matrix to its transpose with `A=randn(5); A=A+A'` (Table 3.1):



```
A =
-0.45 -0.18  1.59 -0.96 -0.54
-0.18 -0.24 -1.04  0.14  0.80
 1.59 -1.04 -2.87 -0.40  1.11
-0.96  0.14 -0.40 -0.26 -1.90
-0.54  0.80  1.11 -1.90  1.64
```

This matrix is symmetric as $a_{ij} = a_{ji}$.

- (b) Find its SVD via `[U,S,V]=svd(A)`

$$\begin{aligned} \mathbf{U} &= \\ &\begin{array}{ccccc} -0.41 & -0.09 & -0.28 & -0.67 & 0.55 \\ 0.25 & -0.11 & -0.05 & 0.53 & 0.80 \\ 0.82 & -0.19 & -0.40 & -0.36 & -0.07 \\ -0.15 & 0.51 & -0.80 & 0.27 & -0.11 \\ -0.27 & -0.83 & -0.36 & 0.25 & -0.22 \end{array} \\ \mathbf{S} &= \\ &\begin{array}{ccccc} 4.28 & 0 & 0 & 0 & 0 \\ 0 & 3.12 & 0 & 0 & 0 \\ 0 & 0 & 1.65 & 0 & 0 \\ 0 & 0 & 0 & 1.14 & 0 \\ 0 & 0 & 0 & 0 & 0.51 \end{array} \\ \mathbf{V} &= \\ &\begin{array}{ccccc} 0.41 & -0.09 & 0.28 & -0.67 & -0.55 \\ -0.25 & -0.11 & 0.05 & 0.53 & -0.80 \\ -0.82 & -0.19 & 0.40 & -0.36 & 0.07 \\ 0.15 & 0.51 & 0.80 & 0.27 & 0.11 \\ 0.27 & -0.83 & 0.36 & 0.25 & 0.22 \end{array} \end{aligned}$$

Observe the second and fourth columns of U and V are identical, and the other columns have opposite signs.

Repeat and observe $\mathbf{u}_j = \pm \mathbf{v}_j$ for all columns j .

■

Why, for symmetric matrices, are the columns of U (almost) always \pm the columns of V ? The answer is connected to the following rearrangement of an SVD. Because $A = USV^T$, post-multiplying by V gives $AV = USV^TV = US$, and then the j th column of $AV = US$ is $A\mathbf{v}_j = \sigma_j \mathbf{u}_j$. Example 4.0.1 observes for symmetric A that $\mathbf{u}_j = \pm \mathbf{v}_j$ (almost always) so this last equation becomes $A\mathbf{v}_j = (\pm \sigma_j)\mathbf{v}_j$. This equation is of the important form $A\mathbf{v} = \lambda\mathbf{v}$. This form is important because it is the mathematical expression of the following geometric question: for what vectors \mathbf{v} does multiplication by A just stretch/shrink \mathbf{v} by some scalar λ ?

Solid modelling Lean with a hand on a table/wall: the force changes depending upon the orientation of the surface. Similarly inside any solid: the internal forces = $A\mathbf{n}$ where \mathbf{n} is the normal unit vector to the internal ‘surface’. Matrix A is always symmetric. To know whether a material will break apart under pulling, or to crumble under compression, we need to know where the extreme forces are. They are found as solutions to $A\mathbf{n} = \lambda\mathbf{n}$ where \mathbf{n} gives the direction and λ the strength of the force.

4.1 Introduction to eigenvalues and eigenvectors

Section Contents

4.1.1	Systematically find eigenvalues and eigenvectors	384
	Compute eigenvalues and eigenvectors	384
	Find eigenvalues and eigenvectors by hand . .	391
4.1.2	Exercises	396

This chapter focuses on some marvellous properties of symmetric matrices. Nonetheless it defines basic concepts in general. Chapter 7 explores analogous properties for general matrices. The marvellously useful properties result from asking for which vectors does multiplication by a given matrix simply stretch or shrink the vector.

Definition 4.1.1. *Let A be a square matrix. A scalar λ is called an **eigenvalue** of A if there is a nonzero vector \mathbf{x} , called an **eigenvector** corresponding to λ , such that $A\mathbf{x} = \lambda\mathbf{x}$.*

Example 4.1.2. Consider the symmetric matrix

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

- (a) Verify an eigenvector is $(1, 0, -1)$. What is the corresponding eigenvalue?
- (b) Verify that $(2, -4, 2)$ is an eigenvector. What is its corresponding eigenvalue?
- (c) Verify that $(1, 2, 1)$ is not an eigenvector.
- (d) Use inspection to guess and verify another eigenvector (not proportional to either of the above two). What is its eigenvalue?

Solution: The simplest approach is to multiply the matrix by the given vector and see what happens.

- (a) For vector $\mathbf{x} = (1, 0, -1)$,

$$A\mathbf{x} = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} = 1 \cdot \mathbf{x}.$$

Hence $(1, 0, -1)$ is an eigenvector of A corresponding to the eigenvalue $\lambda = 1$.

- (b) For vector $\mathbf{x} = (2, -4, 2)$,

$$A\mathbf{x} = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -4 \\ 2 \end{bmatrix} = \begin{bmatrix} 6 \\ -12 \\ 6 \end{bmatrix} = 3 \cdot \mathbf{x}.$$

Hence $(2, -4, 2)$ is an eigenvector of A corresponding to the eigenvalue $\lambda = 3$.

(c) For vector $\mathbf{x} = (1, 2, 1)$,

$$A\mathbf{x} = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \\ -1 \end{bmatrix} \neq \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}.$$

If there was a constant of proportionality (an eigenvalue), then the first component would require the constant $\lambda = -1$ but the second component would require $\lambda = +1$ which is a contradiction. Hence $(1, 2, 1)$ is not an eigenvector of A .

(d) Inspection is useful if it is quick: here one might see that the elements in each row of A sum to the same thing, namely zero, so try vector $\mathbf{x} = (1, 1, 1)$:

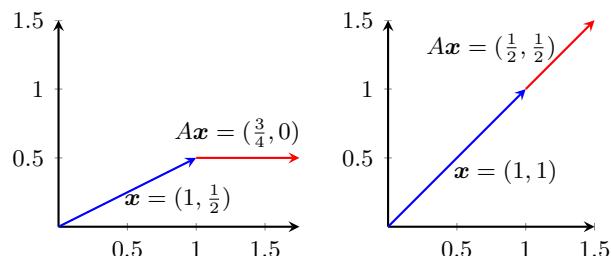
$$A\mathbf{x} = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} = 0 \cdot \mathbf{x}.$$

Hence $(1, 1, 1)$ is an eigenvector of A corresponding to the eigenvalue $\lambda = 0$. ■

Importantly, eigenvectors tell us key directions of a given matrix: the directions in which the multiplication by a matrix is to simply stretch, shrink, or reverse by a factor: the factor being the corresponding eigenvalue. In two dimensional plots we can graphically estimate eigenvectors and eigenvalues.

We might plot a given vector \mathbf{x} and join onto its head the vector $A\mathbf{x}$: if both \mathbf{x} and $A\mathbf{x}$ are aligned in the same direction (or opposite direction), then \mathbf{x} is an eigenvector; if they form a non-trivial angle, then \mathbf{x} is not an eigenvector.

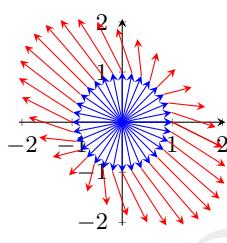
Example 4.1.3. Let the matrix $A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}$. The plot below-left shows the vector $\mathbf{x} = (1, \frac{1}{2})$, and adjoined to its head the matrix-vector product $A\mathbf{x} = (\frac{3}{4}, 0)$: because the two are at an angle, $(1, \frac{1}{2})$ is not an eigenvector.



However, as plotted above-right, for the vector $\mathbf{x} = (1, 1)$ the matrix-vector product $A\mathbf{x} = (\frac{1}{2}, \frac{1}{2})$ and the plot of these vectors head-to-tail illustrates that they are aligned in the same direction. Because of the alignment, $(1, 1)$ is an eigenvector of this matrix. The constant of proportionality is the corresponding eigenvalue: here $A\mathbf{x} = (\frac{1}{2}, \frac{1}{2}) = \frac{1}{2}(1, 1) = \frac{1}{2}\mathbf{x}$ so the eigenvalue is $\lambda = \frac{1}{2}$. ■

As in the next example, we can plot for many directions \mathbf{x} a diagram of vector $A\mathbf{x}$ adjoined head-to-tail to vector \mathbf{x} . Then inspection estimates the eigenvectors and corresponding eigenvalues (Schonefeld 1995).

Example 4.1.4 (graphical eigenvectors one).



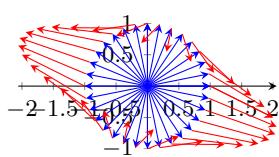
The plot on the left shows (unit) vectors \mathbf{x} (blue), and for some matrix A the corresponding vectors $A\mathbf{x}$ (red) adjoined. Estimate which directions \mathbf{x} are eigenvectors, and for each eigenvector estimate the corresponding eigenvalue.

Solution: We seek vectors \mathbf{x} such that $A\mathbf{x}$ is in the same direction (to graphical accuracy). It appears that vectors at 45° to the axes are the only ones for which $A\mathbf{x}$ is in the same direction as \mathbf{x} :

- the two (blue) vectors $\pm(0.7, 0.7)$ appear to be shrunk to length a half (red) so we estimate two eigenvectors are $\mathbf{x} \approx \pm(0.7, 0.7)$ and the corresponding eigenvalue $\lambda \approx 0.5$;
- the two (blue) vectors $\pm(0.7, -0.7)$ appear to be stretched by a factor about 1.5 (red) so we estimate two eigenvectors are $\mathbf{x} \approx \pm(0.7, -0.7)$ and the corresponding eigenvalue is $\lambda \approx 1.5$;
- and for no other (unit) vector \mathbf{x} is $A\mathbf{x}$ aligned with \mathbf{x} .

Any multiple of these eigenvectors will also be eigenvectors so we may report the directions more simply, perhaps $(1, 1)$ and $(1, -1)$ respectively. ■

Example 4.1.5 (graphical eigenvectors two).



The plot on the left shows (unit) vectors \mathbf{x} (blue), and for some matrix A the corresponding vectors $A\mathbf{x}$ (red) adjoined. Estimate which directions \mathbf{x} are eigenvectors, and for each eigenvector estimate the corresponding eigenvalue.

Solution: We seek vectors \mathbf{x} such that $A\mathbf{x}$ is in the same direction (to graphical accuracy):

- the two (blue) vectors $\pm(0.9, -0.3)$ appear stretched a little by a factor about 1.2 (red) so we estimate eigenvectors are $\mathbf{x} \propto (0.9, -0.3)$ and the corresponding eigenvalue is $\lambda \approx 1.2$;
- the two (blue) vectors $\pm(0.3, 0.9)$ appear shrunk and *reversed* by a factor about 0.4 (red) so we estimate eigenvectors are $\mathbf{x} \propto (0.3, 0.9)$ and the corresponding eigenvalue is $\lambda \approx -0.4$ —negative because the direction is reversed;
- and for no other (unit) vector \mathbf{x} is $A\mathbf{x}$ aligned with \mathbf{x} .

If this matrix arose in the description of forces inside a solid, then the forces would be compressive in directions $\pm(0.3, 0.9)$, and the forces would be (tension) ‘ripping apart’ the solid in directions $\pm(0.9, -0.3)$. ■

Example 4.1.6 (diagonal matrix). The eigenvalues of a (square) diagonal matrix are the entries on the diagonal. Consider an $n \times n$ diagonal matrix

$$D = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix}.$$

Multiply by the standard unit vectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ in turn:

$$D\mathbf{e}_1 = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} d_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = d_1 \mathbf{e}_1;$$

$$D\mathbf{e}_2 = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ d_2 \\ \vdots \\ 0 \end{bmatrix} = d_2 \mathbf{e}_2;$$

$$\vdots$$

$$D\mathbf{e}_n = \begin{bmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ d_n \end{bmatrix} = d_n \mathbf{e}_n.$$

By Definition 4.1.1, each diagonal element d_j is an eigenvalue of the diagonal matrix, and the standard unit vector \mathbf{e}_j is a corresponding eigenvector. ■

Eigenvalues The 3×3 matrix of Example 4.1.2 has three eigenvalues. The 2×2 matrices underlying Examples 4.1.4–4.1.5 both have two eigenvalues. An $n \times n$ diagonal matrix has n eigenvalues. The next section establishes the general pattern that an $n \times n$ *symmetric matrix* generally has n real eigenvalues. However, the eigenvalues of non-symmetric matrices are more complex (in both senses of the word) as explored by Chapter 7.

Eigenvectors It is the direction of eigenvectors that is important. In Example 4.1.2 any nonzero multiple of $(1, -2, 1)$, positive or negative, is also an eigenvector corresponding to eigenvalue $\lambda = 3$. In the diagonal matrices of Example 4.1.6, a straightforward extension of the working shows any nonzero multiple of the standard unit vector e_j is an eigenvector corresponding to the eigenvalue d_j . Let's collect all possible eigenvectors into a subspace.

Theorem 4.1.7. *Let A be a square matrix. A scalar λ is an eigenvalue of A iff the homogeneous linear system $(A - \lambda I)\mathbf{x} = \mathbf{0}$ has nonzero solutions \mathbf{x} . The set of all eigenvectors corresponding to λ , together with the zero vector, is a subspace; the subspace is called the **eigenspace** of λ and is denoted by \mathbb{E}_λ .*

Proof. From Definition 4.1.1, $A\mathbf{x} = \lambda\mathbf{x} \iff A\mathbf{x} - \lambda\mathbf{x} = \mathbf{0} \iff A\mathbf{x} - \lambda I\mathbf{x} = \mathbf{0} \iff (A - \lambda I)\mathbf{x} = \mathbf{0}$. Also, eigenvectors \mathbf{x} must be nonzero, so the homogeneous system $(A - \lambda I)\mathbf{x} = \mathbf{0}$ must have nonzero solutions. Theorem 3.4.11 assures us that the set of solutions to a homogeneous system, such as $(A - \lambda I)\mathbf{x} = \mathbf{0}$ for any given λ , is a subspace. Hence the set of eigenvectors for any given eigenvalue λ , nonzero solutions, together with $\mathbf{0}$, form a subspace. \square

Example 4.1.8. Reconsider the symmetric matrix

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

of Example 4.1.2. Find the eigenspaces \mathbb{E}_1 , \mathbb{E}_3 and \mathbb{E}_0 .

Solution: • The eigenspace \mathbb{E}_1 is the set of solutions of

$$(A - 1I)\mathbf{x} = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 1 & -1 \\ 0 & -1 & 0 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

That is, $-x_2 = 0$, $-x_1 + x_2 - x_3 = 0$ and $-x_2 = 0$. Hence, $x_2 = 0$ and $x_1 = -x_3$. A general solution is $\mathbf{x} = (t, 0, -t)$ so the eigenspace $\mathbb{E}_1 = \{(t, 0, -t) : t \in \mathbb{R}\} = \text{span}\{(1, 0, -1)\}$.

- The eigenspace \mathbb{E}_3 is the set of solutions of

$$(A - 3I)\mathbf{x} = \begin{bmatrix} -2 & -1 & 0 \\ -1 & -1 & -1 \\ 0 & -1 & -2 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

That is, $-2x_1 - x_2 = 0$, $-x_1 - x_2 - x_3 = 0$ and $-x_2 - 2x_3 = 0$. From the first $x_2 = -2x_1$ which substituted into the third gives $2x_3 = -x_2 = 2x_1$. This suggests we try $x_1 = t$, $x_2 = -2t$ and $x_3 = t$; that is, $\mathbf{x} = (t, -2t, t)$. This also satisfies the second equation and so is a general solution. So the eigenspace $\mathbb{E}_3 = \{(t, -2t, t) : t \in \mathbb{R}\} = \text{span}\{(1, -2, 1)\}$.

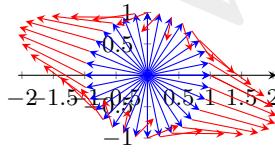
- The eigenspace \mathbb{E}_0 is the set of solutions of

$$(A - 0I)\mathbf{x} = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

That is, $x_1 - x_2 = 0$, $-x_1 + 2x_2 - x_3 = 0$ and $-x_2 + x_3 = 0$. The first and third of these require $x_1 = x_2 = x_3$ which also satisfies the second. Thus a general solution is $\mathbf{x} = (t, t, t)$ so the eigenspace $\mathbb{E}_0 = \{(t, t, t) : t \in \mathbb{R}\} = \text{span}\{(1, 1, 1)\}$.

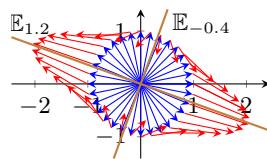
■

Example 4.1.9 (graphical eigenspaces).



The plot on the left shows (unit) vectors \mathbf{x} (blue), and for the matrix A of Example 4.1.5 the corresponding vectors $A\mathbf{x}$ (red) adjoined. Estimate and draw the eigenspaces of matrix A .

Solution:



Example 4.1.5 found directions in which $A\mathbf{x}$ is aligned with \mathbf{x} . Then the corresponding eigenspace is all vectors in the line aligned with that direction, including the opposite direction.

■

Example 4.1.10.

Eigenspaces may be multidimensional. Find the eigenspaces of the diagonal matrix

$$D = \begin{bmatrix} -\frac{1}{3} & 0 & 0 \\ 0 & \frac{3}{2} & 0 \\ 0 & 0 & \frac{3}{2} \end{bmatrix}.$$

Solution: Example 4.1.6 shows this matrix has two distinct eigenvalues $\lambda = -\frac{1}{3}$ and $\lambda = \frac{3}{2}$.

- Eigenvectors corresponding to eigenvalue $\lambda = -\frac{1}{3}$ satisfy

$$(D + \frac{1}{3}I)\mathbf{x} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & \frac{11}{6} & 0 \\ 0 & 0 & \frac{11}{6} \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

Hence $\mathbf{x} = t\mathbf{e}_1$ are eigenvectors, for nonzero t . The eigenspace $\mathbb{E}_{-1/3} = \{t\mathbf{e}_1 : t \in \mathbb{R}\} = \text{span}\{\mathbf{e}_1\}$.

- Eigenvectors corresponding to eigenvalue $\lambda = \frac{3}{2}$ satisfy

$$(D - \frac{3}{2}I)\mathbf{x} = \begin{bmatrix} -\frac{11}{6} & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

Hence $\mathbf{x} = t\mathbf{e}_2 + s\mathbf{e}_3$ are eigenvectors, for nonzero t and s . Then the eigenspace $\mathbb{E}_{3/2} = \{t\mathbf{e}_2 + s\mathbf{e}_3 : t, s \in \mathbb{R}\} = \text{span}\{\mathbf{e}_2, \mathbf{e}_3\}$ is two dimensional.

■

Definition 4.1.11. For a symmetric matrix, the **multiplicity** of an eigenvalue λ is the dimension of the corresponding eigenspace \mathbb{E}_λ .¹

Example 4.1.12. The multiplicity of the various eigenvalues in earlier examples are the following.

4.1.8 Recall that in this example:

- the eigenspace $\mathbb{E}_1 = \text{span}\{(1, 0, -1)\}$ has dimension one, so the multiplicity of eigenvalue $\lambda = 1$ is one;
- the eigenspace $\mathbb{E}_3 = \text{span}\{(1, -2, 1)\}$ has dimension one, so the multiplicity of eigenvalue $\lambda = 3$ is one; and
- the eigenspace $\mathbb{E}_0 = \text{span}\{(1, 1, 1)\}$ has dimension one, so the multiplicity of eigenvalue $\lambda = 0$ is one.

4.1.10 Recall that in this example:

- the eigenspace $\mathbb{E}_{-1/3} = \text{span}\{\mathbf{e}_1\}$ has dimension one, so the multiplicity of eigenvalue $\lambda = -1/3$ is one; and
- the eigenspace $\mathbb{E}_{3/2} = \text{span}\{\mathbf{e}_2, \mathbf{e}_3\}$ has dimension two, so the multiplicity of eigenvalue $\lambda = 3/2$ is two.

■

¹ Section 7.3 discusses that for non-symmetric matrices the dimension of an eigenspace may be less than the multiplicity of an eigenvalue, Theorem 7.3.12, but for symmetric matrices they are the same.

Table 4.1: As well as the Matlab/Octave commands and operations listed in Tables 1.2, 2.3, 3.1, 3.2, and 3.3 we need the eigenvector function.

- $[V, D] = \text{eig}(A)$ computes eigenvectors and the eigenvalues of the $n \times n$ square matrix A .
 - The n eigenvalues of A (repeated according to their multiplicity, Definition 4.1.11) form the diagonal of $n \times n$ square matrix $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$.
 - Corresponding to the j th eigenvalue λ_j , the j th column of $n \times n$ square matrix V is an eigenvector (of unit length).
- $\text{eig}(A)$ by itself just reports, in a vector, the eigenvalues of square matrix A (repeated according to their multiplicity, Definition 4.1.11).
- If the matrix A is a real symmetric matrix, then the eigenvalues and eigenvectors are all real, and the eigenvector matrix V is orthogonal.
If the matrix A is either not symmetric, or is complex valued, then the eigenvalues and eigenvectors may be complex valued.

4.1.1 Systematically find eigenvalues and eigenvectors

Computer packages easily compute eigenvalues and eigenvectors for us. However, sometimes we need to explicitly see dependence upon a parameter so we also explore small systems by hand.

Compute eigenvalues and eigenvectors

Compute in Matlab/Octave. $[V, D] = \text{eig}(A)$ computes eigenvalues and eigenvectors. The eigenvalues are placed in the diagonal of $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. The j th column of V is a unit eigenvector corresponding to the j th eigenvalue λ_j : $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$. If the matrix A is real and symmetric, then V is an orthogonal matrix (Theorem 4.2.18).

Example 4.1.13. Reconsider the symmetric matrix of Example 4.1.2:

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

Use Matlab/Octave to find its eigenvalues and corresponding eigenvectors. Confirm that $AV = VD$ for matrices $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$ and $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, and confirm that the computed V is orthogonal.

Solution: Enter the matrix into Matlab/Octave and execute `eig()`:

```
A=[1 -1 0;-1 2 -1;0 -1 1]
[V,D]=eig(A)
```



The output is

```
A =
 1   -1   0
 -1   2  -1
 0  -1   1

V =
 -0.5774  -0.7071  0.4082
 -0.5774   0.0000 -0.8165
 -0.5774   0.7071  0.4082

D =
 0.0000      0      0
      0  1.0000      0
      0      0  3.0000
```

- The first diagonal element of D is zero (actually 10^{-16}) so eigenvalue $\lambda_1 = 0$. A corresponding eigenvector is the first column of V, namely $v_1 = -0.5774(1, 1, 1)$; since eigenvectors can be scaled by a constant, we could also say an eigenvector is $v_1 = (1, 1, 1)$.
- The second diagonal element of D is one so eigenvalue $\lambda_2 = 1$. A corresponding eigenvector is the second column of V, namely $v_2 = 0.7071(-1, 0, 1)$; we could also say an eigenvector is $v_2 = (-1, 0, 1)$.
- The third diagonal element of D is three so eigenvalue $\lambda_3 = 3$. A corresponding eigenvector is the third column of V, namely $v_3 = 0.4082(1, -2, 1)$; we could also say an eigenvector is $v_3 = (1, -2, 1)$.

Confirm $AV = VD$ simply by computing $A*V-V*D$ and seeing it is zero (to numerical error of circa 10^{-16}):

```
ans =
 5.7715e-17  -1.1102e-16  4.4409e-16
 1.6874e-16   1.2490e-16  0.0000e+00
 -5.3307e-17  -1.1102e-16 -2.2204e-16
```

To verify the computed matrix V is orthogonal (Definition 3.2.35), check $V'*V$ gives the identity:

```
ans =
 1.0000   -0.0000   0.0000
 -0.0000    1.0000   0.0000
 0.0000    0.0000   1.0000
```



Example 4.1.14 (application to vibrations). Consider three masses in a row connected by two springs: on a tiny scale this could represent a molecule of carbon dioxide (CO_2). For simplicity suppose the three

masses are equal, and the spring strengths are equal. Define $y_i(t)$ to be the distance from equilibrium of the i th mass. Newton's law for bodies says the acceleration of the mass, d^2y_i/dt^2 , is proportional to the forces due to the springs. Hooke's law for springs says the force is proportional to the stretching/compression of the springs, $y_2 - y_1$ and $y_3 - y_2$. For simplicity, suppose the constants of proportionality are all one.

- The left mass (y_1) is accelerated by the spring connecting it to the middle mass (y_2); that is, $d^2y_1/dt^2 = y_2 - y_1$.
- The middle mass (y_2) is accelerated by the springs connecting it to the left mass (y_1) and to the right mass (y_3); that is, $d^2y_2/dt^2 = (y_1 - y_2) + (y_3 - y_2) = y_1 - 2y_2 + y_3$.
- The right mass (y_3) is accelerated by the spring connecting it to the middle mass (y_2); that is, $d^2y_3/dt^2 = y_2 - y_3$.

Guess there are solutions oscillating in time, so let's see if we can find solutions $y_i(t) = x_i \cos(ft)$ for some as yet unknown frequency f . Substitute and the three differential equations become

$$\begin{aligned} -f^2x_1 \cos(ft) &= x_2 \cos(ft) - x_1 \cos(ft), \\ -f^2x_2 \cos(ft) &= x_1 \cos(ft) - 2x_2 \cos(ft) + x_3 \cos(ft), \\ -f^2x_3 \cos(ft) &= x_2 \cos(ft) - x_3 \cos(ft). \end{aligned}$$

These are satisfied for all time t only if the coefficients of the cosine are equal on each side of each equation:

$$\begin{aligned} -f^2x_1 &= x_2 - x_1, \\ -f^2x_2 &= x_1 - 2x_2 + x_3, \\ -f^2x_3 &= x_2 - x_3. \end{aligned}$$

Moving the terms on the left to the right, and all terms on the right to the left, this becomes the eigenproblem $A\mathbf{x} = \lambda\mathbf{x}$ for symmetric matrix A of Example 4.1.13 and for eigenvalue $\lambda = f^2$, the square of the as yet unknown frequency. The symmetry of matrix A reflects Newton's law that every action has an equal and opposite reaction: symmetric matrices arise commonly in applications.

Example 4.1.13 tells us that there are three possible eigenvalue and eigenvector solutions for us to interpret.

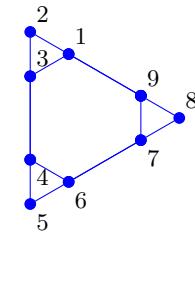
- The eigenvalue $\lambda = 1$ and corresponding eigenvector $\mathbf{x} \propto (-1, 0, 1)$ corresponds to oscillations of frequency $f = \sqrt{\lambda} = \sqrt{1} = 1$. The eigenvector $(-1, 0, 1)$ shows the middle mass is stationary while the outer two masses oscillate in and out in opposition to each other.
- The eigenvalue $\lambda = 3$ and corresponding eigenvector $\mathbf{x} \propto (1, -2, 1)$ corresponds to oscillations of higher frequency $f = \sqrt{\lambda} = \sqrt{3}$. The eigenvector $(1, -2, 1)$ shows the outer two

masses oscillate together, and the middle mass moves opposite to them.

- The eigenvalue $\lambda = 0$ and corresponding eigenvector $\mathbf{x} \propto (1, 1, 1)$ appears as oscillations of zero frequency $f = \sqrt{\lambda} = \sqrt{0} = 0$ which is a static displacement. The eigenvector $(1, 1, 1)$ shows the static displacement is that of all three masses moved all together as a unit.

That these three solutions combine together form a general solution of the system of differential equations is a topic for a course on differential equations. ■

Example 4.1.15 (Sierpinski network). Consider three triangles formed into a triangle (as shown in the margin)—perhaps because triangles make strong structures, or perhaps because of a hierarchical computer/social network. Form an matrix $A = [a_{ij}]$ of ones if node i is connected to node j ; and set the diagonal a_{ii} to be minus the number of other nodes to which node i is connected. The symmetry of the matrix A follows from the symmetry of the connections: construct the matrix, check it is symmetric, and find the eigenvalues and eigenspaces with Matlab/Octave, and their multiplicity. For the computed matrices V and D , check that $AV = VD$ and also that V is orthogonal.



Solution: In Matlab/Octave use

```

A=[-3 1 1 0 0 0 0 0 1
   1 -2 1 0 0 0 0 0 0
   1 1 -3 1 0 0 0 0 0
   0 0 1 -3 1 1 0 0 0
   0 0 0 1 -2 1 0 0 0
   0 0 0 1 1 -3 1 0 0
   0 0 0 0 0 1 -3 1 1
   0 0 0 0 0 0 1 -2 1
   1 0 0 0 0 0 1 1 -3 ]
A-A'
[V,D]=eig(A)

```

To two decimal places so that it fits the page, the computation may give

```

V =
-0.41 0.51 -0.16 -0.21 -0.45 0.18 -0.40 0.06 0.33
 0.00 -0.13 0.28 0.63 0.13 -0.18 -0.58 -0.08 0.33
 0.41 -0.20 -0.49 -0.42 0.32 0.01 -0.36 -0.17 0.33
-0.41 -0.11 0.52 -0.42 0.32 0.01 0.14 -0.37 0.33
-0.00 -0.18 -0.26 0.37 -0.22 0.51 0.36 -0.46 0.33
 0.41 0.53 0.07 0.05 -0.10 -0.51 0.33 -0.23 0.33
-0.41 -0.39 -0.36 0.05 -0.10 -0.51 0.25 0.31 0.33
 0.00 0.31 -0.03 0.16 0.55 0.34 0.22 0.55 0.33
 0.41 -0.33 0.42 -0.21 -0.45 0.18 0.03 0.40 0.33
D =

```

$$\begin{array}{cccccccccc} -5.00 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & -4.30 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -4.30 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -3.00 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -3.00 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & -3.00 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & -0.70 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & -0.70 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & -0.00 \end{array}$$

The five eigenvalues are -5.00 , -4.30 , -3.00 , -0.70 and 0.00 (to two decimal places). Three of the eigenvalues are repeated as a consequence of the additional geometric symmetry in the network. The following are the eigenspaces.

- Corresponding to eigenvalue $\lambda = -5$ are eigenvectors $\mathbf{x} \propto (-0.41, 0, 0.41, -0.41, 0, 0.41, -0.41, 0, 0.41)$; that is, the eigenspace $\mathbb{E}_{-5} = \text{span}\{(-1, 0, 1, -1, 0, 1, -1, 0, 1)\}$. From Definition 4.1.11 the multiplicity of eigenvalue $\lambda = -5$ is one.
- Corresponding to eigenvalue $\lambda = -4.30$ there are two eigenvectors computed by Matlab/Octave. These two eigenvectors are orthogonal (you should check). Because these arise as the solutions of the homogeneous system $(A - \lambda I)\mathbf{x} = \mathbf{0}$, any (nonzero) linear combination of these is also an eigenvector corresponding to the same eigenvalue. That is, the eigenspace

$$\mathbb{E}_{-4.30} = \text{span} \left\{ \begin{bmatrix} 0.51 \\ -0.13 \\ -0.20 \\ -0.11 \\ -0.18 \\ 0.53 \\ -0.39 \\ 0.31 \\ -0.33 \end{bmatrix}, \begin{bmatrix} -0.16 \\ 0.28 \\ -0.49 \\ 0.52 \\ -0.26 \\ 0.07 \\ -0.36 \\ -0.03 \\ 0.42 \end{bmatrix} \right\}.$$

Hence the eigenvalue $\lambda = -4.30$ has multiplicity two.

- Corresponding to eigenvalue $\lambda = -3$ there are three eigenvectors computed by Matlab/Octave. These three eigenvectors are orthogonal (you should check). Thus the eigenspace

$$\mathbb{E}_{-3} = \text{span} \left\{ \begin{bmatrix} -0.21 \\ 0.63 \\ -0.42 \\ -0.42 \\ 0.37 \\ 0.05 \\ 0.05 \\ 0.16 \\ -0.21 \end{bmatrix}, \begin{bmatrix} -0.45 \\ 0.13 \\ 0.32 \\ 0.32 \\ -0.22 \\ -0.10 \\ -0.10 \\ 0.55 \\ -0.45 \end{bmatrix}, \begin{bmatrix} 0.18 \\ -0.18 \\ 0.01 \\ 0.01 \\ 0.51 \\ -0.51 \\ -0.51 \\ 0.34 \\ 0.18 \end{bmatrix} \right\},$$

and so eigenvalue $\lambda = -3$ has multiplicity three.

- Corresponding to eigenvalue $\lambda = -0.70$ there are two eigenvectors computed by Matlab/Octave. These two eigenvectors are orthogonal (you should check). Thus the eigenspace

$$\mathbb{E}_{-0.70} = \text{span} \left\{ \begin{bmatrix} -0.40 \\ -0.58 \\ -0.36 \\ 0.14 \\ 0.36 \\ 0.33 \\ 0.25 \\ 0.22 \\ 0.03 \end{bmatrix}, \begin{bmatrix} 0.06 \\ -0.08 \\ -0.17 \\ -0.37 \\ -0.46 \\ -0.23 \\ 0.31 \\ 0.55 \\ 0.40 \end{bmatrix} \right\},$$

and so eigenvalue $\lambda = -0.70$ has multiplicity two.

- Lastly, corresponding to eigenvalue $\lambda = 0$ are eigenvectors $\mathbf{x} \propto (0.33, 0.33, 0.33, 0.33, 0.33, 0.33, 0.33, 0.33, 0.33)$; that is, the eigenspace $\mathbb{E}_0 = \text{span}\{(1, 1, 1, 1, 1, 1, 1, 1, 1)\}$, and so eigenvalue $\lambda = 0$ has multiplicity one.

Then check $\mathbf{A}*\mathbf{V}-\mathbf{V}*\mathbf{D}$ is zero (2 d.p.),

```
ans =
    0.00  0.00  0.00  0.00  0.00 -0.00 -0.00  0.00 -0.00
    0.00  0.00 -0.00  0.00 -0.00  0.00 -0.00 -0.00  0.00
    0.00  0.00  0.00  0.00 -0.00 -0.00 -0.00  0.00 -0.00
   -0.00  0.00 -0.00 -0.00  0.00 -0.00  0.00 -0.00  0.00
   -0.00 -0.00  0.00 -0.00  0.00 -0.00 -0.00  0.00 -0.00
    0.00  0.00  0.00 -0.00  0.00  0.00 -0.00  0.00  0.00
   -0.00 -0.00  0.00 -0.00 -0.00  0.00 -0.00 -0.00 -0.00
    0.00  0.00  0.00 -0.00  0.00 -0.00  0.00 -0.00  0.00
    0.00  0.00  0.00  0.00 -0.00 -0.00  0.00 -0.00 -0.00
```

and confirm V is orthogonal by checking $V'*V$ is the identity (2 d.p.)

```
ans =
    1.00 -0.00 -0.00  0.00  0.00  0.00 -0.00  0.00  0.00
   -0.00  1.00 -0.00 -0.00  0.00  0.00 -0.00 -0.00 -0.00
   -0.00 -0.00  1.00 -0.00  0.00  0.00 -0.00  0.00  0.00
    0.00 -0.00 -0.00  1.00  0.00 -0.00 -0.00  0.00 -0.00
    0.00  0.00  0.00  0.00  1.00  0.00 -0.00  0.00 -0.00
    0.00  0.00  0.00 -0.00  0.00  1.00 -0.00  0.00 -0.00
   -0.00 -0.00 -0.00 -0.00 -0.00 -0.00  1.00 -0.00  0.00
    0.00 -0.00  0.00  0.00  0.00  0.00 -0.00  1.00  0.00
    0.00 -0.00  0.00 -0.00 -0.00 -0.00  0.00  0.00  1.00
```

Challenge: find the two smallest connected networks that have different connectivity and yet the same eigenvalues (unit strength connections).

In 1966 Mark Kac asked “Can one hear the shape of the drum?” That is, from just knowing the eigenvalues of a network such as the one in Example 4.1.15, can one infer the connectivity of the network? The question for 2D drums was answered “no” in 1992 by Gordon, Webb and Wolpert who constructed two different shaped

2D drums which have the same set of frequencies of oscillation: that is, the same set of eigenvalues.

Why write “the computation may give” in Example 4.1.15? The reason is associated with the duplicated eigenvalues. What is important is the eigenspace. When an eigenvalue of a symmetric matrix is duplicated (or triplicated) there are many choices of eigenvectors that form an orthonormal basis (Definition 3.4.14) of the eigenspace (the same holds for singular vectors of a duplicated singular value). Different algorithms may report different orthonormal bases of the same eigenspace: those given in Example 4.1.15 are just one possibility for each eigenspace.

Theorem 4.1.16. *Consider any $n \times n$ matrix A (not just symmetric). Then $\lambda_1, \lambda_2, \dots, \lambda_m$ are eigenvalues of A with corresponding eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$, for some m (commonly $m = n$), iff $AV = VD$ for diagonal matrix $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$ and $n \times m$ matrix $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_m]$ for non-zero $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$.*

Proof. From Definition 4.1.1, λ_j are eigenvalues and non-zero \mathbf{v}_j are eigenvectors iff $A\mathbf{v}_j = \lambda_j\mathbf{v}_j$. Form into the matrix equation

$$\begin{aligned} & [A\mathbf{v}_1 \ A\mathbf{v}_2 \ \cdots \ A\mathbf{v}_m] = [\lambda_1\mathbf{v}_1 \ \lambda_2\mathbf{v}_2 \ \cdots \ \lambda_m\mathbf{v}_m] \\ \iff & A[\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_m] \\ = & [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_m] \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_m \end{bmatrix} \\ \iff & AV = VD. \end{aligned}$$

□

Example 4.1.17. Use Matlab/Octave to compute eigenvectors and the eigenvalues of (symmetric) matrix

$$A = \begin{bmatrix} 2 & 2 & -2 & 0 \\ 2 & -1 & -2 & -3 \\ -2 & -2 & 4 & 0 \\ 0 & -3 & 0 & 1 \end{bmatrix}$$

Confirm $AV = VD$ for the computed matrices.

Solution: First compute

```
A=[2 2 -2 0
   2 -1 -2 -3
   -2 -2 4 0
   0 -3 0 1]
[V,D]=eig(A)
```

The output is (2 d.p.)



```
V =
-0.23  0.83  0.08  0.50
 0.82  0.01 -0.40  0.42
 0.15  0.52 -0.42 -0.72
 0.51  0.20  0.81 -0.23

D =
-3.80      0      0      0
 0     0.77      0      0
 0      0    2.50      0
 0      0      0   6.53
```

Hence the eigenvalues are (2 d.p.) $\lambda_1 = -3.80$, $\lambda_2 = 0.77$, $\lambda_3 = 2.50$ and $\lambda_4 = 6.53$. From the columns of V , corresponding eigenvectors are (2 d.p.) $v_1 \propto (-0.23, 0.82, 0.15, 0.51)$, $v_2 \propto (0.83, 0.01, 0.52, 0.20)$, $v_3 \propto (0.08, -0.40, -0.42, 0.81)$, and $v_4 \propto (0.50, 0.42, -0.72, -0.23)$.

Then confirm $A*V-V*D$ is zero:

```
ans =
 0.00  0.00  0.00  0.00
 0.00  0.00 -0.00  0.00
-0.00  0.00 -0.00 -0.00
-0.00 -0.00 -0.00  0.00
```

■

Find eigenvalues and eigenvectors by hand

- Recall from previous study (Theorem 3.2.6) that a 2×2 matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ has determinant $\det A = |A| = ad - bc$, and A is not invertible iff $\det A = 0$.
- Similarly, although not justified until Chapter 6, a 3×3 matrix $A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$ has determinant $\det A = |A| = aei + bfg + cdh - ceg - afh - bdi$, and A is not invertible iff $\det A = 0$.

This section shows these two formulas for a determinant are useful for hand calculations on small problems. The formulas are best remembered via the following diagrams where products along the red lines are subtracted from the sum of products along the blue lines, respectively:

$$\begin{array}{c} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \\ \diagup \text{red} \quad \diagdown \text{blue} \\ \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \\ \diagup \text{blue} \quad \diagdown \text{red} \quad \diagup \text{red} \quad \diagdown \text{blue} \\ \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \end{array} \quad (4.1)$$

Chapter 6 extends the determinant to any size matrix, and explores more useful properties, but for now this is the information we need on determinants.

For hand calculation on small matrices the key is the following. By Definition 4.1.1 eigenvalues and eigenvectors are determined from $A\mathbf{x} = \lambda\mathbf{x}$. Rearranging, this equation is equivalent to $(A - \lambda I)\mathbf{x} = \mathbf{0}$. Recall that $(A - \lambda I)\mathbf{x} = \mathbf{0}$ has nonzero solutions \mathbf{x} iff the determinant $\det(A - \lambda I) = 0$. Since eigenvectors must be nonzero, the eigenvalues of a square matrix are precisely the solutions of $\det(A - \lambda I) = 0$. This reasoning leads to the following procedure.

Procedure 4.1.18 (eigenvalues and eigenvectors). *To find by hand eigenvalues and eigenvectors of a (small) square matrix A :*

1. *find all eigenvalues by solving the **characteristic equation** of A , $\det(A - \lambda I) = 0$;*
2. *for each eigenvalue λ , solve $(A - \lambda I)\mathbf{x} = \mathbf{0}$ to find the eigenspace \mathbb{E}_λ ;*
3. *write each eigenspace as the span of a few chosen eigenvectors.*

This procedure applies to general matrices A , as fully established in Section 7.1, but this chapter uses it only for symmetric matrices. Further, this chapter uses it only as a convenient method to illustrate properties by some hand calculation. None of the beautiful theorems of the next Section 4.2 for symmetric matrices are based upon this procedure.

Example 4.1.19. Use Procedure 4.1.18 to find the eigenvalues and eigenvectors of the matrix

$$A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}$$

(this is the matrix illustrated in Examples 4.1.3 and 4.1.4).

Solution: Follow the two steps of Procedure 4.1.18.

(a) Solve $\det(A - \lambda I) = 0$ for the eigenvalues λ . Using (4.1),

$$\det(A - \lambda I) = \begin{vmatrix} 1 - \lambda & -\frac{1}{2} \\ -\frac{1}{2} & 1 - \lambda \end{vmatrix} = (1 - \lambda)^2 - \frac{1}{4} = 0.$$

That is, $(\lambda - 1)^2 = \frac{1}{4}$. Taking account of both square roots this quadratic gives $\lambda - 1 = \pm \frac{1}{2}$; that is, $\lambda = 1 \pm \frac{1}{2} = \frac{1}{2}, \frac{3}{2}$ are the only two eigenvalues.

(b) Consider the two eigenvalues in turn.

i. For eigenvalue $\lambda = \frac{1}{2}$ solve $(A - \lambda I)\mathbf{x} = \mathbf{0}$. That is,

$$(A - \frac{1}{2}I)\mathbf{x} = \begin{bmatrix} 1 - \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & 1 - \frac{1}{2} \end{bmatrix} \mathbf{x} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The first component of this system says $x_1 - x_2 = 0$; that is, $x_2 = x_1$. The second component of this system

says $-x_1 + x_2 = 0$; that is, $x_2 = x_1$ (the same). So a general solution for a corresponding eigenvector is $\mathbf{x} = (1, 1)t$ for any nonzero t . That is, the eigenspace $\mathbb{E}_{1/2} = \text{span}\{(1, 1)\}$.

ii. For eigenvalue $\lambda = \frac{3}{2}$ solve $(A - \lambda I)\mathbf{x} = \mathbf{0}$. That is,

$$(A - \frac{3}{2}I)\mathbf{x} = \begin{bmatrix} 1 - \frac{3}{2} & -\frac{1}{2} \\ -\frac{1}{2} & 1 - \frac{3}{2} \end{bmatrix} \mathbf{x} = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & -\frac{1}{2} \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The first component of this system says $x_1 + x_2 = 0$, as does the second component; that is, $x_2 = -x_1$. So a general solution for a corresponding eigenvector is $\mathbf{x} = (1, -1)t$ for any nonzero t . That is, the eigenspace $\mathbb{E}_{3/2} = \text{span}\{(1, -1)\}$.

■

Example 4.1.20. Use the determinant to confirm that $\lambda = 0, 1, 3$ are the only eigenvalues of the matrix

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

(Example 4.1.8 already found the eigenspaces corresponding to these three eigenvalues.)

Solution: To find all eigenvalues, find all solutions of the characteristic equation $\det(A - \lambda I) = 0$. Using (4.1),

$$\begin{aligned} & \det(A - \lambda I) \\ &= \begin{vmatrix} 1 - \lambda & -1 & 0 \\ -1 & 2 - \lambda & -1 \\ 0 & -1 & 1 - \lambda \end{vmatrix} \\ &= (1 - \lambda)^2(2 - \lambda) + 0 + 0 - 0 - (1 - \lambda) - (1 - \lambda) \\ &= (1 - \lambda)[(1 - \lambda)(2 - \lambda) - 2] \\ &= (1 - \lambda)[2 - 3\lambda + \lambda^2 - 2] \\ &= (1 - \lambda)[-3\lambda + \lambda^2] \\ &= (1 - \lambda)(-3 + \lambda)\lambda. \end{aligned}$$

So the characteristic equation is $(1 - \lambda)(-3 + \lambda)\lambda = 0$ whose only solutions are the three eigenvalues $\lambda = 0, 1, 3$ as previously identified.

■

Example 4.1.21. Use Procedure 4.1.18 to find all eigenvalues and the corresponding eigenspaces of the symmetric matrix

$$A = \begin{bmatrix} -2 & 0 & -6 \\ 0 & 4 & 6 \\ -6 & 6 & -9 \end{bmatrix}.$$

Solution: Follow the steps of Procedure 4.1.18.

- (a) Solve $\det(A - \lambda I) = 0$ for the eigenvalues. Using (4.1),

$$\begin{aligned} & \det(A - \lambda I) \\ &= \begin{vmatrix} -2 - \lambda & 0 & -6 \\ 0 & 4 - \lambda & 6 \\ -6 & 6 & -9 - \lambda \end{vmatrix} \\ &= (-2 - \lambda)(4 - \lambda)(-9 - \lambda) + 0 \cdot 6 \cdot (-6) + (-6) \cdot 0 \cdot 6 \\ &\quad - (-6)(4 - \lambda)(-6) - (-2 - \lambda) \cdot 6 \cdot 6 - 0 \cdot 0 \cdot (-9 - \lambda) \\ &= (2 + \lambda)(4 - \lambda)(9 + \lambda) + 36(-4 + \lambda) + 36(2 + \lambda) \\ &= -\lambda^3 - 7\lambda^2 + 98\lambda \\ &= -\lambda(\lambda^2 - 7\lambda + 98) \\ &= -\lambda(\lambda - 14)(\lambda + 7). \end{aligned}$$

This determinant is zero only for the three eigenvalues $\lambda = 0, -7, 14$.

- (b) Consider the three eigenvalues in turn.

- i. For eigenvalue $\lambda = 0$ solve $(A - \lambda I)\mathbf{v} = \mathbf{0}$. That is,

$$\begin{aligned} (A - 0I)\mathbf{v} &= \begin{bmatrix} -2 & 0 & -6 \\ 0 & 4 & 6 \\ -6 & 6 & -9 \end{bmatrix} \mathbf{v} \\ &= \begin{bmatrix} -2v_1 - 6v_3 \\ 4v_2 + 6v_3 \\ -6v_1 + 6v_2 - 9v_3 \end{bmatrix} = \mathbf{0}. \end{aligned}$$

The first row says $v_1 = -3v_3$, the second row says $v_2 = -\frac{3}{2}v_3$. Substituting these into the left-hand side of the third row gives $-6v_1 + 6v_2 - 9v_3 = 18v_3 - 9v_3 - 9v_3 = 0$ for all v_3 which confirms there are non-zero solutions to form eigenvectors. Eigenvectors may be written in the form $\mathbf{v} = (-3v_3, -\frac{3}{2}v_3, v_3)$; that is, the eigenspace $\mathbb{E}_0 = \text{span}\{(-6, -3, 2)\}$.

- ii. For eigenvalue $\lambda = 7$ solve $(A - \lambda I)\mathbf{v} = \mathbf{0}$. That is,

$$\begin{aligned} (A - 7I)\mathbf{v} &= \begin{bmatrix} -9 & 0 & -6 \\ 0 & -3 & 6 \\ -6 & 6 & -16 \end{bmatrix} \mathbf{v} \\ &= \begin{bmatrix} -9v_1 - 6v_3 \\ -3v_2 + 6v_3 \\ -6v_1 + 6v_2 - 16v_3 \end{bmatrix} = \mathbf{0}. \end{aligned}$$

The first row says $v_1 = -\frac{2}{3}v_3$, the second row says $v_2 = 2v_3$. Substituting these into the left-hand side of the third row gives $-6v_1 + 6v_2 - 16v_3 = 4v_3 + 12v_3 - 16v_3 = 0$ for all v_3 which confirms there are non-zero solutions to form eigenvectors. Eigenvectors may be written in the form $\mathbf{v} = (-\frac{2}{3}v_3, 2v_3, v_3)$; that is, the eigenspace $E_7 = \text{span}\{(-2, 6, 3)\}$.

iii. For eigenvalue $\lambda = -14$ solve $(A - \lambda I)\mathbf{v} = \mathbf{0}$. That is,

$$\begin{aligned}(A + 14I)\mathbf{v} &= \begin{bmatrix} 12 & 0 & -6 \\ 0 & 18 & 6 \\ -6 & 6 & 5 \end{bmatrix} \mathbf{v} \\ &= \begin{bmatrix} 12v_1 - 6v_3 \\ 18v_2 + 6v_3 \\ -6v_1 + 6v_2 + 5v_3 \end{bmatrix} = \mathbf{0}.\end{aligned}$$

The first row says $v_1 = \frac{1}{2}v_3$, the second row says $v_2 = -\frac{1}{3}v_3$. Substituting these into the left-hand side of the third row gives $-6v_1 + 6v_2 + 5v_3 = -3v_3 - 2v_3 + 5v_3 = 0$ for all v_3 which confirms there are non-zero solutions to form eigenvectors. Eigenvectors may be written in the form $\mathbf{v} = (\frac{1}{2}v_3, -\frac{1}{3}v_3, v_3)$; that is, the eigenspace $E_{-14} = \text{span}\{(3, -2, 6)\}$.

■

General matrices may have complex valued eigenvalues and eigenvectors, as seen in the next example, and for good reasons in some applications. One of the key results of the next Section 4.2 is to prove that real symmetric matrices always have real eigenvalues and eigenvectors. There are many applications where this reality is crucial.

Example 4.1.22. Find the eigenvalues and a corresponding eigenvector for

This example aims to recall basic properties of complex numbers as a prelude to the proof of reality of eigenvalues for symmetric matrices.

Solution: the non-symmetric matrix $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$.

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

That is, $\lambda^2 = -1$. Taking square roots we find there are two complex eigenvalues $\lambda = \pm\sqrt{-1} = \pm i$. Despite the appearance of complex numbers, all our arithmetic, algebra and properties continue to hold. Thus we proceed to find complex valued eigenvectors.

(b) Consider the two eigenvalues in turn.

i. For eigenvalue $\lambda = i$ solve $(A - \lambda I)\mathbf{x} = \mathbf{0}$. That is,

$$(A - iI)\mathbf{x} = \begin{bmatrix} -i & 1 \\ -1 & -i \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The first component of this system says $-ix_1 + x_2 = 0$; that is, $x_2 = ix_1$. The second component of this system says $-x_1 - ix_2 = 0$; that is, $x_2 = ix_1$ (the same). So a general corresponding eigenvector is $\mathbf{x} = (1, i)t$ for any nonzero t .

ii. For eigenvalue $\lambda = -i$ solve $(A - \lambda I)\mathbf{x} = \mathbf{0}$. That is,

$$(A + iI)\mathbf{x} = \begin{bmatrix} i & 1 \\ -1 & i \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The first component of this system says $ix_1 + x_2 = 0$; that is, $x_2 = -ix_1$. The second component of this system says $-x_1 + ix_2 = 0$; that is, $x_2 = -ix_1$ (the same). So a general corresponding eigenvector is $\mathbf{x} = (1, -i)t$ for any nonzero t .

■

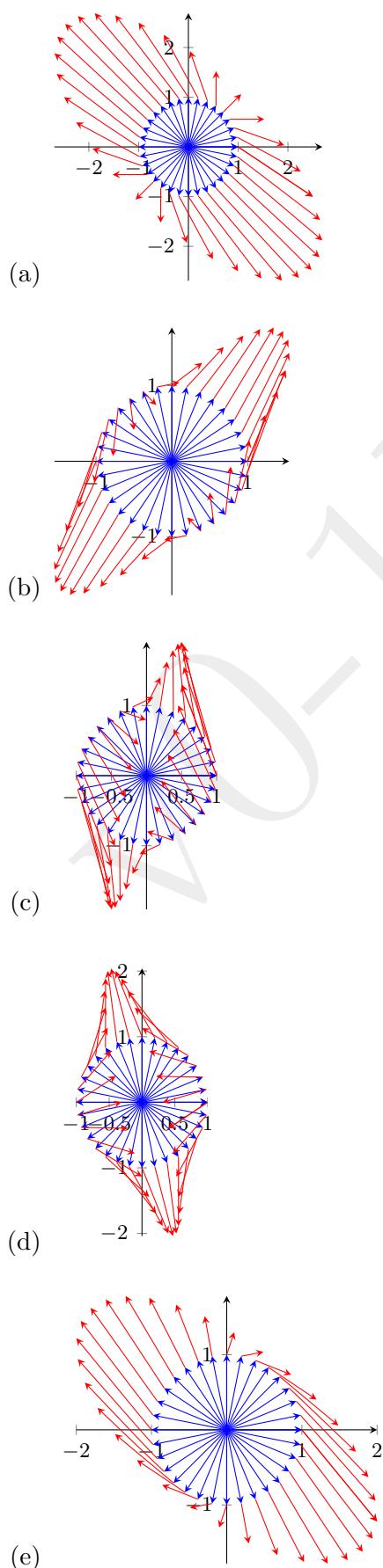
Example 4.1.22 is a problem that might arise using calculus to describe the dynamics of a mass on a spring. Let the displacement of the mass be $y_1(t)$ then Newton's law says the acceleration $d^2y_1/dt^2 \propto -y_1$, the negative of the displacement; for simplicity, let the constant of proportionality be one. Introduce $y_2(t) = dy_1/dt$ then Newton's law becomes $dy_2/dt = -y_1$. Seek solutions of these two first-order differential equations in the form $y_j(t) = x_j e^{\lambda t}$ and the differential equations become $x_2 = \lambda x_1$ and $\lambda x_2 = -x_1$ respectively. Forming into a matrix-vector problem these are

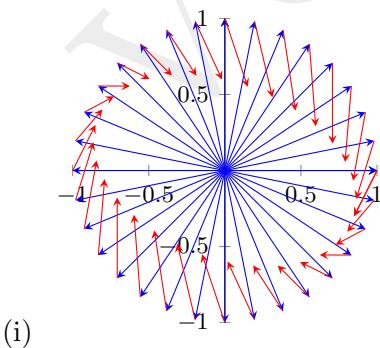
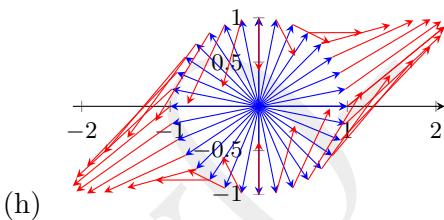
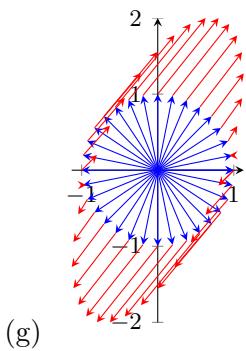
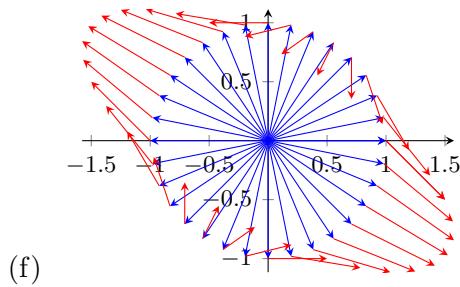
$$\begin{bmatrix} x_2 \\ -x_1 \end{bmatrix} = \lambda \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \iff \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \mathbf{x} = \lambda \mathbf{x}.$$

We need to find the eigenvalues and eigenvectors of the matrix: we derive eigenvalues are $\lambda = \pm\sqrt{-1} = \pm i$. Such complex eigenvalues represent oscillations in time t since, for example, $e^{\lambda t} = e^{it} = \cos t + i \sin t$ by Euler's formula.

4.1.2 Exercises

Exercise 4.1.1. Each plot below shows (unit) vectors \mathbf{x} (blue), and for some matrix the corresponding vectors $A\mathbf{x}$ (red) adjoined. Estimate which directions \mathbf{x} are eigenvectors of matrix A , and for each eigenvector estimate the corresponding eigenvalue.





Exercise 4.1.2. In each case use the matrix-vector product to determine which of the given vectors are eigenvectors of the given matrix? and for each eigenvector what is the corresponding eigenvalue?

$$(a) \begin{bmatrix} 6 & 2 \\ 3 & 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -3 \end{bmatrix}, \begin{bmatrix} 3 \\ -2 \end{bmatrix}, \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{4} \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 \\ -3 \end{bmatrix}$$

$$(b) \begin{bmatrix} -2 & 1 \\ 3 & 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} -2 \\ 2 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -\frac{1}{3} \\ -1 \end{bmatrix}, \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}$$

$$(c) \begin{bmatrix} 3 & 3 & 4 \\ 0 & -4 & 0 \\ -2 & -1 & -6 \end{bmatrix}, \begin{bmatrix} 4 \\ 0 \\ -1 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ 2 \end{bmatrix}, \begin{bmatrix} -2 \\ 6 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} \frac{1}{3} \\ -1 \\ \frac{1}{6} \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 2 \end{bmatrix}$$

$$(d) \begin{bmatrix} 3 & 0 & 0 \\ 2 & -5 & -3 \\ -4 & -2 & 0 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ 2 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 0 \\ 3 \\ 1 \end{bmatrix}, \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ -\frac{1}{3} \end{bmatrix}$$

$$(e) \begin{bmatrix} -2 & 0 & 0 & 0 \\ 1 & -2 & 3 & -2 \\ -5 & -3 & 4 & -2 \\ -5 & 1 & -1 & 0 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \\ 4 \\ 7 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \\ 2 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -2 \\ 2 \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ -1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 2 \\ 1 \end{bmatrix}$$

$$(f) \begin{bmatrix} 5 & -4 & -1 & 0 \\ 1 & 0 & -1 & 0 \\ -5 & 8 & 2 & 5 \\ 4 & -4 & 1 & 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 0 \\ 4 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0 \\ -\frac{3}{5} \end{bmatrix}, \begin{bmatrix} 1 \\ \frac{1}{2} \\ 3 \\ 3 \end{bmatrix}, \begin{bmatrix} 0 \\ \frac{1}{2} \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ -1 \\ 2 \\ 1 \end{bmatrix}$$

Exercise 4.1.3. Use Matlab/Octave function `eig()` to determine the eigenvalues and corresponding eigenspaces of the following symmetric matrices.

$$(a) \begin{bmatrix} -1 & -\frac{4}{5} & \frac{12}{5} & -2 \\ -\frac{4}{5} & \frac{1}{5} & -\frac{8}{5} & -\frac{12}{5} \\ \frac{12}{5} & -\frac{8}{5} & -\frac{11}{5} & -\frac{14}{5} \\ -2 & -\frac{12}{5} & -\frac{14}{5} & 2 \end{bmatrix}$$

$$(b) \begin{bmatrix} \frac{7}{5} & -\frac{2}{5} & 0 & -\frac{2}{5} \\ -\frac{2}{5} & 2 & \frac{2}{5} & 0 \\ 0 & \frac{2}{5} & \frac{13}{5} & \frac{2}{5} \\ -\frac{2}{5} & 0 & \frac{2}{5} & 2 \end{bmatrix}$$

$$(c) \begin{bmatrix} \frac{30}{7} & \frac{16}{7} & \frac{16}{7} & \frac{4}{7} \\ \frac{16}{7} & \frac{30}{7} & \frac{16}{7} & \frac{4}{7} \\ \frac{16}{7} & \frac{16}{7} & \frac{30}{7} & \frac{4}{7} \\ \frac{4}{7} & \frac{4}{7} & \frac{4}{7} & \frac{15}{7} \end{bmatrix}$$

$$(d) \begin{bmatrix} -\frac{36}{7} & -\frac{8}{7} & \frac{20}{7} & \frac{12}{7} \\ -\frac{8}{7} & -\frac{6}{7} & -\frac{12}{7} & \frac{20}{7} \\ \frac{20}{7} & -\frac{12}{7} & -3 & -\frac{32}{7} \\ \frac{12}{7} & \frac{20}{7} & -\frac{32}{7} & -3 \end{bmatrix}$$

$$(e) \begin{bmatrix} -2.6 & -2.7 & 5.2 & 2.1 \\ -2.7 & 4.6 & 9.9 & 5.2 \\ 5.2 & 9.9 & -2.6 & 2.7 \\ 2.1 & 5.2 & 2.7 & 4.6 \end{bmatrix}$$

$$(f) \begin{bmatrix} 1.4 & -7.1 & -0.7 & 6.2 \\ -7.1 & -1.0 & -2.2 & -2.5 \\ -0.7 & -2.2 & -3.4 & -4.1 \\ 6.2 & -2.5 & -4.1 & -1.0 \end{bmatrix}$$

$$(g) \begin{bmatrix} -1 & 1 & 4 & -3 & 1 \\ 1 & 0 & -2 & 1 & 0 \\ 4 & -2 & 1 & -3 & 0 \\ -3 & 1 & -3 & 2 & -1 \\ 1 & 0 & 0 & -1 & -1 \end{bmatrix}$$

$$(h) \begin{bmatrix} 5 & -1 & 3 & 0 & 0 \\ -1 & -1 & 0 & 0 & -2 \\ 3 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & -2 & 1 & 0 & -1 \end{bmatrix}$$

Exercise 4.1.4. For each of the given symmetric matrices, determine all eigenvalues by finding and solving the characteristic equation of the matrix.

$$(a) \begin{bmatrix} 2 & 3 \\ 3 & 2 \end{bmatrix}$$

$$(b) \begin{bmatrix} 6 & \frac{11}{2} \\ \frac{11}{2} & 6 \end{bmatrix}$$

(c)
$$\begin{bmatrix} -5 & 1 \\ 2 & -2 \end{bmatrix}$$

(d)
$$\begin{bmatrix} -5 & 5 \\ 5 & -5 \end{bmatrix}$$

(e)
$$\begin{bmatrix} 5 & -4 \\ -4 & -1 \end{bmatrix}$$

(f)
$$\begin{bmatrix} -2 & \frac{9}{2} \\ \frac{9}{2} & 10 \end{bmatrix}$$

(g)
$$\begin{bmatrix} 6 & 0 & -4 \\ 0 & 6 & 3 \\ -4 & 3 & 6 \end{bmatrix}$$

(h)
$$\begin{bmatrix} -2 & 4 & 6 \\ 4 & 0 & 4 \\ 6 & 4 & -2 \end{bmatrix}$$

(i)
$$\begin{bmatrix} 2 & -3 & -3 \\ -3 & 2 & -3 \\ -3 & -3 & 2 \end{bmatrix}$$

(j)
$$\begin{bmatrix} 4 & -4 & 3 \\ -4 & -2 & 6 \\ 3 & 6 & -8 \end{bmatrix}$$

(k)
$$\begin{bmatrix} 8 & 4 & 2 \\ 4 & 0 & 0 \\ 2 & 0 & 0 \end{bmatrix}$$

(l)
$$\begin{bmatrix} 0 & 0 & -3 \\ 0 & 2 & 0 \\ -3 & 0 & 0 \end{bmatrix}$$

Example 4.1.23. For each symmetric matrix, find the eigenspace of the given ‘eigenvalues’ by hand solution of linear equations, or determine from your solution that the given value cannot be an eigenvalue.

(a)
$$\begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}, 4, -2$$

(b)
$$\begin{bmatrix} 4 & -2 \\ -2 & 7 \end{bmatrix}, 3, 6$$

(c)
$$\begin{bmatrix} -7 & 0 & 2 \\ 0 & -7 & -2 \\ 2 & -2 & -0 \end{bmatrix}, -8, -7, 1$$

(d)
$$\begin{bmatrix} 0 & 6 & -3 \\ 6 & 0 & 7 \\ -3 & 7 & 3 \end{bmatrix}, -6, 4, 9$$

(e)
$$\begin{bmatrix} 0 & -4 & 2 \\ -4 & 1 & -0 \\ 2 & -0 & 1 \end{bmatrix}, -4, 1, 2$$

(f)
$$\begin{bmatrix} 7 & -4 & -2 \\ -4 & 9 & -4 \\ -2 & -4 & 7 \end{bmatrix}, 1, 4, 9$$

Example 4.1.24. For each symmetric matrix, find by hand all eigenvalues and an orthogonal basis for the corresponding eigenspace. What is the multiplicity of each eigenvalue?

(a)
$$\begin{bmatrix} -8 & 3 \\ 3 & 0 \end{bmatrix}$$

(b)
$$\begin{bmatrix} 6 & -5 \\ -5 & 6 \end{bmatrix}$$

(c)
$$\begin{bmatrix} -2 & -2 \\ -2 & -5 \end{bmatrix}$$

(d)
$$\begin{bmatrix} 2 & -3 \\ -3 & -6 \end{bmatrix}$$

(e)
$$\begin{bmatrix} -1 & -3 & -3 \\ -3 & -5 & 3 \\ -3 & 3 & -1 \end{bmatrix}$$

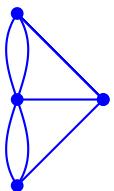
(f)
$$\begin{bmatrix} -5 & 2 & -2 \\ 2 & -2 & 1 \\ -2 & 1 & -10 \end{bmatrix}$$

(g)
$$\begin{bmatrix} 11 & 4 & -2 \\ 4 & 5 & 4 \\ -2 & 4 & 11 \end{bmatrix}$$

(h)
$$\begin{bmatrix} -7 & 2 & 2 \\ 2 & -6 & 0 \\ 2 & 0 & -8 \end{bmatrix}$$

(i)
$$\begin{bmatrix} 6 & 10 & -5 \\ 10 & 13 & -2 \\ -5 & -2 & -2 \end{bmatrix}$$

(j)
$$\begin{bmatrix} 4 & 3 & 1 \\ 3 & -4 & -3 \\ 1 & -3 & 4 \end{bmatrix}$$

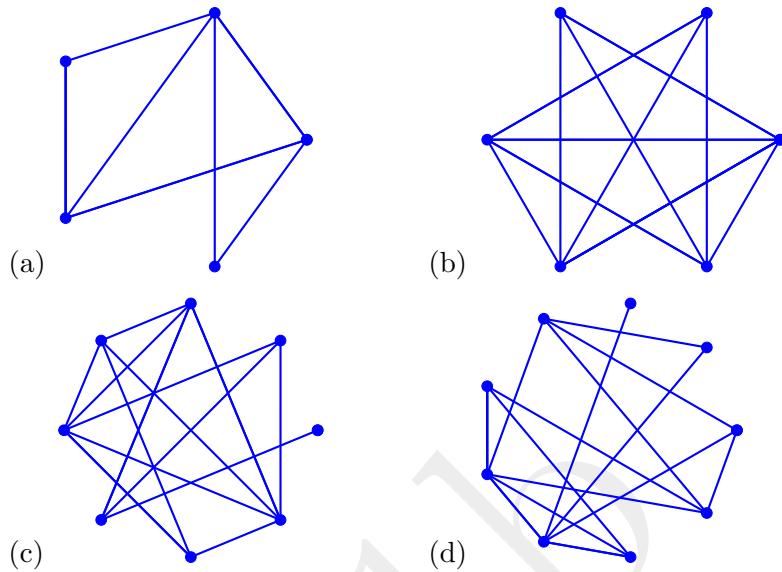


Exercise 4.1.5 (Seven Bridges of Königsberg). The marginal picture shows the abstract graph of the seven bridges of Königsberg during the time of Euler: the small disc nodes represent the islands of Königsberg; the lines between represent the seven different bridges. This abstract graph is famous for its role in founding the theory of such networks, but this exercise addresses an aspect relevant to well used web search software. Number the nodes from 1 to 4. Form the 4×4 symmetric matrix of the number of lines from each node to the other nodes (and zero for the number of lines from a node to itself). Use Matlab/Octave function `eig()` to find the eigenvalues and eigenvectors for this matrix. Analogous to well known web search software, identify the largest eigenvalue and a corresponding eigenvector: then rank the importance of each node in order of the size of the component in the corresponding eigenvector.

Exercise 4.1.6. For each of the following networks:

- label the nodes;
- construct the symmetric adjacency matrix A such that a_{ij} is one if node i is linked to node j , and a_{ij} is zero otherwise (and zero on the diagonal);
- in Matlab/Octave use `eig()` to find all eigenvalues and eigenvectors;
- rank the ‘importance’ of the nodes from the magnitude of their component in the eigenvector of the largest (most positive) eigenvalue.

Although a well known web search engine does much the same thing for web pages, it uses an approximate iterative algorithm more suited to the mind-bogglingly vast size of the internet.



4.2 Beautiful properties for symmetric matrices

Section Contents

4.2.1	Matrix powers maintain eigenvectors	403
4.2.2	Symmetric matrices are orthogonally diagonalisable	408
4.2.3	Change orthonormal basis to classify quadratics	418
	Graph quadratic equations	418
	Simplify quadratic forms	423
4.2.4	Exercises	425

This section starts with two properties for eigenvalues of general matrices, and then proceeds to the special case of real symmetric matrices which have the beautifully useful properties of always having real eigenvalues and orthogonal eigenvectors.

4.2.1 Matrix powers maintain eigenvectors

Recall that section 3.2 introduced the inverse of a matrix (Definition 3.2.2). The first theorem links an eigenvalue of zero to the non-existence of an inverse and hence to more problematic linear equations.

Theorem 4.2.1. *A square matrix is invertible iff zero is not an eigenvalue of the matrix.*

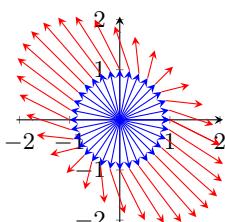
Proof. From Definition 4.1.1, zero is an eigenvalue, $\lambda = 0$, iff $A\mathbf{x} = 0\mathbf{x}$ has nonzero solutions \mathbf{x} ; that is, iff the homogeneous system $A\mathbf{x} = \mathbf{0}$ has nonzero solutions \mathbf{x} . But the Unique Solution Theorems 3.3.21 or 3.4.35 assure us that this occurs iff matrix A is not invertible. Consequently a matrix is invertible iff zero is not an eigenvalue. \square

Example 4.2.2. • The 3×3 matrix of Example 4.1.2 (also 4.1.8, 4.1.13 and 4.1.20) is not invertible as among its eigenvalues of 0, 1 and 3 it has zero as an eigenvalue.

- The plot in the margin shows (unit) vectors \mathbf{x} (blue), and for some matrix A the corresponding vectors $A\mathbf{x}$ (red) adjoined. There are no directions \mathbf{x} for which $A\mathbf{x} = \mathbf{0} = 0\mathbf{x}$. Hence zero cannot be an eigenvalue and the matrix A must be invertible.

Similarly for Example 4.1.5.

- The 3×3 diagonal matrix of Example 4.1.10 has eigenvalues of only $-\frac{1}{3}$ and $\frac{3}{2}$. Since zero is not an eigenvalue, the matrix is invertible.



- The 9×9 matrix of the Sierpinski network in Example 4.1.15 is not invertible as it has zero among its nine eigenvalues.
- The 2×2 matrix of Example 4.1.19 is invertible as its eigenvalues are $\lambda = \frac{1}{2}, \frac{3}{2}$, neither of which are zero. Indeed, the matrix

$$A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}, \quad \text{has inverse } A^{-1} = \begin{bmatrix} \frac{4}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{4}{3} \end{bmatrix}$$

as matrix multiplication confirms.

- The 2×2 non-symmetric matrix of Example 4.1.22 is invertible because zero is not among its eigenvalues of $\lambda = \pm i$. Indeed, the matrix

$$A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad \text{has inverse } A^{-1} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$$

as matrix multiplication confirms. ■

Example 4.2.3. The next theorem considers eigenvalues and eigenvectors of powers of a matrix. Two examples are the following.

- Recall the matrix $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ has eigenvalues $\lambda = \pm i$. The square of this matrix

$$A^2 = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}$$

is diagonal so its eigenvalues are the diagonal elements (Example 4.1.6), namely the only eigenvalue is -1 . Observe that A^2 's eigenvalue, $-1 = (\pm i)^2$, is the square of the eigenvalues of A . That the eigenvalues of A^2 are the square of those of A holds generally.

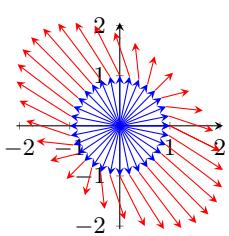
- Also recall matrix

$$A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}, \quad \text{has inverse } A^{-1} = \begin{bmatrix} \frac{4}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{4}{3} \end{bmatrix}.$$

Let's determine the eigenvalues of this inverse. Its characteristic equation (defined in Procedure 4.1.18) is

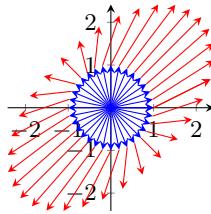
$$\det(A^{-1} - \lambda I) = \begin{vmatrix} \frac{4}{3} - \lambda & \frac{2}{3} \\ \frac{2}{3} & \frac{4}{3} - \lambda \end{vmatrix} = (\frac{4}{3} - \lambda)^2 - \frac{4}{9} = 0.$$

That is, $(\lambda - \frac{4}{3})^2 = \frac{4}{9}$. Taking the square-root of both sides gives $\lambda - \frac{4}{3} = \pm \frac{2}{3}$; that is, the two eigenvalues of the inverse A^{-1} are $\lambda = \frac{4}{3} \pm \frac{2}{3} = 2, \frac{2}{3}$. Observe these eigenvalues



of the inverse are the reciprocals of the eigenvalues $\frac{1}{2}, \frac{3}{2}$ of A . This relation also holds generally.

The marginal pictures illustrates the reciprocal relation graphically: the first picture shows $A\mathbf{x}$ for various \mathbf{x} , the second picture shows $A^{-1}\mathbf{x}$. The eigenvector directions are the same for both matrix and inverse. But in those directions where the matrix stretches, the inverse shrinks, and where the matrix shrinks, the inverse stretches. The relationship is obscure in directions which are not eigenvectors.



Theorem 4.2.4. Let A be a square matrix with eigenvalue λ and corresponding eigenvector \mathbf{x} .

- (a) For any positive integer n , λ^n is an eigenvalue of A^n with corresponding eigenvector \mathbf{x} .
- (b) If A is invertible, then $1/\lambda$ is an eigenvalue of A^{-1} with corresponding eigenvector \mathbf{x} .
- (c) If A is invertible, then for any integer n , λ^n is an eigenvalue of A^n with corresponding eigenvector \mathbf{x} .

Proof. Consider each property in turn.

4.2.4a. Firstly, the result hold for $n = 1$ by Definition 4.1.1, that $A\mathbf{x} = \lambda\mathbf{x}$. Secondly, for the case $n = 2$ consider $A^2\mathbf{x} = (AA)\mathbf{x} = A(A\mathbf{x}) = A(\lambda\mathbf{x}) = \lambda(A\mathbf{x}) = \lambda(\lambda\mathbf{x}) = (\lambda^2)\mathbf{x}$. Hence by definition 4.1.1 λ^2 is an eigenvalue of A^2 corresponding to eigenvector \mathbf{x} . Third, use induction to extend to any power: assume the result for $n = k$ (and proceed to prove it for $n = k + 1$). Consider $A^{k+1}\mathbf{x} = (A^k A)\mathbf{x} = A^k(A\mathbf{x}) = A^k(\lambda\mathbf{x}) = \lambda(A^k\mathbf{x}) = \lambda(\lambda^k\mathbf{x}) = \lambda^{k+1}\mathbf{x}$. Hence by Definition 4.1.1 λ^{k+1} is an eigenvalue of A^{k+1} corresponding to eigenvector \mathbf{x} . By induction the property 4.2.4a. holds for all integer $n \geq 1$.

4.2.4b. For invertible A we know none of the eigenvalues are zero: thus $1/\lambda$ exists. Pre-multiply $A\mathbf{x} = \lambda\mathbf{x}$ by $\frac{1}{\lambda}A^{-1}$ to deduce $\frac{1}{\lambda}A^{-1}A\mathbf{x} = \frac{1}{\lambda}A^{-1}\lambda\mathbf{x}$, which gives $\frac{1}{\lambda}I\mathbf{x} = \frac{1}{\lambda}\lambda A^{-1}\mathbf{x}$, that is, $\frac{1}{\lambda}\mathbf{x} = A^{-1}\mathbf{x}$. Consequently, $1/\lambda$ is an eigenvalue of A^{-1} with corresponding eigenvector \mathbf{x} .

4.2.4c. Proved by Exercise 4.2.11.

□

Example 4.2.5. Recall from Example 4.1.19 that matrix

$$A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}$$

has eigenvalues $1/2$ and $3/2$ with corresponding eigenvectors $(1, 1)$ and $(1, -1)$ respectively. Confirm matrix A^2 has eigenvalues which are these squared, and corresponding to the same eigenvectors.

Solution: Compute

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

Then

$$\begin{aligned} A^2 \begin{bmatrix} 1 \\ 1 \end{bmatrix} &= \begin{bmatrix} \frac{1}{4} \\ \frac{1}{4} \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \\ A^2 \begin{bmatrix} 1 \\ -1 \end{bmatrix} &= \begin{bmatrix} \frac{9}{4} \\ -\frac{9}{4} \end{bmatrix} = \frac{9}{4} \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \end{aligned}$$

and so A^2 has eigenvalues $1/4 = (1/2)^2$ and $9/4 = (3/2)^2$ with the same corresponding eigenvectors $(1, 1)$ and $(1, -1)$ respectively. ■

Example 4.2.6. Given that the matrix

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

has eigenvalues 2 , 1 and -1 with corresponding eigenvectors $(1, 1, 1)$, $(-1, 0, 1)$ and $(1, -2, 1)$ respectively. Confirm matrix A^2 has eigenvalues which are these squared, and corresponding to the same eigenvectors. Given the inverse

$$A^{-1} = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

confirm its eigenvalues are the reciprocals of those of A , and for corresponding eigenvectors.

Solution: • Compute

$$A^2 = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}.$$

Then

$$A^2 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \\ 4 \end{bmatrix} = 4 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix},$$

$$A^2 \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = 1 \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix},$$

$$A^2 \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = 1 \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix},$$

has eigenvalues $4 = 2^2$ and $1 = (\pm 1)^2$ with corresponding eigenvectors $(1, 1, 1)$, and the pair $(-1, 0, 1)$ and $(1, -2, 1)$. Thus here $\text{span}\{(-1, 0, 1), (1, -2, 1)\}$ is the eigenspace of A^2 corresponding to eigenvalue one.

- For the inverse

$$A^{-1} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} = \frac{1}{2} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix},$$

$$A^{-1} \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} = 1 \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix},$$

$$A^{-1} \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \\ -1 \end{bmatrix} = (-1) \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix},$$

has eigenvalues $1/2$, $1 = 1/1$ and $-1 = 1/(-1)$ with corresponding eigenvectors $(1, 1, 1)$, $(-1, 0, 1)$ and $(1, -2, 1)$.

■

Example 4.2.7 (long term age structure). Recall Example 3.1.5 introduced how to use a Leslie matrix to predict the future population of an animal. In the example, letting $\mathbf{x} = (x_1, x_2, x_3)$ be the current number of pups, juveniles, and mature females respectively, then for the Leslie matrix

$$L = \begin{bmatrix} 0 & 0 & 4 \\ \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} \end{bmatrix}$$

the predicted population numbers after a year is $\mathbf{x}' = L\mathbf{x}$, after two years is $\mathbf{x}'' = L\mathbf{x}' = L^2\mathbf{x}$, and so on. Predict what happens after many generations: does the population die out? grow? oscillate?

Solution: Consider what happens after n generations for large n , say $n = 10$ or 100 . The predicted population is $\mathbf{x}^{(n)} = L^n\mathbf{x}$; that is, the matrix L^n transforms the current population to that after n generations. The stretching and/or shrinking of matrix L^n is summarised by its eigenvectors and eigenvalues (section 4.1). By Theorem 4.2.4 the eigenvalues of L^n are λ^n in terms of the

eigenvalues λ of L . By hand (Procedure 4.1.18), the characteristic equation of L is

$$\begin{aligned}\det(L - \lambda I) &= \begin{vmatrix} -\lambda & 0 & 4 \\ \frac{1}{2} & -\lambda & 0 \\ 0 & \frac{1}{3} & \frac{1}{3} - \lambda \end{vmatrix} \\ &= \lambda^2(\frac{1}{3} - \lambda) + 0 + \frac{2}{3} - 0 - 0 - 0 \\ &= (1 - \lambda)(\lambda^2 + \frac{2}{3}\lambda + \frac{2}{3}) \\ &= (1 - \lambda)[(\lambda + \frac{1}{3})^2 + \frac{5}{9}] = 0 \\ \implies \lambda &= 1, (-1 \pm i\sqrt{5})/3.\end{aligned}$$

Such complex valued eigenvalues may arise in real applications when the matrix is not symmetric, as here—the next Theorem 4.2.8 proves such complexities do not arise for symmetric matrices.

But the algebra still works with complex eigenvalues (Chapter 7). Here, the eigenvalues of L^n are λ^n (Theorem 4.2.4) namely $1^n = 1$ for all n and $[(-1 \pm i\sqrt{5})/3]^n$. Because the absolute value $|(-1 \pm i\sqrt{5})/3| = |-1 \pm i\sqrt{5}|/3 = \sqrt{1+5}/3 = \sqrt{6}/3 = \sqrt{2/3} = 0.8165$, then the absolute value of $[(-1 \pm i\sqrt{5})/3]^n$ is 0.8165^n which becomes negligibly small for large n ; for example, $0.8165^{34} \approx 0.001$. Since the eigenvectors of L^n are the same as those of L (Theorem 4.2.4), these negligibly small eigenvalues of L^n imply that any component in the initial population in the direction of the corresponding eigenvectors is shrunk to zero by L^n . For large n , it is only the component in the eigenvector corresponding to eigenvalue $\lambda = 1$ that remains. Find the eigenvector by solving $(L - I)\mathbf{x} = \mathbf{0}$, namely

$$\begin{bmatrix} -1 & 0 & 4 \\ \frac{1}{2} & -1 & 0 \\ 0 & \frac{1}{3} & -\frac{2}{3} \end{bmatrix} \mathbf{x} = \begin{bmatrix} -x_1 + 4x_3 \\ \frac{1}{2}x_1 - x_2 \\ \frac{1}{3}x_2 - \frac{2}{3}x_3 \end{bmatrix} = \mathbf{0}.$$

The first row gives that $x_1 = 4x_3$, the third row that $x_2 = 2x_3$, and the second row confirms these are correct as $\frac{1}{2}x_1 - x_2 = \frac{1}{2}4x_3 - 2x_3 = 0$. Eigenvectors corresponding to $\lambda = 1$ are then of the form $(4x_3, 2x_3, x_3) = (4, 2, 1)x_3$. Because the corresponding eigenvalue of $L^n = 1^n = 1$ the component of \mathbf{x} in this direction remains in $L^n\mathbf{x}$ whereas all other components decay to zero. Thus the model predicts that after many generations the population reaches a steady state of the pups, juveniles, and mature females being in the ratio of 4 : 2 : 1.

■

4.2.2 Symmetric matrices are orthogonally diagonalisable

General matrices may have complex valued eigenvalues (as in Examples 4.1.22 and 4.2.7): that real symmetric matrices always have

real eigenvalues (such as in all matrices of Examples 4.2.5 and 4.2.6) is a special property that often reflects the physical reality of many applications.

To establish the reality of eigenvalues (Theorem 4.2.8), the proof invokes contradiction. The contradiction is to assume a complex valued eigenvalue exists and then prove it cannot. Thus the proof of Theorem 4.2.8 needs to use some complex numbers and some properties of complex numbers. Recall: that any complex number $z = a + bi$ has a complex conjugate $\bar{z} = a - bi$ (denoted by the overbar); a complex number equals its conjugate only if it is real valued (the imaginary part is zero); and properties of complex numbers and operations also hold for complex valued vectors, complex valued matrices, and arithmetic operations with complex matrices and vectors.

Theorem 4.2.8. *Let A be a real symmetric matrix, then the eigenvalues of A are all real.*

Proof. Let λ be any eigenvalue of A with corresponding eigenvector \mathbf{x} ; that is, $A\mathbf{x} = \lambda\mathbf{x}$ for $\mathbf{x} \neq \mathbf{0}$. To establish a contradiction: assume the eigenvalue λ is complex valued, and so correspondingly is the eigenvector \mathbf{x} . First, the complex conjugate $\bar{\lambda}$ must be another eigenvalue of A , corresponding to an eigenvector which is the complex conjugate $\bar{\mathbf{x}}$. To see this just take the complex conjugate of both sides of $A\mathbf{x} = \lambda\mathbf{x}$:

$$\overline{A\mathbf{x}} = \overline{\lambda\mathbf{x}} \implies \bar{A}\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}} \implies A\bar{\mathbf{x}} = \bar{\lambda}\bar{\mathbf{x}}$$

as matrix A is real ($\bar{A} = A$). Second, consider $\mathbf{x}^T A \bar{\mathbf{x}}$ in two ways:

$$\begin{aligned} \mathbf{x}^T A \bar{\mathbf{x}} &= \mathbf{x}^T (A \bar{\mathbf{x}}) = \mathbf{x}^T (\bar{A} \bar{\mathbf{x}}) \quad (\text{as } A \text{ is real}) \\ &= \mathbf{x}^T (\overline{A \bar{\mathbf{x}}}) = \mathbf{x}^T (\overline{\lambda \bar{\mathbf{x}}}) = \mathbf{x}^T (\bar{\lambda} \bar{\mathbf{x}}) = \bar{\lambda} \mathbf{x}^T \bar{\mathbf{x}}; \\ \mathbf{x}^T A \bar{\mathbf{x}} &= (\mathbf{x}^T A) \bar{\mathbf{x}} = (A^T \mathbf{x})^T \bar{\mathbf{x}} = (A \mathbf{x})^T \bar{\mathbf{x}} \quad (\text{by symmetry}) \\ &= (\lambda \mathbf{x})^T \bar{\mathbf{x}} = \lambda \mathbf{x}^T \bar{\mathbf{x}}. \end{aligned}$$

Equating the two ends of this identity gives $\bar{\lambda} \mathbf{x}^T \bar{\mathbf{x}} = \lambda \mathbf{x}^T \bar{\mathbf{x}}$. Rearrange to $\bar{\lambda} \mathbf{x}^T \bar{\mathbf{x}} - \lambda \mathbf{x}^T \bar{\mathbf{x}} = (\bar{\lambda} - \lambda) \mathbf{x}^T \bar{\mathbf{x}} = 0$. Because this product is zero, either $\bar{\lambda} - \lambda = 0$ or $\mathbf{x}^T \bar{\mathbf{x}} = 0$. But we next prove the second is impossible, hence the first must hold; that is, $\bar{\lambda} - \lambda = 0$, equivalently $\bar{\lambda} = \lambda$. Consequently, the eigenvalue λ must be real—it cannot be complex.

Lastly confirm $\mathbf{x}^T \bar{\mathbf{x}} \neq 0$. The nonzero eigenvector \mathbf{x} will be generally complex, say

$$\begin{aligned} \mathbf{x} &= (a_1 + b_1 i, a_2 + b_2 i, \dots, a_n + b_n i) \\ \implies \bar{\mathbf{x}} &= (a_1 - b_1 i, a_2 - b_2 i, \dots, a_n - b_n i). \end{aligned}$$

Then the product

$$\begin{aligned}\mathbf{x}^T \bar{\mathbf{x}} &= [a_1 + b_1 i \ a_2 + b_2 i \ \cdots \ a_n + b_n i] \begin{bmatrix} a_1 - b_1 i \\ a_2 - b_2 i \\ \vdots \\ a_n - b_n i \end{bmatrix} \\ &= (a_1 + b_1 i)(a_1 - b_1 i) + (a_2 + b_2 i)(a_2 - b_2 i) \\ &\quad + \cdots + (a_n + b_n i)(a_n - b_n i) \\ &= (a_1^2 + b_1^2) + (a_2^2 + b_2^2) + \cdots + (a_n^2 + b_n^2) \\ &> 0\end{aligned}$$

since \mathbf{x} is an eigenvector which necessarily is nonzero and so at least one term in the sum is positive. \square

The other property that we have seen graphically for 2D matrices is that the eigenvectors of symmetric matrices are orthogonal. In Example 4.2.3, both the matrices A and A^{-1} in the second part are symmetric and from the marginal illustration their eigenvectors are proportional to $(1, 1)$ and $(-1, 1)$ which are orthogonal directions—they are at right angles in the illustration.

Example 4.2.9. Recall Example 4.1.21 found the 3×3 symmetric matrix

$$\begin{bmatrix} -2 & 0 & -6 \\ 0 & 4 & 6 \\ -6 & 6 & -9 \end{bmatrix}$$

has eigenspaces $\mathbb{E}_0 = \text{span}\{(-6, -3, 2)\}$, $\mathbb{E}_7 = \text{span}\{(-2, 6, 3)\}$ and $\mathbb{E}_{-14} = \text{span}\{(3, -2, 6)\}$. These eigenspaces are orthogonal as seen by the dot products

$$\begin{aligned}(-6, -3, 2) \cdot (-2, 6, 3) &= 12 - 18 + 6 = 0, \\ (-2, 6, 3) \cdot (3, -2, 6) &= -6 - 12 + 18 = 0, \\ (3, -2, 6) \cdot (-6, -3, 2) &= -18 + 6 + 12 = 0.\end{aligned}$$

■

Theorem 4.2.10. Let A be a real symmetric matrix, then for every two distinct eigenvalues of A , any corresponding two eigenvectors are orthogonal.

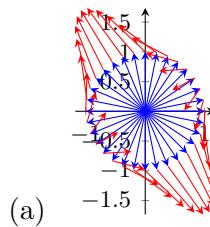
Proof. Let eigenvalues $\lambda_1 \neq \lambda_2$, and let \mathbf{x}_1 and \mathbf{x}_2 be corresponding eigenvectors, respectively; that is, $A\mathbf{x}_1 = \lambda_1\mathbf{x}_1$ and $A\mathbf{x}_2 = \lambda_2\mathbf{x}_2$. Consider $\mathbf{x}_1^T A \mathbf{x}_2$ in two ways:

$$\begin{aligned}\mathbf{x}_1^T A \mathbf{x}_2 &= \mathbf{x}_1^T (A \mathbf{x}_2) = \mathbf{x}_1^T (\lambda_2 \mathbf{x}_2) = \lambda_2 \mathbf{x}_1^T \mathbf{x}_2 = \lambda_2 \mathbf{x}_1 \cdot \mathbf{x}_2; \\ \mathbf{x}_1^T A \mathbf{x}_2 &= \mathbf{x}_1^T A^T \mathbf{x}_2 \quad (\text{as } A \text{ is symmetric}) \\ &= (\mathbf{x}_1^T A^T) \mathbf{x}_2 = (A \mathbf{x}_1)^T \mathbf{x}_2\end{aligned}$$

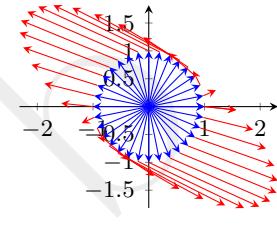
$$= (\lambda_1 \mathbf{x}_1)^T \mathbf{x}_2 = \lambda_1 \mathbf{x}_1^T \mathbf{x}_2 = \lambda_1 \mathbf{x}_1 \cdot \mathbf{x}_2.$$

Equating the two ends of this identity gives $\lambda_2 \mathbf{x}_1 \cdot \mathbf{x}_2 = \lambda_1 \mathbf{x}_1 \cdot \mathbf{x}_2$. Rearrange to $\lambda_2 \mathbf{x}_1 \cdot \mathbf{x}_2 - \lambda_1 \mathbf{x}_1 \cdot \mathbf{x}_2 = (\lambda_2 - \lambda_1)(\mathbf{x}_1 \cdot \mathbf{x}_2) = 0$. Since $\lambda_1 \neq \lambda_2$ it follows that the dot product $\mathbf{x}_1 \cdot \mathbf{x}_2 = 0$. Hence the two eigenvectors are orthogonal. \square

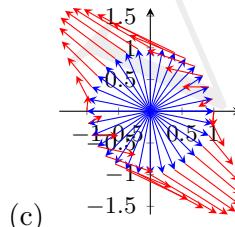
Example 4.2.11. The plots below shows (unit) vectors \mathbf{x} (blue), and for some matrix A the corresponding vectors $A\mathbf{x}$ (red) adjoined. By estimating eigenvectors determine which cases *cannot* be the plot of a real symmetric matrix.



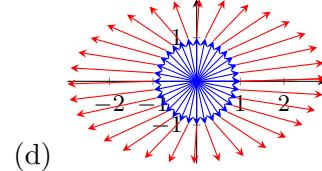
Solution: Estimate eigenvectors $(0.8, 0.5)$ and $(-0.5, 0.8)$ which are orthogonal, so may be a symmetric matrix



Solution: Estimate eigenvectors $(1, 0.1)$ and $(1, -0.3)$ which are not orthogonal, so cannot be from a symmetric matrix



Solution: Estimate eigenvectors $(1, 0.2)$ and $(0.8, -0.7)$ which are not orthogonal, so cannot be from a symmetric matrix



Solution: Estimate eigenvectors $(0.1, -1)$ and $(1, 0.1)$ which are orthogonal, so may be a symmetric matrix

Example 4.2.12. By hand find eigenvectors corresponding to the two distinct eigenvalues of the following matrices and confirm that symmetric matrix A has orthogonal eigenvectors, and that non-symmetric matrix B does not:

$$A = \begin{bmatrix} 1 & \frac{3}{2} \\ \frac{3}{2} & -3 \end{bmatrix}; \quad B = \begin{bmatrix} 0 & -3 \\ -2 & 1 \end{bmatrix}.$$

Solution: • For matrix A , the eigenvalues come from the characteristic equation

$$\begin{aligned}\det(A - \lambda I) &= (1 - \lambda)(-3 - \lambda) - \frac{9}{4} \\ &= \lambda^2 + 2\lambda - \frac{21}{4} = (\lambda + 1)^2 - \frac{25}{4} = 0,\end{aligned}$$

so eigenvalues are $\lambda = -1 \pm \frac{5}{2} = -\frac{7}{2}, \frac{3}{2}$.

– Corresponding to eigenvalue $\lambda = -7/2$, eigenvectors \mathbf{v} satisfy $(A + \frac{7}{2}I)\mathbf{v} = \mathbf{0}$, that is

$$\begin{bmatrix} \frac{9}{2} & \frac{3}{2} \\ \frac{3}{2} & \frac{1}{2} \end{bmatrix} \mathbf{v} = \frac{1}{2} \begin{bmatrix} 9v_1 + 3v_2 \\ 3v_1 + v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix},$$

giving $v_2 = -3v_1$. Eigenvectors must be $\mathbf{v} \propto (1, -3)$.

– Corresponding to eigenvalue $\lambda = 3/2$, eigenvectors \mathbf{v} satisfy $(A - \frac{3}{2}I)\mathbf{v} = \mathbf{0}$, that is

$$\begin{bmatrix} -\frac{1}{2} & \frac{3}{2} \\ \frac{3}{2} & -\frac{9}{2} \end{bmatrix} \mathbf{v} = \frac{1}{2} \begin{bmatrix} -v_1 + 3v_2 \\ 3v_1 - 9v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix},$$

giving $v_1 = 3v_2$. Eigenvectors must be $\mathbf{v} \propto (3, 1)$.

The dot product of the two basis eigenvectors is $(1, -3) \cdot (3, 1) = 3 - 3 = 0$ and hence they are orthogonal.

• For matrix B , the eigenvalues come from the characteristic equation

$$\begin{aligned}\det(B - \lambda I) &= (-\lambda)(1 - \lambda) - 6 \\ &= \lambda^2 - \lambda - 6 = (\lambda - 3)(\lambda + 2) = 0,\end{aligned}$$

so eigenvalues are $\lambda = 3, -2$.

– Corresponding to eigenvalue $\lambda = -2$, eigenvectors \mathbf{v} satisfy $(B + 2I)\mathbf{v} = \mathbf{0}$, that is

$$\begin{bmatrix} 2 & -3 \\ -2 & 3 \end{bmatrix} \mathbf{v} = \begin{bmatrix} 2v_1 - 3v_2 \\ -2v_1 + 3v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix},$$

giving $v_2 = \frac{2}{3}v_1$. Eigenvectors must be $\mathbf{v} \propto (3, 2)$.

– Corresponding to eigenvalue $\lambda = 3$, eigenvectors \mathbf{v} satisfy $(B - 3I)\mathbf{v} = \mathbf{0}$, that is

$$\begin{bmatrix} -3 & -3 \\ -2 & -2 \end{bmatrix} \mathbf{v} = \begin{bmatrix} -3v_1 - 3v_2 \\ -2v_1 - 2v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix},$$

giving $v_1 = -v_2$. Eigenvectors must be $\mathbf{v} \propto (-1, 1)$.

The dot product of the two basis eigenvectors is $(3, 2) \cdot (-1, 1) = -3 + 2 = -1 \neq 0$ and hence they are not orthogonal.

Example 4.2.13. Use Matlab/Octave to compute eigenvectors of the following matrices, and to confirm the eigenvectors are orthogonal for a symmetric matrix.

$$(a) \begin{bmatrix} 0 & 3 & 2 & -1 \\ 0 & 3 & 0 & 0 \\ 3 & 0 & -1 & -1 \\ -3 & 1 & 3 & 0 \end{bmatrix} \quad (b) \begin{bmatrix} -6 & 0 & 1 & 1 \\ 0 & 0 & 2 & 2 \\ 1 & 2 & 2 & -1 \\ 1 & 2 & -1 & -1 \end{bmatrix}$$

Solution: For each matrix, enter the matrix as say A, then execute `[V,D]=eig(A)` to give eigenvectors as the columns of V. Then confirm orthogonality of all pairs of eigenvectors by computing $V'*V$ and confirming the off-diagonal dot products are zero, or confirm the lack of orthogonality if non-zero. (In the case of repeated eigenvalues, Matlab/Octave generates an orthonormal basis for the corresponding eigenspace so the returned matrix V of eigenvectors is still orthogonal for symmetric A.)



(a) The Matlab/Octave code

```
A=[0 3 2 -1
   0 3 0 0
   3 0 -1 -1
   -3 1 3 0]
[V,D]=eig(A)
V'*V
```

gives the following (2 d.p.)

```
V =
-0.49 -0.71  0.41  0.74
 0.00  0.00  0.00  0.34
 0.32 -0.71  0.41  0.57
-0.81 -0.00  0.82 -0.06
D =
-3.00      0      0      0
 0     2.00      0      0
 0      0    0.00      0
 0      0      0    3.00
```

so eigenvectors corresponding to the four distinct eigenvalues are $(-0.49, 0, 0.32, -0.81)$, $(-0.71, 0, -0.71, 0)$, $(0.41, 0, 0.41, 0.82)$ and $(0.74, 0.34, 0.57, -0.6)$. Then $V'*V$ is (2 d.p.)

1.00	0.11	-0.73	-0.13
0.11	1.00	-0.58	-0.93
-0.73	-0.58	1.00	0.49
-0.13	-0.93	0.49	1.00

As the off-diagonal elements are non-zero, the pairs of dot products are non-zero indicating the column vectors are not orthogonal. Hence the matrix A cannot be symmetric.



(b) The Matlab/Octave code

```
B=[-6 0 1 1
   0 0 2 2
   1 2 2 -1
   1 2 -1 -1]
[V,D]=eig(B)
V'*V
```

gives the following (2 d.p.)

```
V =
 0.94 0.32 -0.04 -0.10
 0.13 -0.63 -0.53 -0.55
 -0.17 0.32 0.43 -0.83
 -0.25 0.63 -0.73 -0.08
D =
 -6.45 0 0 0
 0 -3.00 0 0
 0 0 1.11 0
 0 0 0 3.34
```

so eigenvectors corresponding to the four distinct eigenvalues are $(0.94, 0.13, -0.17, -0.25)$, $(0.32, -0.63, 0.32, 0.63)$, $(-0.04, -0.53, 0.43, -0.73)$ and $(-0.10, -0.55, -0.83, -0.08)$. Then $V'*V$ is (2 d.p.)

```
1.00 -0.00 -0.00 -0.00
-0.00 1.00 0.00 -0.00
-0.00 0.00 1.00 -0.00
-0.00 -0.00 -0.00 1.00
```

As the off-diagonal elements are zero, the pairs of dot products are zero indicating the column vectors are orthogonal. The symmetry of this matrix B requires such orthogonality.

■

Recall that when we find eigenvalues by hand for 2×2 or 3×3 matrices we solve a quadratic or cubic characteristic equation, respectively. Thus we get at most two or three eigenvalues, respectively. Further, when we ask Matlab/Octave to compute eigenvalues of an $n \times n$ matrix, it always returns n values.

Theorem 4.2.14. *Let A be an $n \times n$ real symmetric matrix, then A has at most n distinct eigenvalues.*

Proof. Invoke the pigeonhole principle and contradiction. Assume there are more than n distinct eigenvalues, then there would be more than n eigenvectors corresponding to distinct eigenvalues. Theorem 4.2.10 asserts all such eigenvectors are orthogonal. But there cannot be more than n vectors in an orthogonal set in \mathbb{R}^n (Theorem 1.3.20). Hence the assumption is wrong and there cannot be any more than n distinct eigenvalues. \square

The previous theorem establishes there are at most n distinct eigenvalues (for symmetric matrices but it is true more generally). Now we establish that typically there exist n distinct eigenvalues of an $n \times n$ matrix—here symmetric.

Example 4.0.1 started this chapter by observing that in an SVD of a symmetric matrix, $A = USV^T$, the columns of U appear to be (almost) always plus/minus the corresponding columns of V . Exceptions possibly arise in the degenerate cases when two or more singular values are identical. We now prove this close relation between U and V in all non-degenerate cases.

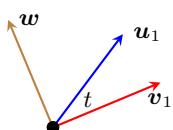
Theorem 4.2.15. *Let A be an $n \times n$ real symmetric matrix with SVD $A = USV^T$. If all the singular values are distinct or zero, $\sigma_1 > \dots > \sigma_r > \sigma_{r+1} = \dots = \sigma_n = 0$, then v_j is an eigenvector of A corresponding to an eigenvalue of either $\lambda_j = +\sigma_j$ or $\lambda_j = -\sigma_j$ (not both).*

This proof modifies parts of the proof of the SVD Theorem 3.3.5 to the specific case of a symmetric matrix.

If non-zero singular values are duplicated, then one can always choose an SVD so the result of this theorem still holds. However, the proof is too involved to give here.

Proof. First, for any zero singular value, $\sigma_j = 0$, then the result is immediate as from the SVD, $AV = US$ the j th column gives $A\mathbf{v}_j = 0\mathbf{u}_j = \mathbf{0} = 0\mathbf{v}_j$ for the nonzero \mathbf{v}_j .

Second, for the singular values $\sigma_j > 0$ use a form of induction allied with a contradiction to prove $\mathbf{u}_j = \pm \mathbf{v}_j$. By contradiction, suppose $\mathbf{u}_1 \neq \pm \mathbf{v}_1$, then we can write $\mathbf{u}_1 = \mathbf{v}_1 \cos t + \mathbf{w} \sin t$ for some vector \mathbf{w} orthogonal to \mathbf{v}_1 ($\mathbf{w} := \text{perp}_{\mathbf{v}_1} \mathbf{u}_1$) and for angle $0 < t < \pi$ (as illustrated in the margin). Multiply by A giving the identity $A\mathbf{u}_1 = A\mathbf{v}_1 \cos t + A\mathbf{w} \sin t$. Now the first column of $AV = US$ gives $A\mathbf{v}_1 = \sigma_1 \mathbf{u}_1$. Also for the symmetric matrix $A = A^T = (USV^T)^T = VS^T U^T = VSU^T$ is an alternative SVD of A : so $AU = VS$ giving in its first column $A\mathbf{u}_1 = \sigma_1 \mathbf{v}_1$. That is, the identity becomes $\sigma_1 \mathbf{v}_1 = \sigma_1 \mathbf{u}_1 \cos t + (A\mathbf{w}) \sin t$. Further, $A\mathbf{w}$ is orthogonal to \mathbf{u}_1 by the proof of the SVD Theorem 3.3.5 (since \mathbf{w} is orthogonal to \mathbf{v}_1). Equate the lengths of both sides: $\sigma_1^2 = \sigma_1^2 \cos^2 t + |A\mathbf{w}|^2 \sin^2 t$ which rearranging implies $(\sigma_1^2 - |A\mathbf{w}|^2) \sin^2 t = 0$. For angles $0 < t < \pi$ this implies $|A\mathbf{w}| = \sigma_1$ for a vector orthogonal to \mathbf{v}_1 which implies the singular value σ_1 is repeated. This contradicts the supposition; hence $\mathbf{u}_1 = \pm \mathbf{v}_1$ (for one of the signs, not both).



Recall the induction proof for the SVD Theorem 3.3.5 (section 3.3.3). Here, since $\mathbf{u}_1 = \pm \mathbf{v}_1$ we can and do choose $\bar{U} = \bar{V}$. Hence $B = \bar{U}^T A \bar{V} = \bar{V}^T A V$ is symmetric as $B^T = V^T A \bar{V}^T = \bar{V}^T A^T V = \bar{V}^T A \bar{V} = B$. Consequently, the same argument applies at all steps in the induction for the proof of an SVD and hence establishes $\mathbf{u}_j = \pm \mathbf{v}_j$ (each for one of the signs, not both).

Third and lastly, from the j th column of $AV = US$, $A\mathbf{v}_j = \sigma_j \mathbf{u}_j = \sigma_j(\pm \mathbf{v}_j) = \lambda_j \mathbf{v}_j$ for eigenvalue λ_j one of $\pm \sigma_j$ but not both. \square

Recall that for any real matrix A an SVD is $A = USV^T$. But specifically for symmetric A , the proof of the previous theorem 4.2.15 identified that the columns of US , $\sigma_j \mathbf{u}_j$, are generally the same as $\lambda_j \mathbf{v}_j$ and hence are the columns of VD where $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. In which case the SVD becomes $A = VDV^T$ and so is intimately connected to the following definition.

Definition 4.2.16. *A real square matrix A is orthogonally diagonalisable if there exists an orthogonal matrix V and a diagonal matrix D such that $V^T AV = D$, equivalently $AV = VD$, equivalently $A = VDV^T$ is a factorisation of A .*

The equivalences in this definition arise immediately from the orthogonality of matrix V (Definition 3.2.35): pre-multiply $V^T AV = D$ by V gives $VV^T AV = AV = VD$; and so on.

Example 4.2.17. (a) Recall from Example 4.2.12 that the symmetric matrix $A = \begin{bmatrix} 1 & \frac{3}{2} \\ \frac{3}{2} & -3 \end{bmatrix}$ has eigenvalues $\lambda = -\frac{7}{2}, \frac{3}{2}$ with corresponding orthogonal eigenvectors $(1, -3)$ and $(3, 1)$. Normalise these eigenvectors to unit length as the columns of the orthogonal matrix

$$\begin{aligned} V &= \begin{bmatrix} \frac{1}{\sqrt{10}} & \frac{3}{\sqrt{10}} \\ -\frac{3}{\sqrt{10}} & \frac{1}{\sqrt{10}} \end{bmatrix} = \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & 3 \\ -3 & 1 \end{bmatrix} \quad \text{then} \\ V^T AV &= \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & -3 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 & \frac{3}{2} \\ \frac{3}{2} & -3 \end{bmatrix} \frac{1}{\sqrt{10}} \begin{bmatrix} 1 & 3 \\ -3 & 1 \end{bmatrix} \\ &= \frac{1}{10} \begin{bmatrix} -\frac{7}{2} & \frac{21}{2} \\ \frac{9}{2} & \frac{3}{2} \end{bmatrix} \begin{bmatrix} 1 & 3 \\ -3 & 1 \end{bmatrix} \\ &= \frac{1}{10} \begin{bmatrix} -35 & 0 \\ 0 & 15 \end{bmatrix} = \begin{bmatrix} -\frac{7}{2} & 0 \\ 0 & \frac{3}{2} \end{bmatrix}. \end{aligned}$$

Hence this matrix is orthogonally diagonalisable.

(b) Recall from Example 4.2.13 that the symmetric matrix

$$B = \begin{bmatrix} -6 & 0 & 1 & 1 \\ 0 & 0 & 2 & 2 \\ 1 & 2 & 2 & -1 \\ 1 & 2 & -1 & -1 \end{bmatrix}$$

has orthogonal eigenvectors computed by Matlab/Octave into the orthogonal matrix V . By additionally computing $V^T * B * V$ we get the following diagonal result (2 d.p.)



```
ans =
-6.45    0.00    0.00    0.00
 0.00   -3.00    0.00   -0.00
 0.00    0.00    1.11   -0.00
-0.00   -0.00   -0.00    3.34
```

and see that this matrix B is orthogonally diagonalisable. ■

These examples of orthogonal diagonalisation involve symmetric matrices. Also, the connection between an SVD and orthogonal matrices was previously discussed only for symmetric matrices. The next theorem establishes that all real symmetric matrices are orthogonally diagonalisable, and vice versa. That is, eigenvectors form an orthogonal set if and only if the matrix is symmetric.

Theorem 4.2.18 (spectral). *Let A be real square matrix. Then matrix A is symmetric iff it is orthogonally diagonalisable.*

Proof. The “if” and the “only if” lead to two parts in the proof.

- If matrix A is orthogonally diagonalisable, then $A = VDV^T$ for orthogonal V and diagonal D (and recall $D^T = D$). Consider

$$A^T = (VDV^T)^T = V^T D^T V^T = VDV^T = A.$$

Consequently the matrix A is symmetric.

- Theorem 4.2.15 establishes the converse for the generic case of distinct singular values. If matrix A is symmetric, then Theorem 4.2.15 asserts an SVD $A = USV^T$ has matrix U such that columns $\mathbf{u}_j = \pm \mathbf{v}_j$. That is, we can write $U = VR$ for diagonal matrix $R = \text{diag}(\pm 1, \pm 1, \dots, \pm 1)$ for appropriately chosen signs. Then by the SVD $A = USV^T = VRSV^T = VDV^T$ for diagonal matrix $D = RS = \text{diag}(\pm \sigma_1, \pm \sigma_2, \dots, \pm \sigma_n)$ for the same pattern of signs. Hence matrix A is orthogonally diagonalisable.

We omit proving the degenerate case when non-zero singular values are repeated. □

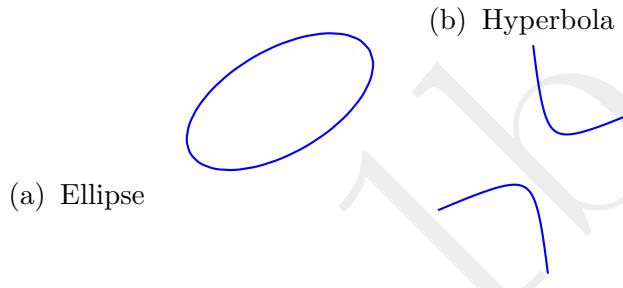
4.2.3 Change orthonormal basis to classify quadratics

An optional subsection which has many uses—although not an application itself as it does not involve real data.

The following preliminary example illustrates the important principle, applicable throughout mathematics, that we often either choose or change to a coordinate system in which the mathematical algebra is simplest.

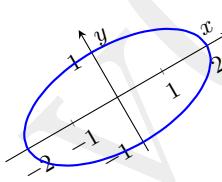
Example 4.2.19 (choose useful coordinates).

Consider the following two quadratic curves. For each curve draw a coordinate system in which the algebraic description of the curve would be most straightforward.



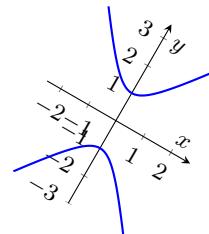
Solution: Among various possibilities are the following.

(a) Ellipse



In this coordinate system the ellipse is algebraically $(x/2)^2 + y^2 = 1$.

(b) Hyperbola



In this coordinate system the hyperbola is algebraically $y^2 = 1 + 2x^2$.

Now let's proceed to see how this geometric idea of choosing good coordinates may be implemented in algebra.

Graph quadratic equations

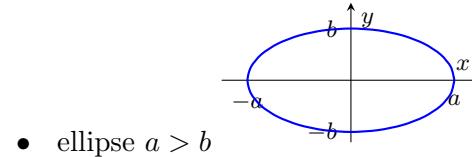
Example 4.2.19 illustrated an ellipse and a hyperbola. These curves are examples of the so-called **conic sections** which arise as solutions of the quadratic equation in two variables, say x and y ,

$$ax^2 + bxy + cy^2 + dx + ey + f = 0 \quad (4.2)$$

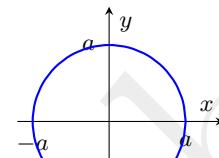
(where a, b, c cannot all be zero). As invoked in the example, the canonical simplest algebraic form of such curves are the following.

The challenge of this subsection is to choose new coordinates so that the quadratic equation (4.2) becomes one of these recognised canonical forms.

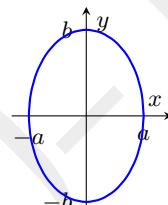
$$\text{Ellipse or circle : } \frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$$



- ellipse $a > b$

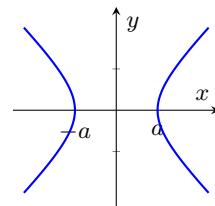


- the circle $a = b$

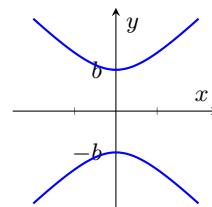


- ellipse $a < b$

$$\text{Hyperbola : } \frac{x^2}{a^2} - \frac{y^2}{b^2} = 1 \text{ or } -\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$$

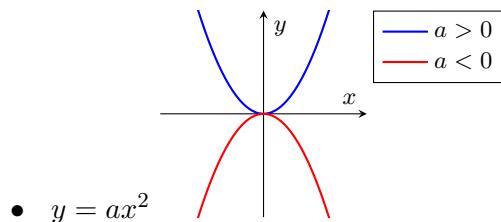


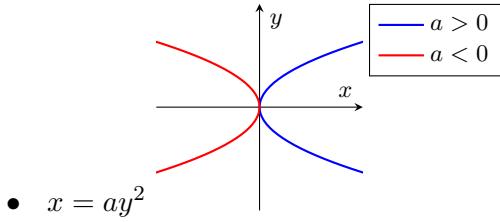
- $\frac{x^2}{a^2} - \frac{y^2}{b^2} = 1$



- $-\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1$

$$\text{Parabola : } y = ax^2 \text{ or } x = ay^2$$





Example 4.2.19 implicitly had two steps: first, we decided upon an orientation for the coordinate axes; second, we decided that the coordinate system should be ‘centred’ in the picture. Algebra follows the same two steps.

Example 4.2.20 (centre coordinates). By shifting coordinates, identify the conic section whose equation is

$$2x^2 + y^2 - 4x + 4y + 2 = 0.$$

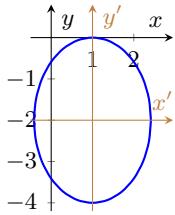
Solution: Group the linear terms with corresponding quadratic powers and seek to rewrite as a perfect square: the equation is

$$\begin{aligned} & (2x^2 - 4x) + (y^2 + 4y) + 2 = 0 \\ \iff & 2(x^2 - 2x) + (y^2 + 4y) = -2 \\ \iff & 2(x^2 - 2x + 1) + (y^2 + 4y + 4) = -2 + 2 + 4 \\ \iff & 2(x - 1)^2 + (y + 2)^2 = 4. \end{aligned}$$

Thus changing to a new (dashed) coordinate system $x' = x - 1$ and $y' = y + 2$, that is choosing the origin of the dashed coordinate system at $(x, y) = (1, -2)$, the quadratic equation becomes

$$2x'^2 + y'^2 = 4, \quad \text{that is } \frac{x'^2}{2} + \frac{y'^2}{4} = 1.$$

In this new coordinate system the equation is that of an ellipse with horizontal axis of half-length $\sqrt{2}$ and vertical axis of half-length $\sqrt{4} = 2$ (as illustrated in the margin). ■



Example 4.2.21 (rotate coordinates). By rotating the coordinate system, identify the conic section whose equation is

$$x^2 + 3xy - 3y^2 - \frac{1}{2} = 0.$$

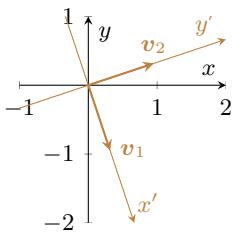
(There are no terms linear in x and y so we do not shift coordinates.)

Solution: The equation contains the product xy . To identify the conic we must eliminate the xy term. To use matrix algebra, and in terms of the vector $\mathbf{x} = (x, y)$, recognise that the quadratic terms may be written as $\mathbf{x}^T A \mathbf{x}$ for symmetric matrix

$$A = \begin{bmatrix} 1 & \frac{3}{2} \\ \frac{3}{2} & -3 \end{bmatrix} \quad \text{as then}$$

$$\begin{aligned}
 \mathbf{x}^T A \mathbf{x} &= \mathbf{x}^T \begin{bmatrix} 1 & \frac{3}{2} \\ \frac{3}{2} & -3 \end{bmatrix} \mathbf{x} \\
 &= [x \ y] \begin{bmatrix} x + \frac{3}{2}y \\ \frac{3}{2}x - 3y \end{bmatrix} \\
 &= x(x + \frac{3}{2}y) + y(\frac{3}{2}x - 3y) \\
 &= x^2 + 3xy - 3y^2.
 \end{aligned}$$

(The matrix form $\mathbf{x}^T A \mathbf{x}$ splits the cross-product term $3xy$ into two equal halves represented by the two off-diagonal elements $\frac{3}{2}$ in matrix A). Suppose we change to some new (dashed) coordinate system with its standard unit vectors \mathbf{v}_1 and \mathbf{v}_2 as illustrated in the margin. The vectors in the plane will be written as the linear combination $\mathbf{x} = \mathbf{v}_1 x' + \mathbf{v}_2 y'$. That is, $\mathbf{x} = V \mathbf{x}'$ for new coordinate vector $\mathbf{x}' = (x', y')$ and matrix $V = [\mathbf{v}_1 \ \mathbf{v}_2]$.



In the new coordinate system, related to the old by $\mathbf{x} = V \mathbf{x}'$, the quadratic terms

$$\mathbf{x}^T A \mathbf{x} = (V \mathbf{x}')^T A (V \mathbf{x}') = \mathbf{x}'^T V^T A V \mathbf{x}' = \mathbf{x}'^T (V^T A V) \mathbf{x}'.$$

Thus choose V to simplify $V^T A V$. Because matrix A is symmetric, Theorem 4.2.18 asserts it is orthogonally diagonalisable (using eigenvectors). Indeed, Example 4.2.17a orthogonally diagonalised this particular matrix A , via its eigenvalues and eigenvectors, using the orthogonal matrix

$$V = [\mathbf{v}_1 \ \mathbf{v}_2] = \begin{bmatrix} \frac{1}{\sqrt{10}} & \frac{3}{\sqrt{10}} \\ -\frac{3}{\sqrt{10}} & \frac{1}{\sqrt{10}} \end{bmatrix}.$$

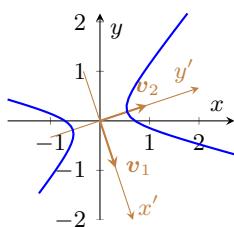
Using this V , in the new dashed coordinate system (illustrated) the quadratic terms in the equation become

$$\mathbf{x}'^T (V^T A V) \mathbf{x}' = \mathbf{x}'^T D \mathbf{x}' = \mathbf{x}'^T \begin{bmatrix} -\frac{7}{2} & 0 \\ 0 & \frac{3}{2} \end{bmatrix} \mathbf{x}' = -\frac{7}{2}x'^2 + \frac{3}{2}y'^2$$

Hence the quadratic equation becomes

$$-\frac{7}{2}x'^2 + \frac{3}{2}y'^2 - \frac{1}{2} = 0 \iff -7x'^2 + 3y'^2 = 1$$

which is the equation of a hyperbola intersecting the y' -axis at $y' = \pm 1/\sqrt{3}$, as illustrated in the margin. ■



Example 4.2.22. Identify the conic section whose equation is

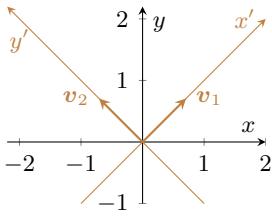
$$x^2 - xy + y^2 + \frac{5}{2\sqrt{2}}x - \frac{7}{2\sqrt{2}}y + \frac{1}{8} = 0.$$

Solution: When there are both the cross-product xy and linear terms, it is easier to first rotate coordinates, and second shift coordinates.

- (a) Rewrite the quadratic terms using vector $\mathbf{x} = (x, y)$ and splitting the cross-product into two equal halves:

$$\begin{aligned} x^2 - xy + y^2 &= x(x - \frac{1}{2}y) + y(-\frac{1}{2}x + y) \\ &= [x \ y] \begin{bmatrix} x - \frac{1}{2}y \\ -\frac{1}{2}x + y \end{bmatrix} \\ &= \mathbf{x}^T \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \\ &= \mathbf{x}^T A \mathbf{x} \quad \text{for matrix } A = \begin{bmatrix} 1 & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{bmatrix}. \end{aligned}$$

Recall that Example 4.1.19 found the eigenvalues of this matrix are $\lambda = \frac{1}{2}, \frac{3}{2}$ with corresponding orthonormal eigenvectors $\mathbf{v}_1 = (1, 1)/\sqrt{2}$ and $\mathbf{v}_2 = (-1, 1)/\sqrt{2}$, respectively. Let's change to a new (dashed) coordinate system (x', y') with \mathbf{v}_1 and \mathbf{v}_2 as its standard unit vectors (as illustrated in the margin). Then throughout the 2D-plane every vector/position



$$\mathbf{x} = \mathbf{v}_1 x' + \mathbf{v}_2 y' = [\mathbf{v}_1 \ \mathbf{v}_2] \begin{bmatrix} x' \\ y' \end{bmatrix} = V \mathbf{x}'$$

for orthogonal matrix

$$V = [\mathbf{v}_1 \ \mathbf{v}_2] = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}.$$

In the new coordinates:

- the quadratic terms

$$\begin{aligned} x^2 - xy + y^2 &= \mathbf{x}^T A \mathbf{x} \\ &= (V \mathbf{x}')^T A (V \mathbf{x}') \\ &= \mathbf{x}'^T V^T A V \mathbf{x}' \\ &= \mathbf{x}'^T \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{3}{2} \end{bmatrix} \mathbf{x}' \quad (\text{as } V^T A V = D) \\ &= \frac{1}{2} x'^2 + \frac{3}{2} y'^2; \end{aligned}$$

- whereas the linear terms

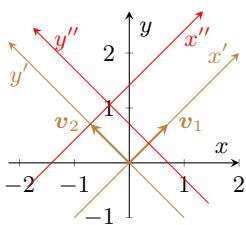
$$\begin{aligned} \frac{5}{2\sqrt{2}}x - \frac{7}{2\sqrt{2}}y &= \left[\frac{5}{2\sqrt{2}} \ -\frac{7}{2\sqrt{2}} \right] \mathbf{x} \\ &= \left[\frac{5}{2\sqrt{2}} \ -\frac{7}{2\sqrt{2}} \right] V \mathbf{x}' \\ &= \left[-\frac{1}{2} \ -3 \right] \mathbf{x}' \\ &= -\frac{1}{2}x' - 3y'; \end{aligned}$$

- so the quadratic equation transforms to

$$\frac{1}{2}x'^2 + \frac{3}{2}y'^2 - \frac{1}{2}x' - 3y' + \frac{1}{8} = 0.$$

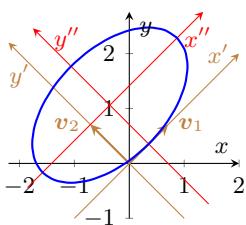
- (b) The second step is to shift coordinates via completing the squares:

$$\begin{aligned} & \frac{1}{2}x'^2 + \frac{3}{2}y'^2 - \frac{1}{2}x' - 3y' + \frac{1}{8} = 0 \\ \iff & \frac{1}{2}(x'^2 - x') + \frac{3}{2}(y'^2 - 2y') = -\frac{1}{8} \\ \iff & \frac{1}{2}(x'^2 - x' + \frac{1}{4}) + \frac{3}{2}(y'^2 - 2y' + 1) = -\frac{1}{8} + \frac{1}{8} + \frac{3}{2} \\ \iff & \frac{1}{2}(x' - \frac{1}{2})^2 + \frac{3}{2}(y' - 1)^2 = \frac{3}{2} \end{aligned}$$



Thus let's change to a new (double dashed) coordinate system $x'' = x' - \frac{1}{2}$ and $y'' = y' - 1$ (equivalently, choose the origin of a new coordinate system to be at $(x', y') = (\frac{1}{2}, 1)$ as illustrated in the margin). In this new coordinate system the quadratic equation becomes

$$\frac{1}{2}x''^2 + \frac{3}{2}y''^2 = \frac{3}{2}, \quad \text{that is } \frac{x''^2}{3} + \frac{y''^2}{1} = 1.$$



In this new coordinate system the equation is that of an ellipse with x'' -axis of half-length $\sqrt{3}$ and y'' -axis of half-length 1 (as illustrated in the margin). ■

Simplify quadratic forms

Definition 4.2.23. A *quadratic form* in variables $\mathbf{x} \in \mathbb{R}^n$ is a function $q : \mathbb{R}^n \rightarrow \mathbb{R}$ that may be written as $q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$ for some real symmetric $n \times n$ matrix A .

Example 4.2.24. (a) The dot product in \mathbb{R}^n is a quadratic form. For all $\mathbf{x} \in \mathbb{R}^n$ consider

$$\mathbf{x} \cdot \mathbf{x} = \mathbf{x}^T \mathbf{x} = \mathbf{x}^T I_n \mathbf{x},$$

which is the quadratic form associated with the identity matrix I_n . ■

Theorem 4.2.25 (principal axes theorem). For every quadratic form, there exists an orthogonal coordinate system that diagonalises the quadratic form. Specifically, for the quadratic form $\mathbf{x}^T A \mathbf{x}$ find the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ and orthonormal eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ of symmetric A , and then in the new coordinate system (y_1, y_2, \dots, y_n) with unit vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ the quadratic form has the **canonical form** $\mathbf{x}^T A \mathbf{x} = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \dots + \lambda_n y_n^2$.

Proof. In the new coordinate system (y_1, y_2, \dots, y_n) the orthonormal vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ (called the **principal axes**) act as the standard unit vectors. Hence any vector $\mathbf{x} \in \mathbb{R}^n$ may be written as a linear combination

$$\mathbf{x} = y_1 \mathbf{v}_1 + y_2 \mathbf{v}_2 + \cdots + y_n \mathbf{v}_n = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = V\mathbf{y}$$

for orthogonal matrix $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$ and vector $\mathbf{y} = (y_1, y_2, \dots, y_n)$. Then the quadratic form

$$\mathbf{x}^\top A \mathbf{x} = (V\mathbf{y})^\top A (V\mathbf{y}) = \mathbf{y}^\top V^\top A V \mathbf{y} = \mathbf{y}^\top D \mathbf{y},$$

since $V^\top A V = D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ by Theorem 4.2.18. Consequently,

$$\mathbf{x}^\top A \mathbf{x} = \mathbf{y}^\top D \mathbf{y} = \lambda_1 y_1^2 + \lambda_2 y_2^2 + \cdots + \lambda_n y_n^2.$$

□

Theorem 4.2.26 (extreme values). *Let A be an $n \times n$ symmetric matrix with eigenvalues $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$ (sorted). Then for all unit vectors $\mathbf{x} \in \mathbb{R}^n$ (that is, $|\mathbf{x}| = 1$), the quadratic form $\mathbf{x}^\top A \mathbf{x}$ has the following properties:*

- (a) $\lambda_1 \leq \mathbf{x}^\top A \mathbf{x} \leq \lambda_n$;
- (b) the minimum of $\mathbf{x}^\top A \mathbf{x}$ is λ_1 , and occurs when \mathbf{x} is a (unit) eigenvector corresponding to λ_1 ;
- (c) the maximum of $\mathbf{x}^\top A \mathbf{x}$ is λ_n , and occurs when \mathbf{x} is a (unit) eigenvector corresponding to λ_n .

Proof. Change to an orthogonal coordinate system that diagonalises the matrix A (Theorem 4.2.25): say $\mathbf{y} = V\mathbf{x}$ for orthogonal matrix V whose columns are orthogonal eigenvectors of A in order so that $D = V^\top A V = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. Then the quadratic form

$$\mathbf{x}^\top A \mathbf{x} = (V\mathbf{y})^\top A (V\mathbf{y}) = \mathbf{y}^\top V^\top A V \mathbf{y} = \mathbf{y}^\top D \mathbf{y}.$$

Since V is orthogonal it preserves lengths (Theorem 3.2.39f) so the unit vector condition $|\mathbf{x}| = 1$ is the same as $|\mathbf{y}| = 1$.

1. To prove the lower bound, consider

$$\begin{aligned} \mathbf{x}^\top A \mathbf{x} &= \mathbf{y}^\top D \mathbf{y} \\ &= \lambda_1 y_1^2 + \lambda_2 y_2^2 + \cdots + \lambda_n y_n^2 \\ &= \lambda_1 y_1^2 + \lambda_1 y_2^2 + \cdots + \lambda_1 y_n^2 \\ &\quad + \underbrace{(\lambda_2 - \lambda_1)y_2^2}_{\geq 0} + \cdots + \underbrace{(\lambda_n - \lambda_1)y_n^2}_{\geq 0} \end{aligned}$$

$$\begin{aligned}
&\geq \lambda_1 y_1^2 + \lambda_1 y_2^2 + \cdots + \lambda_1 y_n^2 \\
&= \lambda_1(y_1^2 + y_2^2 + \cdots + y_n^2) \\
&= \lambda_1 |\mathbf{y}|^2 \\
&= \lambda_1 .
\end{aligned}$$

Similarly for the upper bound (Exercise 4.2.22). Thus $\lambda_1 \leq \mathbf{x}^\top A \mathbf{x} \leq \lambda_n$ for all unit vectors \mathbf{x} .

2. Let \mathbf{v}_1 be a unit eigenvector of A corresponding to the smallest eigenvalue λ_1 ; that is, $A\mathbf{v}_1 = \lambda_1 \mathbf{v}_1$ and $|\mathbf{v}_1| = 1$. Then, setting $\mathbf{x} = \mathbf{v}_1$, the quadratic form

$$\mathbf{x}^\top A \mathbf{x} = \mathbf{v}_1^\top A \mathbf{v}_1 = \mathbf{v}_1^\top \lambda_1 \mathbf{v}_1 = \lambda_1 (\mathbf{v}_1^\top \mathbf{v}_1) = \lambda_1 |\mathbf{v}_1|^2 = \lambda_1 .$$

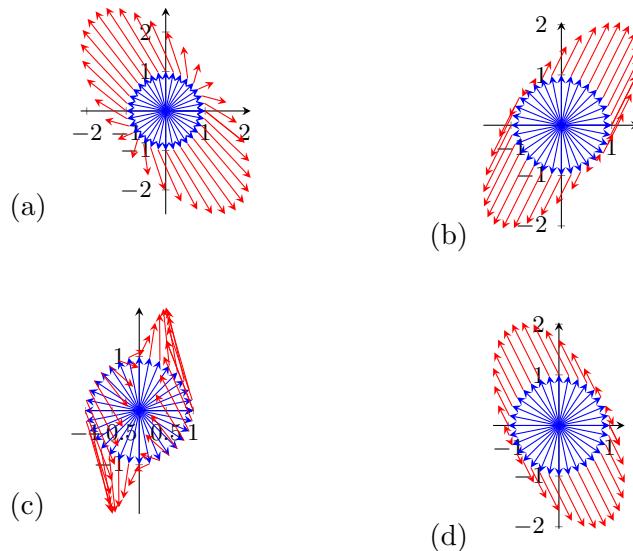
Thus the quadratic form $\mathbf{x}^\top A \mathbf{x}$ takes on the minimum value λ_1 and it occurs when $\mathbf{x} = \mathbf{v}_1$ (at least).

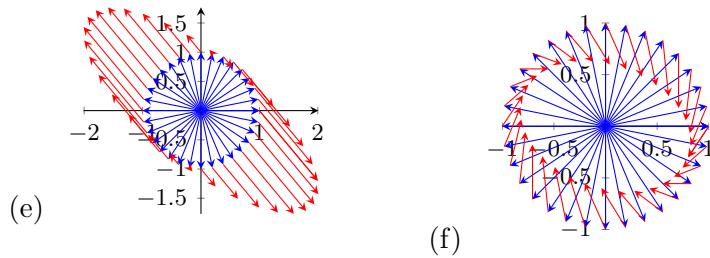
3. Exercise 4.2.22 proves the maximum value occurs.

□

4.2.4 Exercises

Exercise 4.2.1. Each plot below shows (unit) vectors \mathbf{x} (blue), and for some 2×2 matrix A the corresponding vectors $A\mathbf{x}$ (red) adjoined. By assessing whether there are any zero eigenvalues, estimate if the matrix A is invertible or not.





Exercise 4.2.2. For each of the following symmetric matrices: from a hand derivation of the characteristic equation (defined in Procedure 4.1.18), determine whether each matrix has a zero eigenvalue or not, and hence determine whether it is invertible or not.

$$(a) \begin{bmatrix} -1/2 & -3/4 \\ -3/4 & -1/2 \end{bmatrix}$$

$$(b) \begin{bmatrix} 4 & -2 \\ -2 & 1 \end{bmatrix}$$

$$(c) \begin{bmatrix} 0 & -2/5 \\ -2/5 & 3/5 \end{bmatrix}$$

$$(d) \begin{bmatrix} 2 & 1 & -2 \\ 1 & 3 & -1 \\ -2 & -1 & 2 \end{bmatrix}$$

$$(e) \begin{bmatrix} -2 & -1 & 1 \\ -1 & -0 & -1 \\ 1 & -1 & -2 \end{bmatrix}$$

$$(f) \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}$$

$$(g) \begin{bmatrix} -1/2 & 3/2 & 1 \\ 3/2 & -3 & -3/2 \\ 1 & -3/2 & -1/2 \end{bmatrix}$$

$$(h) \begin{bmatrix} 1 & -1 & -1 \\ -1 & -1/2 & 1/2 \\ -1 & 1/2 & -1/2 \end{bmatrix}$$

Exercise 4.2.3. For each of the following matrices, find by hand the eigenvalues and eigenvectors. Using these eigenvectors, confirm that the eigenvalues of the matrix squared are the square of its eigenvalues. If the matrix has an inverse, what are the eigenvalues of the inverse?

$$(a) \ A = \begin{bmatrix} 0 & -2 \\ -2 & 3 \end{bmatrix}$$

$$(b) \ B = \begin{bmatrix} 5/2 & -2 \\ -2 & 5/2 \end{bmatrix}$$

$$(c) \quad C = \begin{bmatrix} 3 & 8 \\ 8 & -9 \end{bmatrix}$$

$$(d) \ D = \begin{bmatrix} -2 & 1 \\ 1 & 14/5 \end{bmatrix}$$

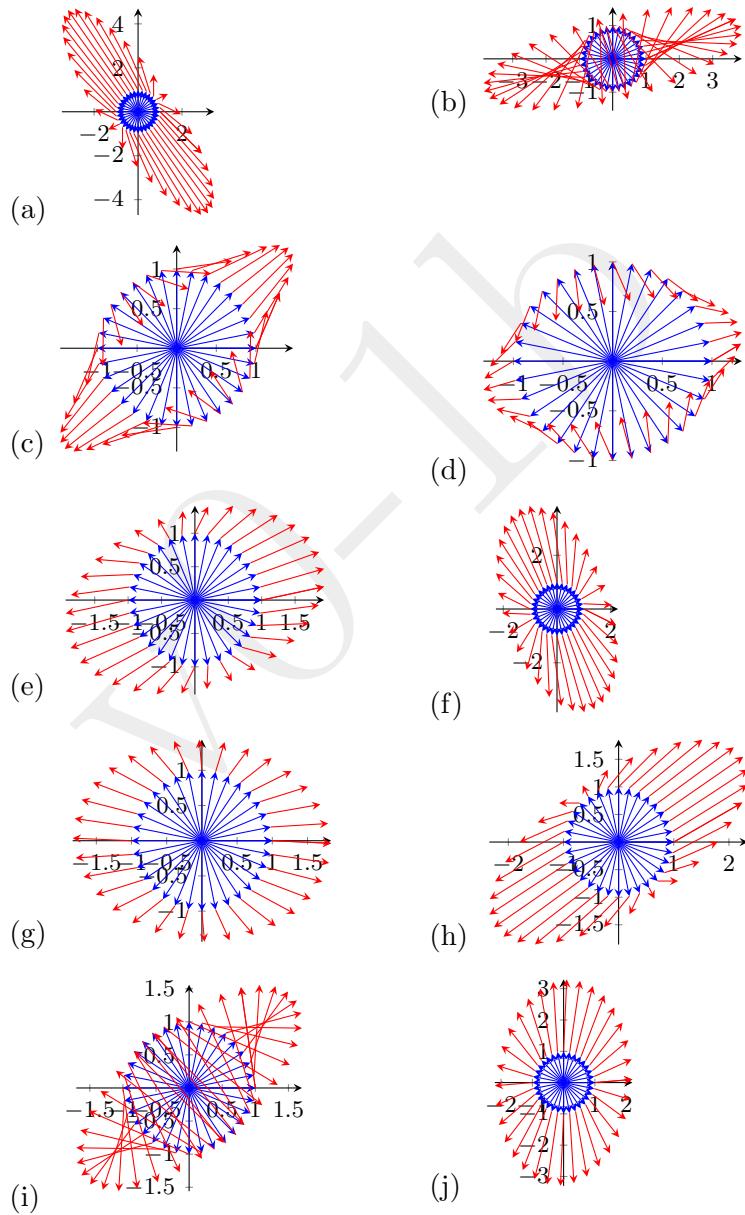
$$(e) \quad E = \begin{bmatrix} -1 & -2 & 0 \\ -2 & 0 & 2 \\ 0 & 2 & 1 \end{bmatrix}$$

$$(f) \quad F = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 0 & -1 \\ 3 & -1 & 2 \end{bmatrix}$$

$$(g) \quad G = \begin{bmatrix} 0 & -1 & -1 \\ -1 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}$$

$$(h) \quad H = \begin{bmatrix} -1 & 3/2 & 3/2 \\ 3/2 & -3 & -1/2 \\ 3/2 & -1/2 & 3 \end{bmatrix}$$

Exercise 4.2.4. Each plot below shows (unit) vectors \mathbf{x} (blue), and for some 2×2 matrix A the corresponding vectors $A\mathbf{x}$ (red) adjoined. For each plot of a matrix A there is a companion plot of the inverse matrix A^{-1} . By roughly estimating eigenvalues and eigenvectors by eye, identify the pairs of plots corresponding to each matrix and its inverse.



Exercise 4.2.5. For the symmetric matrices of Exercise 4.2.3, confirm that eigenvectors corresponding to distinct eigenvalues are orthogonal (Theorem 4.2.10). Show your working.

Exercise 4.2.6. For an $n \times n$ symmetric matrix,

- eigenvectors corresponding to different eigenvalues are orthogonal (Theorem 4.2.10), and

- there are generally n eigenvalues.

Which of the illustrated 2D examples of Exercise 4.1.1 appear to come from symmetric matrices, and which appear to come from non-symmetric matrices?

Exercise 4.2.7. For each of the following *non-symmetric* matrices, confirm that eigenvectors corresponding to distinct eigenvalues are *not* orthogonal. Show and comment on your working. Find eigenvectors by hand for 2×2 and 3×3 matrices, and compute with Matlab/Octave for 3×3 matrices and larger (using `eig()` and `V'*V`).

$$(a) A = \begin{bmatrix} -2 & 2 \\ 3 & -3 \end{bmatrix} \quad (b) B = \begin{bmatrix} 1 & -3 \\ -2 & 2 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} 1 & -2 & 6 \\ 4 & 3 & 2 \\ -3 & 1 & -8 \end{bmatrix} \quad (d) D = \begin{bmatrix} -1 & -2 & 9 \\ 3 & -6 & 3 \\ 0 & 0 & 3 \end{bmatrix}$$

$$(e) E = \begin{bmatrix} 2 & 1 & -2 & 1 \\ 1 & -3 & 4 & -4 \\ 3 & 2 & 4 & -5 \\ -3 & -1 & -3 & 0 \end{bmatrix} \quad (f) F = \begin{bmatrix} -4 & 0 & -5 & 7 \\ 2 & 3 & 1 & -3 \\ 1 & 4 & -2 & 4 \\ 1 & -3 & 4 & 2 \end{bmatrix}$$

$$(g) G = \begin{bmatrix} 1 & 4 & 3 & 1 & -1 \\ 5 & 1 & 6 & -0 & 1 \\ 0 & 4 & -3 & 1 & 4 \\ -3 & 2 & 1 & 4 & -1 \\ -2 & 2 & 2 & 2 & 1 \end{bmatrix} \quad (h) H = \begin{bmatrix} 2 & 0 & 2 & -1 & -1 \\ 2 & 1 & -1 & 2 & -0 \\ 5 & -2 & 6 & 2 & -1 \\ 4 & 0 & -2 & 6 & -5 \\ 2 & 0 & -5 & -3 & -5 \end{bmatrix}$$

Exercise 4.2.8. For the symmetric matrices of Exercise 4.2.3, use Matlab/Octave to compute an SVD (USV^T) of each matrix. Confirm that each column of V is an eigenvector of the matrix (that is, proportional to what the exercise found) and the corresponding singular value is the magnitude of the corresponding eigenvalue (Theorem 4.2.15). Show and discuss your working.

Exercise 4.2.9. To complement the previous exercise, for each of the *non-symmetric* matrices of Exercise 4.2.7, use Matlab/Octave to compute an SVD (USV^T) of each matrix. Confirm that each column of V is *not* an eigenvector of the matrix, and the singular values do *not* appear closely related to the eigenvalues. Show and discuss your working.

Exercise 4.2.10. Let A be an $m \times n$ matrix with SVD $A = USV^T$. Prove that for any $j = 1, 2, \dots, n$, the j th column of V , \mathbf{v}_j , is an eigenvector of the $n \times n$ symmetric matrix $A^T A$ corresponding to the eigenvalue $\lambda_j = \sigma_j^2$ (or $\lambda_j = 0$ if $m < j \leq n$).

Exercise 4.2.11. Prove Theorem part 4.2.4c using parts 4.2.4a and 4.2.4b: that if matrix A is invertible, then for any integer n , λ^n is an eigenvalue of A^n with corresponding eigenvector \mathbf{x} .

Exercise 4.2.12. For each of the following matrices, give reasons as to whether the matrix is orthogonally diagonalisable, and if it is then find an orthogonal matrix V that does so and the corresponding diagonal matrix D . Use Matlab/Octave for the larger matrices.

$$(a) A = \begin{bmatrix} 0 & -2/3 \\ -1 & -1/3 \end{bmatrix}$$

$$(b) B = \begin{bmatrix} 2/3 & 2 \\ 2 & 7/3 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 2 & -1 \\ 1 & -1 & 3 \end{bmatrix}$$

$$(d) D = \begin{bmatrix} 2 & 0 & 2 \\ 1 & -1 & 0 \\ -1 & 0 & -1 \end{bmatrix}$$

$$(e) E = \begin{bmatrix} 0 & -2 & -2 & 0 \\ -2 & 0 & -2 & 0 \\ -2 & -2 & 2 & 2 \\ 0 & 0 & 2 & -2 \end{bmatrix}$$

$$(f) F = \begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & 3 & -2 & -2 \\ 0 & -2 & 0 & 1 \\ 0 & -2 & 1 & 0 \end{bmatrix}$$

$$(g) G = \begin{bmatrix} 3 & 1 & 1 & 1 & -1 \\ -1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 1 & -1 & 1 & -1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$(h) H = \begin{bmatrix} 3 & 0 & 2 & 0 & 0 \\ 0 & 3 & 1 & 0 & -1 \\ 2 & 1 & 1 & 0 & -1 \\ 0 & 0 & 0 & 3 & 2 \\ 0 & -1 & -1 & 2 & 1 \end{bmatrix}$$

Exercise 4.2.13. Let matrix A be invertible and orthogonally diagonalisable. Show that the inverse A^{-1} is orthogonally diagonalisable.

Exercise 4.2.14. Suppose matrices A and B are orthogonally diagonalisable by the same orthogonal matrix V . Show that $AB = BA$ and that the product AB is orthogonally diagonalisable.

Exercise 4.2.15. For each of the given symmetric matrices, say A , find a symmetric matrix X such that $X^2 = A$. That is, find a square-root of the matrix.

$$(a) A = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix}$$

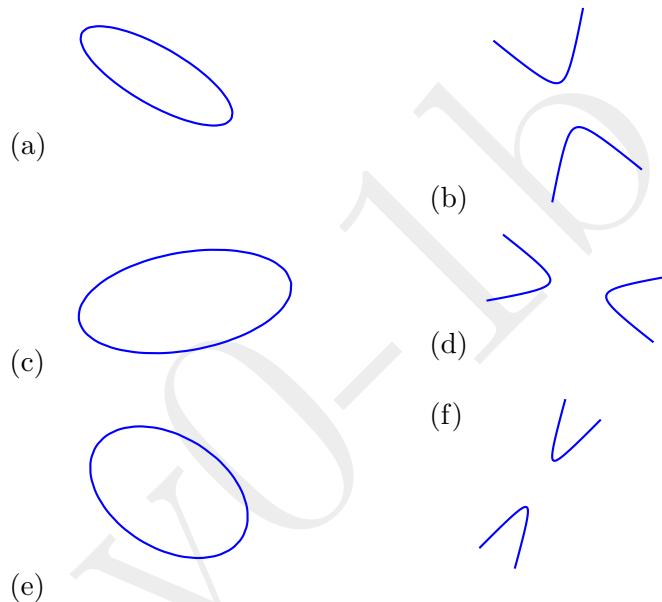
$$(b) B = \begin{bmatrix} 6 & -5 & -5 \\ -5 & 10 & 1 \\ -5 & 1 & 10 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} 2 & 1 & 3 \\ 1 & 2 & 3 \\ 3 & 3 & 6 \end{bmatrix}$$

(d) How many possible answers are there for each of the given matrices? Why?

- (e) For all symmetric matrices A , show that if every eigenvalue of A is non-negative, then there exists a symmetric matrix X such that $A = X^2$.
- (f) Continuing the previous part, how many such matrices X exist? Justify your answer.

Exercise 4.2.16. For each of the following conic sections, draw a pair of coordinate axes for a coordinate system in which the algebraic description of the curves should be simplest.



Exercise 4.2.17. By shifting to a new coordinate axes, find the canonical form of each of the following quadratic equations, and hence describe each curve.

- (a) $-4x^2 + 5y^2 + 4x + 4y - 1 = 0$
- (b) $-4y^2 - 6x - 4y + 2 = 0$
- (c) $5x^2 - y^2 - 2y + 4 = 0$
- (d) $3x^2 + 5y^2 + 6x + 4y + 2 = 0$
- (e) $-2x^2 - 4y^2 + 7x - 2y + 1 = 0$
- (f) $-9x^2 - y^2 - 3y + 6 = 0$
- (g) $-x^2 - 4x - 8y + 4 = 0$
- (h) $-8x^2 - y^2 - 2x + 2y - 3 = 0$

Exercise 4.2.18. By rotating to new coordinate axes, identify each of the following conic sections. Write the quadratic terms in the form $\mathbf{x}^T A \mathbf{x}$, and use the eigenvalues and eigenvectors of matrix A .

- (a) $4xy + 3y^2 - 3 = 0$
- (b) $3x^2 + 8xy - 3y^2 = 0$

(c) $2x^2 - 3xy + 6y^2 - 5 = 0$ (d) $-4x^2 + 3xy - 4y^2 - 2 = 0$

(e) $-4x^2 + 3xy - 4y^2 + 11 = 0$ (f) $3x^2 - 2xy + 3y^2 - 6 = 0$

(g) $-x^2 + 2xy - y^2 + 5 = 0$ (h) $-x^2 - 4xy + 2y^2 - 6 = 0$

Exercise 4.2.19. By rotating and shifting to new coordinate axes, identify each of the following conic sections from its equation.

(a) $-2x^2 - 5xy - 2y^2 - \frac{33}{2}x - 15y - 32 = 0$

(b) $-7x^2 + 3xy - 3y^2 - 52x + \frac{33}{2}y - \frac{381}{4} = 0$

(c) $-4xy - 3y^2 - 18x - 11y + \frac{37}{4} = 0$

(d) $2x^2 - y^2 + 10x + 6y + \frac{11}{2} = 0$

(e) $-2x^2 + 5xy - 2y^2 - \frac{13}{2}x + \frac{5}{2}y + \frac{31}{4} = 0$

(f) $-4xy + 3y^2 + 18x - 13y + \frac{3}{4} = 0$

(g) $6x^2 + 6y^2 - 12x - 42y + \frac{155}{2} = 0$

(h) $2x^2 - 4xy + 5y^2 + 34x - 52y + \frac{335}{2} = 0$

Exercise 4.2.20. For each of the following matrices, say A , consider the quadratic form $q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$. Find coordinate axes, the principal axes, such that the quadratic has the canonical form in the new coordinates y_1, y_2, \dots, y_n . Use eigenvalues and eigenvectors, and use Matlab/Octave for the larger matrices. Over all unit vectors, what is the maximum value of $q(\mathbf{x})$? and what is the minimum value?

(a) $A = \begin{bmatrix} 0 & -5 \\ -5 & 0 \end{bmatrix}$ (b) $B = \begin{bmatrix} 3 & 2 \\ 2 & 0 \end{bmatrix}$

(c) $C = \begin{bmatrix} -1 & 2 & -2 \\ 2 & 0 & 0 \\ -2 & 0 & -2 \end{bmatrix}$ (d) $D = \begin{bmatrix} 2 & 1 & -1 \\ 1 & 3 & 2 \\ -1 & 2 & 3 \end{bmatrix}$

(e) $E = \begin{bmatrix} 6 & 0 & -3 & 1 \\ 0 & -1 & 5 & -3 \\ -3 & 5 & -4 & -7 \\ 1 & -3 & -7 & 0 \end{bmatrix}$ (f) $F = \begin{bmatrix} -5 & -1 & 1 & 2 \\ -1 & 2 & 4 & 1 \\ 1 & 4 & -7 & 7 \\ 2 & 1 & 7 & 2 \end{bmatrix}$

(g) $G = \begin{bmatrix} 1 & 1 & -1 & -6 & -7 \\ 1 & 0 & 3 & 3 & -6 \\ -1 & 3 & 12 & 4 & -4 \\ -6 & 3 & 4 & 1 & -2 \\ -7 & -6 & -4 & -2 & 3 \end{bmatrix}$ (h) $H = \begin{bmatrix} 12 & -3 & -3 & -4 & -6 \\ -3 & 0 & 0 & 2 & -5 \\ -3 & 0 & -4 & 1 & -3 \\ -4 & 2 & 1 & -5 & 2 \\ -6 & -5 & -3 & 2 & 0 \end{bmatrix}$

Exercise 4.2.21. For any given $n \times n$ symmetric matrix A consider the quadratic form $q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}$. For general vectors \mathbf{x} , not necessarily unit vectors, what is the maximum value of $q(\mathbf{x})$ in terms of $|\mathbf{x}|$? and what is the minimum value of $q(\mathbf{x})$? Justify your answer.

Exercise 4.2.22. Complete the proof of Theorem 4.2.26 by detailing the proof of the upper bound, and that the upper bound is achieved for an appropriate unit eigenvector.

Exercise 4.2.23. For a symmetric matrix, discuss the similarities and differences between the SVD and the diagonalisation factorisation, the singular values and the eigenvalues, and the singular vectors and the eigenvectors.

VO-1P

Answers to selected exercises

- 4.1.1b : $\mathbf{v}_1 \propto \pm(-0.5, 0.9)$, $\lambda_1 \approx -0.3$; and $\mathbf{v}_2 \propto \pm(0.6, 0.8)$, $\lambda_2 \approx 1.1$.
- 4.1.1d : $\mathbf{v}_1 \propto \pm(1, -0.2)$, $\lambda_1 \approx -0.7$; and $\mathbf{v}_2 \propto \pm(-0.2, 1)$, $\lambda_2 \approx 1.1$.
- 4.1.1f : $\mathbf{v}_1 \propto \pm(0.5, 0.9)$, $\lambda_1 \approx -0.3$; and $\mathbf{v}_2 \propto \pm(-0.9, 0.5)$, $\lambda_2 \approx 0.8$.
- 4.1.1h : $\mathbf{v}_1 \propto \pm(0, 1)$, $\lambda_1 \approx -0.6$; and $\mathbf{v}_2 \propto \pm(0.9, 0.5)$, $\lambda_2 \approx 1.1$.
- 4.1.2a : Corresponding eigenvalues are: 7, 0, n/a, 7, n/a, n/a
- 4.1.2c : Corresponding eigenvalues are: 2, -5, -4, n/a, -4, n/a
- 4.1.2e : Corresponding eigenvalues are: -2, 0, n/a, 1, 1, n/a
- 4.1.3a : Eigenvalues $-3, -5, 2, 5$.
- 4.1.3c : Eigenvalues $2, 9$, eigenspace \mathbb{E}_2 is 3D.
- 4.1.3e : Eigenvalues $2, -13, 15, 0$.
- 4.1.3g : Eigenvalues $-5.1461, -1.6639, -0.7427, 0.7676, 7.7851$.
- 4.1.4a : $-1, 5$
- 4.1.4c : $-6, -1$
- 4.1.4e : $-3, 7$
- 4.1.4g : $1, 6, 11$
- 4.1.4i : $-4, 5$ (twice)
- 4.1.4k : $-2, 0, 10$
- 4.1.23a : $\mathbb{E}_4 = \text{span}\{(-1, 1)\}$, $\mathbb{E}_{-2} = \text{span}\{(1, 1)\}$
- 4.1.23c : $\mathbb{E}_{-8} = \text{span}\{(2, -2, -1)\}$, $\mathbb{E}_{-7} = \text{span}\{(-1, -1, 0)\}$, $\mathbb{E}_1 = \text{span}\{(1, -1, 4)\}$
- 4.1.23e : $\mathbb{E}_{-4} = \text{span}\{(5, 4, -2)\}$, $\mathbb{E}_1 = \text{span}\{(0, 1, 2)\}$, 2 is not an eigenvalue
- 4.1.24a : $\mathbb{E}_{-9} = \text{span}\{(3, -1)\}$, $\mathbb{E}_1 = \text{span}\{(1, 3)\}$. Both eigenvalues have multiplicity one.
- 4.1.24c : $\mathbb{E}_{-6} = \text{span}\{(1, 2)\}$, $\mathbb{E}_{-1} = \text{span}\{(-2, 1)\}$. Both eigenvalues have multiplicity one.
- 4.1.24e : $\mathbb{E}_{-7} = \text{span}\{(1, 3, -1)\}$, $\mathbb{E}_{-4} = \text{span}\{(1, 0, 1)\}$, $\mathbb{E}_4 = \text{span}\{(-3, 2, 3)\}$. All eigenvalues have multiplicity one.
- 4.1.24g : $\mathbb{E}_1 = \text{span}\{(-1, 2, -1)\}$, $\mathbb{E}_{13} = \text{span}\{(-1, 2, 5), (2, 1, 0)\}$. Eigenvalue $\lambda = 1$ has multiplicity one, whereas $\lambda = 13$ has multiplicity two.
- 4.1.24i : $\mathbb{E}_{-5} = \text{span}\{(5, -2, 7)\}$, $\mathbb{E}_1 = \text{span}\{(1, -1, -1)\}$, $\mathbb{E}_{21} = \text{span}\{(3, 4, -1)\}$. All eigenvalues have multiplicity one.

4.1.5 : Number the nodes as you like, say 1=top, 2=centre, 3=right,

and 4=bottom. Then the matrix is $\begin{bmatrix} 0 & 2 & 1 & 0 \\ 2 & 0 & 1 & 2 \\ 1 & 1 & 0 & 1 \\ 0 & 2 & 1 & 0 \end{bmatrix}$. Eigenvalues

are $-2.86, -0.77, 0.00, 3.63$ (2 d.p.) and an eigenvector corresponding to the largest 3.63 is $(0.46, 0.63, 0.43, 0.46)$. Thus rank the centre left node as the most important, the right node is least important, and the top and bottom nodes equal second importance.

4.2.1b : not invertible

4.2.1d : not invertible

4.2.1f : invertible

4.2.2b : eigenvalues 0, 5 so not invertible.

4.2.2d : eigenvalues 0, 2, 5 so not invertible.

4.2.2f : eigenvalues 1, 4 so invertible.

4.2.2h : eigenvalues $-1, 2$ so invertible.

4.2.3b : Eigenvalues $1/2, 9/2$, and corresponding eigenvectors proportional to $(1, 1), (-1, 1)$. The inverse has eigenvalues $2, 2/9$.

4.2.3d : Eigenvalues $-11/5, 3$, and corresponding eigenvectors proportional to $(-5, 1), (1, 5)$. The inverse has eigenvalues $-5/11, 1/3$.

4.2.3f : Eigenvalues $-2, 1, 5$, and corresponding eigenvectors proportional to $(-1, 1, 1), (1, 2, -1), (1, 0, 1)$. The inverse has eigenvalues $-1/2, 1, 1/5$.

4.2.3h : Eigenvalues $-4, -1/2, 7/2$, and corresponding eigenvectors proportional to $(-3, 5, 1), (3, 2, -1), (1, 0, 3)$. The inverse has eigenvalues $-1/4, -2, 2/7$.

4.2.7a : Eigenvectors proportional to $(1, 1), (-2, 3)$.

4.2.7c : Eigenvectors proportional to $(-2, 3, 1), (2, -2, -1), (-13, 3, 14)$.

4.2.7e : Eigenvectors proportional to $(-.39, .44, .76, -.28), (.58, -.41, -.68, .18), (.22, -.94, .23, .11), (.21, -.53, .53, .62)$ (2 d.p.).

4.2.7g : Eigenvectors proportional to $(.01, -.63, .78, .04, -.05), (-.59, .49, .21, -.46, -.39), (.57, .74, .33, -.01, .15), (-.46, -.46, .07, .62, .43), (-.52, -.07, .32, -.53, .59)$ (2 d.p.).

4.2.12a : Not symmetric, so not orthogonally diagonalisable.

4.2.12c : $V = \begin{bmatrix} -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & 0 & \frac{2}{\sqrt{6}} \end{bmatrix}$, and $D = \text{diag}(1, 2, 4)$.

4.2.12e : $V = \begin{bmatrix} -.50 & -.41 & .71 & -.29 \\ -.50 & -.41 & -.71 & -.29 \\ -.50 & .00 & .00 & .87 \\ .50 & -.82 & .00 & .29 \end{bmatrix}$ (2 d.p.), and $D = \text{diag}(-4, -2, 2, 4)$.

4.2.12g : Not symmetric, so not orthogonally diagonalisable.

4.2.15a : $X = \begin{bmatrix} 3/2 & 1/2 \\ 1/2 & 3/2 \end{bmatrix}$

4.2.15c : $X = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 2 \end{bmatrix}$

4.2.17a : hyperbola, centred $(1/2, -2/5)$, $-\frac{x'^2}{1/5} + \frac{y'^2}{4/25} = 1$

4.2.17c : hyperbola, centred $(0, -1)$, $-\frac{x'^2}{1} + \frac{y'^2}{5} = 1$

4.2.17e : ellipse, centred $(7/4, -1/4)$, $\frac{x'^2}{59/16} + \frac{y'^2}{59/32} = 1$

4.2.17g : parabola, base $(-2, 1)$, $y' = -\frac{1}{8}x'^2$

4.2.18a : With axes $\mathbf{i}' = (-\frac{2}{\sqrt{5}}, \frac{1}{\sqrt{5}})$ and $\mathbf{j}' = (\frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}})$, $-\frac{x'^2}{3} + \frac{y'^2}{3/4} = 1$ is hyperbola.

4.2.18c : In axes $\mathbf{i}' = (\frac{3}{\sqrt{10}}, \frac{1}{\sqrt{10}})$ and $\mathbf{j}' = (-\frac{1}{\sqrt{10}}, \frac{3}{\sqrt{10}})$, $\frac{x'^2}{10/3} + \frac{y'^2}{10/13} = 1$ is ellipse.

4.2.18e : In axes $\mathbf{i}' = (\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}})$ and $\mathbf{j}' = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$, $\frac{x'^2}{2} + \frac{y'^2}{22/5} = 1$ is ellipse.

4.2.18g : In axes $\mathbf{i}' = (\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}})$ and $\mathbf{j}' = (\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}})$, $\frac{x'^2}{5/2} = 1$ is pair of parallel lines.

4.2.19a : hyperbola centred $(-1, -5/2)$ at angle 45°

4.2.19c : hyperbola centred $(4, -9/2)$ at angle -27°

4.2.19e : hyperbola centred $(3/2, 5/2)$ at angle 45°

4.2.19g : circle centred $(1, 7/2)$

4.2.20a : $q = -5y_1^2 + 5y_2^2$, max= 5, min=-5

4.2.20c : $q = -4y_1^2 - y_2^2 + 2y_3^2$, max= 2, min=-4

4.2.20e : (2 d.p.) $q = -10.20y_1^2 - 3.56y_2^2 + 4.81y_3^2 + 9.95y_4^2$, max= 9.95, min=-10.20

4.2.20g : $q = -8.82y_1^2 - 4.71y_2^2 + 3.04y_3^2 + 10.17y_4^2 + 17.03y_5^2$, max= 17.03,
min=-8.82

VO-1b

5 Approximate matrices

Chapter Contents

5.1	Measure changes to matrices	438
5.1.1	Compress images optimally	438
5.1.2	Relate matrix changes to the SVD	443
5.1.3	Principal component analysis	456
5.1.4	Exercises	475
5.2	Regularise linear equations	482
5.2.1	The SVD illuminates regularisation	483
5.2.2	Tikhonov regularisation	498
5.2.3	Exercises	502

This chapter could be studied any time after Chapter 3 to help the transition to more abstract linear algebra. Useful to time as spaced revision of earlier material on the SVD, rank, orthogonality, and so on.

This chapter develops how concepts associated with length and distance applies to matrices as well as vectors. Further courses on Linear Algebra place these in a unifying framework that also encompasses much you see both in solving differential equations (and integral equations) and in problems involving complex numbers (such as those in electrical engineering or quantum physics).

5.1 Measure changes to matrices

Section Contents

5.1.1	Compress images optimally	438
5.1.2	Relate matrix changes to the SVD	443
5.1.3	Principal component analysis	456
	Application to latent semantic indexing	466
5.1.4	Exercises	475

5.1.1 Compress images optimally

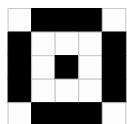
Photographs and other images take a lot of storage. Reducing the amount of storage taken, both for storage and for transmission, is essential. The well-known jpeg format for compressing photographs is incredibly useful: the SVD provides another effective means of compression.

These methods find approximate matrices of the images with the matrices having of various ranks. Recall that a matrix of rank k (Definition 3.3.16) means the matrix has precisely k non-zero singular values, that is, an $m \times n$ matrix

$$\begin{aligned}
 A &= USV^T \\
 &= [\mathbf{u}_1 \ \cdots \ \mathbf{u}_k \ \cdots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_k \\ O_{(m-k) \times k} & & O_{(m-k) \times (n-k)} \end{bmatrix} V^T \\
 &= [\sigma_1 \mathbf{u}_1 \ \cdots \ \sigma_k \mathbf{u}_k \ O_{m \times (n-k)}] \begin{bmatrix} \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_k^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} \\
 &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T.
 \end{aligned}$$

This last form constructs matrix A and *has relatively few components when the rank k is low compared to m and n .*

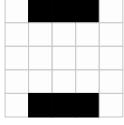
Example 5.1.1. Invent and write down a rank three representation of the 5×5 ‘bulls eye’ matrix (illustrated in the margin)



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

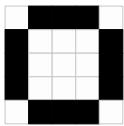
Solution: Here we set all coefficients $\sigma_1 = \sigma_2 = \sigma_3 = 1$ and let the vectors set the magnitude (subsequently, the vectors will be unit vectors and σ_j will set the magnitude).

- (a) Arbitrarily start by addressing together the first and last rows of the image: they can be computed by choosing $\mathbf{u}_1 = (1, 0, 0, 0, 1)$ and $\mathbf{v}_1 = (0, 1, 1, 1, 0)$ as then (as illustrated) a rank one matrix



$$\mathbf{u}_1 \mathbf{v}_1^T = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \end{bmatrix}.$$

- (b) Next choose to address together the first and last columns of the image: they can be computed by choosing $\mathbf{u}_2 = (0, 1, 1, 1, 0)$ and $\mathbf{v}_2 = (1, 0, 0, 0, 1)$ as then (as illustrated) a rank two matrix



$$\begin{aligned} \mathbf{u}_1 \mathbf{v}_1^T + \mathbf{u}_2 \mathbf{v}_2 &= \mathbf{u}_1 \mathbf{v}_1^T + \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 & 1 \end{bmatrix} \\ &= \mathbf{u}_1 \mathbf{v}_1^T + \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 \end{bmatrix}. \end{aligned}$$

- (c) Lastly put the dot in the middle of the image: choose $\mathbf{u}_3 = \mathbf{v}_3 = (0, 0, 1, 0, 0)$ as then (to form the original) a rank three matrix

$$\begin{aligned} &\mathbf{u}_1 \mathbf{v}_1^T + \mathbf{u}_2 \mathbf{v}_2 + \mathbf{u}_3 \mathbf{v}_3 \\ &= \mathbf{u}_1 \mathbf{v}_1^T + \mathbf{u}_2 \mathbf{v}_2 + \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \end{bmatrix} \\ &= \mathbf{u}_1 \mathbf{v}_1^T + \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 \end{bmatrix}.$$

■

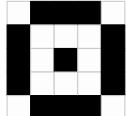
Procedure 5.1.2 (approximate images). *For an image stored as numbers in an $m \times n$ matrix A .*

1. Compute an SVD $A = USV^T$ with $[U, S, V] = \text{svd}(A)$.
2. Choose a desired rank k based upon the singular values (Theorem 5.1.12): typically there will be k ‘large’ singular values and the rest are ‘small’.
3. Then the ‘best’ rank k approximation to the image matrix A is (using the subscript k on the matrix name to denote the rank k approximation)

$$A_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$$

$$= U(:, 1:k) * S(1:k, 1:k) * V(:, 1:k),$$

Example 5.1.3. Use Procedure 5.1.2 to find the ‘best’ rank two and three matrices to approximate the ‘bulls eye’ image matrix (illustrated in the margin)



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

Solution: Enter the matrix into Matlab/Octave and compute an SVD, $A = USV^T$, with $[U, S, V] = \text{svd}(A)$ to find (2 d.p.)

```

U =
-0.47   0.51   0.14   0.71   0.05
-0.35  -0.44   0.43  -0.05   0.71
-0.56  -0.31  -0.77   0.00   0.00
-0.35  -0.44   0.43   0.05  -0.71
-0.47   0.51   0.14  -0.71  -0.05

S =
 2.68      0      0      0      0
      0    2.32      0      0      0
      0      0    0.64      0      0
      0      0      0    0.00      0
      0      0      0      0    0.00

V =
-0.47  -0.51   0.14  -0.68  -0.18
-0.35   0.44   0.43  -0.18   0.68

```



$$\begin{array}{ccccc} -0.56 & 0.31 & -0.77 & 0.00 & -0.00 \\ -0.35 & 0.44 & 0.43 & 0.18 & -0.68 \\ -0.47 & -0.51 & 0.14 & 0.68 & 0.18 \end{array}$$

- For this matrix there are three ‘large’ singular values of 2.68, 2.32 and 0.64, and two ‘small’ singular values of 0.00 (they are precisely zero), thus construct a rank three approximation to the image matrix as

$$A_3 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \sigma_3 \mathbf{u}_3 \mathbf{v}_3^T,$$

computed with `A3=U(:,1:3)*S(1:3,1:3)*V(:,1:3)'`, giving (2 d.p.)

$$\begin{array}{ccccc} A3 = & & & & \\ 0.00 & 1.00 & 1.00 & 1.00 & 0.00 \\ 1.00 & 0.00 & 0.00 & 0.00 & 1.00 \\ 1.00 & 0.00 & 1.00 & 0.00 & 1.00 \\ 1.00 & 0.00 & 0.00 & 0.00 & 1.00 \\ -0.00 & 1.00 & 1.00 & 1.00 & -0.00 \end{array}$$

The rank 3 matrix A_3 exactly reproduces the image matrix A . This exactness is due to the fourth and fifth singular values being precisely zero (to numerical error).

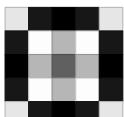
- Alternatively, in the context of some application, we could subjectively decide that there are two ‘large’ singular values of 2.68 and 2.32, and three ‘small’ singular values of 0.64 and 0.00. In such a case, construct a rank two approximation to the image matrix as (illustrated in the margin)

$$A_2 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T,$$

computed with `A2=U(:,1:2)*S(1:2,1:2)*V(:,1:2)'`, giving (2 d.p.)

$$\begin{array}{ccccc} A2 = & & & & \\ -0.01 & 0.96 & 1.07 & 0.96 & -0.01 \\ 0.96 & -0.12 & 0.21 & -0.12 & 0.96 \\ 1.07 & 0.21 & 0.62 & 0.21 & 1.07 \\ 0.96 & -0.12 & 0.21 & -0.12 & 0.96 \\ -0.01 & 0.96 & 1.07 & 0.96 & -0.01 \end{array}$$

This rank two approximation A_2 is indeed roughly the same as the image matrix A , albeit with errors of 20% or so. Subsequent theory confirms that the relative error is characterised by $\sigma_3/\sigma_1 = 0.24$ here.



Euler, 1737, by Vasilij Sokolov

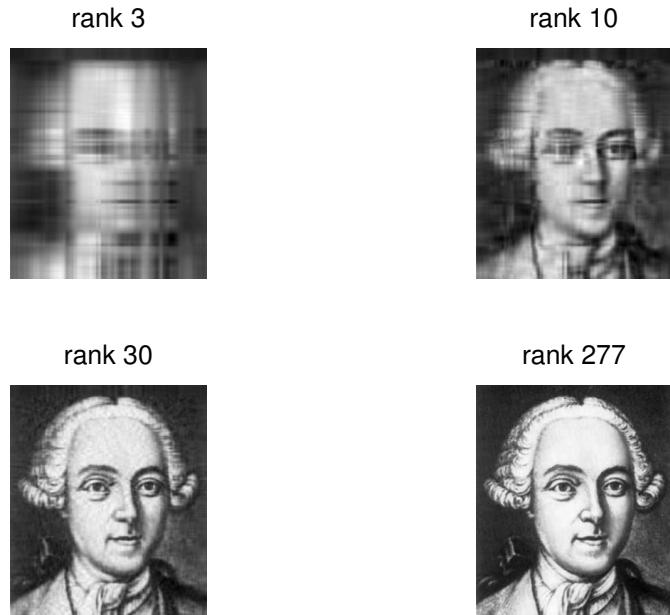


<http://eulerarchive.maa.org/portraits/>

[portraits.html](http://eulerarchive.maa.org/portraits.html) [Sep 2015]

Example 5.1.4. In the margin is a 326×277 greyscale image of Euler at 30 years old. As such the image is coded as 90,302 numbers.

Figure 5.1: four approximate images of Euler ranging from the poor rank 3, via the adequate rank 10, the good rank 30, to the original rank 277.



Let's find a good approximation to the image that uses much fewer numbers, and hence takes less storage. That is, we will effectively compress the image.

Solution: We use an SVD to approximate the image of Euler to controllable levels of approximation. For example, Figure 5.1 shows four approximations to the image of Euler, ranging from the hopeless (labelled “rank 3”) to the original (labelled “rank 277”).

Procedure 5.1.2 is as follows.

- First download the image from the website, and then read the image into Matlab/Octave using

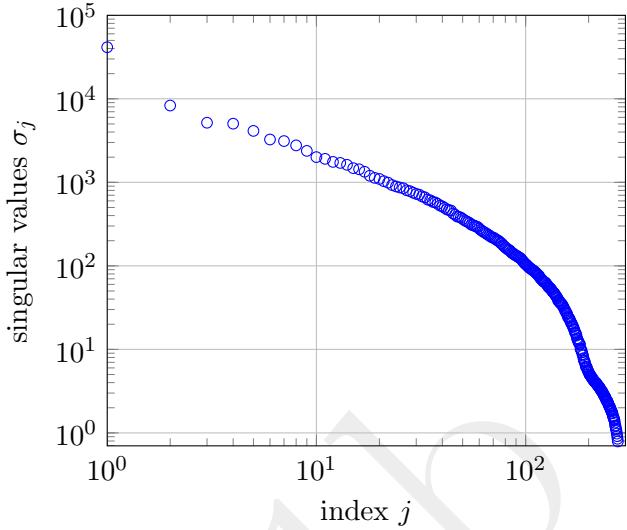
```
rgb=imread('euler1737.png');
A=mean(rgb,3);
```

The `imread` command sets the $326 \times 277 \times 3$ array `rgb` to the red-green-blue values of the image data. Then convert into a grayscale image matrix `A` by averaging the red-green-blue values via the function `mean(rgb,3)` which computes the mean over the third dimension of the array (over the three colours).

- Compute an SVD, $A = USV^T$, of the matrix `A` with the usual command `[U,S,V]=svd(A)`. Here orthogonal `U` is 326×326 , diagonal `S = diag(41 422, 8 309, \dots, 0.79, 0)` is 326×277 , and orthogonal `V` is 277×277 . These matrices are far too big to record in this text.



Figure 5.2: singular values of the image of Euler, 1737.



- Figure 5.2 plots the non-zero singular values from largest to smallest: they cover a range of five orders of magnitude. Choose some number k of singular vectors to use: k is the rank of the approximate images in Figure 5.1. The choice may be guided by the decrease of these singular values: as discussed later, for a say 1% error choose k such that $\sigma_k \approx 0.01\sigma_1$ which from the index j axis of Figure 5.2 is around $k \approx 25$.
- Construct the approximate rank k image

$$\begin{aligned} A_k &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T \\ &= \mathbf{U}(:, 1:k) * \mathbf{S}(1:k, 1:k) * \mathbf{V}(:, 1:k), \end{aligned}$$

- Let's say the rank 30 image of Figure 5.1 is the desired good approximation. To reconstruct it we need 30 singular values $\sigma_1, \sigma_2, \dots, \sigma_{30}$, 30 columns $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_{30}$ of U , 30 columns $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{30}$ of V making a total of

$$30 + 30 \times 326 + 30 \times 277 = 18,120 \text{ numbers.}$$

These 18,120 numbers are much fewer than (one fifth) the $326 \times 277 = 90,302$ numbers of the original image. The SVD provides an effective flexible data compression. ■

5.1.2 Relate matrix changes to the SVD

We need to define what ‘best’ means in the approximation Procedure 5.1.2 and then show the procedure achieves this best. To

Table 5.1: As well as the Matlab/Octave commands and operations listed in Tables 1.2, 2.3, 3.1, 3.2, 3.3, and 3.7 we may invoke these functions.

- `norm(A)` computes the matrix norm of Definition 5.1.5, namely the largest singular value of the matrix A .
Also, and consistent with the matrix norm, recall that `norm(v)` for a vector v computes the length $\sqrt{v_1^2 + v_2^2 + \cdots + v_n^2}$.
 - `scatter(x,y,[] ,c)` draws a 2D scatter plot of points with coordinates in vectors x and y , each point with a colour determined by the corresponding entry of vector c .
Similarly for `scatter3(x,y,z,[] ,c)` but in 3D.
 - `[U,S,V]=svds(A,k)` computes the k largest singular values of the matrix A in the diagonal of $k \times k$ matrix S , and the k columns of U and V are the corresponding singular vectors.
 - `imread('filename')` typically reads an image from a file into an $m \times n \times 3$ array of red-green-blue values. The values are all ‘integers’ in the range $[0, 255]$.
 - `mean(A)` of an $m \times n$ array computes the n elements in the row vector of averages (the arithmetic mean) over each column of A .
Whereas `mean(A,p)` for an ℓ -dimensional array A of dimension $m_1 \times m_2 \times \cdots \times m_\ell$, computes the mean over the p th index to give an array of size $m_1 \times \cdots \times m_{p-1} \times m_{p+1} \times \cdots \times m_\ell$.
 - `std(A)` of an $m \times n$ array computes the n elements in the row vector of the standard deviation over each column of A (close to the root-mean-square from the mean).
 - `csvread('filename')` reads data from a file into a matrix. When each of the m lines in the file is n numbers separated by commas, then the result is an $m \times n$ matrix.
 - `semilogy(x,y,'o')` draws a point plot of y versus x with the vertical axis being logarithmic.
 - `axis` sets some properties of a drawn figure:
 - `axis equal` ensures horizontal and vertical directions are scaled the same—so here there is no distortion of the image;
 - `axis off` means that the horizontal and vertical axes are not drawn—so here the image is unadorned.
-

proceed to understand image compression, we need a measure of the magnitude of matrices and distances between matrices.

In linear algebra we use the double vertical bars, $\|\cdot\|$, to denote the magnitude of a matrix in order to avoid a notational clash with the well-established use of $|\cdot|$ for the determinant of a matrix (Chapter 6).

Definition 5.1.5. *Let A be an $m \times n$ matrix. Define the **matrix norm** (sometimes called the **spectral norm**)*

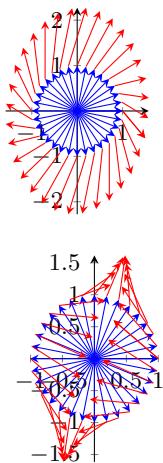
$$\|A\| := \max_{|\mathbf{x}|=1} |A\mathbf{x}|, \quad \text{equivalently } \|A\| = \sigma_1 \quad (5.1)$$

the largest singular value of the matrix A .¹

Proof. The equivalence $\max_{|\mathbf{x}|=1} |A\mathbf{x}| = \sigma_1$ is due to the definition of the largest singular value in the proof of the existence of an SVD (section 3.3.3). \square

Example 5.1.6. The two following 2×2 matrices have the product $A\mathbf{x}$ plotted (red), adjoined to \mathbf{x} (blue), for a complete range of unit vectors \mathbf{x} (as in section 4.1 for eigenvectors). From Definition 5.1.5, the norm of the matrix A is then the length of the longest such plotted $A\mathbf{x}$. For each matrix, use the plot to roughly estimate their norm.

$$(a) A = \begin{bmatrix} 0.5 & 0.5 \\ -0.6 & 1.2 \end{bmatrix}$$



Solution: The longest $A\mathbf{x}$ appear to be near the top and bottom of the plot, and appear to be a little longer than one, so estimate $\|A\| \approx 1.3$.

$$(b) B = \begin{bmatrix} -0.7 & 0.4 \\ 0.6 & 0.5 \end{bmatrix}$$

Solution: Near the top and bottom of the plot, $B\mathbf{x}$ appears to be of length 0.6. But the vectors pointing inwards from the right and left appear longer at about 0.9. So estimate $\|B\| \approx 0.9$. \blacksquare

Example 5.1.7. Consider the 2×2 matrix $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$. Explore products $A\mathbf{x}$ for unit vectors \mathbf{x} , and then find the matrix norm $\|A\|$.

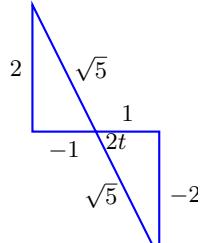
¹ Sometimes this matrix norm is more specifically called a 2-norm and correspondingly denoted by $\|A\|_2$: but not in this book because at other times and places $\|A\|_2$ denotes something slightly different.

- The standard unit vector $\mathbf{e}_2 = (0, 1)$ has $|\mathbf{e}_2| = 1$ and $A\mathbf{e}_2 = (1, 1)$ has length $|A\mathbf{e}_2| = \sqrt{2}$. Since the matrix norm is the maximum of all possible $|A\mathbf{x}|$, so $\|A\| \geq |A\mathbf{e}_2| = \sqrt{2} \approx 1.41$.
- Another unit vector is $\mathbf{x} = (\frac{3}{5}, \frac{4}{5})$. Here $A\mathbf{x} = (\frac{7}{5}, \frac{4}{5})$ has length $\sqrt{49 + 16}/5 = \sqrt{65}/5 \approx 1.61$. Hence the matrix norm $\|A\| \geq |A\mathbf{x}| \approx 1.61$.
- To systematically find the norm, recall all unit vectors in 2D are of the form $\mathbf{x} = (\cos t, \sin t)$. Then

$$\begin{aligned} |A\mathbf{x}|^2 &= |(\cos t, \sin t)|^2 \\ &= (\cos t)^2 + \sin^2 t \\ &= \cos^2 t + 2\cos t \sin t + \sin^2 t + \sin^2 t \\ &= \frac{3}{2} + \sin 2t - \frac{1}{2} \cos 2t. \end{aligned}$$

This length (squared) is maximised (and minimised) for some t determined by calculus. Differentiating with respect to t leads to

$$\frac{d|A\mathbf{x}|^2}{dt} = 2\cos 2t + \sin 2t = 0 \quad \text{for stationary points.}$$



Rearranging determines we require $\tan 2t = -2$. The marginal right-angle triangles illustrate that a stationary point of $|A\mathbf{x}|^2$ occurs for $\sin 2t = \mp 2/\sqrt{5}$ and correspondingly $\cos 2t = \pm 1/\sqrt{5}$ (one gives a minimum and one gives the desired maximum). Substituting these two cases gives

$$\begin{aligned} |A\mathbf{x}|^2 &= \frac{3}{2} + \sin 2t - \frac{1}{2} \cos 2t \\ &= \frac{3}{2} \mp \frac{2}{\sqrt{5}} \mp \frac{1}{2} \frac{1}{\sqrt{5}} \\ &= \frac{1}{2}(3 \mp \sqrt{5}) \\ &= \left(\frac{1 \mp \sqrt{5}}{2}\right)^2. \end{aligned}$$

The plus alternative is the larger so gives the maximum, hence

$$\|A\| = \max_{|\mathbf{x}|=1} |A\mathbf{x}| = \frac{1 + \sqrt{5}}{2} = 1.6180.$$

- Confirm with Matlab/Octave via `svd([1 1;0 1])` which gives the singular values $\sigma_1 = 1.6180$ and $\sigma_2 = 0.6180$. Hence confirming the norm $\|A\| = \sigma_1 = 1.6180$.

Alternatively, see Table 5.1, execute `norm([1 1;0 1])` to compute the norm $\|A\| = 1.6180$.

Example 5.1.8. For larger matrices, Matlab/Octave readily computes the norm either via an SVD or using the `norm` function directly (Table 5.1). Compute the norm of the following matrices.

$$(a) A = \begin{bmatrix} 0.1 & -1.3 & -0.4 & -0.1 & -0.6 \\ 1.9 & 2.4 & -1.8 & 0.2 & 0.8 \\ -0.2 & -0.5 & -0.7 & -2.5 & 1.1 \\ -1.8 & 0.2 & 1.1 & -1.2 & 1.0 \\ -0.0 & 1.2 & 1.1 & -0.1 & 1.7 \end{bmatrix}$$

Solution: Enter the matrix into Matlab/Octave then executing `svd(A)` returns the vector of singular values

$$(4.0175, 3.5044, 2.6568, 0.8571, 0.1618),$$

so $\|A\| = \sigma_1 = 4.0175$. Alternatively, executing `norm(A)` directly gives $\|A\| = 4.0175$.

$$(b) B = \begin{bmatrix} 0 & -2 & -1 & -4 & -5 & 0 \\ 2 & 0 & 1 & -2 & -6 & -2 \\ -2 & 0 & 4 & 2 & 3 & -3 \\ 1 & 2 & -4 & 2 & 1 & 3 \end{bmatrix}$$

Solution: Enter the matrix into Matlab/Octave then executing `svd(B)` returns the vector of singular values

$$(10.1086, 7.6641, 3.2219, 0.8352),$$

so $\|B\| = \sigma_1 = 10.1086$. Alternatively, executing `norm(B)` directly gives $\|B\| = 10.1086$.

■

The Definition 5.1.5 of the magnitude/norm of a matrix may appear a little strange, but, in addition to some marvellously useful properties, it nonetheless has all the familiar properties of a magnitude/length. Recall from Chapter 1 that for vectors:

- $|\mathbf{v}| = 0$ if and only if $\mathbf{v} = \mathbf{0}$ (Theorem 1.1.11);
- $|\mathbf{u} \pm \mathbf{v}| \leq |\mathbf{u}| + |\mathbf{v}|$ (the triangle inequality of Theorem 1.3.13);
- $|t\mathbf{v}| = |t| \cdot |\mathbf{v}|$ (Theorem 1.3.13).

Analogous properties hold for the matrix norm as established in the next theorem.

Theorem 5.1.9 (norm properties). *Let A be an $m \times n$ real matrix:*

- (a) $\|A\| = 0$ if and only if $A = O_{m \times n}$;
- (b) $\|I_n\| = 1$;
- (c) $\|A \pm B\| \leq \|A\| + \|B\|$, for any $m \times n$ matrix B , like a triangle inequality (Theorem 1.3.13c);

- (d) $\|tA\| = |t|\|A\| ;$
- (e) $\|A\| = \|A^T\| ;$
- (f) $\|Q_m A\| = \|A\| = \|AQ_n\|$ for any $m \times m$ orthogonal matrix Q_m and any $n \times n$ orthogonal matrix Q_n ;
- (g) $|Ax| \leq \|A\||x|$ for all $x \in \mathbb{R}^n$, like a Cauchy–Schwarz inequality (Theorem 1.3.13b), as is the following;
- (h) $\|AB\| \leq \|A\|\|B\|$ for any $n \times p$ matrix B .

Proof. Alternative proofs to the following may be invoked (Exercise 5.1.3). Where necessary in the following, let matrix A have the SVD $A = USV^T$.

5.1.9a. If $A = O_{m \times n}$, then from Definition 5.1.5

$$\|A\| = \max_{|x|=1} |Ox| = \max_{|x|=1} |\mathbf{0}| = \max_{|x|=1} 0 = 0 .$$

Conversely, if $\|A\| = 0$, then the largest singular value $\sigma_1 = 0$ (Definition 5.1.5), which implies that all singular values are zero, so the matrix A has an SVD of the form $A = UO_{m \times n}V^T$, which evaluates to $A = O_{m \times n}$.

5.1.9b. From Definition 5.1.5,

$$\|I_n\| = \max_{|x|=1} |Ix| = \max_{|x|=1} |x| = \max_{|x|=1} 1 = 1 .$$

5.1.9c. Using Definition 5.1.5 at the first and last steps:

$$\begin{aligned} \|A \pm B\| &= \max_{|x|=1} |(A \pm B)x| \\ &= \max_{|x|=1} |Ax \pm Bx| \quad (\text{by distributivity}) \\ &\leq \max_{|x|=1} (|Ax| + |Bx|) \quad (\text{by triangle inequality}) \\ &\leq \max_{|x|=1} |Ax| + \max_{|x|=1} |Bx| \\ &= \|A\| + \|B\| . \end{aligned}$$

5.1.9d. Using Definition 5.1.5,

$$\begin{aligned} \|tA\| &= \max_{|x|=1} |(tA)x| \\ &= \max_{|x|=1} |t(Ax)| \quad (\text{by associativity}) \\ &= \max_{|x|=1} |t||Ax| \quad (\text{by Theorem 1.3.13}) \\ &= |t| \max_{|x|=1} |Ax| \\ &= |t|\|A\| . \end{aligned}$$

- 5.1.9e. Recall that then A^T has an SVD $A^T = (USV^T)^T = VS^TU^T$. So the largest singular value of A^T is the same as that of A . Hence $\|A\| = \|A^T\|$.
- 5.1.9f. Recall that multiplication by an orthogonal matrix is a rotation/reflection and so does not change lengths (Theorem 3.2.39f): correspondingly, it also does not change the norm of a matrix as established here.

Now $Q_mA = Q_m(USV^T) = (Q_mU)SV^T$. But Q_mU is an orthogonal matrix (Exercise 3.2.19), so $(Q_mU)SV^T$ is an SVD for Q_mA . From the singular values in S , $\|Q_mA\| = \sigma_1 = \|A\|$.

Also, using 5.1.9e twice: $\|AQ_n\| = \|(AQ_n)^T\| = \|Q_n^T A^T\| = \|A^T\| = \|A\|$.

- 5.1.9g If $\mathbf{x} = \mathbf{0}$, then $|Ax| = |A\mathbf{0}| = |\mathbf{0}| = 0$ whereas $\|A\|\mathbf{x}| = \|A\|\mathbf{0}| = \|A\|0 = 0$, so $|Ax| \leq \|A\|\mathbf{x}|$. Alternatively, if $\mathbf{x} \neq \mathbf{0}$, then we write $\mathbf{x} = \hat{\mathbf{x}}|\mathbf{x}|$ for unit vector $\hat{\mathbf{x}} = \mathbf{x}/|\mathbf{x}|$ so that

$$\begin{aligned} |Ax| &= |A\hat{\mathbf{x}}|\mathbf{x}| \\ &= |A\hat{\mathbf{x}}||\mathbf{x}| \quad (\text{as } |\mathbf{x}| \text{ is a scalar}) \\ &\leq \max_{|\hat{\mathbf{x}}|=1} |A\hat{\mathbf{x}}||\mathbf{x}| \quad (\text{as } \hat{\mathbf{x}} \text{ is a unit vector}) \\ &= \|A\||\mathbf{x}| \quad (\text{by Definition 5.1.5}) \end{aligned}$$

- 5.1.9h See Exercise 5.1.3e.

□

Since the matrix norm has the familiar properties of a measure of magnitude, we use the matrix norm to measure the ‘distance’ between matrices.

Example 5.1.10. (a) Use the matrix norm to estimate the ‘distance’ between matrices

$$B = \begin{bmatrix} -0.7 & 0.4 \\ 0.6 & 0.5 \end{bmatrix} \quad \text{and} \quad C = \begin{bmatrix} -0.2 & 0.9 \\ 0 & 1.7 \end{bmatrix}.$$

Solution: The ‘distance’ is

$$\|C - B\| = \left\| \begin{bmatrix} 0.5 & 0.5 \\ -0.6 & 1.2 \end{bmatrix} \right\| = \|A\| \approx 1.3$$

for the matrix A of Example 5.1.6a and its estimated norm.

- (b) Recall from Example 3.3.2 that the matrix

$$A = \begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix}$$

has an SVD of

$$USV^T = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T.$$

- Find $\|A - B\|$ for the rank one matrix

$$B = \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T = 5\sqrt{2} \begin{bmatrix} -\frac{4}{5} \\ \frac{3}{5} \end{bmatrix} \begin{bmatrix} -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 4 & -4 \\ -3 & 3 \end{bmatrix}.$$

Solution: Let's write matrix

$$\begin{aligned} B &= \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T \\ &= U \begin{bmatrix} 0 & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} V^T. \end{aligned}$$

Then the difference is

$$\begin{aligned} A - B &= U \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} V^T - U \begin{bmatrix} 0 & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} V^T \\ &= U \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 0 \end{bmatrix} V^T. \end{aligned}$$

This is an SVD for $A - B$ with singular values $10\sqrt{2}$ and 0, so by Definition 5.1.5 its norm $\|A - B\| = \sigma_1 = 10\sqrt{2}$.

- Find $\|A - A_1\|$ for the rank one matrix

$$A_1 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T = 10\sqrt{2} \begin{bmatrix} \frac{3}{5} \\ \frac{4}{5} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} = \begin{bmatrix} 6 & 6 \\ 8 & 8 \end{bmatrix}.$$

Solution: Let's write matrix

$$\begin{aligned} A_1 &= \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T \\ &= U \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 0 \end{bmatrix} V^T. \end{aligned}$$

Then the difference is

$$\begin{aligned} A - A_1 &= U \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} V^T - U \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 0 \end{bmatrix} V^T \\ &= U \begin{bmatrix} 0 & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} V^T. \end{aligned}$$

This is an SVD for $A - A_1$ with singular values $5\sqrt{2}$ and 0, albeit out of order, so by Definition 5.1.5 its norm $\|A - A_1\|$ is the largest singular value which here is $5\sqrt{2}$.

Out of these two matrices, A_1 and B , the matrix A_1 is ‘closer’ to A as $\|A - A_1\| = 5\sqrt{2} < 10\sqrt{2} = \|A - B\|$. ■

Example 5.1.11. From Example 5.1.3, recall the ‘bulls eye’ matrix

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

and its rank two and three approximations A_2 and A_3 . Find $\|A - A_2\|$ and $\|A - A_3\|$.

Solution: • Example 5.1.3 found $A_3 = A$ hence $\|A - A_3\| = \|O_5\| = 0$.

- Although $\|A - A_2\|$ is nontrivial, finding it is straightforward using SVDS. Recall that, from the given SVD $A = USV^T$,

$$\begin{aligned} A_2 &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T \\ &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + 0 \mathbf{u}_3 \mathbf{v}_3^T + 0 \mathbf{u}_4 \mathbf{v}_4^T + 0 \mathbf{u}_5 \mathbf{v}_5^T \\ &= U \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} V^T. \end{aligned}$$

Hence the difference

$$\begin{aligned} A - A_2 &= U \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 & 0 \\ 0 & 0 & 0 & \sigma_4 & 0 \\ 0 & 0 & 0 & 0 & \sigma_5 \end{bmatrix} V^T \\ &\quad - U \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} V^T \\ &= U \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 & 0 \\ 0 & 0 & 0 & \sigma_4 & 0 \\ 0 & 0 & 0 & 0 & \sigma_5 \end{bmatrix} V^T. \end{aligned}$$

This is an SVD for $A - A_2$, albeit irregular with the singular values out of order, with singular values of 0, 0, $\sigma_3 = 0.64$,

and $\sigma_4 = \sigma_5 = 0$. The largest singular value gives the norm $\|A - A_2\| = 0.64$ (2 d.p.).

One might further comment that the relative error in the approximate A_2 is $\|A - A_2\|/\|A\| = 0.64/2.68 = 0.24 = 24\%$ (2 d.p.). ■

Theorem 5.1.12 (Eckart–Young). *Let A be an $m \times n$ matrix of rank r with SVD $A = USV^T$. Then for any $k < r$ a closest rank k matrix approximating A , in the matrix norm, is*

$$A_k := US_k V^T = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T \quad (5.2)$$

where $S_k := \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_k, 0, \dots, 0)$. The distance between A and A_k is $\|A - A_k\| = \sigma_{k+1}$.

That is, obtain A_k by ‘setting’ the singular values $\sigma_{k+1} = \cdots = \sigma_r = 0$ from an SVD for A .

Proof. As a prelude, let’s establish the distance between A and A_k . Using their SVDs,

$$\begin{aligned} A - A_k &= USV^T - US_k V^T = U(S - S_k)V^T \\ &= U \text{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_r, 0, \dots, 0)V^T, \end{aligned}$$

and so $A - A_k$ has largest singular value σ_{k+1} . Then from Definition 5.1.5, $\|A - A_k\| = \sigma_{k+1}$.

Use contradiction to prove there is no matrix of rank k closer to A when using $\|\cdot\|$ to measure matrix distances (Trefethen & Bau 1997, p.36). Assume there is some $m \times n$ matrix B with rank $B \leq k$ and closer to A than is A_k , that is, $\|A - B\| < \|A - A_k\|$. First, the Rank Theorem 3.4.32 asserts the null space of B has dimension nullity $B = n - \text{rank } B \geq n - k$ as $\text{rank } B \leq k$. For every $\mathbf{w} \in \text{null } B$, as $B\mathbf{w} = \mathbf{0}$, $A\mathbf{w} = A\mathbf{w} - B\mathbf{w} = (A - B)\mathbf{w}$. Then

$$\begin{aligned} |A\mathbf{w}| &= |(A - B)\mathbf{w}| \\ &\leq \|A - B\| |\mathbf{w}| \quad (\text{by Theorem 5.1.9g}) \\ &< \|A - A_k\| |\mathbf{w}| \quad (\text{by assumption}) \\ &= \sigma_{k+1} |\mathbf{w}| \end{aligned}$$

That is, under the assumption there exists an (at least) $(n - k)$ -dimensional subspace in which $|A\mathbf{w}| < \sigma_{k+1} |\mathbf{w}|$.

Second, consider *any* vector \mathbf{v} in the $(k + 1)$ -dimensional subspace $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k+1}\}$. Say $\mathbf{v} = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \cdots + c_{k+1} \mathbf{v}_{k+1} = V\mathbf{c}$ for any vector of coefficients $\mathbf{c} = (c_1, c_2, \dots, c_{k+1}, 0, \dots, 0) \in \mathbb{R}^n$. Then

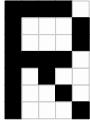
$$|A\mathbf{v}| = |USV^T V\mathbf{c}|$$

$$\begin{aligned}
&= |USe| \quad (\text{as } V^T V = I) \\
&= |Se| \quad (\text{as } U \text{ is orthogonal}) \\
&= |(\sigma_1 c_1, \sigma_2 c_2, \dots, \sigma_{k+1} c_{k+1}, 0, \dots, 0)| \\
&= \sqrt{\sigma_1^2 c_1^2 + \sigma_2^2 c_2^2 + \dots + \sigma_{k+1}^2 c_{k+1}^2} \\
&\geq \sqrt{\sigma_{k+1}^2 c_1^2 + \sigma_{k+1}^2 c_2^2 + \dots + \sigma_{k+1}^2 c_{k+1}^2} \\
&= \sigma_{k+1} \sqrt{c_1^2 + c_2^2 + \dots + c_{k+1}^2} \\
&= \sigma_{k+1} |\mathbf{c}| \\
&= \sigma_{k+1} |V\mathbf{c}| \quad (\text{as } V \text{ is orthogonal}) \\
&= \sigma_{k+1} |\mathbf{v}|.
\end{aligned}$$

That is, there exists a $(k+1)$ -dimensional subspace in which $|A\mathbf{v}| \geq \sigma_{k+1} |\mathbf{v}|$.

Lastly, since the sum of the dimensions of these two subspaces of \mathbb{R}^n is at least $(n-k) + (k+1) > n$, there must be a nonzero vector, say \mathbf{u} , lying in both. So for this \mathbf{u} , simultaneously $|Au| < \sigma_{k+1} |\mathbf{u}|$ and $|Au| \geq \sigma_{k+1} |\mathbf{u}|$. These two deductions contradict each other. Hence the assumption is wrong: there is no rank k matrix more closely approximating A than A_k . \square

Example 5.1.13 (the letter R). In displays with low resolution, letters and numbers are displayed with noticeable pixel patterns: for example, the letter R is pixellated in the margin. Let's see how such pixel patterns are best approximated by matrices of different ranks. (This example is illustrative: it is not a practical image compression since the required singular vectors are more complicated than a small-sized pattern of pixels.)



Solution: Use Procedure 5.1.2. First, form and enter into Matlab/Octave the 7×5 matrix of the pixel pattern as illustrated in the margin

$$R = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

Second, compute an SVD via `[U,S,V]=svd(R)` to find (2 d.p.)

$$\begin{aligned}
\mathbf{U} = & \\
-0.53 & 0.38 -0.00 -0.29 -0.70 -0.06 -0.07 \\
-0.28 & -0.49 0.00 -0.13 0.10 -0.69 -0.42 \\
-0.28 & -0.49 -0.00 -0.13 -0.02 0.72 -0.39 \\
-0.53 & 0.38 -0.00 -0.29 0.70 0.06 0.07 \\
-0.32 & 0.03 -0.71 0.63 -0.00 -0.00 -0.00 \\
-0.32 & 0.03 0.71 0.63 -0.00 0.00 0.00
\end{aligned}$$

```

-0.28 -0.49 -0.00 -0.13 -0.08 -0.02  0.81
S =
 3.47      0      0      0      0
    0    2.09      0      0      0
    0      0    1.00      0      0
    0      0      0    0.75      0
    0      0      0      0    0.00
    0      0      0      0      0
    0      0      0      0      0
V =
-0.73 -0.32  0.00  0.40 -0.45
-0.30  0.36 -0.00 -0.76 -0.45
-0.40  0.37 -0.71  0.07  0.45
-0.40  0.37  0.71  0.07  0.45
-0.24 -0.70 -0.00 -0.50  0.45

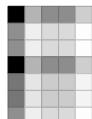
```



The singular values are $\sigma_1 = 3.47$, $\sigma_2 = 2.09$, $\sigma_3 = 1.00$, $\sigma_4 = 0.75$ and $\sigma_5 = 0$. Four successively better approximations to the image are the following.

- The coarsest approximation is $R_1 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T$, that is

$$R_1 = 3.47 \begin{bmatrix} -0.53 \\ -0.28 \\ -0.28 \\ -0.53 \\ -0.32 \\ -0.32 \\ -0.28 \end{bmatrix} \begin{bmatrix} -0.73 & -0.30 & -0.40 & -0.40 & -0.24 \end{bmatrix}^T.$$



Compute with $\mathbf{R1}=\mathbf{U}(:,1)*\mathbf{S}(1,1)*\mathbf{V}(:,1)'$ to find (2 d.p.), as illustrated,

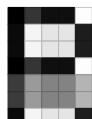
```

R1 =
  1.34   0.55   0.72   0.72   0.44
  0.71   0.29   0.39   0.39   0.24
  0.71   0.29   0.39   0.39   0.24
  1.34   0.55   0.72   0.72   0.44
  0.83   0.34   0.45   0.45   0.27
  0.83   0.34   0.45   0.45   0.27
  0.71   0.29   0.39   0.39   0.24

```

This has difference $\|R - R_1\| = \sigma_2 = 2.09$ which at 60% of σ_1 is large: indeed the letter R is not recognisable.

- The second approximation is $R_2 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T$. Compute via $\mathbf{R2}=\mathbf{U}(:,1:2)*\mathbf{S}(1:2,1:2)*\mathbf{V}(:,1:2)'$ to find (2 d.p.), as illustrated,



```

R2 =
  1.09   0.83   1.02   1.02  -0.11
  1.04  -0.07   0.01   0.01   0.95
  1.04  -0.07   0.01   0.01   0.95

```

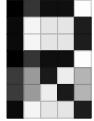
$$\begin{array}{ccccc} 1.09 & 0.83 & 1.02 & 1.02 & -0.11 \\ 0.81 & 0.36 & 0.47 & 0.47 & 0.24 \\ 0.81 & 0.36 & 0.47 & 0.47 & 0.24 \\ 1.04 & -0.07 & 0.01 & 0.01 & 0.95 \end{array}$$

This has difference $\|R - R_2\| = \sigma_3 = 1.00$ which at 29% of σ_1 is large: but one can begin to imagine the letter R in the image.

- The third approximation is

$$R_3 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \sigma_3 \mathbf{u}_3 \mathbf{v}_3^T.$$

Compute with $R3=U(:,1:3)*S(1:3,1:3)*V(:,1:3)'$ to find (2 d.p.), as illustrated,



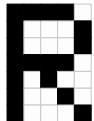
$$\begin{array}{ccccc} R3 = & & & & \\ 1.09 & 0.83 & 1.02 & 1.02 & -0.11 \\ 1.04 & -0.07 & 0.01 & 0.01 & 0.95 \\ 1.04 & -0.07 & 0.01 & 0.01 & 0.95 \\ 1.09 & 0.83 & 1.02 & 1.02 & -0.11 \\ 0.81 & 0.36 & 0.97 & -0.03 & 0.24 \\ 0.81 & 0.36 & -0.03 & 0.97 & 0.24 \\ 1.04 & -0.07 & 0.01 & 0.01 & 0.95 \end{array}$$

This has difference $\|R - R_3\| = \sigma_4 = 0.75$ which at 22% of σ_1 is moderate and one can see the letter R emerging.

- The fourth approximation is

$$R_4 = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_4 \mathbf{u}_4 \mathbf{v}_4^T.$$

Compute with $R4=U(:,1:4)*S(1:4,1:4)*V(:,1:4)'$ to find (2 d.p.), as illustrated,



$$\begin{array}{ccccc} R4 = & & & & \\ 1.00 & 1.00 & 1.00 & 1.00 & -0.00 \\ 1.00 & 0.00 & -0.00 & 0.00 & 1.00 \\ 1.00 & 0.00 & -0.00 & -0.00 & 1.00 \\ 1.00 & 1.00 & 1.00 & 1.00 & -0.00 \\ 1.00 & -0.00 & 1.00 & 0.00 & -0.00 \\ 1.00 & -0.00 & -0.00 & 1.00 & 0.00 \\ 1.00 & 0.00 & -0.00 & -0.00 & 1.00 \end{array}$$

This has difference $\|R - R_4\| = \sigma_5 = 0.00$ and R_4 exactly reproduces R.

■

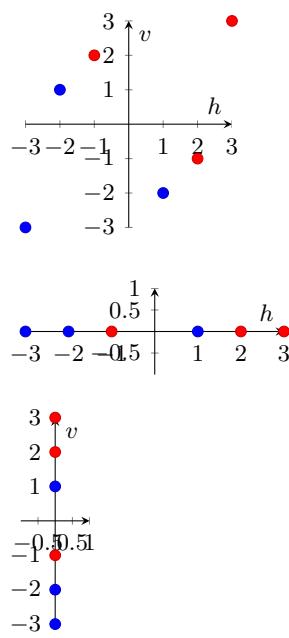
Example 5.1.14. Recall Example 5.1.4 approximated the image of Euler (1737) with various rank k approximates from an SVD of the image. Let the image be denoted by matrix A . From Figure 5.2 the largest singular value of the image is $\|A\| = \sigma_1 \approx 40,000$.

- From Theorem 5.1.12, the rank 3 approximation in Figure 5.1 is a distance $\|A - A_3\| = \sigma_4 \approx 5,000$ (from Figure 5.2) away from the image. That is, image A_3 has a relative error roughly $5000/40000 = 1/8 \approx 12\%$.
 - From Theorem 5.1.12, the rank 10 approximation in Figure 5.1 is a distance $\|A - A_{10}\| = \sigma_{11} \approx 5,000$ (from Figure 5.2) away from the image. That is, image A_{10} has a relative error roughly $2000/40000 = 1/20 = 5\%$.
 - From Theorem 5.1.12, the rank 30 approximation in Figure 5.1 is a distance $\|A - A_{30}\| = \sigma_{31} \approx 800$ (from Figure 5.2) away from the image. That is, image A_{30} has a relative error roughly $800/40000 = 1/50 = 2\%$.
-

5.1.3 Principal component analysis

In its ‘best’ approximation property, Theorem 5.1.12 establishes the effectiveness of an SVD in image compression. Scientists and engineers also use this result for so-called data reduction: often using just a rank two (or three) ‘best’ approximation to high dimensional data, one then plots 2D (or 3D) graphics. Such an approach is often termed a principal component analysis (PCA).

Example 5.1.15 (toy items). Suppose you are given data about six items, three blue and three red. Suppose each item has two measured properties/attributes called h and v as in the following table:



h	v	colour
-3	-3	blue
-2	1	blue
1	-2	blue
-1	2	red
2	-1	red
3	3	red

The item properties/attributes are the points (h, v) in 2D as illustrated in the margin. But humans always prefer simple one dimensional summaries: we do it all the time when we rank sport teams, schools, web pages, and so on.

Challenge: is there a one dimensional summary of these six item’s data that clearly separates the blue from the red? Using just one of the attributes h or v on their own would not suffice:

- using h alone leads to a 1D view where the red and the blue are intermingled as shown in the margin;
- similarly, using v alone leads to a 1D view where the red and the blue are intermingled as shown in the margin.

Solution: Use an SVD to automatically find the best 1D view of the data.

- (a) Enter the 6×2 matrix of data into Matlab/Octave with

```
A=[-3 -3  
-2 1  
1 -2  
-1 2  
2 -1  
3 3 ]
```



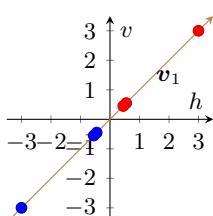
- (b) Then $[U, S, V] = \text{svd}(A)$ computes an SVD, $A = USV^T$, of the data (2 d.p.):

```
U =  
-0.69 0.00 -0.09 0.09 0.14 0.70  
-0.11 -0.50 0.50 -0.50 0.48 -0.08  
-0.11 0.50 0.82 0.18 -0.18 0.01  
0.11 -0.50 0.18 0.82 0.18 -0.01  
0.11 0.50 -0.14 0.14 0.83 -0.09  
0.69 -0.00 0.12 -0.12 0.02 0.70  
  
S =  
6.16 0  
0 4.24  
0 0  
0 0  
0 0  
0 0  
  
V =  
0.71 0.71  
0.71 -0.71
```

- (c) Now what does such an SVD tell us? Recall from the proof of the SVD (section 3.3.3) that $A\mathbf{v}_1 = \sigma_1 \mathbf{u}_1$. Further recall from the proof that \mathbf{v}_1 is the unit vector that maximises $|A\mathbf{v}_1|$ so in some sense it is the direction in which the data in A is most spread out (\mathbf{v}_1 is called the principal vector). We find here (2 d.p.)

$$A\mathbf{v}_1 = \sigma_1 \mathbf{u}_1 = (-4.24, -0.71, -0.71, 0.71, 0.71, 4.24)$$

which neatly separates the blue items (negative) from the red (positive). In essence, the product $A\mathbf{v}_1$ orthogonally projects (Section 3.5.3) the items' (h, v) data onto the subspace $\text{span}\{\mathbf{v}_1\}$ as illustrated in the margin. ■



Although this Example 5.1.15 is just a toy to illustrate concepts, the above steps generalise straightforwardly to be immensely useful on

Table 5.2: part of Edgar Anderson's Iris data, lengths in cm. The measurements come from the flowers of ten each of three different species of Iris.

Sepal length	Sepal width	Petal length	Petal width	Species
4.9	3.0	1.4	0.2	
4.6	3.4	1.4	0.3	
4.8	3.4	1.6	0.2	
5.4	3.9	1.3	0.4	
5.1	3.7	1.5	0.4	Setosa
5.0	3.4	1.6	0.4	
5.4	3.4	1.5	0.4	
5.5	3.5	1.3	0.2	
4.5	2.3	1.3	0.3	
5.1	3.8	1.6	0.2	
6.4	3.2	4.5	1.5	
6.3	3.3	4.7	1.6	
5.9	3.0	4.2	1.5	
5.6	3.0	4.5	1.5	
6.1	2.8	4.0	1.3	Versicolor
6.8	2.8	4.8	1.4	
5.5	2.4	3.7	1.0	
6.7	3.1	4.7	1.5	
6.1	3.0	4.6	1.4	
5.7	2.9	4.2	1.3	
5.8	2.7	5.1	1.9	
4.9	2.5	4.5	1.7	
6.4	2.7	5.3	1.9	
6.5	3.0	5.5	1.8	
5.6	2.8	4.9	2.0	Virginia
6.2	2.8	4.8	1.8	
7.9	3.8	6.4	2.0	
6.3	3.4	5.6	2.4	
6.9	3.1	5.1	2.3	
6.3	2.5	5.0	1.9	

vastly bigger and more challenging data. The next example takes the next step in complexity by introducing how to automatically find a good 2D view of some data in 4D.

Example 5.1.16 (Iris flower data set). Table 5.2 list part of Edgar Anderson's data on the length and widths of sepals and petals of Iris flowers.² There are three species of Irises in the data. The data is 4D: each instance of thirty Iris flowers is characterised by

² <http://archive.ics.uci.edu/ml/datasets/Iris> gives the full dataset (Lichman 2013).

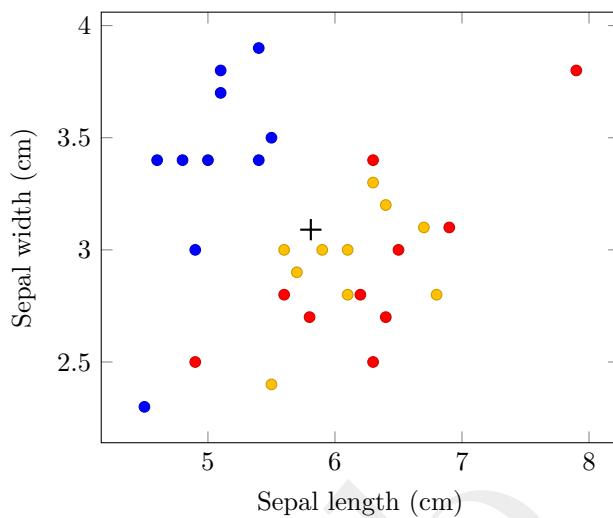


Figure 5.3: scatter plot of sepal widths versus lengths for Edgar Anderson’s iris data of Table 5.2: blue, Setosa; brown, Versicolor; red, Virginia. The black + marks the mean sepal width and length.

the four lengths. Our challenge is to plot a 2D picture of this data in such a way that separates the flowers as best as possible. For high-D data (although 4D is not really that high), simply plotting one characteristic against another is rarely useful. For example, Figure 5.3 plots the attributes of sepal widths versus sepal lengths and results in the three species being intermingled together rather than reasonably separated. Our aim is to instead plot Figure 5.4 which successfully separates the three species.

Solution: Use an SVD to find a best low-rank view of the data.

- (a) Enter the 30×5 matrix of Iris data into Matlab/Octave with a complete version of

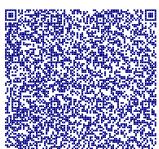
```
iris=[  
4.9 3.0 1.4 0.2 1  
4.6 3.4 1.4 0.3 1  
...  
6.3 2.5 5.0 1.9 3  
]
```

where the fifth column of 1, 2, 3 corresponds to the species Setosa, Versicolor or Virginia, respectively. Then a scatter plot such as Figure 5.3 may be drawn with the command

```
scatter(iris(:,1),iris(:,2),[],iris(:,5))
```

The above `scatter(x,y,[],s)` plots a scatter plot of points with colour depending upon `s` which here corresponds to each different species.

- (b) If we were on a walk to a scenic lookout to get a view of the countryside, then the scenic lookout would be in the



countryside: it is no good going to a lookout a long way away from the scene we wish to view. Correspondingly, to best view a dataset we typically look it at from the very centre of the data, namely its mean. That is, here we use an SVD of the data matrix only after subtracting the mean of each attribute so that the SVD analyses the variations from the mean. Here the mean Iris sepal length and width is 5.81 cm and 3.09 cm (the black + in Figure 5.3), and the mean petal length and width is 3.69 cm and 1.22 cm. In Matlab/Octave execute the following to form a matrix A of the variations from the mean, and compute an SVD:

```
meaniris=mean(iris(:,1:4))
A=iris(:,1:4)-ones(30,1)*meaniris
[U,S,V]=svd(A)
```

The resulting SVD is (2 d.p.)

```
U = ...
S =
    10.46      0      0      0
        0    2.86      0      0
        0      0    1.47      0
        0      0      0    0.85
    ...
V =
    0.34    0.72   -0.56   -0.20
   -0.07    0.65    0.74    0.14
    0.87   -0.17    0.14    0.45
    0.36   -0.15    0.33   -0.86
```

where a ... indicates information that is not directly of interest.

- (c) As justified shortly, the two most important components of a flower's shape are those in the directions of \mathbf{v}_1 and \mathbf{v}_2 (called the two principal vectors). Because \mathbf{v}_1 and \mathbf{v}_2 are orthonormal, the first component for each Iris flower is $\mathbf{x} = A\mathbf{v}_1$ and the second component for each is $\mathbf{y} = A\mathbf{v}_2$. The beautiful Figure 5.4 is a scatter plot of the components of \mathbf{y} versus the components of \mathbf{x} that untangles the three species. Obtain Figure 5.4 in Matlab/Octave with the command

```
scatter(A*V(:,1),A*V(:,2),[],iris(:,5))
```

Figure 5.4 shows our SVD based analysis largely separates the three species using these two different combinations of the flowers' attributes.



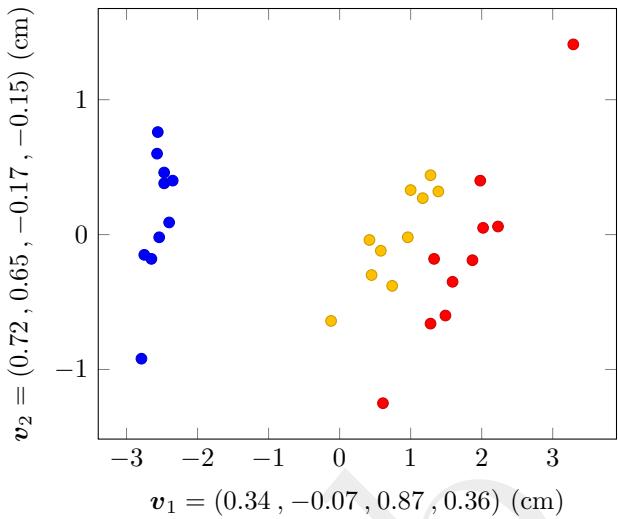


Figure 5.4: best 2D scatter plot of Edgar Anderson's iris data: blue, Setosa; brown, Versicolor; red, Virginica.

Transpose the usual mathematical convention Perhaps you noticed that the previous Example 5.1.16 flips our usual mathematical convention that vectors are column vectors. The example uses row vectors of the four attributes of each flower: Table 5.2 lists that the first Iris Setosa flower has a row vector of attributes $[4.9 \ 3.0 \ 1.4 \ 0.2]$ (cm) corresponding to the sepal length and width, and the petal length and width, respectively. Similarly, the last Virginia Iris flower has row vector of attributes of $[46.3 \ 2.5 \ 5.0 \ 1.9]$ (cm), and the mean vector is the row vector $[5.81 \ 3.09 \ 3.69 \ 1.22]$ (cm). The reason for this mathematical transposition is that throughout science and engineering, data results are most often presented as rows of different instances of flowers, animals, clients or experiments, and each row contains the list of characteristic measured or derived properties/attributes. Table 5.2 has this most common structure. Thus in this sort of application, the mathematics we do needs to reflect this most common structure. Hence many vectors in this subsection are row vectors.

Definition 5.1.17 (principal components). *Given a $m \times n$ data matrix A (usually with zero mean when averaged over all rows), with SVD $A = USV^T$ the j th column \mathbf{v}_j of V is called the j th **principal vector** and the vector $\mathbf{x}_j := A\mathbf{v}_j$ is called the j th **principal components** of the data matrix A .*

Now what does an SVD tell us for 2D plots of data? We know A_2 is the best rank two approximation to the data matrix A (Theorem 5.1.12). That is, if we are only to plot two components, those two components are best to come from A_2 . Recall from (5.2) that

$$A_2 = US_2V^T = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T = (\sigma_1 \mathbf{u}_1) \mathbf{v}_1^T + (\sigma_2 \mathbf{u}_2) \mathbf{v}_2^T.$$

That is, in this best rank two approximation of the data, the row vector of attributes of the i th Iris are the linear combination of row vectors $(\sigma_1 u_{i1})\mathbf{v}_1^T + (\sigma_2 u_{i2})\mathbf{v}_2^T$. The vectors \mathbf{v}_1 and \mathbf{v}_2 are orthonormal vectors so we treat them as the horizontal and vertical unit vectors of a scatter plot. That is, $x_i = \sigma_1 u_{i1}$ and $y_i = \sigma_2 u_{i2}$ are horizontal and vertical coordinates of the i th Iris in the best 2D plot. Consequently, in Matlab/Octave we draw a scatter plot of the components of vectors $\mathbf{x} = \sigma_1 \mathbf{u}_1$ and $\mathbf{y} = \sigma_2 \mathbf{u}_2$ (Figure 5.4).

Theorem 5.1.18. *Using the matrix norm to measure ‘best’, the best k -dimensional summary of the $m \times n$ data matrix A (usually of zero mean) are the first k principal components in the directions of the first k principal vectors.*

Proof. Let $A = USV^T$ be an SVD of matrix A . For any $k < \text{rank } A$, Theorem 5.1.12 establishes that

$$A_k := US_k V^T = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T$$

is the best rank k approximation to A in the matrix norm. Letting matrix $U = [u_{ij}]$, write the i th row of A_k as $(\sigma_1 u_{i1})\mathbf{v}_1^T + (\sigma_2 u_{i2})\mathbf{v}_2^T + \cdots + (\sigma_k u_{ik})\mathbf{v}_k^T$ and hence the transpose of each row of A_k lies in the k D subspace $\text{span}\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$. This establishes that these are principal vectors.

Since $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is an orthonormal set, we now use them as standard unit vectors of a coordinate system for the k D subspace. From the above linear combination, the components of the i th data point approximation in this subspace coordinate system are $\sigma_1 u_{i1}, \sigma_2 u_{i2}, \dots, \sigma_k u_{ik}$. That is, the j th coordinate for all data points, the principal components, is $\sigma_j \mathbf{u}_j$. By post-multiplying the SVD $A = USV^T$ by orthogonal V , recall that $AV = US$ which written in terms of columns is

$$[Av_1 \ Av_2 \ \cdots \ Av_r] = [\sigma_1 \mathbf{u}_1 \ \sigma_2 \mathbf{u}_2 \ \cdots \ \sigma_r \mathbf{u}_{mr}]$$

where $r = \text{rank } A$. Consequently, the vector $\sigma_j \mathbf{u}_j$ of j th coordinates in the subspace are equal to Av_j , the principal components. \square

Example 5.1.19 (wine recognition). From the Lichman (2013) repository³ download the data file `wine.data` and its description file `wine.names`. The wine data has 178 rows of different wine samples, and 14 columns of attributes of which the first column is the cultivar class number and the remaining 13 columns are the amounts of different chemicals measured in the wine. Question: is there a two-dimensional view of these chemical measurements that largely separates the cultivars?

Solution: Use and SVD to find the best two-dimensional, rank two, view of the data.

³ <http://archive.ics.uci.edu/ml/datasets/Wine>

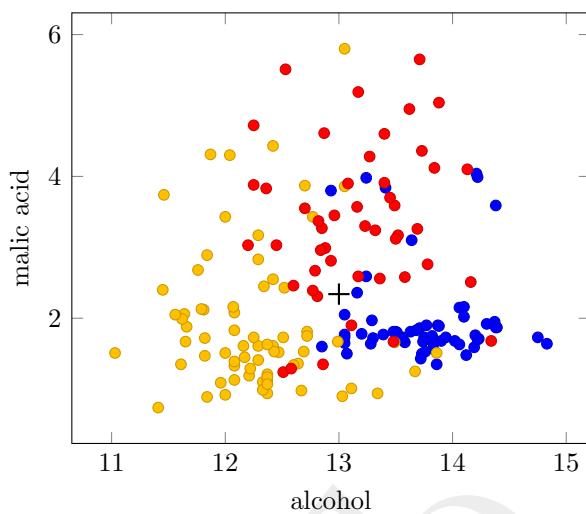


Figure 5.5: for the wine data of Example 5.1.19, a plot of the measured malic acid versus measured alcohol, and coloured depending upon the cultivar, shows these measurements alone cannot effectively discriminate between the cultivars.

- (a) Read in the 178×14 matrix of data into Matlab/Octave with the commands

```
wine=csvread('wine.data')
[m,n]=size(wine)
scatter(wine(:,2),wine(:,3),[],wine(:,1))
```

The scatter plot, Figure 5.5, shows that if we just plot the first two chemicals, alcohol and malic acid, then the three cultivars are inextricably intermingled. Our aim is to automatically draw Figure 5.6 in which the three cultivars are largely separated.

- (b) To find the principal components of the wine chemicals it is best to remove the mean with

```
meanw=mean(wine(:,2:14))
A=wine(:,2:14)-ones(m,1)*meanw;
```

where the `mean(X)` computes the mean/average of each column of `X`.

But now a further issue arises: the values in the columns are of widely different magnitudes; moreover, each column has different physical units (in contrast, the Iris flower measurements were all cm). In practice we *must not* mix together quantities with different physical units. The general rule, after making each column zero mean, is to scale each column by dividing by its standard deviation, equivalently by its root-mean-square. This scaling does two practically useful things:

- since the standard deviation measures the spread of



data in a column, it has the same physical units as the column of data, so dividing by it renders the results dimensionless, and so suitable for mixing with other scaled columns;

- also the spread of data in each column is now comparable to each other, namely around about size one, instead of some columns being of the order of one-tenths and other columns being in the hundreds.

Consequently, form the 178×13 matrix to analyse by commands

```
meanw=mean(wine(:,2:14))
stdw=std(wine(:,2:14))
A=(wine(:,2:14)-ones(m,1)*meanw)*diag(1./stdw);
```

where the `std(X)` computes the standard deviation of each column of `X`.

- (c) Now compute and use an SVD $A = USV^T$. But for low rank approximations we only ever use the first few singular values and first few singular vectors. Thus it is pointless computing a full SVD which here has 178×178 matrix U and 13×13 matrix V .⁴ Consequently, use `[U,S,V]=svds(A,4)` to economically compute only the first four singular values and singular vectors (change the four to suit your purpose) to find (2 d.p.)

```
U =
S =
    28.86      0      0      0
        0    21.02      0      0
        0      0   16.00      0
        0      0      0   12.75
V =
   -0.14    0.48   -0.21   -0.02
    0.25    0.22    0.09    0.54
    0.00    0.32    0.63   -0.21
    0.24   -0.01    0.61    0.06
   -0.14    0.30    0.13   -0.35
   -0.39    0.07    0.15    0.20
   -0.42   -0.00    0.15    0.15
    0.30    0.03    0.17   -0.20
   -0.31    0.04    0.15    0.40
    0.09    0.53   -0.14    0.07
   -0.30   -0.28    0.09   -0.43
   -0.38   -0.16    0.17    0.18
```

⁴ Yes, on modern computers this is here done within a millisecond. But for modern datasets with thousands to billions of rows a full SVD is infeasible so let's see how to analyse modern large datasets.

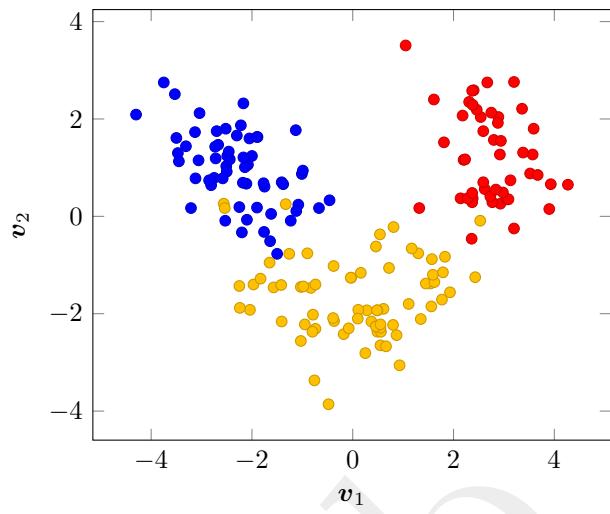


Figure 5.6: for the wine data of Example 5.1.19, a plot of the first two principal components almost entirely separates the three cultivars.

$-0.29 \quad 0.36 \quad -0.13 \quad -0.23$

where the \dots indicates we do not here need to know U .

- (d) Recall that the columns of this orthogonal V are the principal vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_4$, and the j th principal components of the data are $\mathbf{x}_j = \mathbf{A}\mathbf{v}_j$. We form a 2D plotted view of the data, Figure 5.6, by drawing a scatter plot of the first two principal components with

```
scatter(A*V(:,1),A*V(:,2),[],wine(:,1))
```

Figure 5.6 shows these two principal components do an amazingly good job of almost completely disentangling the three wine cultivars (use `scatter3` to explore the first three principal components).

■

The previous three examples develop the following procedure for ‘best’ viewing data in low dimensions. The procedure is automatic. However, additional information about the data or preferred results would lead to modifications.

Procedure 5.1.20 (principal component analysis). *Consider the case when you have data values consisting of n attributes for each of m instances, and the aim is to find a good k -dimensional summary/view of the data.*

1. Form/enter the $m \times n$ data matrix B .
2. Scale the data matrix B to form $m \times n$ matrix A :

- (a) usually make each column have zero mean by subtracting its mean \bar{b}_j , algebraically $\mathbf{a}_j = \mathbf{b}_j - \bar{b}_j$;
- (b) but ensure each column has the same ‘physical dimensions’, often by dividing by the standard deviation s_j of each column, algebraically $\mathbf{a}_j = (\mathbf{b}_j - \bar{b}_j)/s_j$.

Compute $\mathbf{A}=(\mathbf{B}-\text{ones}(\mathbf{m},1)\text{mean}(\mathbf{B}))*\text{diag}(1./\text{std}(\mathbf{B}))$ in Matlab/Octave.*

3. Economically compute an SVD for the best rank k approximation to the scaled data matrix with $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svds}(\mathbf{A}, k)$.
4. Then the j th column of \mathbf{V} is the j th principal vector, and the principal components are the entries of the $m \times k$ matrix $\mathbf{A} * \mathbf{V}$.

Courses on multivariate statistics prove that, for any (usually zero mean) data matrix A , the first k principal vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ are orthogonal unit vectors that *maximise the total variance* in the principal components $\mathbf{x}_j = \mathbf{Av}_j$; that is, that maximise $|\mathbf{x}_1|^2 + |\mathbf{x}_2|^2 + \dots + |\mathbf{x}_k|^2$. Indeed, this maximisation of the variance is closely tied to the constructive proof of the existence of SVDS (section 3.3.3) which successively maximises $|\mathbf{Av}|$ subject to \mathbf{v} being orthonormal to the singular/principal vectors already determined. Consequently, when data is approximated in the space of the first k principal vectors, then the data is the most spread out it can be in k -dimensions.

Application to latent semantic indexing

This ability to retrieve relevant information based upon meaning rather than literal term usage is the main motivation for using LSI [latent semantic indexing].

(Berry et al. 1995, p.579)

Searching for information based upon word matching results in surprisingly poor retrieval of relevant documents (Berry et al. 1995, §5.5). Instead, the so-called latent semantic indexing improves retrieval by replacing individual words with nearness of word vectors derived via the singular value decomposition. This section introduces latent semantic indexing via a very small example.

The Society for Industrial and Applied Mathematics (SIAM) reviews many mathematical books. In 2015 six of those books had the following titles:

1. Introduction to Finite and Spectral Element Methods using MATLAB
2. Iterative Methods for Linear Systems: Theory and Applications

3. Singular Perturbations: Introduction to System Order Reduction Methods with Applications
4. Risk and Portfolio Analysis: Principles and Methods
5. Stochastic Chemical Kinetics: Theory and Mostly Systems Biology Applications
6. Quantum Theory for Mathematicians

Consider the capitalised words. For those words that appear in more than one title, let's form a word vector (Example 1.1.7) for each title, then use principal components to summarise these six books on a 2D plane. This task is part of what is called latent semantic indexing (Berry et al. 1995). (We should also count words that are used only once, but this example omits for simplicity.)

Follow the principal component analysis Procedure 5.1.20.

1. First find the set of words that are used more than once. Ignoring pluralisation, they are: Application, Introduction, Method, System, Theory. The corresponding word vector for each book title is then the following:

- $\mathbf{w}_1 = (0, 1, 1, 0, 0)$ Introduction to Finite and Spectral Element Methods using MATLAB
- $\mathbf{w}_2 = (1, 0, 1, 1, 1)$ Iterative Methods for Linear Systems: Theory and Applications
- $\mathbf{w}_3 = (1, 1, 1, 1, 0)$ Singular Perturbations: Introduction to System Order Reduction Methods with Applications
- $\mathbf{w}_4 = (0, 0, 1, 0, 0)$ Risk and Portfolio Analysis: Principles and Methods
- $\mathbf{w}_5 = (1, 0, 0, 1, 1)$ Stochastic Chemical Kinetics: Theory and Mostly Systems Biology Applications
- $\mathbf{w}_6 = (0, 0, 0, 0, 1)$ Quantum Theory for Mathematicians

2. Second, form the data matrix with $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_6$ as rows (not columns). We could remove the mean word vector, but choose not to: here the position of each book title relative to an empty title (the origin) is interesting. There is no need to scale each column as each column has the same ‘physical’ dimensions, namely a word count. The data matrix of word vectors is then

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



3. Third, to compute a representation in the 2D plane, principal components uses, as an orthonormal basis, the singular vectors corresponding to the two largest singular values. So compute the economical SVD with $[U, S, V] = \text{svds}(A, 2)$ giving (2 d.p.)

```

U = ...
S =
    3.14      0
        0   1.85
V =
    +0.52  -0.20
    +0.26  +0.52
    +0.50  +0.57
    +0.52  -0.20
    +0.37  -0.57

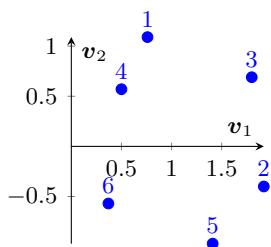
```

4. Columns of V are word vectors in the 5D space of counts of Application, Introduction, Method, System, and Theory. The two given columns of $V = [v_1 \ v_2]$ are the two orthonormal principal vectors:

- the first v_1 , from its largest components, mainly identifies the overall direction of Application, Method and System;
- whereas the second v_2 , from its largest positive and negative components, mainly distinguishes Introduction and Method from Theory.

The corresponding principal components are the entries of the 6×2 matrix

$$AV = \begin{bmatrix} 0.76 & 1.09 \\ 1.92 & -0.40 \\ 1.80 & 0.69 \\ 0.50 & 0.57 \\ 1.41 & -0.97 \\ 0.37 & -0.57 \end{bmatrix} :$$



for each of the six books, the book title has components in the two principal directions given by the corresponding row in this product. We plot the six books on a 2D plane with the Matlab/Octave command

```
scatter(A*V(:,1), A*V(:,2), [], 1:6)
```

to produce a picture like that in the margin. The SVD analysis nicely distributes the books.

The above procedure would approximate the original word vector

data, formed into a matrix, by the following rank two matrix (2 d.p.)

$$A_2 = US_2V^T = \begin{bmatrix} 0.18 & 0.77 & 1.01 & 0.18 & -0.33 \\ 1.08 & 0.29 & 0.74 & 1.08 & 0.95 \\ 0.80 & 0.82 & 1.30 & 0.80 & 0.28 \\ 0.15 & 0.43 & 0.58 & 0.15 & -0.14 \\ 0.93 & -0.14 & 0.16 & 0.93 & 1.08 \\ 0.31 & -0.20 & -0.14 & 0.31 & 0.46 \end{bmatrix}.$$

The largest components in each row do correspond to the ones in the original word vector matrix A . However, in this application we work with the representation in the low dimensional, 2D, subspace spanned by the first two principal vectors v_1 and v_2 .

Angles measure similarity Recall that Example 1.3.7 introduced using the dot product to measure the similarity between word vectors. We could use the dot product in the 5D space of the word vectors to find the ‘angles’ between the book titles. However, we know that the 2D view just plotted is the ‘best’ 2D summary of the book titles, so we could more economically estimate the angle between book titles using just the 2D summary.

Example 5.1.21. What is the ‘angle’ between the first two listed books?

- Introduction to Finite and Spectral Element Methods using MATLAB
- Iterative Methods for Linear Systems: Theory and Applications

Solution: Find the angle two ways.

- (a) First, the corresponding 5D word vectors are $w_1 = (0, 1, 1, 0, 0)$ and $w_2 = (1, 0, 1, 1, 1)$, with lengths $|w_1| = \sqrt{2}$ and $|w_2| = \sqrt{4} = 2$. The dot product then determines

$$\cos \theta = \frac{\mathbf{w}_1 \cdot \mathbf{w}_2}{|\mathbf{w}_1| |\mathbf{w}_2|} = \frac{0 + 0 + 1 + 0 + 0}{2\sqrt{2}} = 0.3536.$$

Hence the angle $\theta = 69.30^\circ$.

- (b) Secondly, estimate the angle using the 2D view. For these two books the principal component vectors are $(0.76, 1.09)$ and $(1.92, -0.40)$, respectively, with lengths 1.33 and 1.96 (2 d.p.). The dot product gives

$$\cos \theta \approx \frac{(0.76, 1.09) \cdot (1.92, -0.40)}{1.33 \cdot 1.96} = \frac{1.02}{2.61} = 0.39.$$

Hence the angle $\theta \approx 67^\circ$ which is effectively the same as the first exact calculation.

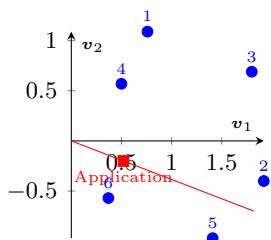
Because of the relatively large ‘angle’ between these two book titles, we deduce that the two books are quite dissimilar.

■

We can also use the 2D plane to economically measure similarity between the book titles and any other title or words of interest.

Example 5.1.22. Let's ask which of the six books is 'closest' to a book about Applications.

Solution: The word Application has word vector $\mathbf{w} = (1, 0, 0, 0, 0)$. So we could do some computations in the original 5D space of word vectors finding precise angles between this word vector and the word vectors of all titles. Alternatively, let's draw a picture in 2D. The Application word vector \mathbf{w} projects onto the 2D plane of principal components by computing $\mathbf{w} \cdot \mathbf{v}_1 = \mathbf{w}^T \mathbf{v}_1$ and $\mathbf{w} \cdot \mathbf{v}_2 = \mathbf{w}^T \mathbf{v}_2$, that is, $\mathbf{w}^T V$. Here the Application word vector $\mathbf{w} = (1, 0, 0, 0, 0)$, so $\mathbf{w}^T V = [0.52 \ -0.20]$, as plotted in the margin. Which of the six books makes the smallest angle with the line through $(0.52, -0.20)$? Visually, books 2 and 5 are closest, and book 2 appears to have slightly smaller angle to the line than book 5. On this data, we deduce that closest to "Application" is book 2: "Iterative Methods for Linear Systems: Theory and Applications"



■

Search for information from more books Berry et al. (1995) reviewed the application of the SVD to the problem of searching for information. Let's explore this further with more data, albeit still very restricted. Berry et al. (1995) listed some mathematical books including the following fourteen titles.

1. a Course on Integral Equations
2. Automatic Differentiation of Algorithms: Theory, Implementation, and Application
3. Geometrical Aspects of Partial Differential Equations
4. Introduction to Hamiltonian Dynamical Systems and the n-Body Problem
5. Knapsack Problems: Algorithms and Computer Implementations
6. Methods of Solving Singular Systems of Ordinary Differential Equations
7. Nonlinear Systems
8. Ordinary Differential Equations
9. Oscillation Theory of Delay Differential Equations

10. Pseudodifferential Operators and Nonlinear Partial Differential Equations
11. Sinc Methods for Quadrature and Differential Equations
12. Stability of Stochastic Differential Equations with Respect to Semi-Martingales
13. the Boundary Integral Approach to Static and Dynamic Contact Problems
14. the Double Mellin–Barnes Type Integrals and their Applications to Convolution Theory

Principal component analysis summarises and relates these titles.

Follow Procedure 5.1.20.

1. The significant (capitalised) words which appear more than once in these titles (ignoring pluralisation) are the fourteen words

Algorithm, Application, Differential/tion,
 Dynamic/al, Equation, Implementation, Integral,
 Method, Nonlinear, Ordinary, Partial, Problem,
 System, and Theory. (5.3)

With this dictionary of significant words, the titles have the following word vectors.

- $\mathbf{w}_1 = (0, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0)$ a Course on Integral Equations
 - $\mathbf{w}_2 = (1, 1, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1)$ Automatic Differentiation of Algorithms: Theory, Implementation, and Application
 - ...
 - $\mathbf{w}_{14} = (0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1)$ the Double Mellin–Barnes Type Integrals and their Applications to Convolution Theory
2. Form the 14×14 data matrix with the word count for each

title in rows

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

Each row corresponds to a book title, and each column corresponds to a word.



3. To compute a representation of the titles in 3D space, principal components uses, as an orthonormal basis, the singular vectors corresponding to the three largest singular values. So in Matlab/Octave compute the economical SVD with $[U, S, V] = \text{svds}(A, 3)$ giving (2 d.p.)

$$\begin{aligned} U &= \dots \\ S &= \begin{bmatrix} 4.20 & 0 & 0 \\ 0 & 2.65 & 0 \\ 0 & 0 & 2.36 \end{bmatrix} \\ V &= \begin{bmatrix} 0.07 & 0.40 & 0.14 \\ 0.07 & 0.38 & 0.25 \\ 0.65 & 0.00 & 0.15 \\ 0.01 & 0.23 & -0.46 \\ 0.64 & -0.21 & -0.07 \\ 0.07 & 0.40 & 0.14 \\ 0.06 & 0.30 & -0.18 \\ 0.19 & -0.09 & -0.12 \\ 0.10 & -0.05 & -0.11 \\ 0.19 & -0.09 & -0.12 \\ 0.17 & -0.09 & 0.02 \\ 0.02 & 0.40 & -0.50 \\ 0.12 & 0.05 & -0.48 \\ 0.16 & 0.41 & 0.32 \end{bmatrix} \end{aligned}$$

4. The three columns of V are word vectors in the 14D space of counts of the dictionary words (5.3) Algorithm, Application, Differential, Dynamic, Equation, Implementation, Integral, Method, Nonlinear, Ordinary, Partial, Problem, System, and Theory.

- The first column \mathbf{v}_1 of V , from its largest components, mainly identifies the two most common words of Differential and Equation.
- The second column \mathbf{v}_2 of V , from its largest components, identifies books with Algorithms, Applications, Implementations, Problems, and Theory.
- The third column \mathbf{v}_3 of V , from its largest components, largely distinguishes Dynamics, Problems and Systems, from Differential and Theory.

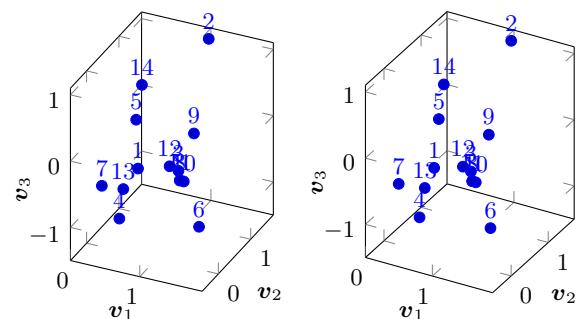
The corresponding principal components are the entries of the 14×3 matrix (2 d.p.)

$$AV = \begin{bmatrix} 0.70 & 0.09 & -0.25 \\ 1.02 & 1.59 & 1.00 \\ 1.46 & -0.29 & 0.10 \\ 0.16 & 0.67 & -1.44 \\ 0.16 & 1.19 & -0.22 \\ 1.78 & -0.34 & -0.64 \\ 0.22 & -0.00 & -0.58 \\ 1.48 & -0.29 & -0.04 \\ 1.45 & 0.21 & 0.40 \\ 1.56 & -0.34 & -0.01 \\ 1.48 & -0.29 & -0.04 \\ 1.29 & -0.20 & 0.08 \\ 0.10 & 0.92 & -1.14 \\ 0.29 & 1.09 & 0.39 \end{bmatrix}.$$

Each of the fourteen books is represented in 3D space by the corresponding row of these coordinates. Plot these books in Matlab/Octave with

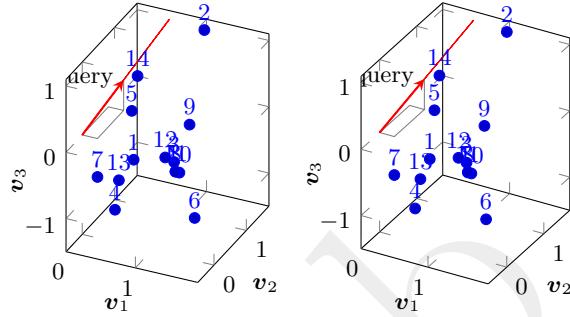
```
scatter3(A*V(:,1),A*V(:,2),A*V(:,3),[],1:14)
```

as shown below in stereo.



There is a cluster of five books near the front along the \mathbf{v}_1 -axis (numbered 3, 8, 10, 11 and 12, their focus is Differential Equations), the other nine are spread out.

Queries Suppose we search for books on *Application and Theory*. In our dictionary (5.3), the corresponding word vector for this search is $\mathbf{w} = (1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1)$. Project this query into the 3D space of principal components with the product $\mathbf{w}^T V$ which evaluates to the query vector $\mathbf{q} = (0.22, 0.81, 0.46)$ whose direction is added to the picture as shown below.



Books 2 and 14 appear close to the direction of the query vector and so should be returned as a match: these books are no surprise as both their titles have both *Application* and *Theory* in their titles. But the above plot also suggests Book 5 is near to the direction of the query vector, and so is also worth considering despite not having either of the search words in its title! The power of this latent semantic indexing is that it extracts additional titles that are relevant to the query yet share no common words with the query—as commented at the start of this section.

The angles between the query vector and the book title 3D vectors confirm the graphical appearance claimed above. Recall that the dot product determines the angle between vectors (Theorem 1.3.4).

- From the second row of the above product AV , Book 2 has the principal component vector $(1.02, 1.59, 1.00)$ which has length 2.14. Consequently, it is at small angle 15° to the 3D query vector $\mathbf{q} = (0.22, 0.81, 0.46)$, of length $|\mathbf{q}| = 0.96$, because its cosine

$$\cos \theta = \frac{(1.02, 1.59, 1.00) \cdot \mathbf{q}}{2.14 \cdot 0.96} = 0.97.$$

- Similarly, Book 14 has the principal component vector $(0.29, 1.09, 0.39)$ which has length 1.20. Consequently, it is at small angle 10° to the 3D query vector $\mathbf{q} = (0.22, 0.81, 0.46)$ because its cosine

$$\cos \theta = \frac{(0.29, 1.09, 0.39) \cdot \mathbf{q}}{1.20 \cdot 0.96} = 0.99.$$

- Whereas Book 5 has the principal component vector $(0.16, 1.19, -0.22)$ which has length 1.22. Consequently, it is at moderate angle 40° to the 3D query vector $\mathbf{q} = (0.22, 0.81, 0.46)$ because its cosine

$$\cos \theta = \frac{(0.16, 1.19, -0.22) \cdot \mathbf{q}}{1.20 \cdot 0.96} = 0.76.$$

Such a significant cosine suggests that Book 5 is also of interest.

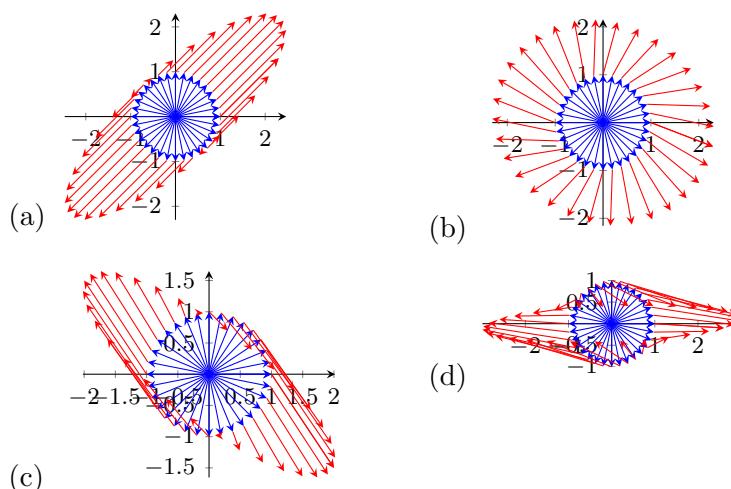
If we were to compute the angles in the original 14D space of the full dictionary (5.3), then the title of Book 5 would be orthogonal to the query, because it has no words in common, and so Book 5 would not be flagged as of interest. The principal component analysis reduces the dimensionality to those relatively few directions that are important, and it is in these important directions that the title of Book 5 appears promising for the query.

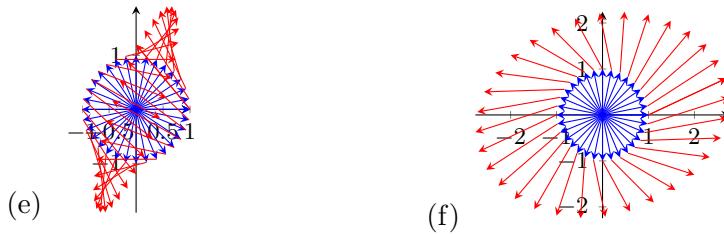
- All the other book titles have angles greater than 62° and so are significantly less related to the query.

Latent semantic indexing in practice This application of principal components to analysing a few book titles is purely indicative. In practice one would analyse the many thousands of words used throughout hundreds or thousands of documents. Moreover, one would be interested in not just plotting the documents in a 2D plane or 3D space, but in representing the documents in say a 70D space of seventy principal components. Berry et al. (1995) reviews how such statistically derived principal word vectors are a more robust indicator of meaning than individual terms. Hence this SVD analysis of documents becomes an effective way of retrieving information from a search without requiring the results actually match any of the words in the search request—the results just need to match cognate words.

5.1.4 Exercises

Exercise 5.1.1. For some 2×2 matrices A the following plots adjoin the product $A\mathbf{x}$ to \mathbf{x} for a complete range of unit vectors \mathbf{x} . Use each plot to roughly estimate the norm of the underlying matrix for that plot.





Exercise 5.1.2. For the following matrices, use few unit vectors \mathbf{x} to determine how the matrix-vector product varies with \mathbf{x} . Using calculus to find the appropriate maximum, find from definition the norm of the matrices (hint: all norms here are integers).

$$(a) \begin{bmatrix} 2 & -3 \\ 0 & 2 \end{bmatrix}$$

$$(b) \begin{bmatrix} -4 & -1 \\ -1 & -4 \end{bmatrix}$$

$$(c) \begin{bmatrix} 2 & -1 \\ 1 & -2 \end{bmatrix}$$

$$(d) \begin{bmatrix} 2 & -2 \\ 2 & 1 \end{bmatrix}$$

$$(e) \begin{bmatrix} 5 & 2 \\ 2 & 2 \end{bmatrix}$$

$$(f) \begin{bmatrix} 6 & -1 \\ 4 & 6 \end{bmatrix}$$

$$(g) \begin{bmatrix} 2 & 4 \\ -7 & -2 \end{bmatrix}$$

$$(h) \begin{bmatrix} 2 & -4 \\ 1 & -2 \end{bmatrix}$$

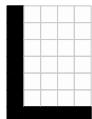
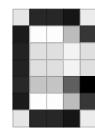
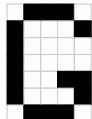
Exercise 5.1.3. Many properties of a norm may be proved in other ways than those given for Theorem 5.1.9.

- (a) Use an SVD factorisation of I_n to prove Theorem 5.1.9b.
- (b) Use the existence of an SVD factorisation of A to prove Theorem 5.1.9d.
- (c) Use the definition of a norm as a maximum to prove Theorem 5.1.9f.
- (d) Use the existence of an SVD factorisation of A to prove Theorem 5.1.9g.
- (e) Prove Theorem 5.1.9h using 5.1.9g and the definition of a norm as a maximum.

Exercise 5.1.4. Let $m \times n$ matrix A have in each row and column at most one non-zero element. Argue that there exists an SVD which establishes that the norm $\|A\| = \max_{i,j} |a_{ij}|$.

Exercise 5.1.5.

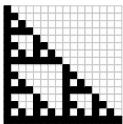
The margin shows a 7×5 pixel image of the letter G. Compute an SVD of the pixel image. By inspecting various rank approximations from this SVD, determine the rank of the approximation to G shown to the right.



Exercise 5.1.6. Write down two different rank two representations of the pixel image of the letter L, as shown in the margin. Compute SVD representations of the letter L. Compare and comment on the various representations.

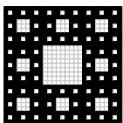
Exercise 5.1.7 (Sierpinski triangle). Whereas mostly we deal with the smooth geometry of lines, planes and curves, the subject of fractal geometry recognises that there is much in the world around us that has a rich fractured structure: from clouds, and rocks, to the cardiovascular structure within each of us. The Sierpinski triangle, illustrated in the margin, is a simple fractal. Generate such fractals using recursion as in the following Matlab/Octave code⁵: the recursion is that the next generation image A is computed from the previous generation A .

```
A=1
A=[A 0*A;A A]
A=[A 0*A;A A]
A=[A 0*A;A A]
A=[A 0*A;A A]
imagesc(1-A)
colormap('gray')
axis equal, axis off
```

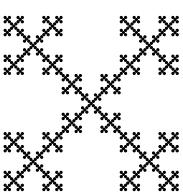


- (a) Add code to this recursion to compute and print the singular values of each generation of the Sierpinski triangle image. What do you conjecture about the number of distinct singular values as a function of generation number k ? Test your conjecture for more iterations in the recursion.
- (b) Returning to the 16×16 Sierpinski triangle formed after four iterations, use an SVD to form the best rank five approximation of the Sierpinski triangle (as illustrated in the margin, it has a beautiful structure). Comment on why such a rank five approximation may be a reasonable one to draw. What is the next rank for an approximation that is reasonable to draw? Justify.
- (c) Modify the code to compute and draw the 256×256 image of the Sierpinski triangle. Use an SVD to generate and draw the best rank nine approximation to the image.

Exercise 5.1.8 (Sierpinski carpet). The Sierpinski carpet is another fractal easily generated by recursion (as illustrated in the margin after three iterations in the recursion).



- (a) Modify the recursive Matlab/Octave code of the previous Exercise 5.1.7 to generate such an image.
- (b) For a range of generations of the image, compute the singular values and comment on the apparent patterns in the singular values.



Exercise 5.1.9 (another fractal). Illustrated in the margin is another fractal easily generated by recursion.

⁵ Many students would know that a for-loop would more concisely compute the recursion.

- (a) Modify the recursive Matlab/Octave code of Exercise 5.1.7 to generate such an image.
- (b) For a range of generations of the image, compute the singular values and comment on the apparent patterns in the singular values.

Exercise 5.1.10.

Ada Lovelace (1815–52)



This is an image of Countess Ada Lovelace: the first computer programmer, she invented, developed and wrote programs for Charles Babbage's analytical engine. Download the 249×178 image. Using an SVD draw various rank approximations to the image. Using the matrix norm to measures errors, what is the smallest rank to reproduce the image to an error of 5%? and of 1%?

http://www.maa.org/sites/default/files/images/upload_library/46/Portraits/Lovelace_Ada.jpg [Sep 2015]

Exercise 5.1.11 (hiding information project). Steganographic methods embed secret messages in images. Some methods are based upon the singular value decomposition (Gorodetski et al. 2001, e.g.). Perhaps the easiest method is to use the unimportant small singular values. For example, Exercise 5.1.10 indicates a rank 32 approximation of the image of Ada Lovelace accurately reproduces the image. Using the SVD $A = USV^T$, let's keep the singular values $\sigma_1, \sigma_2, \dots, \sigma_{32}$ the same to produce a good image. The unimportant singular values $\sigma_{33}, \sigma_{34}, \sigma_{35}, \dots$ can hide a small message in the image. Suppose the message is forty binary digits such as

0000001000000100000110001000011010001111

The scheme is to set the unimportant singular values recursively

$$\sigma_{32+k} = \begin{cases} 0.99\sigma_{31+k}, & \text{if } k\text{th digit is one,} \\ 0.96\sigma_{31+k}, & \text{if } k\text{th digit is zero.} \end{cases}$$

These ratios achieve two things: first, the singular values are decreasing as is necessary to maintain the order of the singular vectors in the standard computation of an SVD; second the ratios are sufficiently different to be robustly detected in an image.

- Download the 249×178 image A .
- Compute an SVD, $A = USV^T$.

- Change the singular values matrix S to S' with the first 32 singular values unchanged, and the next 40 singular values encoding the message. The Matlab/Octave function `cumprod()` neatly computes the recursive product.
- Using the first 72 columns of U and V , compute and draw the new image $B = \text{round}(US'V^T)$ (as shown in the margin it is essentially the same as the original), where `round()` is the Matlab/Octave function that rounds the real values to the nearest integer (greyscale values should be in the range zero to 255).



It is this image that contains the hidden message.

- To check the hidden message is recoverable, compute an SVD of the new image, $B = QDR^T$. Compare the singular values in D , say $\delta_1, \delta_2, \dots, \delta_{72}$, with those of S' : comment on the effect of the rounding in the computation of B . Invent and test Matlab/Octave commands to extract the hidden message: perhaps use `diff(log())` to undo much of the recursive product in the singular values.
- Report on all code, its role, and the results.

Exercise 5.1.12. As in Example 5.1.22, consider the 2D plot of the six books. Add to the plot the word vector corresponding to a query for books relevant to “Introduction”. Using angle to measure closeness, which book is closest to “Introduction”? Confirm by computing the angles, in the 2D plane, between “Introduction” and all six books.

Exercise 5.1.13. As in Example 5.1.22, consider the 2D plot of the six books. Add to the plot the word vector corresponding to a query for books relevant to “Application and Method”. Using angle to measure closeness, which book is closest to “Application and Method”? Confirm by computing, in the 2D plane, the angles between “Application and Method” and all six books.

Exercise 5.1.14. Reconsider the word vectors of the group of fourteen mathematical books listed in the text. Instead of computing a representation of these books in 3D space, here use principal component analysis to compute a representation in a 2D plane.

- (a) Plot the titles of the fourteen books in the 2D plane of the two principal vectors.
- (b) Add the vector corresponding to the query of *Differential and Equation*: which books appear to have a small angle to this query in this 2D representation?
- (c) Add the vector corresponding to the query of *Application and Theory*: which books appear to have a small angle to this query in this 2D representation?

Table 5.3: twenty user reviews of bathrooms in a major chain of hotels. The data comes from the Opinosis Opinion/Review in the UCI Machine Learning Repository.

- The room was not overly big, but clean and very comfortable beds, a great shower and very clean bathrooms
- The second room was smaller, with a very inconvenient bathroom layout, but at least it was quieter and we were able to sleep
- Large comfortable room, wonderful bathroom
- The rooms were nice, very comfy bed and very clean bathroom
- Bathroom was spacious too and very clean
- The bathroom only had a single sink, but it was very large
- The room was a standard but nice motel room like any other, bathroom seemed upgraded if I remember
- The room was quite small but perfectly formed with a super bathroom
- You could eat off the bathroom floor it was so clean
- The bathroom door does the same thing, making the bathroom seem slightly larger
- bathroom spotless and nicely appointed
- The rooms are exceptionally clean and also the bathrooms
- The bathroom was clean and the bed was comfy
- They provide you with great aveda products in the bathroom
- Also, the bathroom was a bit dirty , brown water came out of the bath tub faucet initially and the sink wall by the toilet was dirty
- If your dog tends to be a little disruptive or on the noisy side, there is a bathroom fan that you can keep on to make noise
- The bathroom was big and clean as well
- Also, the bathrooms were quite well set up, with a separate toilet shower to basin, so whilst one guest is showering another can use the basin
- The bathroom was marble and we had luxurious bathrobes and really, every detail attended to
- It was very clean, had a beautiful bathroom, and was comfortable

Exercise 5.1.15. The discussion on latent semantic indexing only considered queries which were “and”, such as a search for books relevant to *Application and Theory*. What if we wanted to search for books relevant to *either Application or Theory*? Discuss how such an “or” query might be phrased in terms of the angle between vectors and a multi-D subspace.

Exercise 5.1.16. Table 5.3 lists twenty short reviews about bathrooms of a major chain of hotels.⁶ There are about 17 meaningful words

⁶ From <http://archive.ics.uci.edu/ml>

common to more than one review: create a list of such words. Then form corresponding word vectors for each review. Use an SVD to best plot these reviews in 2D. Discuss any patterns in the results.

VO-1b

5.2 Regularise linear equations

Section Contents

5.2.1	The SVD illuminates regularisation	483
5.2.2	Tikhonov regularisation	498
5.2.3	Exercises	502

Singularity is almost invariably a clue.

Sherlock Holmes, in The Boscombe Valley Mystery, by Sir Arthur Conan Doyle, 1892

Often we need to approximate matrices involved in solving linear equations. This is especially so when the matrix itself comes from experimental measurements and so is subject to errors. We do not want such errors to affect results. By avoiding division with small singular values, the procedure developed in this section avoids unwarranted magnification of errors. Sometimes such error magnification is disastrous, so avoiding it is essential.

Example 5.2.1. Suppose from some experiment we to solve the linear equations

$$0.5x + 0.3y = 1 \quad \text{and} \quad 1.1x + 0.7y = 2,$$

where all the coefficients on *both* the left-hand sides and the right-hand sides are determined from experimental measurements. In particular, suppose they are measured to errors ± 0.05 . Solve the equations.

Solution: • Following Procedure 2.2.4 in Matlab/Octave, form the matrix and the right-hand side

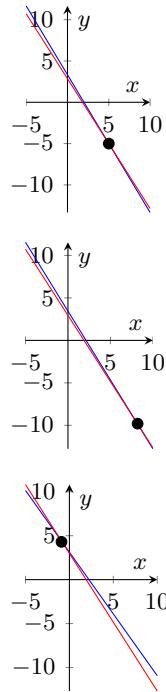
$$\begin{aligned} A &= [0.5 \ 0.3; 1.1 \ 0.7] \\ b &= [1.0; 2.0] \end{aligned}$$

Then check the condition number with `rcond(A)` to find it is 0.007 which previously we would call only just outside the ‘good’ range (Procedure 2.2.4). So proceed to compute the solution with `A\b` to find $(x, y) = (5, -5)$ (as illustrated in the margin).

- Is this solution reasonable? No. Not when the matrix itself has errors. Let’s perturb the matrix A by amounts consistent with its experimental error of ± 0.05 and explore the predicted solutions (the first two illustrated):

$$A = \begin{bmatrix} 0.47 & 0.29 \\ 1.06 & 0.68 \end{bmatrix} \implies \mathbf{x} = \mathbf{A}\backslash\mathbf{b} = (8.2, -9.8);$$

$$A = \begin{bmatrix} 0.45 & 0.32 \\ 1.05 & 0.67 \end{bmatrix} \implies \mathbf{x} = \mathbf{A}\backslash\mathbf{b} = (-0.9, 4.3);$$



$$A = \begin{bmatrix} 0.46 & 0.31 \\ 1.06 & 0.73 \end{bmatrix} \implies \mathbf{x} = \mathbf{A} \setminus \mathbf{b} = (15.3, -19.4).$$

For equally valid matrices A , the predicted solutions are all over the place!

- If the matrix itself has errors, then we must reconsider Procedure 2.2.4. The SVD empowers a sensible resolution. Compute an SVD of the matrix, $A = USV^T$, with $[U, S, V] = \text{svd}(A)$ to find (2 d.p.)



$$A = \begin{bmatrix} -0.41 & -0.91 \\ -0.91 & 0.41 \end{bmatrix} \begin{bmatrix} 1.43 & 0 \\ 0 & 0.01 \end{bmatrix} \begin{bmatrix} -0.85 & -0.53 \\ -0.53 & 0.85 \end{bmatrix}^T$$

Because the matrix A has errors ± 0.05 , the small singular value of 0.01 might as well be zero: it is zero to experimental error. The appropriate solution algorithm is then Procedure 3.5.3 for inconsistent equations (not Procedure 2.2.4). Thus (2 d.p.)

- (a) $\mathbf{z} = U^T \mathbf{b} = (-2.23, -0.10)$;
 - (b) due to the approximately zero singular value, neglect $z_2 = -0.10$ as an error, and solve $\begin{bmatrix} 1.43 & 0 \\ 0 & 0 \end{bmatrix} \mathbf{y} = \begin{bmatrix} -2.23 \\ 0 \end{bmatrix}$ to deduce $\mathbf{y} = (-1.56, y_2)$;
 - (c) consequently, we find reasonable solutions are $\mathbf{x} = V\mathbf{y} = (1.32, 0.83) + y_2(-0.53, 0.85)$.
- To choose between this infinitude of solutions, extra information must be provided by the context/application/modelling. For example, often one prefers the solution of smallest length/magnitude, obtained by setting $y_2 = 0$ (Theorem 3.5.9); that is, $\mathbf{x}_{\text{smallest}} = (1.32, 0.83)$.

■

5.2.1 The SVD illuminates regularisation

I think it is much more interesting to live with uncertainty than to live with answers that might be wrong.

Richard Feynman

Procedure 5.2.2 (approximate linear equations). *Suppose the system of linear equation $A\mathbf{x} = \mathbf{b}$ arises from experiment where both the $m \times n$ matrix A and the right-hand side vector \mathbf{b} are subject to experimental error. Suppose the expected error in the matrix entries are of size e (here “e” denotes “error”, not the exponential e)*

1. When forming the matrix A and vector \mathbf{b} , scale the data so that

- all $m \times n$ components in A have the same physical units, and they are of roughly the same size; and
- similarly for \mathbf{b} .

Determine the error e corresponding to this matrix A .

2. Compute an SVD $A = USV^T$.
3. Choose ‘rank’ k to be the number of singular values bigger than the error e ; that is, $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k > e > \sigma_{k+1} \geq \dots \geq 0$. Then the rank k approximation to A is

$$\begin{aligned} A_k &:= US_k V^T \\ &= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \dots + \sigma_k \mathbf{u}_k \mathbf{v}_k^T \\ &= \mathbf{U}(:, 1:k) * \mathbf{S}(1:k, 1:k) * \mathbf{V}(:, 1:k)^T. \end{aligned}$$

We usually do not construct A_k as we only need its SVD to solve the system.

4. Solve the approximating linear equation $A_k \mathbf{x} = \mathbf{b}$ as in Theorems 3.5.5–3.5.9 (often as an inconsistent set of equations). Usually use the SVD $A_k = US_k V^T$.
5. Among all the solutions allowed, choose the ‘best’ according to the needs of the application: often the smallest solution overall; or equally often a solution with the most zero components.

That is, treat as zero any singular values smaller than the expected error in the matrix entries. For example, in computation on modern computers (with nearly sixteen significant decimal digits accuracy) any singular value smaller than $10^{-8}\sigma_1$ should be treated as zero. Certainly, any smaller than $10^{-16}\sigma_1$ must be treated as zero.

The final step in Procedure 5.2.2 arises because in many cases an infinite number of possible solutions are derived. The linear algebra cannot presume which is best for your application. Consequently, you will have to be aware of the freedom, and make a choice based on extra information from your particular application.

- For example, in a CT scan such as Example 3.5.12 one would usually prefer the greyest result in order to avoid diagnosing artifices.
- For another example, in the data mining task of fitting curves or surfaces to data, one would instead usually prefer a curve or surface with fewest non-zero coefficients.

Such extra information from the application is essential.

Example 5.2.3. For the following matrices A and right-hand side vectors \mathbf{b} , solve $A\mathbf{x} = \mathbf{b}$. But suppose the matrix entries come from experiments and are only known to within errors ± 0.05 , solve $A'\mathbf{x}' = \mathbf{b}$ for some matrices A' which approximate A to this error. Finally,

use an SVD to find a general solution consistent with the error in matrix A . Report to two decimal places.

$$(a) \quad A = \begin{bmatrix} -0.2 & -0.6 & 1.8 \\ 0.0 & 0.2 & -0.4 \\ -0.3 & 0.7 & 0.3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -0.5 \\ 0.1 \\ -0.2 \end{bmatrix}$$



Solution: Enter the matrix and vector into Matlab/Octave and solve with $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$ to determine $\mathbf{x} = (0.06, -0.13, -0.31)$. To within the experimental error of ± 0.05 the following matrices approximate A : that is, they might have been what was measured for A . Then Matlab/Octave $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$ gives the corresponding equally valid solutions.

- $A' = \begin{bmatrix} -0.16 & -0.58 & 1.83 \\ 0.01 & 0.16 & -0.45 \\ -0.28 & 0.74 & 0.30 \end{bmatrix}$ gives $\mathbf{x}' = \begin{bmatrix} 0.85 \\ 0.12 \\ -0.16 \end{bmatrix}$
- $A'' = \begin{bmatrix} -0.22 & -0.62 & 1.77 \\ 0.01 & 0.17 & -0.42 \\ -0.26 & 0.66 & 0.26 \end{bmatrix}$ gives $\mathbf{x}'' = \begin{bmatrix} 0.42 \\ -0.04 \\ -0.24 \end{bmatrix}$

There are major differences between these equally valid solutions \mathbf{x} , \mathbf{x}' and \mathbf{x}'' . The problem is that, relative to the experimental error, there is a small singular value in the matrix A . We must use an SVD to find all solutions consistent with the experimental error. Consequently, compute $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{A})$ to find (2 d.p.)

```

 $\mathbf{U} =$ 
    -0.97    0.03   -0.23
     0.22   -0.10   -0.97
    -0.05   -0.99    0.09

 $\mathbf{S} =$ 
    1.96      0      0
     0    0.82      0
     0      0    0.02

 $\mathbf{V} =$ 
    0.11    0.36    0.93
    0.30   -0.90    0.31
   -0.95   -0.25    0.20
  
```

The singular value 0.02 is less than the error ± 0.05 so is effectively zero. Hence we solve the system as if this singular value is zero; that is, as if matrix A has rank two. Compute the smallest consistent solution with the three steps $\mathbf{z}=\mathbf{U}(:,1:2)' * \mathbf{b}$, $\mathbf{y}=\mathbf{z} ./ \text{diag}(\mathbf{S}(1:2,1:2))$, and $\mathbf{x}=\mathbf{V}(:,1:2) * \mathbf{y}$. Then add an arbitrary multiple of the last column of \mathbf{V} to determine a general solution $\mathbf{x} = (0.10, -0.11, -0.30) + t(0.93, 0.31, 0.20)$.

$$(b) \quad A = \begin{bmatrix} -1.1 & 0.1 & 0.7 & -0.1 \\ 0.1 & -0.1 & 1.2 & -0.6 \\ 0.8 & -0.2 & 0.4 & -0.8 \\ 0.8 & 0.1 & -2.0 & 1.0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -1.1 \\ -0.1 \\ 1.1 \\ 0.8 \end{bmatrix}$$

Solution: Enter the matrix and vector into Matlab/Octave and solve with $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$ to determine $\mathbf{x} = (0.61, -0.64, -0.65, -0.93)$.

To within experimental error the following matrices approximate A , and Matlab/Octave $\mathbf{x}=\mathbf{A}\backslash\mathbf{b}$ gives the corresponding solutions.

$$\bullet \quad A' = \begin{bmatrix} -1.10 & 0.11 & 0.67 & -0.08 \\ 0.08 & -0.10 & 1.17 & -0.59 \\ 0.75 & -0.21 & 0.39 & -0.83 \\ 0.79 & 0.08 & -1.96 & 0.98 \end{bmatrix}, \quad \mathbf{x}' = \begin{bmatrix} 0.64 \\ -0.40 \\ -0.64 \\ -0.95 \end{bmatrix}$$

$$\bullet \quad A'' = \begin{bmatrix} -1.10 & 0.08 & 0.66 & -0.14 \\ 0.08 & -0.09 & 1.22 & -0.58 \\ 0.77 & -0.18 & 0.39 & -0.78 \\ 0.75 & 0.11 & -2.01 & 1.04 \end{bmatrix}, \quad \mathbf{x}'' = \begin{bmatrix} 0.87 \\ 1.09 \\ -0.58 \\ -1.09 \end{bmatrix}$$

There are significant differences, mainly in the second component, between these equally valid solutions \mathbf{x} , \mathbf{x}' and \mathbf{x}'' . The problem is that, relative to the experimental error, there is a small singular value in the matrix A . We must use an SVD to find all solutions consistent with the experimental error: compute $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{A})$ to find (2 d.p.)

$$\mathbf{U} = \begin{bmatrix} -0.33 & 0.59 & -0.36 & 0.64 \\ -0.43 & -0.31 & 0.71 & 0.46 \\ -0.16 & -0.74 & -0.59 & 0.27 \\ 0.82 & -0.07 & 0.12 & 0.55 \end{bmatrix}$$

$$\mathbf{S} = \begin{bmatrix} 2.89 & 0 & 0 & 0 \\ 0 & 1.50 & 0 & 0 \\ 0 & 0 & 0.26 & 0 \\ 0 & 0 & 0 & 0.02 \end{bmatrix}$$

$$\mathbf{V} = \begin{bmatrix} 0.30 & -0.88 & 0.35 & 0.09 \\ 0.04 & 0.15 & 0.09 & 0.98 \\ -0.85 & -0.08 & 0.52 & 0.00 \\ 0.43 & 0.43 & 0.78 & -0.16 \end{bmatrix}$$

The singular value 0.02 is less than the error ± 0.05 so is effectively zero. Hence solve the system as if this singular value is zero; that is, as if matrix A has rank three. Compute the smallest consistent solution with $\mathbf{z}=\mathbf{U}(:,1:3)' * \mathbf{b}$, $\mathbf{y}=\mathbf{z} ./ \text{diag}(\mathbf{S}(1:3,1:3))$, and $\mathbf{x}=\mathbf{V}(:,1:3) * \mathbf{y}$. Then add an arbitrary multiple of the last column of \mathbf{V} to determine a general solution $\mathbf{x} = (0.65, -0.22, -0.65, -1.00) + t(0.09, 0.98, 0, -0.16)$.

That the second component of $(0.09, 0.98, 0, -0.16)$ is the largest corresponds to the second component in each of \mathbf{x} , \mathbf{x}'



and \mathbf{x}'' being the most sensitive, as seen in the above three cases.

Both of these examples gave an infinite number of solutions which are equally valid as far as the linear algebra is concerned. In each example, more information from an application would be needed to choose which of the infinity of solutions is preferred. ■

Most often the singular values are spread over a wide range of orders of magnitude. In such cases an assessment of the errors in the matrix is crucial in what one reports as a solution. The following artificial example illustrates the range.

Example 5.2.4 (various errors). The matrix

$$A = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \frac{1}{8} \\ \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \frac{1}{8} & \frac{1}{9} \end{bmatrix}$$

is an example of a so-called Hilbert matrix. Explore the effects of various assumptions about possible errors in A upon the solution to $A\mathbf{x} = \mathbf{1}$ where $\mathbf{1} := (1, 1, 1, 1, 1)$.

Solution: Enter the matrix A into Matlab/Octave with `A=hilb(5)` for the above 5×5 Hilbert matrix, and enter the right-hand side with `b=ones(5,1)`.

- First assume there is insignificant error in A (there is always the base error of 10^{-15} in computation). Then Procedure 2.2.4 finds that although the reciprocal of the condition number `rcond(A) ≈ 10-6` is bad, the unique solution to $A\mathbf{x} = \mathbf{1}$, obtained via `x=A\b`, is

$$\mathbf{x} = (5, -120, 630, -1120, 630).$$



- Second suppose the errors in A are roughly 10^{-5} . This level of error is a concern as `rcond ≈ 10-6` so errors would be magnified by 10^6 in a direct solution of $A\mathbf{x} = \mathbf{1}$. Here we explore when all errors are in A and none in the right-hand side vector $\mathbf{1}$. To explore, adopt Procedure 5.2.2.

- (a) Find an SVD $A = USV^T$ via `[U,S,V]=svd(A)` (2 d.p.)

$U =$

$$\begin{array}{ccccc} -0.77 & 0.60 & -0.21 & 0.05 & 0.01 \\ -0.45 & -0.28 & 0.72 & -0.43 & -0.12 \\ -0.32 & -0.42 & 0.12 & 0.67 & 0.51 \\ -0.25 & -0.44 & -0.31 & 0.23 & -0.77 \end{array}$$

$$\begin{aligned}
 S &= \begin{matrix} -0.21 & -0.43 & -0.57 & -0.56 & 0.38 \\ 1.57 & 0 & 0 & 0 & 0 \\ 0 & 0.21 & 0 & 0 & 0 \\ 0 & 0 & 0.01 & 0 & 0 \\ 0 & 0 & 0 & 0.00 & 0 \\ 0 & 0 & 0 & 0 & 0.00 \end{matrix} \\
 V &= \begin{matrix} -0.77 & 0.60 & -0.21 & 0.05 & 0.01 \\ -0.45 & -0.28 & 0.72 & -0.43 & -0.12 \\ -0.32 & -0.42 & 0.12 & 0.67 & 0.51 \\ -0.25 & -0.44 & -0.31 & 0.23 & -0.77 \\ -0.21 & -0.43 & -0.57 & -0.56 & 0.38 \end{matrix}
 \end{aligned}$$

More informatively, the singular values have the following wide range of magnitudes,

$$\begin{aligned}
 \sigma_1 &= 1.57, & \sigma_2 &= 0.21, & \sigma_3 &= 1.14 \cdot 10^{-2}, \\
 \sigma_4 &= 3.06 \cdot 10^{-4}, & \sigma_5 &= 3.29 \cdot 10^{-6}.
 \end{aligned}$$

- (b) Because the assumed error 10^{-5} satisfies $\sigma_4 > 10^{-5} > \sigma_5$ the matrix A is effectively of rank four, $k = 4$.
- (c) Solving the system $A\mathbf{x} = USV^T\mathbf{x} = \mathbf{1}$ as rank four, in the least square sense, Procedure 3.5.3 gives (2 d.p.)
 - i. $\mathbf{z} = U^T\mathbf{1} = (-2.00, -0.97, -0.24, -0.04, 0.00)$,
 - ii. neglect the fifth component of \mathbf{z} as an error and obtain the first four components of \mathbf{y} via $\mathbf{y} = \mathbf{z}(1:4) ./ \text{diag}(S(1:4, 1:4))$ so that

$$\mathbf{y} = (-1.28, -4.66, -21.43, -139.69, y_5),$$

- iii. then the smallest, least square, solution determined with $\mathbf{x} = V(:, 1:4) * \mathbf{y}$ is

$$\mathbf{x} = (-3.82, 46.78, -93.41, -23.53, 92.27),$$

and a general solution includes the arbitrary multiple y_5 of the last column of V to be

$$\begin{aligned}
 \mathbf{x} &= (-3.82, 46.78, -93.41, -23.53, 92.27) \\
 &\quad + y_5(0.01, -0.12, 0.51, -0.77, 0.38).
 \end{aligned}$$

- Third suppose the errors in A are roughly 10^{-3} . Re-adopt Procedure 5.2.2.
 - (a) Use the same SVD, $A = USV^T$.
 - (b) Because the assumed error 10^{-3} satisfies $\sigma_3 > 10^{-3} > \sigma_4$ the matrix A is effectively of rank three, $k = 3$.

- (c) Solving the system $A\mathbf{x} = USV^T\mathbf{x} = \mathbf{1}$ as rank three, in the least square sense, Procedure 3.5.3 gives (2 d.p.) the same \mathbf{z} , and the same first three components in

$$\mathbf{y} = (-1.28, -4.66, -21.43, y_4, y_5),$$

then the smallest, least square, solution determined with $\mathbf{x} = \mathbf{V}(:, 1:3) * \mathbf{y}$ is

$$\mathbf{x} = (2.76, -13.66, -0.19, 9.03, 14.38),$$

and a general solution includes the arbitrary multiples of the last columns of V to be

$$\begin{aligned}\mathbf{x} &= (2.76, -13.66, -0.19, 9.03, 14.38) \\ &\quad + y_4(0.05, -0.43, 0.67, 0.23, -0.56) \\ &\quad + y_5(0.01, -0.12, 0.51, -0.77, 0.38).\end{aligned}$$

- Lastly suppose the errors in A are roughly 0.05. Re-adopt Procedure 5.2.2.
 - (a) Use the same SVD, $A = USV^T$.
 - (b) Because the assumed error 0.05 satisfies $\sigma_2 > 0.05 > \sigma_3$ the matrix A is effectively of rank two, $k = 2$.
 - (c) Solving the system $A\mathbf{x} = USV^T\mathbf{x} = \mathbf{1}$ as rank two, in the least square sense, Procedure 3.5.3 gives (2 d.p.) the same \mathbf{z} , and the same first two components in

$$\mathbf{y} = (-1.28, -4.66, y_3, y_4, y_5),$$

then the smallest, least square, solution determined with $\mathbf{x} = \mathbf{V}(:, 1:2) * \mathbf{y}$ is

$$\mathbf{x} = (-1.83, 1.85, 2.39, 2.39, 2.27),$$

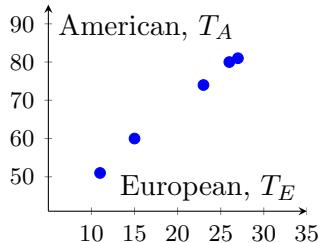
and a general solution includes the arbitrary multiples of the last columns of V to be

$$\begin{aligned}\mathbf{x} &= (-1.83, 1.85, 2.39, 2.39, 2.27) \\ &\quad + y_3(-0.21, 0.72, 0.12, -0.31, -0.57) \\ &\quad + y_4(0.05, -0.43, 0.67, 0.23, -0.56) \\ &\quad + y_5(0.01, -0.12, 0.51, -0.77, 0.38).\end{aligned}$$

The level of error makes a major difference in the qualitative nature of allowable solutions: here from a unique solution through to a three parameter family of equally valid solutions. To appropriately solve systems of linear equations we must know the level of error.



Example 5.2.5 (translating temperatures). Recall Example 2.2.10 attempts to fit a quartic polynomial to observations (plotted in the margin) of the relation between Celsius and Fahrenheit temperature. The attempt failed because `rcond` is too small. Let's try again now that we can cater for matrices with errors. Recall the data between temperatures reported by a European and an American are the following:



$$\begin{array}{c|ccccc} T_E & 15 & 26 & 11 & 23 & 27 \\ \hline T_A & 60 & 80 & 51 & 74 & 81 \end{array}$$

Example 2.2.10 attempts to fit the data with the quartic polynomial

$$T_A = c_1 + c_2 T_E + c_3 T_E^2 + c_4 T_E^3 + c_5 T_E^4,$$

and deduced the following system of equations for the coefficients

$$\begin{bmatrix} 1 & 15 & 225 & 3375 & 50625 \\ 1 & 26 & 676 & 17576 & 456976 \\ 1 & 11 & 121 & 1331 & 14641 \\ 1 & 23 & 529 & 12167 & 279841 \\ 1 & 27 & 729 & 19683 & 531441 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \end{bmatrix} = \begin{bmatrix} 60 \\ 80 \\ 51 \\ 74 \\ 81 \end{bmatrix}.$$

In order to find a robust solution, here let's approximate both the matrix and the right-hand side vector because both the matrix and the vector come from real data with errors of about up to $\pm 0.5^\circ$.

Solution: Now invoke Procedure 5.2.2 to approximate the system of linear equations, and solve the approximate problem.

- (a) There is a problem in approximating the matrix: the columns are of wildly different sizes. In contrast, our mathematical analysis treats all columns the same. The problem is that each column comes from different powers of temperatures. To avoid the problem we must scale the temperature data for the matrix. The simplest is to divide by a typical temperature. That is, instead of seeking a fit in terms of powers of T_E , we seek a fit in powers of $T_E/20^\circ$ as 20 degrees is a typical temperature in the data. Here we fit the data with the quartic polynomial

$$T_A = c_1 + c_2 \frac{T_E}{20} + c_3 \left(\frac{T_E}{20}\right)^2 + c_4 \left(\frac{T_E}{20}\right)^3 + c_5 \left(\frac{T_E}{20}\right)^4,$$

which gives the following system for the coefficients (2 d.p.)

$$\begin{bmatrix} 1 & 0.75 & 0.56 & 0.42 & 0.32 \\ 1 & 1.30 & 1.69 & 2.20 & 2.86 \\ 1 & 0.55 & 0.30 & 0.17 & 0.09 \\ 1 & 1.15 & 1.32 & 1.52 & 1.75 \\ 1 & 1.35 & 1.82 & 2.46 & 3.32 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \\ c_5 \end{bmatrix} = \begin{bmatrix} 60 \\ 80 \\ 51 \\ 74 \\ 81 \end{bmatrix}.$$

Now all the components of the matrix are roughly the same size, as required.

There is no need to scale the right-hand side vector as all components are all roughly the same size, they are all simply ‘American temperatures’.

In script construct the scaled matrix and right-hand side vector with



```
te=[15;26;11;23;27]
ta=[60;80;51;74;81]
tes=te/20
A=[ones(5,1) tes tes.^2 tes.^3 tes.^4]
```

- (b) Compute an SVD, $A = USV^T$, with $[U, S, V] = \text{svd}(A)$ to get (2 d.p.)

```
U =
-0.16   0.64   0.20   0.72  -0.12
-0.59  -0.13  -0.00  -0.15  -0.78
-0.10   0.67  -0.59  -0.45   0.05
-0.42   0.23   0.68  -0.42   0.36
-0.66  -0.28  -0.39   0.29   0.49

S =
    7.26      0      0      0      0
      0    1.44      0      0      0
      0      0    0.21      0      0
      0      0      0    0.02      0
      0      0      0      0    0.00

V =
-0.27   0.78  -0.49  -0.27   0.09
-0.32   0.39   0.36   0.66  -0.42
-0.40   0.09   0.55  -0.09   0.72
-0.50  -0.17   0.25  -0.62  -0.52
-0.65  -0.44  -0.51   0.32   0.14
```

- (c) Now choose the effective rank of the matrix to be the number of singular values bigger than the error. Here recall that the temperatures in the matrix have been divided by 20° . Hence the errors of roughly $\pm 0.5^\circ$ in each temperature becomes roughly $\pm 0.5/20 = \pm 0.025$ in the scaled components in the matrix. There are three singular values larger than the error 0.025, so the matrix effectively has rank three. The two singular values less than the error 0.025 are effectively zero. That is, although it is not necessary to construct, we

approximate the matrix A by (2 d.p.)

$$A_3 = US_3V^T = \begin{bmatrix} 1 & 0.74 & 0.56 & 0.43 & 0.31 \\ 1 & 1.30 & 1.69 & 2.20 & 2.86 \\ 1 & 0.56 & 0.30 & 0.16 & 0.09 \\ 1 & 1.16 & 1.32 & 1.52 & 1.75 \\ 1 & 1.35 & 1.82 & 2.46 & 3.32 \end{bmatrix} :$$

the differences between this and A are only ± 0.01 , so matrix A_3 is indeed close to A .

(d) Solve the equations as if matrix A has rank three.

i. Find $\mathbf{z} = U'\mathbf{b}$ via $\mathbf{z}=\mathbf{U}'*\mathbf{ta}$ to find

$$\begin{aligned} \mathbf{z} = \\ -146.53 \\ 56.26 \\ 0.41 \\ 1.06 \\ -0.69 \end{aligned}$$

As matrix A has effective rank of three, we approximate the right-hand side data by neglecting the last two components in this \mathbf{z} . That the last two components in \mathbf{z} are small compared to the others indicates this neglect is a reasonable approximation.

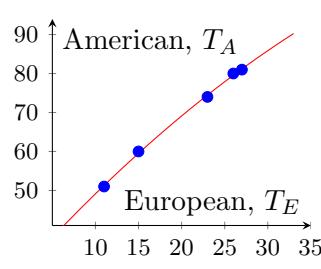
ii. Find \mathbf{y} by solving $S\mathbf{y} = \mathbf{z}$ as a rank three system via $\mathbf{y}=\mathbf{z}(1:3) ./ \text{diag}(\mathbf{S}(1:3,1:3))$ to find (2 d.p.)

$$\mathbf{y} = (-20.19, 39.15, 1.95, y_4, y_5).$$

The smallest solution would be obtained by setting $y_4 = y_5 = 0$.

iii. Finally determine the coefficients $\mathbf{c} = V\mathbf{y}$ with command $\mathbf{c}=\mathbf{V}(:,1:3)*\mathbf{y}$ and then add arbitrary multiples of the remaining columns of V to obtain the general solution (2 d.p.)

$$\mathbf{c} = \begin{bmatrix} 35.09 \\ 22.50 \\ 12.77 \\ 3.99 \\ -5.28 \end{bmatrix} + y_4 \begin{bmatrix} -0.27 \\ 0.66 \\ -0.09 \\ -0.62 \\ 0.32 \end{bmatrix} + y_5 \begin{bmatrix} 0.09 \\ -0.42 \\ 0.72 \\ -0.52 \\ 0.14 \end{bmatrix}.$$



(e) Obtain the solution with smallest coefficients by setting $y_4 = y_5 = 0$. This would fit the data with the quartic polynomial

$$T_A = 35.09 + 22.50 \frac{T_E}{20} + 12.77 \left(\frac{T_E}{20} \right)^2 + 3.99 \left(\frac{T_E}{20} \right)^3 - 5.28 \left(\frac{T_E}{20} \right)^4.$$

But choosing the polynomial with smallest coefficients has little meaning in this application. Surely we prefer a polynomial with fewer terms, fewer non-zero coefficients. Surely we would prefer, say, the quadratic

$$T_A = 25.93 + 49.63 \frac{T_E}{20} - 6.45 \left(\frac{T_E}{20} \right)^2,$$

as plotted in the margin.

In this application, let's use the freedom in y_4 and y_5 to set two of the coefficients to zero in the quartic. Since $c_4 = 3.99$ and $c_5 = -5.28$ are the smallest coefficients, and because they correspond to the highest powers in the quartic, it is natural to choose to make them both zero. Let's redo the last step in the procedure.⁷

The last step in the procedure is to solve $V^T \mathbf{c} = \mathbf{y}$ for \mathbf{c} . With the last two components of \mathbf{c} set to zero, from the computed SVD this is the system of equations (2 d.p.)

$$\begin{bmatrix} -0.27 & -0.32 & -0.40 & -0.50 & -0.65 \\ 0.78 & 0.39 & 0.09 & -0.17 & -0.44 \\ -0.49 & 0.36 & 0.55 & 0.25 & -0.51 \\ -0.27 & 0.66 & -0.09 & -0.62 & 0.32 \\ 0.09 & -0.42 & 0.72 & -0.52 & 0.14 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -20.19 \\ 39.15 \\ 1.95 \\ y_4 \\ y_5 \end{bmatrix},$$

where y_4 and y_5 can be anything for equally good solutions. Considering only the first three rows of this system, and using the zeros in \mathbf{c} , this system becomes

$$\begin{bmatrix} -0.27 & -0.32 & -0.40 \\ 0.78 & 0.39 & 0.09 \\ -0.49 & 0.36 & 0.55 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} -20.19 \\ 39.15 \\ 1.95 \end{bmatrix}.$$

This is a basic system of three equations for three unknowns. Since the matrix is the first three rows and columns of V^T and the right-hand side is the three components of \mathbf{y} already computed, we solve the equation by

- checking the condition number, `rcond(V(1:3,1:3))` is 0.05 which is good, and
- then `c=V(1:3,1:3)'\y` determines the coefficients $(c_1, c_2, c_3) = (25.93, 49.63, -6.45)$ (2 d.p.).

That is, a just as good polynomial fit, consistent with errors in the data, is the simpler quadratic polynomial

$$T_A = 25.93 + 49.63 \frac{T_E}{20} - 6.45 \left(\frac{T_E}{20} \right)^2,$$

⁷ Alternatively, one could redo all the linear algebra to seek a quadratic from the outset rather than a quartic. The two alternative answers for a quadratic are not the same, but they are nearly the same. The small differences in the answers are because one modifies the matrix by recognising its errors, and the other does not. © AJ Roberts, ORCID:0000-0001-8950-1552, July 26, 2016

as plotted previously in the margin.

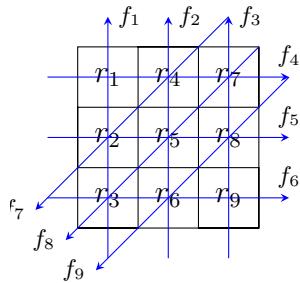
■

Occam's razor: Non sunt multiplicanda entia sine necessitate [Entities must not be multiplied beyond necessity]

John Punch (1639)

Example 5.2.6.

Recall that Exercise 3.5.15 introduced extra ‘diagonal’ measurements into a 2D CT-scan. As shown in the margin, the 2D region is divided into a 3×3 grid of nine blocks. Then measurements taken of the X-rays not absorbed along the shown nine paths: three horizontal, three vertical, and three diagonal. Suppose the measured fractions of X-ray energy are $\mathbf{f} = (0.048, 0.081, 0.042, 0.020, 0.106, 0.075, 0.177, 0.181, 0.105)$. Use an SVD to find the ‘grayest’ transmission factors consistent with the measurements and likely errors.



```

A=[1 1 1 0 0 0 0 0 0
  0 0 0 1 1 1 0 0 0
  0 0 0 0 0 0 1 1 1
  1 0 0 1 0 0 1 0 0
  0 1 0 0 1 0 0 1 0
  0 0 1 0 0 1 0 0 1
  0 1 0 1 0 0 0 0 0
  0 0 1 0 1 0 1 0 0
  0 0 0 0 0 1 0 1 0]
b=log([0.048
       0.081
       0.042
       0.020
       0.106
       0.075
       0.177
       0.181
       0.105])
[U,S,V]=svd(A)
z=U'*b
y=z(1:7)./diag(S(1:7,1:7))
x=V(:,1:7)*y
r=reshape(exp(A\b),3,3)
r7=reshape(exp(x),3,3)

```

Solution: Nine X-ray measurements are made through the body where f_1, f_2, \dots, f_9 denote the fraction of energy in the measurements relative to the power of the X-ray beam. Thus we need to solve nine equations for the nine unknown transmission factors:

$$r_1 r_2 r_3 = f_1, \quad r_4 r_5 r_6 = f_2, \quad r_7 r_8 r_9 = f_3,$$

$$\begin{aligned} r_1r_4r_7 &= f_4, & r_2r_5r_8 &= f_5, & r_3r_6r_9 &= f_6, \\ r_2r_4 &= f_7, & r_3r_5r_7 &= f_8, & r_6r_8 &= f_9. \end{aligned}$$

Turn such nonlinear equations into linear equations by taking the logarithm (to any base, but here say the natural logarithm to base e) of both sides of all equations:

$$r_i r_j r_k = f_l \iff (\log r_i) + (\log r_j) + (\log r_k) = (\log f_l).$$

That is, letting new unknowns $x_i = \log r_i$ and new right-hand sides $b_i = \log f_i$, we aim to solve a system of nine linear equations for nine unknowns:

$$\begin{aligned} x_1 + x_2 + x_3 &= b_1, & x_4 + x_5 + x_6 &= b_2, & x_7 + x_8 + x_9 &= b_3, \\ x_1 + x_4 + x_7 &= b_4, & x_2 + x_5 + x_8 &= b_5, & x_3 + x_6 + x_9 &= b_6, \\ x_2 + x_4 &= b_7, & x_3 + x_5 + x_7 &= b_8, & x_6 + x_8 &= b_9. \end{aligned}$$

These forms the matrix-vector system $A\mathbf{x} = \mathbf{b}$ for 9×9 matrix

$$A = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 0 \end{bmatrix}, \quad \mathbf{b} = \log \begin{bmatrix} .048 \\ .081 \\ .042 \\ .020 \\ .106 \\ .075 \\ .177 \\ .181 \\ .105 \end{bmatrix} = \begin{bmatrix} -3.04 \\ -2.51 \\ -3.17 \\ -3.91 \\ -2.24 \\ -2.59 \\ -1.73 \\ -1.71 \\ -2.25 \end{bmatrix}.$$

Implement Procedure 5.2.2.



- (a) Here there is no need to scale the vector \mathbf{b} as all entries are also the same. There is no need to scale the matrix A as all entries mean the same, namely simply zero or one depending upon whether beam passes through the a pixel square. However, the entries of A are in error in two ways.

- A diagonal beam has a path through a pixel that is up to 41% longer than horizontal or vertical beam—which is not accounted for. Further, a beam has finite width so it will also pass through part of some off-diagonal pixels—which is not represented.
- Similarly, a horizontal or vertical beam has finite width and may underrepresent the sides of the pixels it goes through, and/or involve parts of neighbouring pixels—neither effect is represented.

Consequently the entries in the matrix A could easily have error of 0.5.

- (b) Compute an SVD, $A = USV^T$, via `[U,S,V]=svd(A)` (2 d.p.)

```

U =
  0.33 -0.41  0.27  0.21  0.54  0.23 -0.29  0.13 -0.41
  0.38 -0.00 -0.47  0.36 -0.45 -0.19 -0.00  0.32 -0.41
  0.33  0.41  0.27 -0.57 -0.08  0.23  0.29  0.13 -0.41
  0.33 -0.41  0.27 -0.21 -0.54  0.23 -0.29  0.13  0.41
  0.38 -0.00 -0.47 -0.36  0.45 -0.19 -0.00  0.32  0.41
  0.33  0.41  0.27  0.57  0.08  0.23  0.29  0.13  0.41
  0.23 -0.41 -0.29 -0.00  0.00  0.30  0.58 -0.52 -0.00
  0.41  0.00  0.33  0.00  0.00 -0.74 -0.00 -0.43  0.00
  0.23  0.41 -0.29 -0.00 -0.00  0.30 -0.58 -0.52  0.00

S =
  2.84      0      0      0      0      0      0      0      0
    0  2.00      0      0      0      0      0      0      0
    0      0  1.84      0      0      0      0      0      0
    0      0      0  1.73      0      0      0      0      0
    0      0      0      0  1.73      0      0      0      0
    0      0      0      0      0  1.51      0      0      0
    0      0      0      0      0      0  1.00      0      0
    0      0      0      0      0      0      0  0.51      0
    0      0      0      0      0      0      0      0  0.00

V =
  0.23 -0.41  0.29  0.00  0.00  0.30 -0.58  0.52 -0.00
  0.33 -0.41 -0.27 -0.08  0.57  0.23  0.29 -0.13 -0.41
  0.38  0.00  0.47  0.45  0.36 -0.19 -0.00 -0.32  0.41
  0.33 -0.41 -0.27  0.08 -0.57  0.23  0.29 -0.13  0.41
  0.41 -0.00 -0.33  0.00 -0.00 -0.74 -0.00  0.43 -0.00
  0.33  0.41 -0.27  0.54 -0.21  0.23 -0.29 -0.13 -0.41
  0.38  0.00  0.47 -0.45 -0.36 -0.19  0.00 -0.32 -0.41
  0.33  0.41 -0.27 -0.54  0.21  0.23 -0.29 -0.13  0.41
  0.23  0.41  0.29  0.00  0.00  0.30  0.58  0.52  0.00

```

- (c) Here choose the rank of the matrix to be effectively seven as two of the nine singular values, namely 0.51 and 0.00, are less than about the size of the expected error, roughly 0.5.
- (d) Use the rank seven SVD to solve the approximate system as in Procedure 3.5.3.

i. Find $\mathbf{z} = \mathbf{U}^T \mathbf{b}$ via $\mathbf{z} = \mathbf{U}' * \mathbf{b}$ to find

```

z =
-7.63
 0.27
-0.57
 0.42
 0.64
-1.95
 0.64
-0.40
-0.01

```

ii. Neglect the last two rows in solving $S_7 \mathbf{y} = \mathbf{z}$ to find via $\mathbf{y} = \mathbf{z}(1:7) ./ \text{diag}(\mathbf{S}(1:7, 1:7))$ that the first seven components of \mathbf{y} are

```
y =
-2.69
0.14
-0.31
0.24
0.37
-1.29
0.64
```

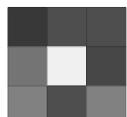
The last two components of \mathbf{y} , y_8 and y_9 , are free variables.

- iii. Obtain a particular solution to $V^T \mathbf{x} = \mathbf{y}$, the one of smallest magnitude, by setting $y_8 = y_9 = 0$ and determining \mathbf{x} from $\mathbf{x} = V(:, 1:7) * \mathbf{y}$ to get the smallest solution

```
x =
-1.53
-0.78
-0.67
-1.16
-0.05
-1.18
-1.16
-1.28
-0.68
```

Obtain other equally valid solutions, in the context of the identified error in matrix A , by adding arbitrary multiples of the last two columns of V .

- (e) Here we aim to make predictions from the CT-scan. The ‘best’ solution in this application is the one with least artificial features. The smallest magnitude \mathbf{x} seems to reasonably implement this criteria. Thus use the above particular \mathbf{x} to determine the transmission factors, $r_i = \exp(x_i)$. Here use $\mathbf{r} = \text{reshape}(\exp(\mathbf{x}), 3, 3)$ to compute and form into the 3×3 array of pixels



0.22	0.31	0.31
0.46	0.95	0.28
0.51	0.31	0.51

as illustrated with `colormap(gray), imagesc(r)`

The CT-scan identifies that there is a significant ‘hole’ in the middle of the body being scanned.

■

5.2.2 Tikhonov regularisation

This optional extension connects to much established practice that graduates may encounter.

Regularisation of poorly-posed linear equations is a widely used practical necessity. Many people have invented alternative techniques. Many have independently re-invented techniques. Perhaps the most common is the so-called Tikhonov regularisation. This section introduces and discusses Tikhonov regularisation.

In statistics, the method is known as ridge regression, and with multiple independent discoveries, it is also variously known as the Tikhonov–Miller method, the Phillips–Twomey method, the constrained linear inversion method, and the method of linear regularization.

Wikipedia (2015)

Definition 5.2.7. *In seeking to solve the poorly-posed system $A\mathbf{x} = \mathbf{b}$ for $m \times n$ matrix A , a **Tikhonov regularisation** is the system $(A^T A + \alpha^2 I_n) \mathbf{x} = A^T \mathbf{b}$ for some chosen regularisation parameter value $\alpha > 0$.*⁸

Example 5.2.8. Use Tikhonov regularisation to solve Example 5.2.1:

$$0.5x + 0.3y = 1 \quad \text{and} \quad 1.1x + 0.7y = 2,$$

Solution: Here the matrix and right-hand side vector are

$$A = \begin{bmatrix} 0.5 & 0.3 \\ 1.1 & 0.7 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}.$$

A Tikhonov regularisation, $(A^T A + \alpha^2 I_n) \mathbf{x} = A^T \mathbf{b}$, is then the system

$$\begin{bmatrix} 1.46 + \alpha^2 & 0.92 \\ 0.92 & 0.58 + \alpha^2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 2.7 \\ 1.7 \end{bmatrix}.$$

Choose regularisation parameter α to be roughly the error: here the error is ± 0.05 so let's choose $\alpha = 0.1$ ($\alpha^2 = 0.01$). Enter into Matlab/Octave with

```
A=[0.5 0.3;1.1 0.7]
b=[1.0;2.0]
At=(A'*A+0.01*eye(2))
bt=A'*b
rcondA=rcond(At)
x=At\bt
```

to find the Tikhonov regularised solution is $\mathbf{x} = (1.39, 0.72)$.⁹ This solution is reasonably close to the smallest solution found by the

⁸ Some will notice that a Tikhonov regularisation is closely connected to the so-called normal equation $A^T A \mathbf{x} = A^T \mathbf{b}$. Tikhonov regularisation shares some of the practical limitations of the normal equation.

⁹ Interestingly, $\text{rcond} = 0.003$ for the Tikhonov system which is worse than $\text{rcond}(A)$. The regularisation only works because pre-multiplying by A^T puts both sides in the row space of A (except for numerical error and the small $\alpha^2 I$ factor).



SVD which is $(1.32, 0.83)$. However, Tikhonov regularisation gives no hint of the reasonable general solutions found by the SVD approach of Example 5.2.1.

Change the regularisation parameter to $\alpha = 0.01$ and $\alpha = 1$ and see that both degrade the Tikhonov solution.

■

Do not apply Tikhonov regularisation blindly as it does introduce biases. The following example illustrates the bias.

Example 5.2.9. Recall the example at the start of Section 3.5.1 where my scales variously reported my weight in kg as 84.8, 84.1, 84.7 and 84.4. To best estimate my weight x we rewrote the problem in matrix-vector form

$$Ax = \mathbf{b}, \quad \text{namely} \quad \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} x = \begin{bmatrix} 84.8 \\ 84.1 \\ 84.7 \\ 84.4 \end{bmatrix}.$$

A Tikhonov regularisation of this inconsistent system is

$$\left([1 \ 1 \ 1 \ 1] \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \alpha^2 \right) x = [1 \ 1 \ 1 \ 1] \begin{bmatrix} 84.8 \\ 84.1 \\ 84.7 \\ 84.4 \end{bmatrix}.$$

That is, $(4 + \alpha^2)x = 338$ kg with solution $x = 84.5/(1 + \alpha^2/4)$ kg. This Tikhonov answer is biased because it is systematically below the average 84.5 kg. For small Tikhonov parameter α it is only a small bias, but even so such a bias is unpleasant.

■

Example 5.2.10. Use Tikhonov regularisation to solve $A\mathbf{x} = \mathbf{b}$ for the matrix and vector of Example 5.2.3a.

Solution: Here the matrix and right-hand side vector are

$$A = \begin{bmatrix} -0.2 & -0.6 & 1.8 \\ 0.0 & 0.2 & -0.4 \\ -0.3 & 0.7 & 0.3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} -0.5 \\ 0.1 \\ -0.2 \end{bmatrix}.$$

A Tikhonov regularisation, $(A^T A + \alpha^2 I_3)\mathbf{x} = A^T \mathbf{b}$, is then the system

$$\begin{bmatrix} 0.13 + \alpha^2 & -0.09 & -0.45 \\ -0.09 & 0.89 + \alpha^2 & -0.95 \\ -0.45 & -0.95 & 3.49 + \alpha^2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 0.16 \\ 0.18 \\ -1.00 \end{bmatrix}.$$

Choose regularisation parameter α to be roughly the error: here the error is ± 0.05 so let's choose $\alpha = 0.1$ ($\alpha^2 = 0.01$). Enter into and solve with Matlab/Octave via

```

A=[-0.2 -0.6 1.8
  0.0 0.2 -0.4
  -0.3 0.7 0.3 ]
b=[-0.5;0.1;-0.2]
At=(A'*A+0.01*eye(3))
bt=A'*b
rcondA=rcond(At)
x=At\bt

```



which then finds the Tikhonov regularised solution is $\mathbf{x} = (0.10, -0.11, -0.30)$. To two decimal places this is the same as the smallest solution found by an SVD. However, Tikhonov regularisation gives no hint of the reasonable general solutions found by the SVD approach of Example 5.2.3a.

Theorem 5.2.11 (Tikhonov regularisation). *Solving the Tikhonov regularisation, with parameter α , of $A\mathbf{x} = \mathbf{b}$ is equivalent to finding the smallest, least square, solution of the system $\tilde{A}\mathbf{x} = \mathbf{b}$ where the matrix \tilde{A} is obtained from A by replacing each of its non-zero singular values σ_i by $\tilde{\sigma}_i := \sigma_i + \alpha^2/\sigma_i$.*

Proof. Let $m \times n$ matrix A have SVD $A = USV^T$. First, the left-hand side matrix in a Tikhonov regularisation is

$$\begin{aligned}
A^T A + \alpha^2 I_n &= (USV^T)^T USV^T + \alpha^2 I_n VV^T \\
&= VS^T U^T USV^T + \alpha^2 V I_n V^T \\
&= VS^T SV^T + V(\alpha^2 I_n) V^T \\
&= V(S^T S + \alpha^2 I_n) V^T,
\end{aligned}$$

whereas the right-hand side is $A^T \mathbf{b} = (USV^T)^T \mathbf{b} = VS^T U^T \mathbf{b}$. Corresponding to the variables used in previous procedures, let $\mathbf{z} = U^T \mathbf{b} \in \mathbb{R}^m$ and as yet unknown $\mathbf{y} = V^T \mathbf{x} \in \mathbb{R}^n$. Then equating the above two sides, and premultiplying by the orthogonal V^T , means the Tikhonov regularisation is equivalent to solving $(S^T S + \alpha^2 I_n) \mathbf{y} = S^T \mathbf{z}$ for \mathbf{y} .

Second, suppose rank $A = r$ so that the singular value matrix

$$S = \begin{bmatrix} \sigma_1 & \cdots & 0 & & \\ \vdots & \ddots & \vdots & O_{r \times (n-r)} & \\ 0 & \cdots & \sigma_r & & \\ & & & O_{(m-r) \times r} & O_{(m-r) \times (n-r)} \end{bmatrix}$$

(where the bottom right zero block contains all the zero singular values). Consequently, the equivalent Tikhonov regularisation,

$(S^T S + \alpha^2 I_n) \mathbf{y} = S^T \mathbf{z}$, becomes

$$\begin{bmatrix} \sigma_1^2 + \alpha^2 & \cdots & 0 & O_{r \times (n-r)} \\ \vdots & \ddots & \vdots & \\ 0 & \cdots & \sigma_r^2 + \alpha^2 & \\ O_{(n-r) \times r} & & & \alpha^2 I_{n-r} \end{bmatrix} \mathbf{y} = \begin{bmatrix} \sigma_1 z_1 \\ \vdots \\ \sigma_r z_r \\ \mathbf{0}_{n-r} \end{bmatrix}.$$

Dividing each of the first r rows by the corresponding non-zero singular value, $\sigma_1, \sigma_2, \dots, \sigma_r$, the equivalent system is

$$\begin{bmatrix} \sigma_1 + \alpha^2/\sigma_1 & \cdots & 0 & O_{r \times (n-r)} \\ \vdots & \ddots & \vdots & \\ 0 & \cdots & \sigma_r + \alpha^2/\sigma_r & \\ O_{(n-r) \times r} & & & \alpha^2 I_{n-r} \end{bmatrix} \mathbf{y} = \begin{bmatrix} z_1 \\ \vdots \\ z_r \\ \mathbf{0}_{n-r} \end{bmatrix},$$

with solution

- $y_i = z_i / (\sigma_i + \alpha^2/\sigma_i)$ for $i = 1, \dots, r$, and
- $y_i = 0$ for $i = r+1, \dots, n$ (since $\alpha^2 > 0$).

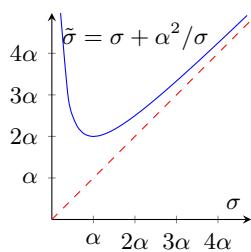
This establishes that solving the Tikhonov system is equivalent to performing the SVD Procedure 3.5.3 for the least square solution to $A\mathbf{x} = \mathbf{b}$ but with two changes in Step 3:

- for $i = 1, \dots, r$ divide by $\tilde{\sigma}_i := \sigma_i + \alpha^2/\sigma_i$ instead of the true singular value σ_i (the margin plots $\tilde{\sigma}$ versus σ), and
- for $i = r+1, \dots, n$ set $y_i = 0$ to obtain the smallest possible solution (Theorem 3.5.9).

Thus Tikhonov regularisation of $A\mathbf{x} = \mathbf{b}$ is equivalent to finding the smallest, least square, solution of the system $\tilde{A}\mathbf{x} = \mathbf{b}$. \square

There is another reason to be careful when using Tikhonov regularisation. Yes, it gives a nice, neat, unique solution. However, it does not hint that there may be an infinite number of equally good nearby solutions (as found through Procedure 5.2.2). Among those equally good nearby solutions may be ones that you prefer in your application.

Choose a good regularisation parameter



- One strategy to choose the regularisation parameter α is that the effective change in the matrix, from A to \tilde{A} , should be about the size of errors expected in A .¹⁰ Since changes in the matrix are largely measured by the singular values we need to consider the relation between $\tilde{\sigma} = \sigma + \alpha^2/\sigma$ and σ . From

¹⁰This strategic choice is sometimes called the discrepancy principle (Kress 2015, §7).

the marginal graph the small singular values are changed by a lot, but these are the ones for which we want $\tilde{\sigma}$ large in order give a ‘least square’ approximation. Significantly, the marginal graph also shows that singular values larger than α change by less than α . Thus the parameter α should not be much larger than the expected error in the elements of the matrix A .

- Another consideration is the effect of regularisation upon errors in the right-hand side vector. The condition number of A may be very bad. However, as the marginal graph shows the smallest $\tilde{\sigma} \geq 2\alpha$. Thus, in the regularised system the condition number of the effective matrix \tilde{A} is approximately $\sigma_1/(2\alpha)$. We need to choose the regularisation parameter α large enough so that $\frac{\sigma_1}{2\alpha} \times (\text{relative error in } \mathbf{b})$ is an acceptable relative error in the solution \mathbf{x} (Theorem 3.3.24). It is only when the regularisation parameter α is big enough that the regularisation will be effective in finding a least square approximation.

5.2.3 Exercises

Exercise 5.2.1. For each of the following matrices, say A , and right-hand side vectors, say \mathbf{b}_1 , solve $A\mathbf{x} = \mathbf{b}_1$. But suppose the matrix entries come from experiments and are only known to within errors ± 0.05 . Thus within experimental error the given matrices A' and A'' may be the ‘true’ matrix A . Solve $A'\mathbf{x}' = \mathbf{b}_1$ and $A''\mathbf{x}'' = \mathbf{b}_1$ and comment on the results. Finally, use an SVD to find a general solution consistent with the error in the matrix.

$$(a) A = \begin{bmatrix} -1.3 & -0.4 \\ 0.7 & 0.2 \end{bmatrix}, \mathbf{b}_1 = \begin{bmatrix} 2.4 \\ -1.3 \end{bmatrix}, A' = \begin{bmatrix} -1.27 & -0.43 \\ 0.71 & 0.19 \end{bmatrix}, \\ A'' = \begin{bmatrix} -1.27 & -0.38 \\ 0.66 & 0.22 \end{bmatrix}.$$

$$(b) B = \begin{bmatrix} -1.8 & -1.1 \\ -0.2 & -0.1 \end{bmatrix}, \mathbf{b}_2 = \begin{bmatrix} -0.7 \\ -0.1 \end{bmatrix}, B' = \begin{bmatrix} -1.81 & -1.13 \\ -0.24 & -0.12 \end{bmatrix}, \\ B'' = \begin{bmatrix} -1.81 & -1.13 \\ -0.18 & -0.1 \end{bmatrix}.$$

$$(c) C = \begin{bmatrix} 0.8 & -0.1 \\ -1.0 & 0.1 \end{bmatrix}, \mathbf{b}_3 = \begin{bmatrix} 0.2 \\ -0.3 \end{bmatrix}, C' = \begin{bmatrix} 0.81 & -0.07 \\ -1.01 & 0.06 \end{bmatrix}, \\ C'' = \begin{bmatrix} 0.79 & -0.08 \\ -1.03 & 0.09 \end{bmatrix}.$$

$$(d) D = \begin{bmatrix} 0.0 & 0.5 & -0.5 \\ 0.6 & 0.5 & 0.9 \\ 0.6 & 1.3 & 0.0 \end{bmatrix}, \mathbf{b}_4 = \begin{bmatrix} -1.4 \\ 0.4 \\ -1.9 \end{bmatrix}, \\ D' = \begin{bmatrix} -0.02 & 0.49 & -0.49 \\ 0.58 & 0.54 & 0.9 \\ 0.61 & 1.34 & -0.02 \end{bmatrix}, D'' = \begin{bmatrix} -0.04 & 0.52 & -0.48 \\ 0.64 & 0.52 & 0.87 \\ 0.57 & 1.33 & 0.04 \end{bmatrix}.$$

$$(e) \quad E = \begin{bmatrix} 0.6 & -0.8 & -0.2 \\ -0.9 & 1.0 & 1.2 \\ -0.9 & 0.9 & 1.4 \end{bmatrix}, \quad \mathbf{b}_5 = \begin{bmatrix} 1.1 \\ -3.7 \\ -4.1 \end{bmatrix},$$

$$E' = \begin{bmatrix} 0.57 & -0.78 & -0.23 \\ -0.91 & 0.99 & 1.22 \\ -0.93 & 0.9 & 1.39 \end{bmatrix},$$

$$E'' = \begin{bmatrix} 0.56 & -0.77 & -0.21 \\ -0.87 & 1.01 & 1.22 \\ -0.87 & 0.9 & 1.39 \end{bmatrix}.$$

$$(f) \quad F = \begin{bmatrix} 0.1 & -1.0 & 0.0 \\ 2.1 & -0.2 & -0.5 \\ 0.0 & -1.6 & 0.0 \end{bmatrix}, \quad \mathbf{b}_6 = \begin{bmatrix} -0.2 \\ 1.6 \\ -0.5 \end{bmatrix},$$

$$F' = \begin{bmatrix} 0.1 & -0.98 & -0.04 \\ 2.11 & -0.17 & -0.47 \\ -0.04 & -1.62 & -0.01 \end{bmatrix}, \quad F'' = \begin{bmatrix} 0.14 & -0.96 & 0.01 \\ 2.13 & -0.23 & -0.47 \\ 0.0 & -1.57 & -0.02 \end{bmatrix}.$$

$$(g) \quad G = \begin{bmatrix} 1.0 & -0.3 & 0.3 & -0.4 \\ 1.8 & 0.5 & 0.1 & 0.2 \\ 0.2 & -0.3 & 1.3 & -0.6 \\ 0.0 & 0.5 & 1.2 & 0.0 \end{bmatrix}, \quad \mathbf{b}_7 = \begin{bmatrix} 2.0 \\ 1.6 \\ 1.4 \\ -0.2 \end{bmatrix},$$

$$G' = \begin{bmatrix} 0.98 & -0.3 & 0.31 & -0.44 \\ 1.8 & 0.54 & 0.06 & 0.21 \\ 0.24 & -0.33 & 1.27 & -0.58 \\ 0.01 & 0.52 & 1.23 & -0.01 \end{bmatrix},$$

$$G'' = \begin{bmatrix} 1.03 & -0.32 & 0.33 & -0.36 \\ 1.82 & 0.49 & 0.08 & 0.16 \\ 0.2 & -0.31 & 1.33 & -0.64 \\ 0.0 & 0.49 & 1.22 & 0.0 \end{bmatrix}.$$

$$(h) \quad H = \begin{bmatrix} -0.9 & -0.5 & -0.3 & -0.4 \\ -0.1 & 0.1 & -0.2 & 0.8 \\ -1.0 & 0.4 & -1.1 & 0.6 \\ 1.0 & 2.2 & -1.0 & -0.1 \end{bmatrix}, \quad \mathbf{b}_8 = \begin{bmatrix} 0.4 \\ 0.3 \\ 0.2 \\ -2.0 \end{bmatrix},$$

$$H' = \begin{bmatrix} -0.88 & -0.52 & -0.33 & -0.41 \\ -0.11 & 0.13 & -0.17 & 0.78 \\ -0.96 & 0.44 & -1.12 & 0.61 \\ 0.98 & 2.19 & -0.99 & -0.13 \end{bmatrix},$$

$$H'' = \begin{bmatrix} -0.86 & -0.49 & -0.29 & -0.37 \\ -0.06 & 0.14 & -0.18 & 0.83 \\ -0.96 & 0.38 & -1.11 & 0.58 \\ 1.01 & 2.21 & -1.04 & -0.13 \end{bmatrix}.$$

Exercise 5.2.2. Recall Example 5.2.4 explores the effective rank of the 5×5 Hilbert matrix depending upon a supposed level of error. Similarly, explore the effective rank of the 7×7 Hilbert matrix (`hilb(7)` in Matlab/Octave) depending upon supposed levels of error in the matrix. What levels of error in the components would give what effective rank of the matrix?

Exercise 5.2.3. Recall Exercise 2.2.12 considered the inner four planets in the solar system. The exercise fitted fit a quadratic polynomial to the orbital period $T = c_1 + c_2R + c_3R^2$ as a function of distance R using the data of Table 2.4. In view of the bad condition number, $\text{rcond} = 6 \cdot 10^{-6}$, revisit the task with the more powerful techniques of this section. Use the data for Mercury, Venus and Earth to fit the quadratic and predict the period for Mars. Discuss how the bad condition number is due to the failure in Exercise 2.2.12 of scaling the data in the matrix.

Exercise 5.2.4. Recall Exercise 3.5.16 used a 4×4 grid of pixels in the computed tomography of a CT-scan. Redo this exercise recognising that the entries in matrix A have errors up to roughly 0.5. Discuss any change in the prediction.

Exercise 5.2.5. Reconsider each of the matrix-vector systems you explored in Exercise 5.2.1. Also solve each system using Tikhonov regularisation; for example, in the first system solve $A\mathbf{x} = \mathbf{b}_1$, $A'\mathbf{x}' = \mathbf{b}_1$ and $A''\mathbf{x}'' = \mathbf{b}_1$. Discuss why \mathbf{x} , \mathbf{x}' and \mathbf{x}'' are all reasonably close to the smallest solution of those obtained via an SVD.

Exercise 5.2.6. Recall that Example 5.2.4 explores the effective rank of the 5×5 Hilbert matrix depending upon a supposed level of error. Here do the alternative and solve the system $A\mathbf{x} = \mathbf{1}$ via Tikhonov regularisation using a wide range of various regularisation parameters α . Comment on the relation between the solutions obtained for various α and those obtained in the example for the various presumed error—perhaps plot the components of \mathbf{x} versus parameter α (on a log-log plot).

Exercise 5.2.7. Recall Example 5.2.6 used a 3×3 grid of pixels in the computed tomography of a CT-scan. Redo this example with Tikhonov regularisation recognising that the entries in matrix A have errors up to roughly 0.5. Discuss the relation between the solution of Example 5.2.6 that of Tikhonov regularisation.

Exercise 5.2.8. Recall Exercise 3.5.16 used a 4×4 grid of pixels in the computed tomography of a CT-scan. Redo this exercise with Tikhonov regularisation recognising that the entries in matrix A have errors up to roughly 0.5. Discuss any change in the prediction.

Answers to selected exercises

5.1.1b : 1.4

5.1.1d : 2.3

5.1.1f : 1.9

5.1.2b : 5

5.1.2d : 3

5.1.2f : 8

5.1.2h : 5

5.1.5 : Rank three.

5.1.8a : Use $A = [A \ A \ A; A \ 0*A \ A; A \ A \ A]$

5.1.9a : Use $A = [A \ 0*A \ A; 0*A \ A \ 0*A; A \ 0*A \ A]$

5.1.10 : ranks 6, 24

5.1.13 : Book 3: angles $35^\circ, 32^\circ, 1^\circ, 29^\circ, 54^\circ, 76^\circ$.

5.1.14c : Books 14, 4, 5, 13 and 2 are at angles $1^\circ, 3^\circ, 9^\circ, 10^\circ$ and 16° , respectively.

5.2.1b : $\mathbf{x} = (1.0, -1.0)$, $\mathbf{x}' = (0.54, -0.24)$, $\mathbf{x}'' = (1.92, -2.46)$, $\mathbf{x} = (0.28, 0.17) + t(-0.52, 0.85)$ (2 d.p.).

5.2.1d : $\mathbf{x} = (1.6, -2.2, 0.6)$, $\mathbf{x}' = (1.31, -2.0, 0.8)$, $\mathbf{x}'' = (-0.28, -1.35, 1.47)$, $\mathbf{x} = (-0.09, -1.43, 1.32) + t(0.85, -0.39, -0.36)$ (2 d.p.).

5.2.1f : $\mathbf{x} = (1.12, 0.31, 1.4)$, $\mathbf{x}' = (0.59, 0.3, -0.84)$, $\mathbf{x}'' = (0.77, 0.32, -0.08)$, $\mathbf{x} = (0.75, 0.3, -0.18) + t(-0.23, -0.01, -0.97)$ (2 d.p.).

5.2.1h : $\mathbf{x} = (-0.76, -0.23, 0.69, 0.48)$, $\mathbf{x}' = (-1.36, 0.33, 1.35, 0.43)$, $\mathbf{x}'' = (-0.09, -0.92, -0.17, 0.47)$, $\mathbf{x} = (-0.38, -0.63, 0.2, 0.47) + t(0.52, -0.54, -0.67, -0.02)$ (2 d.p.).

5.2.4 : The matrix has effective rank of eleven. Pixel ten is still the most absorbing. The corner pixels are the most affected.

6 Determinants distinguish matrices

Chapter Contents

6.1	Geometry underlies determinants	507
6.1.1	Exercises	518
6.2	Laplace expansion theorem for determinants	523
6.2.1	Exercises	543

Although much of the theoretical role of determinants is usurped by the SVD, nonetheless, determinants aid in establishing forthcoming properties of eigenvalues and eigenvectors, and empower graduates to connect to much extant practice.

Recall from previous study (section 4.1.1, e.g.)

- a 2×2 matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ has determinant $\det A = |A| = ad - bc$, and matrix A is invertible iff $\det A \neq 0$;
- a 3×3 matrix $A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$ has determinant $\det A = |A| = aei + bfg + cdh - ceg - afh - bdi$, and matrix A is invertible iff $\det A \neq 0$.

These two formulas for a determinant are best remembered via the following diagrams where products along the red lines are subtracted from the products along the blue lines, respectively:

$$\begin{array}{c} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \end{array} \quad \begin{array}{l} \text{red lines} \\ \text{blue lines} \end{array} \quad (6.1)$$

This chapter extends these determinants to any size matrix, and explores more useful properties—especially those properties useful for understanding and developing the general eigenvalue problems and applications of Chapter 7.

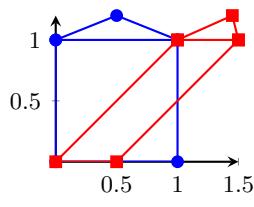
6.1 Geometry underlies determinants

Section Contents

6.1.1 Exercises	518
---------------------------	-----

Sections 3.2.2, 3.2.3 and 3.6 introduced that multiplication by a matrix transforms areas and volumes. Determinants give precisely how much a square matrix transforms areas and volumes.

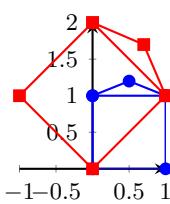
Example 6.1.1. Consider the square matrix $A = \begin{bmatrix} \frac{1}{2} & 1 \\ 0 & 1 \end{bmatrix}$. Use matrix multiplication to find the image of the unit square under the transformation by A . How much is the area of the unit square scaled up/down? Compare with the determinant.



Solution: Consider the corner points of the unit square, under multiplication by A : $(0, 0) \mapsto (0, 0)$, $(1, 0) \mapsto (\frac{1}{2}, 0)$, $(0, 1) \mapsto (0, 1)$, and $(1, 1) \mapsto (\frac{3}{2}, 1)$, as shown in the marginal picture (the ‘roof’ is only plotted to uniquely identify the sides). The resultant parallelogram has area of $\frac{1}{2}$ as its base is $\frac{1}{2}$ and its height is 1. This parallelogram area is the same as the determinant since here (6.1) gives $\det A = \frac{1}{2} \cdot 1 - 0 \cdot 1 = \frac{1}{2}$.

■

Example 6.1.2. Consider the square matrix $B = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}$. Use matrix multiplication to find the image of the unit square under the transformation by B . How much is the unit area scaled up/down? Compare with the determinant.

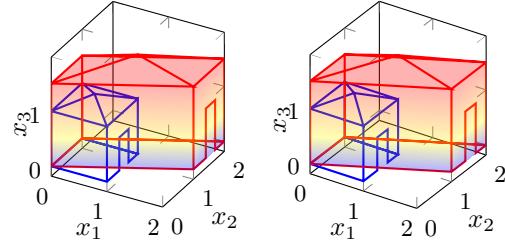


Solution: Consider the corner points of the unit square, under multiplication by B : $(0, 0) \mapsto (0, 0)$, $(1, 0) \mapsto (-1, 1)$, $(0, 1) \mapsto (1, 1)$, and $(1, 1) \mapsto (0, 2)$, as shown in the marginal picture. Through multiplication by the matrix B , the unit square is expanded, rotated and reflected. The resultant square has area of 2 as its sides are all of length $\sqrt{2}$. This area has the same magnitude as the determinant since here (6.1) gives $\det B = (-1) \cdot 1 - 1 \cdot 1 = -2$.

■

Example 6.1.3. Consider the square matrix $A = \begin{bmatrix} 2 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & \frac{3}{2} \end{bmatrix}$. Use matrix multiplication to find the image of the unit cube under the transformation by A . How much is the volume of the unit cube scaled up/down? Compare with the determinant.

Solution: Consider the corner points of the unit cube, under multiplication by A : $(0, 0, 0) \mapsto (0, 0, 0)$, $(1, 0, 0) \mapsto (2, 1, 0)$, $(0, 1, 0) \mapsto (0, 1, 0)$, $(0, 0, 1) \mapsto (0, 0, \frac{3}{2})$, and so on, as shown below (in stereo):

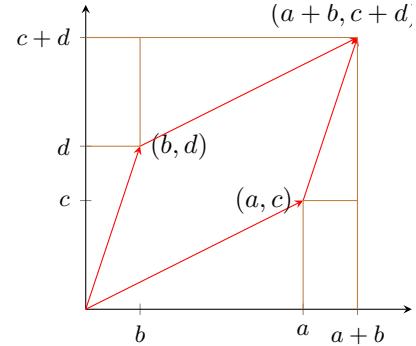


Through multiplication by the matrix A , the unit cube is deformed to a parallelepiped. The resultant parallelepiped has volume of 3 as it has height $\frac{3}{2}$ and the parallelogram base has area $2 \cdot 1$. This volume is the same as the matrix determinant since (6.1) gives $\det A = 2 \cdot 1 \cdot \frac{3}{2} + 0 + 0 - 0 - 0 - 0 = 3$.

■

Determinants determine area transformation Consider multiplication by the general 2×2 matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$.

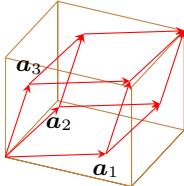
Under multiplication by this matrix A the unit square becomes the parallelogram shown with corners at the origin, (a, c) , (b, d) and $(a+b, c+d)$. Let's determine the area of the parallelogram by that of the containing rectangle less the two small rectangles and the four small triangles. The two small rectangles have the same area, namely bc . The two small triangles on the left and the right also have the same area, namely $\frac{1}{2}bd$. The two small triangles on the top and the bottom similarly have the same area, namely $\frac{1}{2}ac$. Thus, under multiplication by matrix A the image of the unit square is the parallelogram with



$$\begin{aligned} \text{area} &= (a+b)(c+d) - 2bc - 2 \cdot \frac{1}{2}bd - 2 \cdot \frac{1}{2}ac \\ &= ac + ad + bc + bd - 2bc - bd - ac \\ &= ad - bc = \det A. \end{aligned}$$

This picture is the case when the matrix does not also reflect the image: if the matrix also reflects, as in Example 6.1.2, then the determinant is the negative of the area. In either case, the

area of the unit square after transforming by the matrix A is the magnitude $|\det A|$.



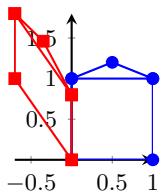
Analogous geometric arguments relate determinants of 3×3 matrices with transformations of volumes. Under multiplication by a 3×3 matrix $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \mathbf{a}_3]$, the image of the unit cube is a parallelepiped with edges \mathbf{a}_1 , \mathbf{a}_2 and \mathbf{a}_3 as illustrated in the margin. By computing the volumes of various rectangular boxes, prisms and tetrahedra (or by other methods), the volume of such a parallelepiped can be expressed as the 3×3 determinant formula (6.1).

In higher dimensions we want the determinant to behave analogously and so next define the determinant to do so. We use the terms **n D-cube** to generalise a square and cube to n dimensions (\mathbb{R}^n), **n D-volume** to generalise the notion of area and volume to n dimensions, and so on. When the dimension of the space is unspecified, then we may say **hyper-cube**, **hyper-volume**, and so on.

Definition 6.1.4. Let A be an $n \times n$ square matrix, and let C be the unit n D-cube in \mathbb{R}^n . Consider the image C' in \mathbb{R}^n of C by the transformation $\mathbf{x} \mapsto A\mathbf{x}$. Define the **determinant** of A , denoted either $\det A$ or $|A|$ such that:

- the magnitude $|\det A|$ is the n D-volume of C' ; and
- the sign of $\det A$ to be negative iff the transformation reflects the orientation of the n D-cube.

Example 6.1.5. Roughly estimate the determinant of the matrix that transforms the unit square to the parallelogram as shown in the margin.



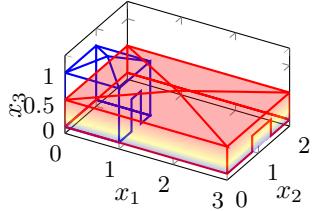
Solution: The image is a parallelogram with a vertical base of length about 0.8 and a horizontal height of about 0.7 so the area of the image is about $0.8 \times 0.7 = 0.56 \approx 0.6$. But the image has been reflected as one cannot rotate and stretch to get the image (remember the origin is fixed under matrix multiplication): thus the determinant must be negative. Our estimate for the matrix determinant is -0.6 . ■

Basic properties of a determinant follow direct from Definition 6.1.4.

Theorem 6.1.6. (a) Let D be an $n \times n$ diagonal matrix. The determinant of D is the product of the diagonal entries: $\det D = d_{11}d_{22} \cdots d_{nn}$.
 (b) An orthogonal matrix Q has $\det Q = \pm 1$ (only one alternative, not both), and $\det Q = \det(Q^T)$.
 (c) $\det(kA) = k^n \det A$ for any scalar k .

Proof. Use Definition 6.1.4.

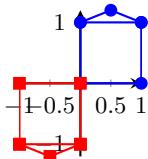
- 6.1.6a. The unit n D-cube in \mathbb{R}^n has edges e_1, e_2, \dots, e_n (the unit vectors). Multiplying each of these edges by the diagonal matrix $D = \text{diag}(d_{11}, d_{22}, \dots, d_{nn})$ maps the unit n D-cube to a n D-rectangle with edges $d_{11}e_1, d_{22}e_2, \dots, d_{nn}e_n$ (as illustrated in the margin). Being a n D-rectangle with all edges orthogonal, its n D-volume is the product of the length of the sides; that is, $|\det D| = |d_{11}| \cdot |d_{22}| \cdots |d_{nn}|$. The n D-cube is reflected only if there are an odd number of negative diagonal elements, hence the sign of the determinant is such that $\det D = d_{11}d_{22} \cdots d_{nn}$.



- 6.1.6b. Multiplication by an orthogonal matrix Q is a rotation and/or reflection as it preserves all lengths and angles (Theorem 3.2.39f). Hence it preserves n D-volumes. Consequently the image of the unit n D-cube under multiplication by Q has the same volume of one; that is, $|\det Q| = 1$. The sign of $\det Q$ characterises whether multiplication by Q has a reflection.

When Q is orthogonal then so is Q^T (Theorem 3.2.39d). Hence $\det Q^T = \pm 1$. Multiplication by Q involves a reflection iff its inverse of multiplication by Q^T involves the reflection back again. Hence the signs of the two determinants must be the same: that is, $\det Q = \det Q^T$.

- 6.1.6c. Let the matrix A transform the unit n D-cube to a n D-parallelepiped C' that has n D-volume $|\det A|$. Multiplication by the matrix (kA) then forms a n D-parallelepiped which is $|k|$ -times bigger than C' in every direction. In \mathbb{R}^n its n D-volume is then $|k|^n$ -times bigger; that is, $|\det(kA)| = |k|^n |\det A| = |k^n \det A|$. If the scalar k is negative then the orientation of the image is reversed (reflected) only in odd n dimensions; that is, the sign of the determinant is multiplied by $(-1)^n$. (For example, the unit square shown in the margin is transformed through multiplication by $(-1)I_2$ and the effect is the same as rotation by 180° , without any reflection as $(-1)^2 = 1$.) Hence for all real k , the orientation is such that $\det(kA) = k^n \det A$.



□

Example 6.1.7. The determinant of the $n \times n$ identity matrix is one: that is, $\det I_n = 1$. We justify this result

- either because it is a diagonal matrix and hence its determinant is the product of the diagonal entries (Theorem 6.1.6a) here all ones,
- or because multiplication by the identity does not change the unit n D-cube and so does not change its n D-volume (Definition 6.1.4).

The determinant of $-I_n$ is $(-1)^n$ either by the product of the diagonals (Theorem 6.1.6a), or because of the scalar multiplication Theorem 6.1.6c. ■

Example 6.1.8. Use (6.1) to compute the determinant of the orthogonal matrix

$$Q = \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix}.$$

Then use Theorem 6.1.6 to deduce the determinants of the following matrices:

$$\begin{bmatrix} \frac{1}{3} & \frac{2}{3} & -\frac{2}{3} \\ -\frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix}, \quad \begin{bmatrix} -1 & 2 & -2 \\ -2 & -2 & -1 \\ 2 & -1 & -2 \end{bmatrix}, \quad \begin{bmatrix} \frac{1}{6} & -\frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{6} \\ -\frac{1}{3} & \frac{1}{6} & \frac{1}{3} \end{bmatrix}.$$

Solution:

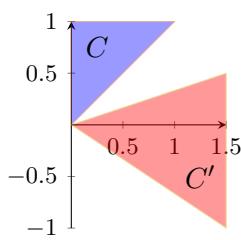
- Using (6.1) the determinant

$$\begin{aligned} \det Q &= \frac{1}{3} \frac{2}{3} \frac{2}{3} + (-\frac{2}{3}) \frac{1}{3} (-\frac{2}{3}) + \frac{2}{3} \frac{2}{3} \frac{1}{3} \\ &\quad - \frac{2}{3} \frac{2}{3} (-\frac{2}{3}) - \frac{1}{3} \frac{1}{3} \frac{1}{3} - (-\frac{2}{3}) \frac{2}{3} \frac{2}{3} \\ &= \frac{1}{27} (4 + 4 + 4 + 8 - 1 + 8) = 1. \end{aligned}$$

- The first matrix is Q^T for which $\det Q^T = \det Q = 1$.
- The second matrix is minus three times Q so, being 3×3 matrices, its determinant is $(-3)^3 \det Q = -27$.
- The third matrix is half of Q so, being 3×3 matrices, its determinant is $(\frac{1}{2})^3 \det Q = \frac{1}{8}$.

■

A consequence of Theorem 6.1.6c is that a determinant characterises the transformation of any sized hyper-cube. Consider the transformation by a matrix A of an n D-cube of side length k ($k \geq 0$). The n D-cube has edges $k\mathbf{e}_1, k\mathbf{e}_2, \dots, k\mathbf{e}_n$. The transformation results in an n D-parallelepiped with edges $A(k\mathbf{e}_1), A(k\mathbf{e}_2), \dots, A(k\mathbf{e}_n)$, which by commutativity and associativity (Theorem 3.1.18d) are the same edges as $(kA)\mathbf{e}_1, (kA)\mathbf{e}_2, \dots, (kA)\mathbf{e}_n$. That is, the resulting n D-parallelepiped is the same as applying matrix kA to the unit n D-cube, and so must have n D-volume $k^n |\det A|$. Crucially, this property that matrix multiplication multiplies all sizes of hyper-cubes by the determinant holds for all other shapes and sizes, not just hyper-cubes. Let's see an specific example, before proving the general theorem.



Example 6.1.9. Multiplication by some specific matrix transforms the (blue) triangle C to the (red) triangle C' as shown in the margin. By finding the ratio of the areas, estimate the magnitude of the determinant of the matrix.

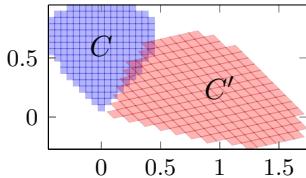
Solution: The (blue) triangle C has vertical height one and horizontal base one, so has area 0.5. The mapped (red) triangle C' has vertical base of 1.5 and horizontal height of 1.5 so its area is $\frac{1}{2} \times 1.5 \times 1.5 = 1.25$. The mapped area is thus $1.25/0.5 = 2.5$ times bigger than the initial area; hence the determinant of the transformation matrix has magnitude $|\det A| = 2.5$.

We cannot determine the sign of the determinant as we do not know about the orientation of C' relative to C . ■

Theorem 6.1.10. Consider any bounded smooth nD -volume C in \mathbb{R}^n and its image C' after multiplication by $n \times n$ matrix A . Then

$$\det A = \pm \frac{nD\text{-volume of } C'}{nD\text{-volume of } C}$$

with the negative sign when matrix A changes the orientation.



Proof. A proof is analogous to integration in calculus (Hannah 1996, p.402). In 2D, as drawn in the margin, we divide a given region C into many small squares of side length k , each of area k^2 : each of these transforms to a small parallelogram of area $k^2|\det A|$ (by Theorem 6.1.6c), then the sum of the transformed areas is just $|\det A|$ times the original area of C .

In n -dimensions, divide a given region C into many small nD -cubes of side length k , each of nD -volume k^n : each of these transforms to a small nD -parallelepiped of nD -volume $k^n|\det A|$ (by Theorem 6.1.6c), then the sum of the transformed nD -volume is just $|\det A|$ times the nD -volume of C . □

A more rigorous proof would involve upper and lower sums for the original and transformed regions, and also explicit restrictions to regions where these upper and lower sums converge to a unique nD -volume.

This property of transforming general areas and volumes also establishes the next crucial property of determinants, namely that the determinant of a product is the product of the determinants: $\det(AB) = \det(A)\det(B)$ for all matrices A and B (of the same size).

Example 6.1.11. Recall the two 2×2 matrices of Examples 6.1.1 and 6.1.2:

$$A = \begin{bmatrix} \frac{1}{2} & 1 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}.$$

Check that the determinant of their product is the product of their determinants.

Solution: First, Examples 6.1.1 and 6.1.2 computed $\det A = \frac{1}{2}$ and $\det B = -2$. Thus the product of their determinants is $\det(A) \det(B) = -1$.

Secondly, calculate the matrix product

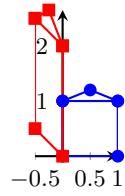
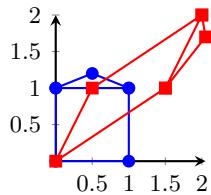
$$AB = \begin{bmatrix} \frac{1}{2} & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{3}{2} \\ 1 & 1 \end{bmatrix},$$

whose action upon multiplication is illustrated in the margin. By (6.1), $\det(AB) = \frac{1}{2} \cdot 1 - \frac{3}{2} \cdot 1 = -1$, as required.

Thirdly, we should also check the other product (as the question does not specify the order of the product):

$$BA = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} -\frac{1}{2} & 0 \\ \frac{1}{2} & 2 \end{bmatrix},$$

whose action upon multiplication is illustrated in the margin. By (6.1), $\det(BA) = -\frac{1}{2} \cdot 2 - 0 \cdot \frac{1}{2} = -1$, as required. ■



Theorem 6.1.12. For any two $n \times n$ matrices A and B , $\det(AB) = \det(A) \det(B)$. Further, for $n \times n$ matrices A_1, A_2, \dots, A_ℓ , $\det(A_1 A_2 \cdots A_\ell) = \det(A_1) \det(A_2) \cdots \det(A_\ell)$.

Proof. Consider the unit n D-cube C , its image C' upon transforming by B , and the image C'' after transforming C' by A . That is, each edge e_j of cube C is mapped to the edge Be_j of C' , which is in turn mapped to edge $A(Be_j)$ of C'' . By Definition 6.1.4, C' has (signed) n D-volume $\det B$. Theorem 6.1.10 implies C'' has (signed) n D-volume $\det A$ times that of C' ; that is, C'' has (signed) n D-volume $\det(A) \det(B)$.

By associativity, the j th edge $A(Be_j)$ of C'' is the same as $(AB)e_j$ and so C'' is the image of C under the transformation by matrix (AB) . Consequently, the (signed) n D-volume of C'' is alternatively given by $\det(AB)$. These two expressions for the n D-volume of C'' must be equal: that is, $\det(AB) = \det(A) \det(B)$.

Exercise 6.1.13 uses induction to then prove the second statement in the theorem that $\det(A_1 A_2 \cdots A_\ell) = \det(A_1) \det(A_2) \cdots \det(A_\ell)$. □

Example 6.1.13. (a) Confirm the product rule for determinants, Theorem 6.1.12, for the product

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

Solution: Although the determinant of the left-hand matrix is $(-3)(-3) - 3(-2) = 9 - (-6) = 15$, we cannot confirm the product rule because it does not apply: the matrices on the right-hand side are not square matrices and so do not have determinants.

(b) Given $\det A = 2$ and $\det B = \pi$, what is $\det(AB)$?

Solution: Strictly, there is no answer as we do not know that matrices A and B are of the same size. However, if we are *additionally* given that A and B are the same size, then Theorem 6.1.12 gives $\det(AB) = \det(A)\det(B) = 2\pi$. ■

Example 6.1.14. Use the product theorem to help find the determinant of matrix

$$C = \begin{bmatrix} 45 & -15 & 30 \\ -2\pi & \pi & 2\pi \\ \frac{1}{9} & \frac{2}{9} & -\frac{1}{3} \end{bmatrix}.$$

Solution: One route is to observe that there is a common factor in each row of the matrix so it may be factored as

$$C = \begin{bmatrix} 15 & 0 & 0 \\ 0 & \pi & 0 \\ 0 & 0 & \frac{1}{9} \end{bmatrix} \begin{bmatrix} 3 & -1 & 2 \\ -2 & 1 & 2 \\ 1 & 2 & -3 \end{bmatrix}.$$

The first matrix, being diagonal, has determinant that is the product of its diagonal elements (Theorem 6.1.6a) so its determinant = $15\pi\frac{1}{9} = \frac{5}{3}\pi$. The second matrix, from (6.1), has determinant = $-9 - 2 - 8 - 2 - 12 + 6 = -27$. Theorem 6.1.12 then gives $\det C = \frac{5}{3}\pi \cdot (-27) = -45\pi$. ■

Theorem 6.1.15. For any square matrix A , $\det A = \det(A^T)$.

Proof. Use an SVD of the matrix, say $A = USV^T$. Then

$$\begin{aligned} \det A &= \det(USV^T) \\ &= \det(U)\det(S)\det(V^T) \quad (\text{by Theorem 6.1.12 twice}) \end{aligned}$$

$$\begin{aligned}
&= \det(U^T)(\sigma_1\sigma_2 \cdots \sigma_n) \det(V) \quad (\text{by Theorem 6.1.6}) \\
&= \det(U^T) \det(S^T) \det(V) \\
&= \det(V) \det(S^T) \det(U^T) \quad (\text{by scalar commutativity}) \\
&= \det(VS^TU^T) \quad (\text{by Theorem 6.1.12 twice}) \\
&= \det[(USV^T)^T] = \det(A^T).
\end{aligned}$$

□

Example 6.1.16. Example 6.1.13a determined that $\det \begin{bmatrix} -3 & -2 \\ 3 & -3 \end{bmatrix} = 15$.

By (6.1), its transpose has determinant

$$\det \begin{bmatrix} -3 & 3 \\ -2 & -3 \end{bmatrix} = (-3)^2 - 3(-2) = 9 + 6 = 15.$$

The determinants are the same. ■

Example 6.1.17. A general 3×3 matrix $A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$ has determinant

$\det A = |A| = aei + bfg + cdh - ceg - afh - bdi$. Its transpose,

$$A^T = \begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix}, \text{ from the rule (6.1)}$$

$$\begin{bmatrix} a & d & g \\ b & e & h \\ c & f & i \end{bmatrix} \xrightarrow{\text{rule (6.1)}} \begin{bmatrix} a & b \\ d & e \\ g & h \end{bmatrix}$$

has determinant

$$\det A^T = aei + dhc + gbh - gec - ahf - dbi = \det A.$$

■

Example 6.1.18. Recall Example 3.3.3 showed that the following matrix has the given SVD:

$$\begin{aligned}
A &= \begin{bmatrix} -4 & -2 & 4 \\ -8 & -1 & -4 \\ 6 & 6 & 0 \end{bmatrix} \\
&= \begin{bmatrix} \frac{1}{3} & -\frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{2}{3} & \frac{1}{3} \\ -\frac{2}{3} & \frac{1}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 3 \end{bmatrix} \begin{bmatrix} -\frac{8}{9} & -\frac{1}{9} & -\frac{4}{9} \\ -\frac{4}{9} & \frac{4}{9} & \frac{7}{9} \\ -\frac{1}{9} & -\frac{8}{9} & \frac{4}{9} \end{bmatrix}^T.
\end{aligned}$$

Use this SVD to find the magnitude $|\det A|$.

Solution: Given the SVD $A = USV^T$, the Product Theorem 6.1.12 gives $\det A = \det(U) \det(S) \det(V^T)$.

- $\det U = \pm 1$ by Theorem 6.1.6b as U is an orthogonal matrix.
- Using the Transpose Theorem 6.1.15, $\det(V^T) = \det V = \pm 1$ as V is orthogonal.
- Since $S = \text{diag}(12, 6, 3)$ is diagonal, Theorem 6.1.6a asserts its determinant is the product of the diagonal elements; that is, $\det S = 12 \cdot 6 \cdot 3 = 216$.

Consequently $\det A = (\pm 1)216(\pm 1) = \pm 216$, so $|\det A| = 216$. ■

Theorem 6.1.19. *For any $n \times n$ square matrix A , the magnitude of its determinant $|\det A| = \sigma_1\sigma_2 \cdots \sigma_n$, the product of all its singular values.*

Proof. Consider an SVD of the matrix $A = USV^T$. Theorems 6.1.6 and 6.1.12 empowers the following identities:

$$\begin{aligned}\det A &= \det(USV^T) \\ &= \det(U)\det(S)\det(V^T) \quad (\text{by Theorem 6.1.12}) \\ &= (\pm 1)\det(S)(\pm 1) \quad (\text{by Theorem 6.1.6b}) \\ &= \pm \det S \quad (\text{and since } S \text{ is diagonal}) \\ &= \pm \sigma_1\sigma_2 \cdots \sigma_n. \quad (\text{by Theorem 6.1.6a})\end{aligned}$$

Hence $|\det A| = \sigma_1\sigma_2 \cdots \sigma_n$. □

Example 6.1.20. Confirm Theorem 6.1.19 for the matrix $A = \begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix}$ of Example 3.3.2.

Solution: Example 3.3.2 gave the SVD

$$\begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix} = \begin{bmatrix} \frac{3}{5} & -\frac{4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{bmatrix} \begin{bmatrix} 10\sqrt{2} & 0 \\ 0 & 5\sqrt{2} \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}^T,$$

so Theorem 6.1.19 asserts $|\det A| = 10\sqrt{2} \cdot 5\sqrt{2} = 100$. Using (6.1) directly, $\det A = 10 \cdot 11 - 2 \cdot 5 = 110 - 10 = 100$ which agrees with the product of the singular values. ■

Example 6.1.21. Use an SVD of the following matrix to find the magnitude of its determinant:

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

Solution: Although an SVD exists for this matrix and so we could form the product of its singular values, the concept of a determinant only applies to square matrices and so there is no such thing as a determinant for this 3×4 matrix A . The task is meaningless.

■

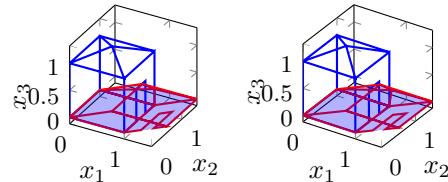
One of the main reasons for studying determinants is to establish when solutions to linear equations may exist or not (albeit only applicable when there are n linear equations in n unknowns). One example lies in finding eigenvalues by hand (section 4.1.1) where we solve $\det(A - \lambda I) = 0$.

Recall that for 2×2 and 3×3 matrices we commented that a matrix is invertible only when its determinant is non-zero. Theorem 6.1.23 establishes this in general. The geometric reason for this connection between invertibility and determinants is that when a determinant is zero the action of multiplying by the matrix ‘squashes’ the unit n D-cube into a n D-parallelepiped of zero thickness. Such extreme squashing cannot be uniquely undone.

Example 6.1.22. Consider multiplication by the matrix

$$A = \begin{bmatrix} 1 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

whose effect on the unit cube is illustrated below:



As illustrated, this matrix squashes the unit cube onto the x_1x_2 -plane ($x_3 = 0$). Consequently the resultant volume is zero and so $\det A = 0$. Because many points in 3D space are squashed onto the same point in the $x_3 = 0$ plane, the action of the matrix cannot be undone. Hence the matrix is not invertible. That the matrix is not invertible and its determinant is zero is not a coincidence. ■

Theorem 6.1.23. *A square matrix A is invertible iff $\det A \neq 0$. If matrix A is invertible, then $\det(A^{-1}) = 1/(\det A)$.*

Proof. First, Theorem 6.1.19 establishes that $\det A \neq 0$ iff all the singular values of square matrix A are non-zero, which by Theorem 3.4.35d is iff matrix A is invertible.

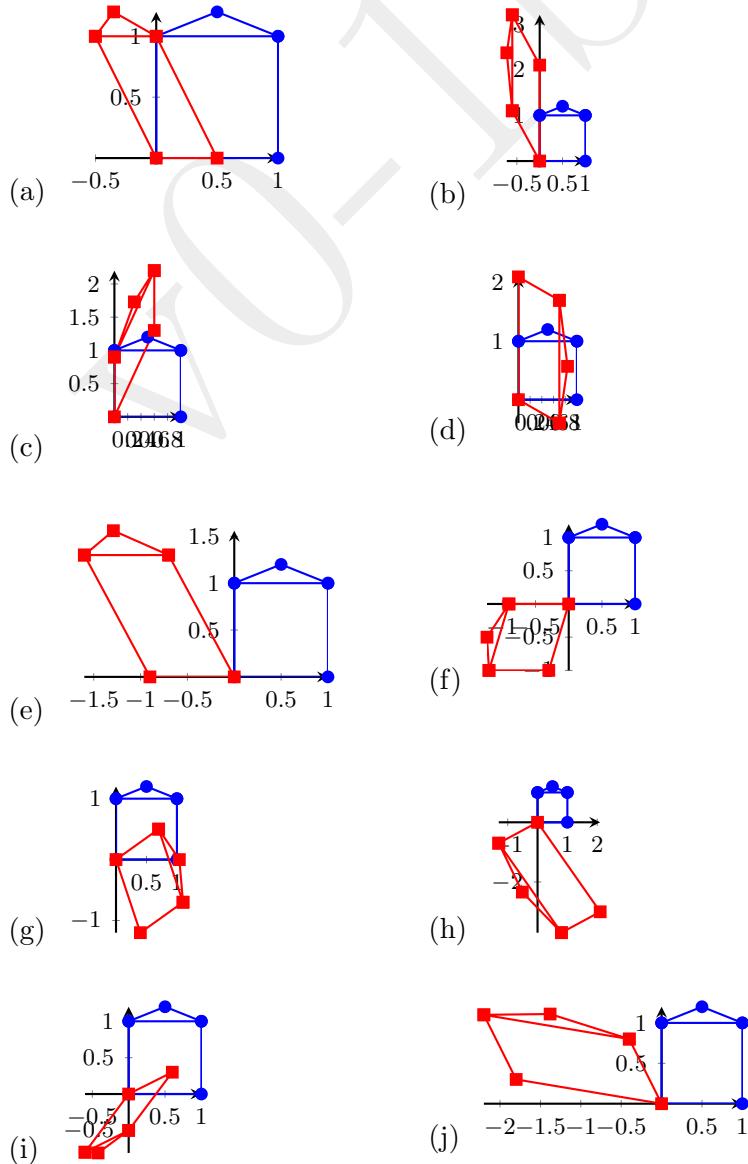
Second, as matrix A is invertible, an inverse A^{-1} exists such that $AA^{-1} = I_n$. Then the product of determinants

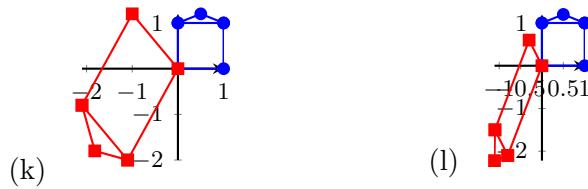
$$\begin{aligned}\det(A)\det(A^{-1}) &= \det(AA^{-1}) \quad (\text{by Theorem 6.1.12}) \\ &= \det I_n \quad (\text{from } AA^{-1} = I_n) \\ &= 1 \quad (\text{by Example 6.1.7})\end{aligned}$$

For an invertible matrix A , $\det A \neq 0$; hence dividing by $\det A$ gives $\det(A^{-1}) = 1/\det A$. \square

6.1.1 Exercises

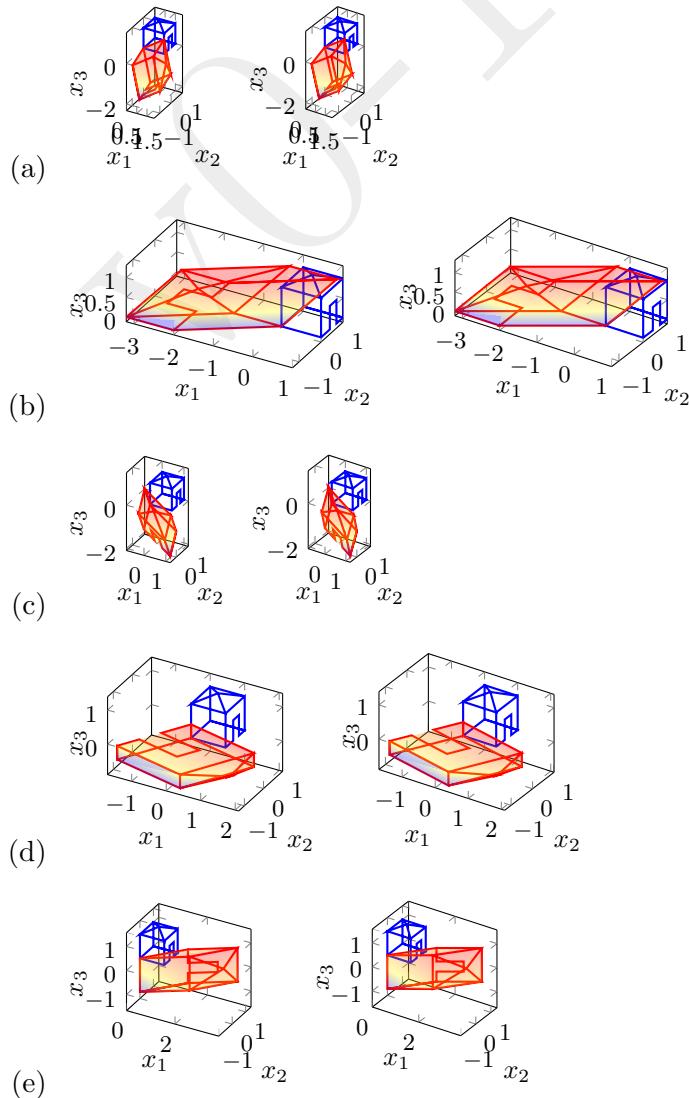
Exercise 6.1.1. For each of the given illustrations of a linear transformation of the unit square, ‘guesstimate’ by eye the determinant of the matrix of the transformation (estimate to say 33% or so).

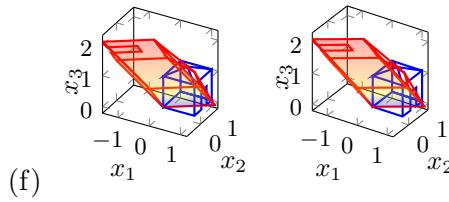




Exercise 6.1.2. For each of the transformations illustrated in Exercise 6.1.1, estimate the matrix of the linear transformation (to within say 10%). Then use formula (6.1) to estimate the determinant of your matrix and confirm Exercise 6.1.1.

Exercise 6.1.3. For each of the following stereo illustrations of a linear transformation of the unit cube, estimate the matrix of the linear transformation (to within say 20%). Then use formula (6.1) to estimate the determinant of the matrix of the transformation.





Exercise 6.1.4. For each of the following matrices, use (6.1) to find all the values of k for which the matrix is *not* invertible.

$$(a) A = \begin{bmatrix} 0 & 6 - 2k \\ -2k & -4 \end{bmatrix}$$

$$(b) B = \begin{bmatrix} 3k & 4 - k \\ -4 & 0 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} 2 & 0 & -2k \\ 4k & 0 & -1 \\ 0 & k & -4 + 3k \end{bmatrix}$$

$$(d) D = \begin{bmatrix} 2 & -2 - 4k & -1 + k \\ -1 - k & 0 & -5k \\ 0 & 0 & 4 + 2k \end{bmatrix}$$

$$(e) E = \begin{bmatrix} -1 - 2k & 5 & 1 - k \\ 0 & -2 & 0 \\ 0 & 0 & -7 + k \end{bmatrix}$$

$$(f) F = \begin{bmatrix} k & 0 & -3 - 3k \\ 3 & 7 & 3k \\ 0 & 2k & 2 - k \end{bmatrix}$$

Exercise 6.1.5. Find the determinants of each of the following matrices.

$$(a) \begin{bmatrix} -3 & 0 & 0 & 0 \\ 0 & -4 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix}$$

$$(b) \begin{bmatrix} -3 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$(c) \begin{bmatrix} -1/2 & 0 & 0 & 0 \\ 0 & 3/2 & 0 & 0 \\ 0 & 0 & -5/2 & 0 \\ 0 & 0 & 0 & -1/2 \end{bmatrix}$$

$$(d) \begin{bmatrix} 5/6 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 7/6 & 0 \\ 0 & 0 & 0 & 2/3 \end{bmatrix}$$

$$(e) \begin{bmatrix} 1/3 & 0 & -4/3 & -1 \\ -1/3 & 1/3 & 5/3 & 4/3 \\ 0 & -1/3 & 5/3 & -1 \\ -1/3 & -1/3 & 8/3 & -1/3 \end{bmatrix}$$

given $\det \begin{bmatrix} 1 & 0 & -4 & -3 \\ -1 & 1 & 5 & 4 \\ 0 & -1 & 5 & -3 \\ -1 & -1 & 8 & -1 \end{bmatrix} = -8$

$$(f) \begin{bmatrix} -1 & 3/2 & -1/2 & 1/2 \\ -2 & -1/2 & -1 & -3/2 \\ -0 & -3/2 & -5/2 & 1/2 \\ -0 & -1/2 & 3 & 3/2 \end{bmatrix}$$

given $\det \begin{bmatrix} 2 & -3 & 1 & -1 \\ 4 & 1 & 2 & 3 \\ 0 & 3 & 5 & -1 \\ 0 & 1 & -6 & -3 \end{bmatrix} = -524$

$$(g) \begin{bmatrix} -0 & -2/3 & -2 & -4/3 \\ -0 & 1/3 & -1/3 & 1/3 \\ 1 & -1/3 & -5/3 & 2/3 \\ -7/3 & 1/3 & 4/3 & 2/3 \end{bmatrix}$$

given $\det \begin{bmatrix} 0 & 2 & 6 & 4 \\ 0 & -1 & 1 & -1 \\ -3 & 1 & 5 & -2 \\ 7 & -1 & -4 & -2 \end{bmatrix} = 246$

$$(h) \begin{bmatrix} -12 & -16 & -4 & 12 \\ -4 & 8 & -4 & -16 \\ -0 & -4 & -12 & 4 \\ 4 & -4 & -8 & 4 \end{bmatrix}$$

given $\det \begin{bmatrix} 3 & 4 & 1 & -3 \\ 1 & -2 & 1 & 4 \\ 0 & 1 & 3 & -1 \\ -1 & 1 & 2 & -1 \end{bmatrix} = -34$

Exercise 6.1.6. Use Theorems 3.2.21 and 6.1.6a to prove that for any diagonal square matrix D , $\det(D^{-1}) = 1/\det D$ provided $\det D \neq 0$.

Exercise 6.1.7. For each pair of following matrices, by computing in full using (6.1) confirm $\det(AB) = \det(BA) = \det(A)\det(B)$. Show your working.

$$(a) A = \begin{bmatrix} -2 & 3/2 \\ -1/2 & 0 \end{bmatrix}, B = \begin{bmatrix} 0 & -1/2 \\ 1/2 & -1 \end{bmatrix}$$

$$(b) A = \begin{bmatrix} -1 & 1 \\ -7/2 & -3/2 \end{bmatrix}, B = \begin{bmatrix} 1/2 & -5/2 \\ -1/2 & 3/2 \end{bmatrix}$$

$$(c) A = \begin{bmatrix} 4 & -1 \\ 4 & -3 \end{bmatrix}, B = \begin{bmatrix} -3/2 & 0 \\ -1 & 3/2 \end{bmatrix}$$

$$(d) A = \begin{bmatrix} 0 & -1 \\ 1/2 & 0 \end{bmatrix}, B = \begin{bmatrix} 2 & 0 \\ 1 & -2 \end{bmatrix}$$

$$(e) A = \begin{bmatrix} 0 & -1/2 & 1/2 \\ 2 & 0 & 0 \\ 0 & 1/2 & 0 \end{bmatrix}, B = \begin{bmatrix} -1 & -1 & 0 \\ 0 & 0 & -1 \\ -2 & 0 & 0 \end{bmatrix}$$

$$(f) A = \begin{bmatrix} 0 & 0 & -2 \\ 0 & 1 & 0 \\ -1/2 & 0 & 0 \end{bmatrix}, B = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

$$(g) \quad A = \begin{bmatrix} -2 & -1/2 & 0 \\ 0 & 0 & 2 \\ 0 & -1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} -1 & -1 & 0 \\ 0 & 0 & -1 \\ 0 & -2 & 0 \end{bmatrix}$$

$$(h) \quad A = \begin{bmatrix} 1 & 2 & -3/2 \\ 0 & 3/2 & 0 \\ -1 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 0 \\ -4 & 2 & 0 \\ -1 & 0 & -5 \end{bmatrix}$$

Exercise 6.1.8. Given that $\det(AB) = \det(A)\det(B)$ for any two square matrices of the same size, prove that $\det(AB) = \det(BA)$ (despite $AB \neq BA$ in general).

Exercise 6.1.9. Given that $n \times n$ matrices A and B have $\det A = 3$ and $\det B = -5$, determine the following determinants (if possible).

(a) $\det(AB)$

(b) $\det(B^2)$

(c) $\det(A^4)$

(d) $\det(A + B)$

(e) $\det(A^{-1}B)$

(f) $\det(2A)$

(g) $\det(B^T/2)$

(h) $\det(B^T B)$

Exercise 6.1.10. Let A and P be square matrices of the same size, and let matrix P be invertible. Prove that $\det(P^{-1}AP) = \det A$.

Exercise 6.1.11. Suppose square matrix A satisfies $A^2 = A$ (called idempotent). Determine all possible values of $\det A$. Invent and verify a nontrivial example of a idempotent matrix.

Exercise 6.1.12. Suppose a square matrix A satisfies $A^p = O$ for some integer exponent $p \geq 2$ (called nilpotent). Determine all possible values of $\det A$. Invent and verify a nontrivial example of a nilpotent matrix.

Exercise 6.1.13. Recall that $\det(AB) = \det(A)\det(B)$ for any two square matrices of the same size. For $n \times n$ matrices A_1, A_2, \dots, A_ℓ , use induction to prove $\det(A_1 A_2 \cdots A_\ell) = \det(A_1) \det(A_2) \cdots \det(A_\ell)$ for all integer $\ell \geq 2$.

Exercise 6.1.14. To complement the algebraic argument of Theorem 6.1.23, use a geometric argument based upon the transformation of n D-volumes to establish that $\det(A^{-1}) = 1/(\det A)$ for an invertible matrix A .

Exercise 6.1.15. Suppose square matrix $A = \begin{bmatrix} P & O \\ O & Q \end{bmatrix}$ for some square matrices P and Q , and appropriately sized zero matrices O . Give a geometric argument justifying that $\det A = (\det P)(\det Q)$.

6.2 Laplace expansion theorem for determinants

Section Contents

6.2.1 Exercises	543
---------------------------	-----

This section develops a so-called row/column algebraic expansion for determinants. This expansion is useful for theoretical purposes. But there are vastly more efficient ways of computing determinants than using a row/column expansion. In Matlab/Octave one may invoke `det(A)` to compute the determinant of a matrix. You may find this function useful for checking the results of some examples and exercises. However, just like computing an inverse, computing the determinant is expensive and error prone. In medium to large scale problems avoid computing the determinant, something else is almost always better.

Nonetheless, a row/column algebraic expansion for a determinant is useful for small matrix problems, as well as its beautiful theoretical uses. We start with examples of row properties that underpin a row/column algebraic expansion.

Example 6.2.1 (Theorem 6.2.5a). Example 6.1.22 argued geometrically that the determinant is zero for the matrix

$$A = \begin{bmatrix} 1 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

Confirm this determinant algebraically.

Solution: Using (6.1), $\det A = 1 \cdot \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 1 \cdot 0 + 0 \cdot 0 \cdot 0 - 0 \cdot \frac{1}{2} \cdot 0 - 1 \cdot 1 \cdot 0 - \frac{1}{2} \cdot 0 \cdot 0 = 0$. (There is a zero from the last row in every term.)

■

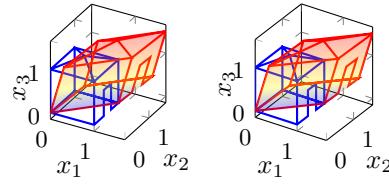
Example 6.2.2 (Theorem 6.2.5b). Consider the matrix with two identical rows,

$$A = \begin{bmatrix} 1 & \frac{1}{2} & \frac{1}{5} \\ 1 & \frac{1}{2} & \frac{1}{5} \\ 0 & \frac{1}{2} & 1 \end{bmatrix}.$$

Confirm algebraically that its determinant is zero. Give a geometric reason for why its determinant has to be zero.

Solution: Using (6.1), $\det A = 1 \cdot \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot \frac{1}{5} \cdot 0 + \frac{1}{5} \cdot 1 \cdot \frac{1}{2} - \frac{1}{5} \cdot 0 - 1 \cdot \frac{1}{5} \cdot \frac{1}{2} - \frac{1}{2} \cdot 1 \cdot 1 = \frac{1}{2} + \frac{1}{10} - \frac{1}{10} - \frac{1}{2} = 0$.

Geometrically, consider the image of the unit cube under multiplication by A illustrated in stereo below.



Because the first two rows of A are identical the first two components of $A\mathbf{x}$ are always identical and hence all points are mapped onto the plane $x_1 = x_2$. The image of the cube thus has zero thickness and hence zero volume. By Definition 6.1.4, $\det A = 0$.

■

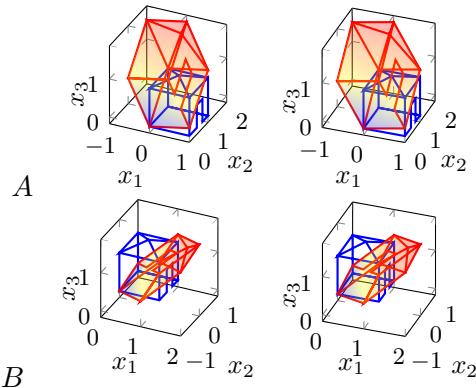
Example 6.2.3 (Theorem 6.2.5c). Consider the two matrices with two rows swapped:

$$A = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & 1 \\ \frac{1}{5} & \frac{1}{2} & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 1 \\ 1 & -1 & 0 \\ \frac{1}{5} & \frac{1}{2} & 1 \end{bmatrix}$$

Confirm algebraically that their determinants are the negative of each other. Give a geometric reason why this should be so.

Solution: Using (6.1), $\det A = 1 \cdot 1 \cdot 1 + (-1) \cdot 1 \cdot \frac{1}{5} + 0 \cdot 0 \cdot \frac{1}{2} - 0 \cdot 1 \cdot \frac{1}{5} - 1 \cdot 1 \cdot \frac{1}{2} - (-1) \cdot 0 \cdot 1 = 1 - \frac{1}{5} - \frac{1}{2} = \frac{3}{10}$. Using (6.1), $\det B = 0 \cdot (-1) \cdot 1 + 1 \cdot 0 \cdot \frac{1}{5} + 1 \cdot 1 \cdot \frac{1}{2} - 1 \cdot (-1) \cdot \frac{1}{5} - 0 \cdot 0 \cdot \frac{1}{2} - 1 \cdot 1 \cdot 1 = \frac{1}{2} + \frac{1}{5} - 1 = -\frac{3}{10} = -\det A$.

Geometrically, since the first two rows in A and B are swapped that means that multiplying by the matrix as in $A\mathbf{x}$ and $B\mathbf{x}$ has the first two components swapped. Hence $A\mathbf{x}$ and $B\mathbf{x}$ are always the reflection of each other in the plane $x_1 = x_2$. Consequently, the images of the unit cubes under multiplication by A and by B are the reflection of each other in the plane $x_1 = x_2$, as illustrated below, and so the determinants must be the negative of each other.



■

Example 6.2.4 (Theorem 6.2.5d). Compute the determinant of the matrix

$$B = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 1 & 1 \\ 2 & 5 & 10 \end{bmatrix}.$$

Compare B with matrix A given in Example 6.2.3, and compare their determinants.

Solution: Using (6.1), $\det B = 1 \cdot 1 \cdot 10 + (-1) \cdot 1 \cdot 2 + 0 \cdot 0 \cdot 5 - 0 \cdot 1 \cdot 2 - 1 \cdot 1 \cdot 5 - (-1) \cdot 0 \cdot 10 = 10 - 2 - 5 = 3$. Matrix B is the same as matrix A except the third row is a factor of ten times bigger. Correspondingly, $\det B = 3 = 10 \times \frac{3}{10} = 10 \det A$.

■

The above four examples are specific cases of the four general properties established as the four parts of the following theorem.

Theorem 6.2.5 (row and column properties of determinants). *Let A be an $n \times n$ matrix.*

- (a) *If A has a zero row or column, then $\det A = 0$.*
- (b) *If A has two identical rows or columns, then $\det A = 0$.*
- (c) *Let B be obtained by interchanging two rows or columns of A , then $\det B = -\det A$.*
- (d) *Let B be obtained by multiplying any one row or column of A by a scalar k , then $\det B = k \det A$.*

Proof. We establish the properties for matrix rows. Then the same property holds for the columns because $\det(A^T) = \det(A)$ (Theorem 6.1.15).

6.2.5d Suppose row i of matrix A is multiplied by k to give a new matrix B . Let the diagonal matrix

$$D := \text{diag}(1, \dots, 1, k, 1, \dots, 1),$$

with the factor k being in the i th row and column. Then $\det D = 1 \times \dots \times 1 \times k \times 1 \times \dots \times 1 = k$ by Theorem 6.1.6a. Because multiplication by D multiplies the i th row by the factor k and leaves everything else the same, $B = DA$. Equate determinants of both sides and use the product Theorem 6.1.12: $\det B = \det(DA) = \det(D) \det(A) = k \det A$.

6.2.5a This arises from the case $k = 0$ of property 6.2.5d. Then $A = DA$ because multiplying the i th row of A by $k = 0$ maintains that the row is zero. Consequently, $\det A = \det(DA) = \det(D) \det(A) = 0 \det(A)$ and hence $\det A = 0$.

6.2.5c Suppose rows i and j of matrix A are swapped to form matrix B . Let the matrix E be the identity except with rows i and j swapped:

$$E := \begin{bmatrix} \ddots & & & & \\ & 0 & & 1 & \\ & & \ddots & & \\ & 1 & & 0 & \\ & & & & \ddots \\ i & & & j & \end{bmatrix} \begin{array}{l} \text{row } i \\ \text{row } j \end{array}$$

where the diagonal dots \ddots denote diagonals of ones, and all other entries of E are zero. Then $B = EA$ as multiplication by E copies row i into row j and vice versa. Equate determinants of both sides and use Theorem 6.1.12: $\det B = \det(EA) = \det(E) \det(A)$.

To find $\det E$ observe that $EE^T = E^2 = I_n$ so E is orthogonal and hence $\det E = \pm 1$ by Theorem 6.1.6b. Geometrically, multiplication by E is a simple reflection in the $(n-1)$ D-plane $x_i = x_j$ hence its determinant must be negative, so $\det E = -1$. Consequently, $\det B = \det(E) \det(A) = -\det A$.

6.2.5b Suppose rows i and j of matrix A are identical. Using the matrix E to swap these two identical rows results in the same matrix: that is, $A = EA$. Take determinants of both sides: $\det A = \det(EA) = \det(E) \det(A)$. Since $\det E = -1$ it follows that $\det A = -\det A$. Zero is the only number that equals its negative: thus $\det A = 0$.

□

Example 6.2.6. You are given that $\det A = -9$ for the matrix

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

Use Theorem 6.2.5 to find the determinant of the following matrices, giving reasons.

$$(a) \begin{bmatrix} 0 & 2 & 3 & 0 & 4 \\ -2 & 2 & -2 & 0 & -3 \\ 4 & -2 & -4 & 0 & 0 \\ 2 & -1 & -4 & 0 & 2 \\ 5 & 4 & 3 & 0 & -5 \end{bmatrix}$$

Solution: $\det = 0$ as the fourth column is all zeros.

$$(b) \begin{bmatrix} 0 & 2 & 3 & 1 & 4 \\ -2 & 2 & -2 & 0 & -3 \\ 2 & -1 & -4 & 2 & 2 \\ 4 & -2 & -4 & 1 & 0 \\ 5 & 4 & 3 & -2 & -5 \end{bmatrix}$$

Solution: $\det = -\det A = +9$ as the 3rd and 4th rows are swapped.

$$(c) \begin{bmatrix} 0 & 2 & 3 & 1 & 4 \\ -2 & 2 & -2 & 0 & -3 \\ 4 & -2 & -4 & 1 & 0 \\ 2 & -1 & -4 & 2 & 2 \\ -2 & 2 & -2 & 0 & -3 \end{bmatrix}$$

Solution: $\det = 0$ as the 2nd and 5th rows are identical.

$$(d) \begin{bmatrix} 0 & 1 & 3 & 1 & 4 \\ -2 & 1 & -2 & 0 & -3 \\ 4 & -1 & -4 & 1 & 0 \\ 2 & -\frac{1}{2} & -4 & 2 & 2 \\ 5 & 2 & 3 & -2 & -5 \end{bmatrix}$$

Solution: $\det = \frac{1}{2} \det A = -\frac{9}{2}$ as the 2nd column is half that of A

$$(e) \begin{bmatrix} -2 & 2 & -6 & 0 & -3 \\ 4 & -2 & -12 & 1 & 0 \\ 2 & -1 & -12 & 2 & 2 \\ 0 & 2 & 9 & 1 & 4 \\ 5 & 4 & 9 & -2 & -5 \end{bmatrix}$$

Solution: $\det = 3(-\det A) = 27$ as this matrix is A with 1st and 4th rows swapped, and the 3rd column multiplied by 3.

$$(f) \begin{bmatrix} 0 & 3 & 3 & 1 & 4 \\ -2 & 0 & -2 & 0 & -5 \\ 5 & -1 & -4 & 1 & 0 \\ 2 & -1 & -4 & 2 & 2 \\ 5 & 4 & 6 & -2 & -5 \end{bmatrix}$$

Solution: Cannot answer as none of these row and column operations on A appear to give this matrix.

Example 6.2.7. Without evaluating the determinant, use Theorem 6.2.5 to establish that the determinant equation

$$\begin{vmatrix} 1 & x & y \\ 1 & 2 & 3 \\ 1 & 4 & 5 \end{vmatrix} = 0 \quad (6.2)$$

is the equation of the straight line in the xy -plane through the two points $(2, 3)$ and $(4, 5)$.

Solution: First, equation (6.2) is linear in variables x and y as at most one of x and y occur in each term of the determinant expansion:

$$\begin{bmatrix} 1 & x & y \\ 1 & 2 & 3 \\ 1 & 4 & 5 \end{bmatrix} \begin{matrix} 1 & x \\ 1 & 2 \\ 1 & 4 \end{matrix}$$

Since the equation is linear in x and y , any solution set of (6.2) must be a straight line. Second, equation (6.2) is satisfied when $(x, y) = (2, 3)$ or when $(x, y) = (4, 5)$ as then two rows in the determinant are identical, and so Theorem 6.2.5b assures us the determinant is zero. Thus the solution straight line passes through the two required points.

Let's evaluate the determinant to check: equation (6.2) becomes $1 \cdot 2 \cdot 5 + x \cdot 3 \cdot 1 + y \cdot 1 \cdot 4 - y \cdot 2 \cdot 1 - 1 \cdot 3 \cdot 4 - x \cdot 1 \cdot 5 = -2 - 2x + 2y = 0$. That is, $y = x + 1$ which does indeed pass through $(2, 3)$ and $(4, 5)$.

■

Example 6.2.8. Without evaluating the determinant, use Theorem 6.2.5 to establish that the determinant equation

$$\begin{vmatrix} x & y & z \\ -1 & -2 & 2 \\ 3 & 5 & 2 \end{vmatrix} = 0$$

is, in xyz -space, the equation of the plane through the origin and the two points $(-1, -2, 2)$ and $(3, 5, 2)$.

Solution: As in the previous example, the determinant is linear in x , y and z , so the solutions must be those of one linear equation, namely a plane (Section 1.3.4).

- The solutions include the origin since when $x = y = z = 0$ the first row of the matrix is zero, hence the determinant is zero, and the equation satisfied.
- The solutions include the points $(-1, -2, 2)$ and $(3, 5, 2)$ since when (x, y, z) are either, then two rows in the determinant are identical, so the determinant is zero, and the equation satisfied.

Hence the solutions are the points in the plane passing through the origin and $(-1, -2, 2)$ and $(3, 5, 2)$.

■

The next step in developing a general ‘formula’ for a determinant is the special class of matrices for which one column or row is zero except for one element.

Example 6.2.9. Find the determinant of $A = \begin{bmatrix} -2 & -1 & -1 \\ 1 & -3 & -2 \\ 0 & 0 & 2 \end{bmatrix}$ which has two zeros in its last row.

Solution: Let's use two different arguments, both illustrating the next theorem and proof.

- Using (6.1),

$$\begin{aligned} \det A &= (-2)(-3)2 + (-1)(-2) \cdot 0 + (-1)1 \cdot 0 \\ &\quad - (-1)(-3)0 - (-2)(-2)0 - (-1)1 \cdot 2 \\ &= 2[(-2)(-3) - (-1)1] = 2 \cdot 7 = 14. \end{aligned}$$

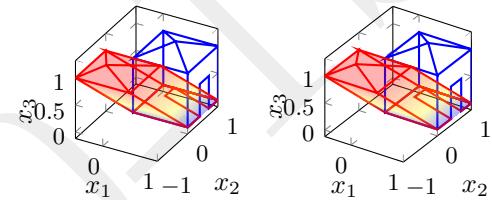
Observe $\det A = 2 \times \det \begin{bmatrix} -2 & -1 \\ 1 & -3 \end{bmatrix}$ which is the expression to be generalised.

- Alternatively we use the product rule for determinants. Recognise that the matrix A may be factored to $A = FB$ where

$$F = \begin{bmatrix} 1 & 0 & -\frac{1}{2} \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} -2 & -1 & 0 \\ 1 & -3 & 0 \\ 0 & 0 & 2 \end{bmatrix};$$

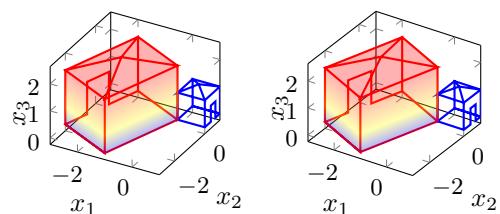
just multiply out and see that the last column of F ‘fills in’ the last column of A from that of B . Consider the geometry of the two transformations arising from multiplication by F and by B .

- Multiplication by F shears the unit cube as illustrated below.



Thus the volume of the unit cube after multiplication by F is the square base of area one, times the height of one, which is a volume of one. Consequently $\det F = 1$ by Definition 6.1.4.

- As illustrated below, multiplication by B has two components.



Firstly, the 2×2 top-left sub-matrix, being bordered by zeros, maps the unit square in the x_1x_2 -plane into the base parallelogram in the x_1x_2 -plane. The shape of the parallelogram is determined by the top-left sub-matrix $\begin{bmatrix} -2 & -1 \\ 1 & -3 \end{bmatrix}$ (to be called the A_{33} minor) acting on the unit square. Thus the area of the parallelogram is $\det A_{33} = (-2)(-3) - (-1)1 = 7$. Secondly, the 2 in the bottom corner of B stretches objects vertically by a factor of 2 to form a parallelepiped of height 2. Thus the volume of the parallelepiped is $2 \cdot \det A_{33} = 2 \cdot 7 = 14$, which is $\det B$ by Definition 6.1.4.

By the product Theorem 6.1.12, $\det A = \det(FB) = \det(F) \det(B) = 1 \cdot 14 = 14$. ■

Theorem 6.2.10 (almost zero row/column). *Let A be an $n \times n$ matrix.*

*Define the (i, j) th **minor** A_{ij} to be the $(n - 1) \times (n - 1)$ square matrix obtained from A by omitting the i th row and j th column. If, except for the entry a_{ij} , the i th row (or j th column) is all zero, then*

$$\det A = (-1)^{i+j} a_{ij} \det A_{ij}. \quad (6.3)$$

The pattern of signs in this formula, $(-1)^{i+j}$, is

$$\begin{array}{cccccc} + & - & + & - & + & \cdots \\ - & + & - & + & - & \cdots \\ + & - & + & - & + & \cdots \\ - & + & - & + & - & \cdots \\ + & - & + & - & + & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{array}$$

Proof. We establish the determinant formula (6.3) for matrix rows, then the same result holds for the columns because $\det(A^T) = \det(A)$ (Theorem 6.1.15). First, if the entry $a_{ij} = 0$, then the whole i th row (or j th column) is zero and so $\det A = 0$ by Theorem 6.2.5a. Also, the expression $(-1)^{i+j} a_{ij} \det A_{ij} = 0$ as $a_{ij} = 0$. Consequently, the identity $\det A = (-1)^{i+j} a_{ij} \det A_{ij}$ holds. The rest of this proof addresses the case $a_{ij} \neq 0$.

Second, consider the special case when the last row *and* last column of matrix A is all zero except for $a_{nn} \neq 0$; that is,

$$A = \begin{bmatrix} A_{nn} & \mathbf{0} \\ \mathbf{0}^T & a_{nn} \end{bmatrix}$$

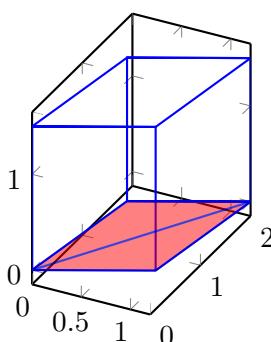
for the minor A_{nn} and $\mathbf{0} \in \mathbb{R}^{n-1}$. Recall Definition 6.1.4: the image of the n D-cube under multiplication by the matrix A is the image of the $(n - 1)$ D-cube under multiplication by A_{nn} extended orthogonally a length a_{nn} in the orthogonal direction e_n (as illustrated in the margin in 3D). The volume of the n D-image is thus $a_{nn} \times (\text{volume of the } (n - 1)\text{D-image})$. Consequently, $\det A = a_{nn} \det A_{nn}$.

Third, consider the special case when the last row of matrix A is all zero except for $a_{nn} \neq 0$; that is,

$$A = \begin{bmatrix} A_{nn} & \mathbf{a}'_n \\ \mathbf{0}^T & a_{nn} \end{bmatrix}$$

for the minor A_{nn} , and where $\mathbf{a}'_n = (a_{1n}, a_{2n}, \dots, a_{n-1,n})$. Define the two $n \times n$ matrices

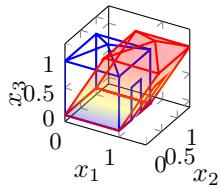
$$F := \begin{bmatrix} I_{n-1} & \mathbf{a}'_n / a_{nn} \\ \mathbf{0}^T & 1 \end{bmatrix} \quad \text{and} \quad B := \begin{bmatrix} A_{nn} & \mathbf{0} \\ \mathbf{0}^T & a_{nn} \end{bmatrix}.$$



Then $A = FB$ since

$$\begin{aligned} FB &= \begin{bmatrix} I_{n-1} & \mathbf{a}'_n/a_{nn} \\ \mathbf{0}^T & 1 \end{bmatrix} \begin{bmatrix} A_{nn} & \mathbf{0} \\ \mathbf{0}^T & a_{nn} \end{bmatrix} \\ &= \begin{bmatrix} I_{n-1}A_{nn} + \mathbf{a}'_n/a_{nn}\mathbf{0}^T & I_{n-1}\mathbf{0} + \mathbf{a}'_n/a_{nn} \cdot a_{nn} \\ \mathbf{0}^TA_{nn} + 1 \cdot \mathbf{0}^T & \mathbf{0}^T\mathbf{0} + 1 \cdot a_{nn} \end{bmatrix} \\ &= \begin{bmatrix} A_{nn} & \mathbf{a}'_n \\ \mathbf{0}^T & a_{nn} \end{bmatrix} = A. \end{aligned}$$

By Theorem 6.1.12, $\det A = \det(FB) = \det(F)\det(B)$. From the previous part $\det B = a_{nn}\det A_{nn}$, so we just need to determine $\det F$. As illustrated for 3D in the margin, the action of matrix F on the unit n D-cube is that of a simple shear keeping the $(n-1)$ D-cube base unchanged (due to the identity I_{n-1} in F). Since the height orthogonal to the $(n-1)$ D-cube base is unchanged (due to the one in the bottom corner of F), the action of multiplying by F leaves the volume of the unit n D-cube unchanged at one. Hence $\det F = 1$. Thus $\det A = 1 \det(B) = a_{nn}\det A_{nn}$ as required.



Fourth, suppose row i of matrix A is all zero except for entry a_{ij} . Swap rows i and $i+1$, then swap rows $i+1$ and $i+2$, and so on until the original row i is in the last row, and the order of all other rows are unchanged: this takes $(n-i)$ row swaps which changes the sign of the determinant $(n-i)$ times, that is, multiplies it by $(-1)^{n-i}$. Then swap columns j and $j+1$, then swap columns $j+1$ and $j+2$, and so on until the original column j is in the last column: this takes $(n-j)$ column swaps which change the determinant by a factor $(-1)^{n-j}$. The resulting matrix, say C , has the form

$$C = \begin{bmatrix} A_{ij} & \mathbf{a}'_j \\ \mathbf{0}^T & a_{ij} \end{bmatrix}$$

for \mathbf{a}'_j denoting the j th column of A with the i th entry omitted. Since matrix C has the form addressed in the first part, we know $\det C = a_{ij}\det A_{ij}$. From the row and column swapping, $\det A = (-1)^{n-i}(-1)^{n-j}\det C = (-1)^{2n-i-j}\det C = (-1)^{-(i+j)}\det C = (-1)^{i+j}\det C = (-1)^{i+j}a_{ij}\det A_{ij}$. \square

Example 6.2.11. Use Theorem 6.2.10 to evaluate the determinant of the following matrices.

$$(a) \begin{bmatrix} -3 & -3 & -1 \\ -3 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

Solution: There are two zeros in the bottom row so the determinant is

$$(-1)^6 2 \det \begin{bmatrix} -3 & -3 \\ -3 & 2 \end{bmatrix} = 2(-6 - 9) = -30.$$

$$(b) \begin{bmatrix} 2 & -1 & 7 \\ 0 & 3 & 0 \\ 2 & 2 & 5 \end{bmatrix}$$

Solution: There are two zeros in the middle row so the determinant is

$$(-1)^4 3 \det \begin{bmatrix} 2 & 7 \\ 2 & 5 \end{bmatrix} = 3(10 - 14) = -12.$$

$$(c) \begin{bmatrix} 2 & 4 & 3 \\ 8 & 0 & -1 \\ -5 & 0 & -2 \end{bmatrix}$$

Solution: There are two zeros in the middle column so the determinant is

$$(-1)^3 4 \det \begin{bmatrix} 8 & -1 \\ -5 & -2 \end{bmatrix} = -4(-16 - 5) = 84.$$

$$(d) \begin{bmatrix} 2 & 1 & 3 \\ 0 & -2 & -3 \\ 0 & 2 & 4 \end{bmatrix}$$

Solution: There are two zeros in the first column so the determinant is

$$(-1)^2 2 \det \begin{bmatrix} -2 & -3 \\ 2 & 4 \end{bmatrix} = 2(-8 + 6) = -4.$$

■

Example 6.2.12. Use Theorem 6.2.10 to evaluate the determinant of the so-called triangular matrix

$$A = \begin{bmatrix} 2 & -2 & 3 & 1 & 0 \\ 0 & 2 & -1 & -1 & -7 \\ 0 & 0 & 5 & -2 & -9 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 3 \end{bmatrix}$$

Solution: The last row is all zero except for the last element, so

$$\det A = (-1)^{10} 3 \det \begin{bmatrix} 2 & -2 & 3 & 1 \\ 0 & 2 & -1 & -1 \\ 0 & 0 & 5 & -2 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

(as the last row is zero except the last)

$$= 3 \cdot (-1)^8 1 \det \begin{bmatrix} 2 & -2 \\ 0 & 2 \\ 0 & 5 \end{bmatrix}$$

(as the last row is zero except the last)

$$= 3 \cdot 1 \cdot (-1)^6 5 \det \begin{bmatrix} 2 & -2 \\ 0 & 2 \end{bmatrix}$$

(as the last row is zero except the last)

$$= 3 \cdot 1 \cdot 5(-1)^4 2 \det [2]$$

$$= 3 \cdot 1 \cdot 5 \cdot 2 \cdot 2 = 60.$$

■

The relative simplicity of finding the determinant in Example 6.2.12 indicates that there is something special and memorable about matrices with zeros in the lower-left ‘triangle’. There is, as expressed by the following definition and theorem.

Definition 6.2.13. A **triangular matrix** is a square matrix where all entries are zero either to the lower-left of the diagonal or to the upper-right: ¹

- an *upper triangular matrix* has the form (although any of the a_{ij} may also be zero)

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n-1} & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n-1} & a_{2n} \\ \vdots & 0 & \ddots & \vdots & \vdots \\ 0 & \vdots & \ddots & a_{n-1n-1} & a_{n-1n} \\ 0 & 0 & \cdots & 0 & a_{nn} \end{bmatrix};$$

- a *lower triangular matrix* has the form (although any of the a_{ij} may also be zero)

$$\begin{bmatrix} a_{11} & 0 & \cdots & 0 & 0 \\ a_{21} & a_{22} & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ a_{n-11} & a_{n-12} & \cdots & a_{n-1n-1} & 0 \\ a_{n1} & a_{n2} & \cdots & a_{nn-1} & a_{nn} \end{bmatrix}.$$

Any square diagonal matrix is also an upper triangular matrix, and also a lower triangular matrix. Thus the following theorem encompasses square diagonal matrices and so generalises Theorem 6.1.6a.

Theorem 6.2.14 (triangular matrix). Let A be an $n \times n$ triangular matrix. The determinant of A is the product of the diagonal entries, $\det A = a_{11}a_{22} \cdots a_{nn}$.

Proof. A little induction proves the determinant of a triangular matrix is the product of its diagonal entries: only consider upper triangular matrices as transposing the matrix caters for lower triangular matrices.

First, for 1×1 matrices the result is trivial. The results is also straightforward for 2×2 matrices since the determinant

$$\begin{vmatrix} a_{11} & a_{12} \\ 0 & a_{22} \end{vmatrix} = a_{11}a_{22} - 0a_{12} = a_{11}a_{22}$$

which is the product of the diagonal entries as required.

Second, assume the property for $(n-1) \times (n-1)$ matrices. Also the upper triangular matrix A has the form

$$A = \begin{bmatrix} A_{nn} & \mathbf{a}'_n \\ \mathbf{0}^T & a_{nn} \end{bmatrix}.$$

¹ From time-to-time, some people may call an upper triangular matrix either a right triangular or an upper-right triangular matrix. Correspondingly, from time-to-time, some people may call a lower triangular matrix either a left triangular or a lower-left triangular matrix.

for $(n - 1) \times (n - 1)$ minor A_{nn} . Theorem 6.2.10 establishes $\det A = a_{nn} \det A_{nn}$. Since the minor A_{nn} is upper triangular and $(n - 1) \times (n - 1)$, by assumption $\det A_{nn} = a_{11}a_{22} \cdots a_{n-1,n-1}$. Consequently, $\det A = a_{nn} \det A_{nn} = a_{nn}a_{11}a_{22} \cdots a_{n-1,n-1}$. Induction then establishes the theorem. \square

Example 6.2.15. Find the determinant of those of the following matrices which are triangular.

$$(a) \begin{bmatrix} -1 & -1 & -1 & -5 \\ 0 & -4 & 1 & 4 \\ 0 & 0 & 7 & 0 \\ 0 & 0 & 0 & -3 \end{bmatrix}$$

Solution: This is upper triangular, and its determinant is $(-1) \cdot (-4) \cdot 7 \cdot (-3) = -84$.

$$(b) \begin{bmatrix} -3 & 0 & 0 & 0 \\ -4 & 2 & 0 & 0 \\ -1 & 1 & 1 & 0 \\ -2 & -3 & 7 & -1 \end{bmatrix}$$

Solution: This is lower triangular, and its determinant is $(-3) \cdot 2 \cdot 1 \cdot (-1) = 6$.

$$(c) \begin{bmatrix} -4 & 0 & 0 & 0 \\ 2 & -2 & 0 & 0 \\ -5 & -3 & -2 & 0 \\ -2 & 5 & -2 & 0 \end{bmatrix}$$

Solution: This is lower triangular, and its determinant is zero as it has a column of zeros.

$$(d) \begin{bmatrix} 0.2 & 0 & 0 & 0 \\ 0 & 1.1 & 0 & 0 \\ 0 & 0 & -0.5 & 0 \\ 0 & 0 & 0 & 0.9 \end{bmatrix}$$

Solution: This diagonal matrix is both upper and lower triangular, and its determinant is $0.2 \cdot 1.1 \cdot (-0.5) \cdot 0.9 = -0.099$.

$$(e) \begin{bmatrix} 1 & -1 & 1 & -3 \\ 0 & 0 & 0 & -5 \\ 0 & 0 & -3 & -4 \\ 0 & -2 & 1 & -2 \end{bmatrix}$$

Solution: This is not triangular, so we do not have to compute its determinant. Nonetheless, if we swap the 2nd and 4th rows, then the result is the upper triangular

$$\begin{bmatrix} 1 & -1 & 1 & -3 \\ 0 & -2 & 1 & -2 \\ 0 & 0 & -3 & -4 \\ 0 & 0 & 0 & -5 \end{bmatrix}$$

and its determinant is $1 \cdot (-2) \cdot (-3) \cdot (-5) = -30$. But the row swap changes the sign so the determinant of the original matrix is $-(-30) = 30$.

$$(f) \begin{bmatrix} 0 & 0 & 0 & -3 \\ 0 & 0 & 2 & -4 \\ 0 & -1 & 4 & -1 \\ -6 & 1 & 5 & 1 \end{bmatrix}$$

Solution: This is *not* triangular, so we do not have to compute its determinant. Nonetheless, if we swap the 1st and 4th rows, and the 2nd and 3rd rows, then the result is

the lower triangular $\begin{bmatrix} -3 & 0 & 0 & 0 \\ -4 & 2 & 0 & 0 \\ -1 & 4 & -1 & 0 \\ 1 & 5 & 1 & -6 \end{bmatrix}$ and its determinant is

$(-3) \cdot 2 \cdot (-1) \cdot (-6) = -36$. But each row swap changes the sign so the determinant of the original matrix is $(-1)^2(-36) = 36$.

$$(g) \begin{bmatrix} -1 & 0 & 0 & 1 \\ -2 & 0 & 0 & 0 \\ 2 & -2 & -1 & -2 \\ -1 & 0 & 4 & 2 \end{bmatrix}$$

Solution: This is not triangular, so we do not have to compute its determinant. Nonetheless, if we swap the 2nd and 4th columns, the 1st and 2nd rows, and the 3rd and 4th rows,

then the result is the lower triangular $\begin{bmatrix} -2 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -1 & 2 & 4 & 0 \\ 2 & -2 & -1 & -2 \end{bmatrix}$ and

its determinant is $(-2) \cdot 1 \cdot 4 \cdot (-2) = 16$. But each row and column swap changes the sign so the determinant of the original matrix is $(-1)^3 16 = -16$.

■

The above case of triangular matrices is a short detour from the main development of this section which is to derive a formula for determinants in general. The following two examples introduce the next property we need before establishing a general formula for determinants.

Example 6.2.16. Let's rewrite the explicit formulas (6.1) for 2×2 and 3×3 determinants explicitly as the sum of simpler determinants.

- Recall that the 2×2 determinant

$$\begin{aligned} \begin{vmatrix} a & b \\ c & d \end{vmatrix} &= ad - bc \\ &= (ad - 0c) + (0d - bc) \\ &= \begin{vmatrix} a & 0 \\ c & d \end{vmatrix} + \begin{vmatrix} 0 & b \\ c & d \end{vmatrix}. \end{aligned}$$

That is, the original determinant is the same as the sum of two determinants, each with a zero in the first row and the other

row unchanged. Equivalently, the first row is decomposed as $[a \ b] = [a \ 0] + [0 \ b]$, while the other row is unchanged.

- Recall from (6.1) that the 3×3 determinant

$$\begin{aligned} \begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} &= aei + bfg + cdh - ceg - afh - bdi \\ &= +aei + 0fg + 0dh - 0eg - afh - 0di \\ &\quad + 0ei + bfg + 0dh - 0eg - 0fh - bdi \\ &\quad + 0ei + 0fg + cdh - ceg - 0fh - 0di \\ &= \begin{vmatrix} a & 0 & 0 \\ d & e & f \\ g & h & i \end{vmatrix} + \begin{vmatrix} 0 & b & 0 \\ d & e & f \\ g & h & i \end{vmatrix} + \begin{vmatrix} 0 & 0 & c \\ d & e & f \\ g & h & i \end{vmatrix}. \end{aligned}$$

That is, the original determinant is the same as the sum of three determinants, each with two zeros in the first row and the other rows unchanged. Equivalently, the first row is decomposed as $[a \ b \ c] = [a \ 0 \ 0] + [0 \ b \ 0] + [0 \ 0 \ c]$, while the other rows are unchanged.

This sort of rearrangement of a determinant makes progress because then Theorem 6.2.10 helps by finding the determinant of the resultant matrices that have almost all zero rows. ■

Example 6.2.17. A 2×2 example of a more general summation property is furnished by the determinant of matrix $A = \begin{bmatrix} a_{11} & b_1 + c_1 \\ a_{21} & b_2 + c_2 \end{bmatrix}$.

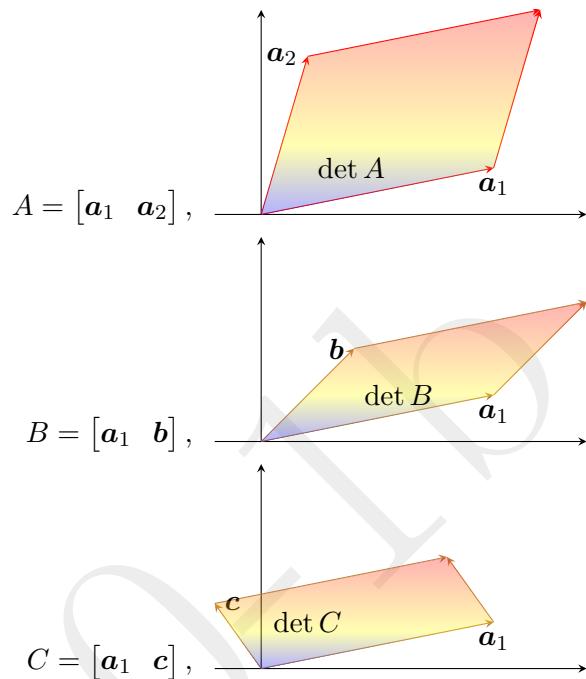
$$\begin{aligned} \det A &= a_{11}(b_2 + c_2) - a_{21}(b_1 + c_1) \\ &= a_{11}b_2 + a_{11}c_2 - a_{21}b_1 - a_{21}c_1 \\ &= (a_{11}b_2 - a_{21}b_1) + (a_{11}c_2 - a_{21}c_1) \\ &= \det \begin{bmatrix} a_{11} & b_1 \\ a_{21} & b_2 \end{bmatrix} + \det \begin{bmatrix} a_{11} & c_1 \\ a_{21} & c_2 \end{bmatrix} \\ &= \det B + \det C, \end{aligned}$$

where matrices B and C have the same first column as A , and their second columns add up to that of A . ■

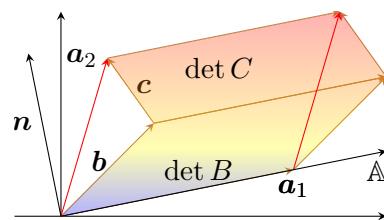
Theorem 6.2.18 (sum formula). *Let A , B and C be $n \times n$ matrices. If A , B and C are identical except for their i th column, and that the i th column of A is the sum of the i th columns of B and C , then $\det A = \det B + \det C$. Further, the same property holds when “column” is replaced by “row” throughout.*

Proof. We establish the theorem for matrix columns. Then the same results holds for the rows because $\det(A^T) = \det(A)$ (Theorem 6.1.15). As a prelude to the general geometric proof, consider

the 2×2 case and the second column (as also established algebraically by Example 6.2.17). Write the matrices in terms of their column vectors, and draw the determinant parallelogram areas as shown below: let



The matrices A , B and C all have the same first column \mathbf{a}_1 , whereas the second columns satisfy $\mathbf{a}_2 = \mathbf{b} + \mathbf{c}$ by the condition of the theorem. Because these parallelograms have common side \mathbf{a}_1 we can stack the area for $\det C$ on top of that for $\det B$, and because $\mathbf{a}_2 = \mathbf{b} + \mathbf{c}$ the top edge of the stack matches that for the area $\det A$, as shown below



The base of the stacked shape lies on the line \mathbb{A} , and let \mathbf{n} denote the orthogonal/normal direction (as shown). Because the shape has the same cross-section in lines parallel to \mathbb{A} , its area is the area of the base times the height of the stacked shape in the direction \mathbf{n} . But this is precisely the same height and base as the area for $\det A$, hence $\det A = \det B + \det C$.

A general proof for the last column uses the same diagrams, albeit schematically. Let matrices

$$A = [A' \ \mathbf{a}_n], \quad B = [A' \ \mathbf{b}], \quad C = [A' \ \mathbf{c}],$$

where the $n \times (n - 1)$ matrix $A' = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_{n-1}]$ is common to all three, and where the three last columns satisfy $\mathbf{a}_n = \mathbf{b} + \mathbf{c}$. Consider the n D-parallelepipeds whose n D-volumes are the three determinants, as before. Because these n D-parallelepipeds have common base of the $(n - 1)$ D-parallelepiped formed by the columns of A' , we can and do stack the n D-volume for $\det C$ on top of that for $\det B$, and because $\mathbf{a}_n = \mathbf{b} + \mathbf{c}$ the top $(n - 1)$ D-parallelepiped of the stack matches that for the n D-volume $\det A$, as shown schematically before. The base of the stacked shape lies on the subspace $\mathbb{A} = \text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_{n-1}\}$, and let \mathbf{n} denote the orthogonal/normal direction to \mathbb{A} (as shown schematically). Because the shape has the same cross-section parallel to \mathbb{A} , its n D-volume is the $(n - 1)$ D-volume of the base $(n - 1)$ D-parallelepiped times the height of the stacked shape in the direction \mathbf{n} . But this is precisely the same height and base as the n D-volume for $\det A$, hence $\det A = \det B + \det C$.

Lastly, when it is the j th column for which $\mathbf{a}_j = \mathbf{b} + \mathbf{c}$ and all others columns are identical, then swap column j with column n in all matrices. Theorem 6.2.5c asserts the signs of the three determinants are changed by this swapping. The above proof for the last column case then assures us $(-\det A) = (-\det B) + (-\det C)$; that is, $\det A = \det B + \det C$, as required. \square

The sum formula Theorem 6.2.18 leads to the common way to compute determinants by hand for matrices larger than 3×3 , albeit not generally practical for matrices significantly larger.

Example 6.2.19. Use Theorems 6.2.18 and 6.2.10 to evaluate the determinant of matrix

$$A = \begin{bmatrix} -2 & 1 & -1 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{bmatrix}.$$

Solution: Write the first row of A as the sum

$$\begin{aligned} [-2 \ 1 \ -1] &= [-2 \ 0 \ 0] + [0 \ 1 \ -1] \\ &= [-2 \ 0 \ 0] + [0 \ 1 \ 0] + [0 \ 0 \ -1]. \end{aligned}$$

Then using Theorem 6.2.18 twice, the determinant

$$\begin{aligned} &\left| \begin{array}{ccc} -2 & 1 & -1 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{array} \right| \\ &= \left| \begin{array}{ccc} -2 & 0 & 0 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{array} \right| + \left| \begin{array}{ccc} 0 & 1 & -1 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{array} \right| \\ &= \left| \begin{array}{ccc} -2 & 0 & 0 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{array} \right| + \left| \begin{array}{ccc} 0 & 1 & 0 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{array} \right| + \left| \begin{array}{ccc} 0 & 0 & -1 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{array} \right| \end{aligned}$$

Each of these last three matrices has the first row zero except for one element, so Theorem 6.2.10 applies to each of the three determinants to give

$$\begin{aligned} & \begin{vmatrix} -2 & 1 & -1 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{vmatrix} \\ &= (-1)^2(-2) \begin{vmatrix} -6 & -1 \\ 1 & 0 \end{vmatrix} + (-1)^3(1) \begin{vmatrix} 1 & -1 \\ 2 & 0 \end{vmatrix} + (-1)^4(-1) \begin{vmatrix} 1 & -6 \\ 2 & 1 \end{vmatrix} \\ &= (-2) \cdot 1 - (1) \cdot 2 + (-1) \cdot 13 = -17 \end{aligned}$$

upon using the well-known formula (6.1) for the three 2×2 determinants.

Alternatively, we could have used any row or column instead of the first row. For example, let's use the last column as it usefully already has a zero entry: write the last column of matrix A as $(-1, -1, 0) = (-1, 0, 0) + (0, -1, 0)$, then by Theorem 6.2.18 the determinant

$$\begin{aligned} \begin{vmatrix} -2 & 1 & -1 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{vmatrix} &= \begin{vmatrix} -2 & 1 & -1 \\ 1 & -6 & 0 \\ 2 & 1 & 0 \end{vmatrix} + \begin{vmatrix} -2 & 1 & 0 \\ 1 & -6 & -1 \\ 2 & 1 & 0 \end{vmatrix} \\ &\quad (\text{so by Theorem 6.2.10}) \\ &= (-1)^4(-1) \begin{vmatrix} 1 & -6 \\ 2 & 1 \end{vmatrix} + (-1)^5(-1) \begin{vmatrix} -2 & 1 \\ 2 & 1 \end{vmatrix} \\ &= (-1) \cdot 13 - (-1) \cdot (-4) = -17, \end{aligned}$$

as before. ■

Theorem 6.2.20 (Laplace expansion theorem). *Consider an $n \times n$ matrix $A = [a_{ij}]$ ($n \geq 2$). Recall the (i, j) th minor A_{ij} to be the $(n-1) \times (n-1)$ matrix obtained from A by omitting the i th row and j th column. Then the determinant of A can be computed via expansion in any row i or any column j as, respectively,*

$$\begin{aligned} \det A &= (-1)^{i+1} a_{i1} \det A_{i1} + (-1)^{i+2} a_{i2} \det A_{i2} \\ &\quad + \cdots + (-1)^{i+n} a_{in} \det A_{in} \\ &= (-1)^{j+1} a_{1j} \det A_{1j} + (-1)^{j+2} a_{2j} \det A_{2j} \\ &\quad + \cdots + (-1)^{j+n} a_{nj} \det A_{nj}. \end{aligned} \tag{6.4}$$

Proof. We establish the expansion for matrix rows: then the same property holds for the columns because $\det(A^T) = \det(A)$ (Theorem 6.1.15). First prove the expansion for a first row expansion,

and then second for any row. So first use the sum Theorem 6.2.18 ($n - 1$) times to deduce

$$\begin{aligned}
 & \left| \begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right| \\
 &= \left| \begin{array}{cccc} a_{11} & 0 & \cdots & 0 \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right| + \left| \begin{array}{cccc} 0 & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right| \\
 &\quad \vdots \\
 &= \left| \begin{array}{cccc} a_{11} & 0 & \cdots & 0 \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right| + \left| \begin{array}{cccc} 0 & a_{12} & \cdots & 0 \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right| \\
 &\quad + \cdots + \left| \begin{array}{cccc} 0 & 0 & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right|
 \end{aligned}$$

As each of these n determinants has the first row zero except for one element, Theorem 6.2.10 applies to give

$$\begin{aligned}
 & \left| \begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{array} \right| \\
 &= (-1)^2 a_{11} \left| \begin{array}{ccc} a_{22} & \cdots & a_{2n} \\ \vdots & \ddots & \vdots \\ a_{n2} & \cdots & a_{nn} \end{array} \right| + (-1)^3 a_{12} \left| \begin{array}{ccc} a_{21} & \cdots & a_{2n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{array} \right| \\
 &\quad + \cdots + (-1)^{n+1} a_{1n} \left| \begin{array}{ccc} a_{21} & a_{22} & \cdots & a_{2,n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{n,n-1} \end{array} \right| \\
 &= (-1)^2 a_{11} \det A_{11} + (-1)^3 a_{12} \det A_{12} \\
 &\quad + \cdots + (-1)^{n+1} a_{1n} \det A_{1n},
 \end{aligned}$$

which is the case $i = 1$ of formula (6.4).

Second, for the general i th row expansion, let a new matrix B be obtained from A by swapping the i th row up $(i - 1)$ times to form the first row of B and leaving the other rows from A in the same order. Then the elements $b_{ij} = a_{ij}$, and also the minors $B_{1j} = A_{ij}$. Apply formula (6.4) to the first row of B (just proved) to give

$$\det B = (-1)^2 b_{11} \det B_{11} + (-1)^3 b_{12} \det B_{12}$$

$$\begin{aligned}
& + \cdots + (-1)^{n+1} b_{1n} \det B_{1n} \\
& = (-1)^2 a_{i1} \det A_{i1} + (-1)^3 a_{i2} \det A_{i2} \\
& \quad + \cdots + (-1)^{n+1} a_{in} \det A_{in}.
\end{aligned}$$

But by Theorem 6.2.5c each of the $(i - 1)$ row swaps in forming B changes the sign of the determinant: hence

$$\begin{aligned}
\det A &= (-1)^{i-1} \det B \\
&= (-1)^{i-1+2} a_{i1} \det A_{i1} + (-1)^{i-1+3} a_{i2} \det A_{i2} \\
&\quad + \cdots + (-1)^{i-1+n+1} a_{in} \det A_{in} \\
&= (-1)^{i+1} a_{i1} \det A_{i1} + (-1)^{i+2} a_{i2} \det A_{i2} \\
&\quad + \cdots + (-1)^{i+n} a_{in} \det A_{in},
\end{aligned}$$

as required. \square

Example 6.2.21. Use the Laplace expansion (6.4) to find the determinant of the following matrices.

$$(a) \begin{bmatrix} 0 & 2 & 1 & 2 \\ -1 & 2 & -1 & -2 \\ 1 & 2 & -1 & -1 \\ 0 & -1 & -1 & 1 \end{bmatrix}$$

Solution: The first column has two zeros, so expand in the first column:

$$\begin{aligned}
\det &= (-1)^3(-1) \det \begin{bmatrix} 2 & 1 & 2 \\ 2 & -1 & -1 \\ -1 & -1 & 1 \end{bmatrix} \\
&\quad + (-1)^4(1) \det \begin{bmatrix} 2 & 1 & 2 \\ 2 & -1 & -2 \\ -1 & -1 & 1 \end{bmatrix} \\
&\quad \text{(using (6.1) for these } 3 \times 3 \text{ matrices)} \\
&= (-2 + 1 - 4 - 2 - 2 - 2) \\
&\quad + (-2 + 2 - 4 - 2 - 4 - 2) \\
&= -23.
\end{aligned}$$

$$(b) \begin{bmatrix} -3 & -1 & 1 & 0 \\ -2 & 0 & -2 & 0 \\ -3 & -2 & 0 & 0 \\ 1 & -2 & 0 & 3 \end{bmatrix}$$

Solution: The last column has three zeros, so expand in the last column:

$$\begin{aligned}
\det &= (-1)^8(3) \det \begin{bmatrix} -3 & -1 & 1 \\ -2 & 0 & -2 \\ -3 & -2 & 0 \end{bmatrix} \\
&\quad \text{(expand in the middle row (say) due to its zero)}
\end{aligned}$$

$$\begin{aligned}
&= 3 \left\{ (-1)^3(-2) \det \begin{bmatrix} -1 & 1 \\ -2 & 0 \end{bmatrix} \right. \\
&\quad \left. + (-1)^5(-2) \det \begin{bmatrix} -3 & -1 \\ -3 & -2 \end{bmatrix} \right\} \\
&= 3 \{2(0+2) + 2(6-3)\} \\
&= 30.
\end{aligned}$$

■

The Laplace expansion is generally too computationally expensive for all but small matrices. The reason is that computing the determinant of an $n \times n$ matrix with the Laplace expansion generally takes $n!$ operations (the next Theorem 6.2.23), and the factorial $n! = n(n-1)\cdots 3 \cdot 2 \cdot 1$ grows very quickly even for medium n . Even for just a 20×20 matrix the Laplace expansion has over two quintillion terms ($2 \cdot 10^{18}$). Exceptional matrices are those with lots of zeros, such as triangular matrices (Theorem 6.2.14). In any case, remember that except for theoretical purposes there is rarely need to compute a medium to large determinant.

Example 6.2.22. The determinant of a 3×3 matrix has $3! = 6$ terms, each a product of three factors: diagram (6.1) gives the determinant

$$\begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = aei + bfg + cdh - ceg - afh - bdi.$$

Further, observe that within each term the factors come from different rows and columns. For example, a never appears in a term with the entries b, c, d or g (the elements from either the same row or the same column). Similarly, f never appears in a term with the entries d, e, c or i . ■

Theorem 6.2.23. *The determinant of any $n \times n$ matrix expands to the sum (\pm) of $n!$ terms, where each term is a product of n factors such that each factor comes from different rows and columns of the matrix.*

Proof. Use induction on the size of the matrix. First, the properties holds for 1×1 matrices as $\det [a_{11}] = a_{11}$ is one term of one factor from the only row and column of the matrix.

Second, assume the determinant of any $(n-1) \times (n-1)$ matrix may be written the sum (\pm) of $(n-1)!$ terms, where each term is a product of $(n-1)$ factors such that each factor comes from different rows and columns. Consider any $n \times n$ matrix A . By the Laplace Expansion Theorem 6.2.20, $\det A$ may be written as the sum (\pm)

of n terms of the form $a_{ij} \det A_{ij}$. By induction assumption, the $(n-1) \times (n-1)$ minors A_{ij} have determinants with $(n-1)!$ terms, each of $(n-1)$ factors and so the n terms in a Laplace Expansion of $\det A$ expands to $n(n-1)! = n!$ terms, each term being of n factors through the multiplication by the entry a_{ij} . Further, recall the minor A_{ij} is obtained from A by omitting row i and column j , and so the minor has no elements from the same row or column as a_{ij} . Consequently, each term in the determinant only has factors from different rows and columns, as required. By induction the theorem holds for all n . \square

6.2.1 Exercises

Exercise 6.2.1. In each of the following, the determinant of a matrix is given. Use Theorem 6.2.5 on the row and column properties of a determinant to find the determinant of the other four listed matrices. Give reasons for your answers.

$$(a) \det \begin{bmatrix} -2 & 1 & -4 \\ -2 & -1 & 2 \\ -2 & 5 & -1 \end{bmatrix} = 60$$

$$\text{i. } \begin{bmatrix} 1 & 1 & -4 \\ 1 & -1 & 2 \\ 1 & 5 & -1 \end{bmatrix}$$

$$\text{ii. } \begin{bmatrix} -2 & 1 & -4 \\ -2 & 1 & -4 \\ -2 & 5 & -1 \end{bmatrix}$$

$$\text{iii. } \begin{bmatrix} -2 & 1 & -4 \\ -0.2 & -0.1 & 0.2 \\ -2 & 5 & -1 \end{bmatrix}$$

$$\text{iv. } \begin{bmatrix} -2 & 1 & -4 \\ -2 & 5 & -1 \\ -2 & -1 & 2 \end{bmatrix}$$

$$(b) \det \begin{bmatrix} -1 & -1 & 4 & -6 \\ 4 & -2 & -2 & -1 \\ 0 & -3 & -1 & -4 \\ 3 & 2 & 1 & 1 \end{bmatrix} = 72$$

$$\text{i. } \begin{bmatrix} 0 & -3 & -1 & -4 \\ 4 & -2 & -2 & -1 \\ -1 & -1 & 4 & -6 \\ 3 & 2 & 1 & 1 \end{bmatrix}$$

$$\text{ii. } \begin{bmatrix} -1 & -1 & 4 & -6 \\ 4 & -2 & -2 & -1 \\ 0 & 0 & 0 & 0 \\ 3 & 2 & 1 & 1 \end{bmatrix}$$

$$\text{iii. } \begin{bmatrix} -1 & -1 & 4 & -6 \\ 2 & -1 & -1 & -1/2 \\ 0 & -3 & -1 & -4 \\ 3 & 2 & 1 & 1 \end{bmatrix}$$

$$\text{iv. } \begin{bmatrix} -1 & -1 & 4 & -6 \\ -2 & 4 & -2 & -1 \\ -3 & 0 & -1 & -4 \\ 2 & 3 & 1 & 1 \end{bmatrix}$$

$$(c) \det \begin{bmatrix} 2 & -3 & 2 & -3 \\ 0 & -1 & -1 & -2 \\ 2 & 1 & -2 & -3 \\ -4 & -1 & -4 & 0 \end{bmatrix} = 16$$

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

ii. $\begin{bmatrix} 1 & 12 & 2 & -3 \\ 0 & 4 & -1 & -2 \\ 1 & -4 & -2 & -3 \\ -2 & 4 & -4 & 0 \end{bmatrix}$

iii. $\begin{bmatrix} 4 & 4 & -6 & -6 \\ 0 & -1 & -1 & -2 \\ 2 & -2 & 1 & -3 \\ -4 & -4 & -1 & 0 \end{bmatrix}$

iv. $\begin{bmatrix} 0 & -1 & -0.5 & -2 \\ 2 & -3 & 1 & -3 \\ 2 & 1 & -1 & -3 \\ -4 & -1 & -2 & 0 \end{bmatrix}$

(d) $\det \begin{bmatrix} 0.3 & -0.1 & -0.1 & 0.4 \\ 0.2 & 0.3 & 0 & 0.1 \\ 0.1 & -0.1 & -0.3 & -0.2 \\ -0.1 & -0.2 & 0.4 & 0.2 \end{bmatrix} = 0.01$

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

ii. $\begin{bmatrix} 0.3 & 0.2 & 0 & 2 \\ 0.2 & -0.6 & 0 & 0.5 \\ 0.1 & 0.2 & 0 & -1 \\ -0.1 & 0.4 & 0 & 1 \end{bmatrix}$

iii. $\begin{bmatrix} 0.2 & 0.3 & 0 & 0.1 \\ 0.1 & -0.1 & -0.3 & -0.2 \\ -0.1 & -0.2 & 0.4 & 0.2 \\ 0.3 & -0.1 & -0.1 & 0.4 \end{bmatrix}$

iv. $\begin{bmatrix} 0.3 & 0.4 & -0.1 & 0.4 \\ 0.2 & 0.1 & 0 & 0.1 \\ 0.1 & -0.2 & -0.3 & -0.2 \\ -0.1 & 0.2 & 0.4 & 0.2 \end{bmatrix}$

Exercise 6.2.2. Recall Example 6.2.7. For each pair of given points, (x_1, y_1) and (x_2, y_2) , evaluate the determinant in the equation

$$\det \begin{bmatrix} 1 & x & y \\ 1 & x_1 & y_1 \\ 1 & x_2 & y_2 \end{bmatrix} = 0$$

to find an equation for the straight line through the two given points. Show your working.

(a) $(-3, -6), (2, 3)$

(b) $(3, -2), (-3, 0)$

(c) $(1, -4), (-3, 1)$

(d) $(-1, 0), (-2, 1)$

(e) $(6, 1), (2, -1)$

(f) $(3, -8), (7, -2)$

Exercise 6.2.3. Using mainly the properties of Theorem 6.2.5 detail an argument that the following determinant equations each give an equation for the line through two given points (x_1, y_1) and (x_2, y_2) .

(a) $\det \begin{bmatrix} 1 & x_1 & y_1 \\ 1 & x & y \\ 1 & x_2 & y_2 \end{bmatrix} = 0$

(b) $\det \begin{bmatrix} 1 & 1 & 1 \\ x_2 & x_1 & x \\ y_2 & y_1 & y \end{bmatrix} = 0$

Exercise 6.2.4. Recall Example 6.2.8. For each pair of given points, (x_1, y_1, z_1) and (x_2, y_2, z_2) , evaluate the determinant in the equation

$$\det \begin{bmatrix} x & y & z \\ x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \end{bmatrix} = 0$$

to find an equation for the plane that passes through the two given points and the origin. Show your working.

- (a) $(-1, -1, -3), (3, -5, -1)$ (b) $(0, 0, -2), (4, -4, 0)$

- (c) $(-1, 2, 2), (-1, -3, 2)$ (d) $(4, -2, 0), (-3, -4, -1)$

- (e) $(-4, -1, 2), (-3, -2, 2)$ (f) $(2, 2, 3), (2, 1, 4)$

Exercise 6.2.5. Using mainly the properties of Theorem 6.2.5 detail an argument that the following determinant equations each give an equation for the plane passing through the origin and the two given points (x_1, y_1, z_1) and (x_2, y_2, z_2) .

$$(a) \det \begin{bmatrix} x_1 & y_1 & z_1 \\ x_2 & y_2 & z_2 \\ x & y & z \end{bmatrix} = 0 \quad (b) \det \begin{bmatrix} x_2 & x & x_1 \\ y_2 & y & y_1 \\ z_2 & z & z_1 \end{bmatrix} = 0$$

Exercise 6.2.6. Prove Theorems 6.2.5a, 6.2.5d and 6.2.5b using basic geometric arguments about the transformation of the unit n D-cube.

Exercise 6.2.7. Use Theorem 6.2.5 to prove that if a square matrix A has two non-zero rows proportional to each other, then $\det A = 0$. Why does it immediately follow that (instead of rows) if the matrix has two non-zero columns proportional to each other, then $\det A = 0$.

Exercise 6.2.8. Use Theorem 6.2.10, and then (6.1), to evaluate the following determinants. Show your working.

$$(a) \det \begin{bmatrix} 6 & 1 & 1 \\ -1 & 3 & -8 \\ -6 & 0 & 0 \end{bmatrix} \quad (b) \det \begin{bmatrix} 4 & 8 & 0 \\ 3 & -2 & 0 \\ -1 & -1 & -3 \end{bmatrix}$$

$$(c) \det \begin{bmatrix} 0 & 0 & 3 \\ -1 & -3 & -3 \\ -3 & -5 & 2 \end{bmatrix} \quad (d) \det \begin{bmatrix} -4 & 0 & -5 \\ 1 & -7 & -1 \\ 4 & 0 & 4 \end{bmatrix}$$

$$(e) \det \begin{bmatrix} 2 & -4 & -2 & -2 \\ 0 & 1 & -3 & -2 \\ -2 & 0 & 0 & 0 \\ 5 & -8 & 1 & 7 \end{bmatrix} \quad (f) \det \begin{bmatrix} 0 & -5 & 0 & 0 \\ -7 & 2 & 2 & 1 \\ 1 & -2 & -2 & -5 \\ 6 & 8 & -2 & 0 \end{bmatrix}$$

$$(g) \det \begin{bmatrix} 0 & 2 & -4 & 3 \\ -6 & 6 & -2 & 0 \\ -1 & -8 & 4 & 0 \\ 2 & -2 & -1 & 0 \end{bmatrix} \quad (h) \det \begin{bmatrix} -3 & 4 & -1 & -1 \\ 4 & 8 & 1 & 6 \\ 0 & 7 & 0 & 0 \\ 2 & -6 & 1 & 2 \end{bmatrix}$$

Exercise 6.2.9. Use the triangular matrix Theorem 6.2.14, as well as the row/column properties of Theorem 6.2.5, to find the determinants of each of the following matrices. Show your argument.

$$\begin{array}{ll} (a) \begin{bmatrix} -6 & -4 & -7 & 2 \\ 0 & -2 & -1 & 1 \\ 0 & 0 & -4 & 1 \\ 0 & 0 & 0 & -2 \end{bmatrix} & (b) \begin{bmatrix} 2 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 \\ -5 & -1 & 2 & 0 \\ 2 & 4 & -1 & 1 \end{bmatrix} \\ (c) \begin{bmatrix} 0 & 0 & -6 & -6 \\ -2 & -2 & -2 & 1 \\ 0 & 0 & 0 & -4 \\ 0 & 0 & -1 & 7 \end{bmatrix} & (d) \begin{bmatrix} 0 & 0 & -2 & 6 \\ 0 & 0 & 0 & 1 \\ -7 & -6 & 8 & 2 \\ 0 & -2 & -4 & 1 \end{bmatrix} \\ (e) \begin{bmatrix} 0 & 0 & 8 & 0 \\ -5 & -6 & 6 & -1 \\ 0 & 0 & -5 & 6 \\ 0 & -6 & -4 & 3 \end{bmatrix} & (f) \begin{bmatrix} 0 & 0 & 7 & 0 \\ 0 & -3 & 7 & -4 \\ 0 & 7 & -4 & 0 \\ 1 & 2 & -1 & -3 \end{bmatrix} \\ (g) \begin{bmatrix} 0 & 0 & 1 & 8 & 5 \\ 6 & 1 & -5 & -8 & -1 \\ 0 & 0 & 0 & 0 & 5 \\ 0 & -1 & -6 & -5 & 4 \\ 0 & 0 & 0 & -1 & -8 \end{bmatrix} & (h) \begin{bmatrix} -6 & 0 & 0 & 0 & 0 \\ -4 & -3 & 0 & 0 & 0 \\ 0 & -4 & 0 & 4 & 0 \\ 0 & 1 & -5 & 12 & 0 \\ -2 & -1 & -5 & -2 & 5 \end{bmatrix} \end{array}$$

Exercise 6.2.10. Given that the determinant

$$\begin{vmatrix} a & b & c \\ d & e & f \\ g & h & i \end{vmatrix} = 6,$$

find the following determinants. Give reasons.

$$\begin{array}{ll} (a) \begin{vmatrix} 3a & b & c \\ 3d & e & f \\ 3g & h & i \end{vmatrix} & (b) \begin{vmatrix} a & b & c/2 \\ -d & -e & -f/2 \\ g & h & i/2 \end{vmatrix} \\ (c) \begin{vmatrix} d & e & f \\ a & b & c \\ g & h & i \end{vmatrix} & (d) \begin{vmatrix} a+d & b+e & c+f \\ d & e & f \\ g & h & i \end{vmatrix} \\ (e) \begin{vmatrix} a & b & c-a \\ d & e & f-d \\ g & h & i-g \end{vmatrix} & (f) \begin{vmatrix} a & b & c \\ d & e & f \\ a+2g & b+2h & c+2i \end{vmatrix} \end{array}$$

$$(g) \begin{vmatrix} d & e & f \\ g+a & h+b & i+c \\ a & b & c \end{vmatrix} \quad (h) \begin{vmatrix} a-3g & b & c \\ d/3-f & e/3 & f/3 \\ g-3i & h & i \end{vmatrix}$$

Exercise 6.2.11. Consider a general 3×3 matrix $A = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix}$. Derive a first column Laplace expansion (Theorem 6.2.20) of the 3×3 determinant, and rearrange to show it is the same as the determinant formula (6.1).

Exercise 6.2.12. Use the Laplace expansion (Theorem 6.2.20) to find the determinant of the following matrices. Use rows or columns with many zeros. Show your working.

$$(a) \begin{bmatrix} 1 & 4 & 0 & 0 \\ 0 & 2 & -1 & 0 \\ -3 & 0 & -3 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \quad (b) \begin{bmatrix} -4 & -3 & 0 & -3 \\ -2 & 0 & -1 & 0 \\ -1 & 2 & 4 & 2 \\ 0 & 0 & 0 & -4 \end{bmatrix}$$

$$(c) \begin{bmatrix} -4 & 0 & -2 & 2 \\ 0 & 1 & 5 & -2 \\ 4 & 0 & 0 & 4 \\ -2 & 4 & 4 & 5 \end{bmatrix} \quad (d) \begin{bmatrix} 3 & -1 & 5 & 0 \\ 1 & -2 & 0 & 2 \\ -1 & 3 & -6 & 2 \\ 0 & -2 & 0 & -3 \end{bmatrix}$$

$$(e) \begin{bmatrix} 0 & -1 & 0 & 0 & 6 \\ -4 & 0 & 1 & -1 & 0 \\ 0 & 0 & 3 & 0 & -4 \\ 5 & 0 & -4 & 0 & 0 \\ 3 & -1 & 0 & -3 & -6 \end{bmatrix} \quad (f) \begin{bmatrix} 0 & 3 & -7 & 2 & 0 \\ 0 & -4 & 0 & -6 & 0 \\ 0 & 1 & 3 & 0 & 0 \\ -3 & 0 & -4 & -3 & -3 \\ 0 & -6 & 0 & 0 & 5 \end{bmatrix}$$

$$(g) \begin{bmatrix} 0 & 0 & -2 & 0 & -6 \\ 0 & -3 & 0 & 1 & 0 \\ -4 & -5 & 2 & 0 & 2 \\ 0 & 2 & 0 & 3 & -3 \\ 2 & 0 & 0 & -1 & 0 \end{bmatrix} \quad (h) \begin{bmatrix} 0 & 0 & -2 & 7 & 4 \\ 0 & 4 & -4 & 1 & 0 \\ 2 & -2 & 0 & 1 & 0 \\ 3 & 2 & 0 & 2 & 0 \\ 0 & -2 & 0 & 0 & 1 \end{bmatrix}$$

Exercise 6.2.13. For each of the following matrices, use the Laplace expansion (Theorem 6.2.20) to find all the values of k for which the matrix is *not* invertible. Show your working.

$$(a) \begin{bmatrix} 3 & -2k & -1 & -1-2k \\ 0 & 2 & 1 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & -2 & -5 & 3+2k \end{bmatrix}$$

$$(b) \begin{bmatrix} -1 & -2 & 0 & 2k \\ 0 & 0 & 5 & 0 \\ 0 & 2 & -1+k & -k \\ -2+k & 1+2k & 4 & 0 \end{bmatrix}$$

$$(c) \begin{bmatrix} -1+k & 0 & -2+3k & -1+k \\ -6 & 1+2k & -3 & k \\ 0 & 3k & -4 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$(d) \begin{bmatrix} 0 & 3 & 2 & 2k \\ -3k & 3 & 0 & 0 \\ 2+k & 4 & 2 & -1+k \\ 0 & -2 & 3 & 4k \end{bmatrix}$$

$$(e) \begin{bmatrix} 0 & 0 & -4 & 0 & 2+k \\ 0 & 0 & 5 & 0 & -1 \\ 0 & 0 & 0 & 1-2k & -3-2k \\ 1+2k & -1 & 3 & 1-4k & -4 \\ -1+2k & 0 & 0 & -1 & -2 \end{bmatrix}$$

$$(f) \begin{bmatrix} k & 1 & 0 & -5-k & -1-k \\ -4+6k & -1 & 1+3k & -3-5k & 0 \\ 0 & 2 & 0 & -2 & -5 \\ 0 & 0 & 0 & 0 & 2-k \\ -2k & 1+k & 0 & -2 & 3 \end{bmatrix}$$

Exercise 6.2.14. Using Theorem 6.2.23 and the properties of Theorem 6.2.5, detail an argument that the following determinant equation generally forms an equation for the plane passing through the three given points (x_1, y_1, z_1) , (x_2, y_2, z_2) and (x_3, y_3, z_3) :

$$\det \begin{bmatrix} 1 & x & y & z \\ 1 & x_1 & y_1 & z_1 \\ 1 & x_2 & y_2 & z_2 \\ 1 & x_3 & y_3 & z_3 \end{bmatrix} = 0.$$

Exercise 6.2.15. Using Theorem 6.2.23 and the properties of Theorem 6.2.5, detail an argument that the following determinant equation generally forms an equation for the parabola passing through the three given points (x_1, y_1) , (x_2, y_2) and (x_3, y_3) :

$$\det \begin{bmatrix} 1 & x & x^2 & y \\ 1 & x_1 & x_1^2 & y_1 \\ 1 & x_2 & x_2^2 & y_2 \\ 1 & x_3 & x_3^2 & y_3 \end{bmatrix} = 0.$$

Exercise 6.2.16. Using Theorem 6.2.23 and the properties of Theorem 6.2.5, detail an argument that the equation

$$\det \begin{bmatrix} 1 & x & x^2 & y & y^2 & xy \\ 1 & x_1 & x_1^2 & y_1 & y_1^2 & x_1 y_1 \\ 1 & x_2 & x_2^2 & y_2 & y_2^2 & x_2 y_2 \\ 1 & x_3 & x_3^2 & y_3 & y_3^2 & x_3 y_3 \\ 1 & x_4 & x_4^2 & y_4 & y_4^2 & x_4 y_4 \\ 1 & x_5 & x_5^2 & y_5 & y_5^2 & x_5 y_5 \end{bmatrix} = 0$$

generally forms an equation for the conic section passing through the five given points (x_i, y_i) , $i = 1, \dots, 5$.

Answers to selected exercises6.1.1b : $\det \approx 1.26$ 6.1.1d : $\det \approx -1.47$ 6.1.1f : $\det \approx -0.9$ 6.1.1h : $\det \approx -5.4$ 6.1.1j : $\det \approx -1.32$ 6.1.1l : $\det \approx -1.11$ 6.1.3b : $\det \approx 1.6$ 6.1.3d : $\det \approx 1.0$ 6.1.3f : $\det \approx 2.3$

6.1.4b : 4

6.1.4d : $-\frac{1}{2}, -1, -2$ 6.1.4f : $0, -\frac{1}{6}, -4$

6.1.5b : 3

6.1.5d : $-35/54$ 6.1.5f : $-131/4$ 6.1.5h : -8704

6.1.9b : 25

6.1.9d : Unknowable on the given information.

6.1.9f : $3 \cdot 2^n$

6.1.9h : 25

6.2.1b : $-72, 0, 36, -72$ 6.2.1d : $100, -0.1, -0.01, 0$ 6.2.2b : $-2(x + 3y + 3) = 0$ 6.2.2d : $-(x + y + 1) = 0$ 6.2.2f : $2(-3x + 2y + 25) = 0$ 6.2.4b : $8(x + y) = 0$ 6.2.4d : $2(x + 2y - 11z) = 0$ 6.2.4f : $5x - 2y - 2z = 0$

6.2.8b : 96

6.2.8d : -28

6.2.8f : 100

- 6.2.8h : -42
- 6.2.9b : 12
- 6.2.9d : -28
- 6.2.9f : 196
- 6.2.9h : 1800
- 6.2.10b : -3
- 6.2.10d : 6
- 6.2.10f : 12
- 6.2.10h : 2
- 6.2.12b : 140
- 6.2.12d : -137
- 6.2.12f : 1080
- 6.2.12h : 236
- 6.2.13b : 0 , 3/4
- 6.2.13d : 0 , -1
- 6.2.13f : -1/3 , 0 , -9 , 2

7 Eigenvalues and eigenvectors in general

Chapter Contents

7.1	Find eigenvalues and eigenvectors of matrices	559
7.1.1	A characteristic equation gives eigenvalues	559
7.1.2	Repeated eigenvalues are sensitive	574
7.1.3	Application: discrete dynamics of populations	577
7.1.4	Extension: Connect SVDs to eigen-problems	592
7.1.5	Exercises	595
7.2	Linear independent vectors may form a basis	605
7.2.1	Linearly (in)dependent sets	606
7.2.2	Form a basis for subspaces	616
7.2.3	Exercises	631
7.3	Diagonalisation identifies the transformation	637
7.3.1	Solve systems of differential equations	648
7.3.2	Exercises	659

Population modelling Suppose two species of animals interact: how do their populations evolve in time? Let $y(t)$ and $z(t)$ be the number of female animals in each of the species at time t in years (biologists usually just count females in population models as females usually determine reproduction). Modelling might deduce the populations interact according to the rule that the population one year later is $y(t+1) = 2y(t) - 4z(t)$ and $z(t+1) = -y(t) + 2z(t)$: that is, if it was not for the other species, then for each species the number of females would both double every year (since then $y(t+1) = 2y(t)$ and $z(t+1) = 2z(t)$); but the other species decreases each of these growths via the $-4z(t)$ and $-y(t)$ terms.

Question: can we find special solutions in the form $(y, z) = \mathbf{x}\lambda^t$ for some constant λ ? Substitute $y = x_1\lambda^t$ and $z = x_2\lambda^t$ into the equations:

$$\begin{aligned} & y(t+1) = 2y(t) - 4z(t), \quad z(t+1) = -y(t) + 2z(t) \\ \iff & x_1\lambda^{t+1} = 2x_1\lambda^t - 4x_2\lambda^t, \quad x_2\lambda^{t+1} = -x_1\lambda^t + 2x_2\lambda^t \\ \iff & 2x_1 - 4x_2 = \lambda x_1, \quad -x_1 + 2x_2 = \lambda x_2 \end{aligned}$$

after dividing by the factor λ^t (assuming constant λ is non-zero). Forming as a matrix-vector equation these last two equations are

$$\begin{bmatrix} 2 & -4 \\ -1 & 2 \end{bmatrix} \mathbf{x} = \lambda \mathbf{x}.$$

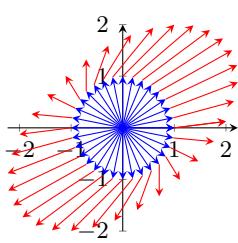
That is, this substitution $(y, z) = \mathbf{x}\lambda^t$ shows the question about solutions of the population equations reduces to solving $A\mathbf{x} = \lambda\mathbf{x}$, called an eigen-problem.

The linear algebra of the eigen-problem will empower us to predict that the general solution for the population is, in terms of two constants c_1 and c_2 , that one species has female population $y(t) = 2c_1 4^t + 2c_2$ whereas the second species has female population $z(t) = -c_1 4^t + c_2$.

The basic eigen-problem Recall from Section 4.1 that the eigen-problem equation $A\mathbf{x} = \lambda\mathbf{x}$ is just asking can we find directions \mathbf{x} such that matrix A acting on \mathbf{x} , resulting in $A\mathbf{x}$, is in the same direction as \mathbf{x} , namely $\lambda\mathbf{x}$ for some proportionality constant λ . Now $\mathbf{x} = \mathbf{0}$ is always a solution of the equation $A\mathbf{x} = \lambda\mathbf{x}$. Consequently, we are only interested in those values of the eigenvalue λ when non-zero solutions for the eigenvector \mathbf{x} exist (as it is the directions which are of interest). Rearranging the equation $A\mathbf{x} = \lambda\mathbf{x}$ as the homogeneous system $(A - \lambda I)\mathbf{x} = \mathbf{0}$, let's invoke properties of linear equations to solve the eigen-problem.

- Procedure 4.1.18 establishes that one way to find the eigenvalues λ is to solve the characteristic equation $\det(A - \lambda I) = 0$.
- Then for each eigenvalue, solving the homogeneous system $(A - \lambda I)\mathbf{x} = \mathbf{0}$ gives corresponding eigenvectors \mathbf{x} .
- The set of eigenvectors for a given eigenvalue forms a subspace called the eigenspace \mathbb{E}_λ (Theorem 4.1.7).

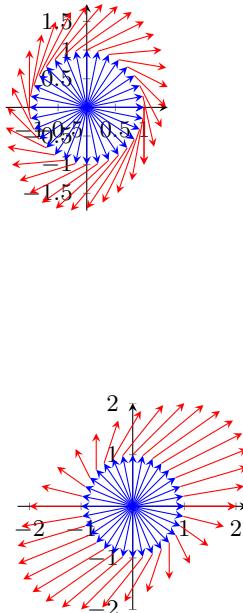
Three general difficulties in eigen-problems Recall that Section 4.1 introduced one way to visually estimate eigenvectors and eigenvalues of a given matrix A (Schonefeld 1995). The method is to plot many unit vectors \mathbf{x} , and at the end of each \mathbf{x} to adjoin the vector $A\mathbf{x}$. Since eigenvectors satisfy $A\mathbf{x} = \lambda\mathbf{x}$ for some scalar eigenvalue λ , we visually identify eigenvectors as those \mathbf{x} which point in the same (or opposite) direction to $A\mathbf{x}$. Let's use this approach to identify three general difficulties.



1. In this first picture, for matrix $A = \begin{bmatrix} 1 & 1 \\ \frac{1}{8} & 1 \end{bmatrix}$, the eigenvectors appear to be in directions $\mathbf{x}_1 \approx \pm(0.9, 0.3)$ and $\mathbf{x}_2 \approx \pm(0.9, -0.3)$ corresponding to eigenvalues $\lambda_1 \approx 1.4$ and $\lambda_2 \approx 0.6$. (Recall that scalar multiples of an eigenvector are always

also eigenvectors, §4.1, so we always see \pm pairs of eigenvectors in these pictures.) These pairs of eigenvectors are not orthogonal, not at right angles—as happens for symmetric matrices (Theorem 4.2.10). The lack of orthogonality in general means we soon generalise the concept of orthogonal sets of vectors to the new concept of linearly independent sets (Section 7.2).

2. In this second case, for $A = \begin{bmatrix} 0 & 1 \\ -1 & \frac{1}{2} \end{bmatrix}$, there appear to be no eigenvectors at all. No eigenvectors and eigenvalues is the answer if we require real answers. However, in most applications we find it sensible to have complex valued eigenvalues and eigenvectors (Section 7.1). So although we cannot see them graphically, for this matrix there are two complex eigenvalues and two families of complex eigenvectors (analogous to those found in Example 4.1.22). ¹



3. In this third case, for $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, there appears to be only the vectors $\mathbf{x} = \pm(1, 0)$, aligned along the horizontal axis, for which $A\mathbf{x} = \lambda\mathbf{x}$. Whereas for symmetric matrices there were always two pairs, here we only appear to have one pair of eigenvectors (Theorem 7.3.12). Such degeneracy occurs for matrices on the border between reality and complexity.

The first problem of the general lack of orthogonality of the eigenvectors is most clearly seen in the case of triangular matrices (Definition 6.2.13). The reason is linked to Theorem 6.2.14 that the determinant of a triangular matrix is simply the product of its diagonal entries.

Example 7.0.1. Find algebraically the eigenvalues and eigenvectors of the triangular matrix $A = \begin{bmatrix} 2 & 1 \\ 0 & 3 \end{bmatrix}$.

Solution: Recall Procedure 4.1.18.

- Find all eigenvalues by solving the characteristic equation $\det(A - \lambda I) = 0$. Here $\det(A - \lambda I) = \det \begin{bmatrix} 2 - \lambda & 1 \\ 0 & 3 - \lambda \end{bmatrix}$ which being a triangular matrix has determinant that is the product of the diagonals, namely $\det(A - \lambda I) = (2 - \lambda)(3 - \lambda)$. This determinant is zero only for eigenvalues $\lambda = 2$ or 3 . These are the diagonal entries in the triangular matrix.
- For each eigenvalue, find corresponding eigenvectors by solving the system $(A - \lambda I)\mathbf{x} = \mathbf{0}$.

¹ In this second case the vectors $A\mathbf{x}$ all appear to be pointing clockwise. Such a consistent rotation in $A\mathbf{x}$ is characteristic of matrices with complex valued eigenvalues and eigenvectors.

- For $\lambda = 2$, the system is $\begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} \mathbf{x} = \mathbf{0}$ which requires $x_2 = 0$. That is, all eigenvectors are $x_1(1, 0)$.
- For $\lambda = 3$, the system is $\begin{bmatrix} -1 & 1 \\ 0 & 0 \end{bmatrix} \mathbf{x} = \mathbf{0}$ which requires $x_1 = x_2$. That is, all eigenvectors are $x_2(1, 1)$.

The eigenvectors corresponding to the different eigenvalues are not orthogonal as their dot product $(1, 0) \cdot (1, 1) = 1 + 0 = 1 \neq 0$: instead they are at 45° to each other.

■

Theorem 7.0.2 (triangular matrices). *The diagonal entries of a triangular matrix are the only eigenvalues of the matrix. The corresponding eigenvectors of distinct eigenvalues are generally not orthogonal.*

Proof. We detail only the case of upper triangular matrices as the argument is similar for lower triangular matrices. First establish that the diagonal entries are eigenvalues, and second prove there are no others. Let λ be any value in the diagonal of the matrix A , and let k be the smallest index such that $a_{kk} = \lambda$ (this ‘smallest’ caters for duplicated diagonal values). Let’s construct an eigenvector in the form $\mathbf{x} = (x_1, x_2, \dots, x_{k-1}, 1, 0, \dots, 0)$. Set $x_k = 1$ and $x_j = 0$ for $j > k$. Then set

$$\begin{aligned} x_{k-1} &= -a_{k-1,k}x_k/(a_{k-1,k-1} - \lambda), \\ x_{k-2} &= -(a_{k-2,k}x_k + a_{k-2,k-1}x_{k-1})/(a_{k-2,k-2} - \lambda), \\ &\vdots \\ x_1 &= -(a_{1,k}x_k + a_{1,k-1}x_{k-1} + \dots + a_{1,2}x_2)/(a_{1,1} - \lambda). \end{aligned}$$

Since k is the smallest index for which $\lambda = a_{k,k}$ none of the above expressions involve divisions by zero, and so all are well defined. Rearranging the above equations shows that this vector \mathbf{x} satisfies, for $\lambda = a_{k,k}$,

$$\begin{aligned} (a_{1,1} - \lambda)x_1 + a_{1,2}x_2 + \dots + a_{1,k-1}x_{k-1} + a_{1,k}x_k &= 0, \\ (a_{2,2} - \lambda)x_2 + \dots + a_{2,k-1}x_{k-1} + a_{2,k}x_k &= 0, \\ &\vdots \\ (a_{k-1,k-1} - \lambda)x_{k-1} + a_{k-1,k}x_k &= 0, \\ (a_{k,k} - \lambda)x_k &= 0; \end{aligned}$$

that is, $(A - \lambda I)\mathbf{x} = \mathbf{0}$. Rearranging gives $A\mathbf{x} = \lambda\mathbf{x}$ for non-zero eigenvector \mathbf{x} and corresponding eigenvalue $\lambda = a_{kk}$.

Second, there can be no other eigenvalues. Every eigenvalue has to have non-trivial solutions, eigenvectors \mathbf{x} , to $(A - \lambda I)\mathbf{x} = \mathbf{0}$, which by Theorem 6.1.23 requires $\det(A - \lambda I) = 0$. But matrix $(A - \lambda I)$

is upper triangular, as A is upper triangular, so Theorem 6.2.14 asserts the determinant

$$\det(A - \lambda I) = (a_{1,1} - \lambda)(a_{2,2} - \lambda) \cdots (a_{n,n} - \lambda) = 0$$

iff the eigenvalue λ is one of the diagonal elements of A .

As an example of the non-orthogonality of eigenvectors, consider the two eigenvalues $\lambda_1 = a_{1,1}$ and $\lambda_2 = a_{2,2}$ with corresponding eigenvectors $\mathbf{x}_1 = (1, 0, 0, \dots, 0)$ and $\mathbf{x}_2 = (-a_{1,2}/(a_{1,1} - a_{2,2}), 1, 0, \dots, 0)$. Then the dot product $\mathbf{x}_1 \cdot \mathbf{x}_2 = -a_{1,2}/(a_{1,1} - a_{2,2}) \neq 0$ in general (the dot product is zero only when $a_{1,2} = 0$). Since the dot product is generally non-zero, \mathbf{x}_1 and \mathbf{x}_2 are generally not orthogonal. Similarly for other pairs of eigenvectors corresponding to distinct eigenvalues. \square

Example 7.0.3. Use Theorem 7.0.2 to find the eigenvalues, corresponding eigenvectors, and corresponding eigenspaces, of the following triangular matrices.

$$(a) A = \begin{bmatrix} -3 & 2 & 0 \\ 0 & -4 & 2 \\ 0 & 0 & 4 \end{bmatrix}$$

Solution: Matrix A is upper triangular so read off the eigenvalues from the diagonal to be -3 and ± 4 .

- For $\lambda = -3$, and by inspection, all eigenvectors are proportional to $(1, 0, 0)$. Hence eigenspace $\mathbb{E}_{-3} = \text{span}\{(1, 0, 0)\}$.
- For $\lambda = -4$ we need to solve

$$(A + 4I)\mathbf{x} = \begin{bmatrix} 1 & 2 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 8 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

By inspection an eigenvector must be of the form $(x_1, 1, 0)$. And the first line of the system then asserts $x_1 + 2 = 0$. Hence eigenvectors are proportional to $(-2, 1, 0)$. That is, the eigenspace $\mathbb{E}_{-4} = \text{span}\{(-2, 1, 0)\}$.

- For $\lambda = +4$ we need to solve

$$(A - 4I)\mathbf{x} = \begin{bmatrix} -7 & 2 & 0 \\ 0 & -8 & 2 \\ 0 & 0 & 0 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

Consider eigenvectors of the form $(x_1, x_2, 1)$. The second line asserts $-8x_2 + 2 = 0$, that is $x_2 = \frac{1}{4}$. The first line asserts $-7x_1 + 2x_2 = 0$, that is $x_1 = \frac{2}{7}x_2 = \frac{1}{14}$. Hence eigenvectors are proportional to $(\frac{1}{14}, \frac{1}{4}, 1)$. That is, the eigenspace $\mathbb{E}_4 = \text{span}\{(\frac{1}{14}, \frac{1}{4}, 1)\}$.

$$(b) \quad B = \begin{bmatrix} 3 & 0 & 0 & 0 \\ -2 & -4 & 0 & 0 \\ -3 & 1 & 0 & 0 \\ 0 & 0 & -3 & 1 \end{bmatrix}$$

Solution: Matrix B is lower triangular so read the eigenvalues from the diagonal to be 3, -4, 0 and 1.

- For $\lambda = 1$, by inspection all eigenvectors are of the form $(0, 0, 0, 1)$. Hence eigenspace $\mathbb{E}_1 = \text{span}\{(0, 0, 0, 1)\}$.
- For $\lambda = 0$, seek an eigenvector $(0, 0, 1, x_4)$ then the last line of the system

$$(B - 0I)\mathbf{x} = \begin{bmatrix} 3 & 0 & 0 & 0 \\ -2 & -4 & 0 & 0 \\ -3 & 1 & 0 & 0 \\ 0 & 0 & -3 & 1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

requires $-3 + x_4 = 0$. Hence eigenvectors are proportional to $(0, 0, 1, 3)$. That is, the eigenspace $\mathbb{E}_0 = \text{span}\{(0, 0, 1, 3)\}$.

- For $\lambda = -4$, seek an eigenvector $(0, 1, x_3, x_4)$ then the third line of the system

$$(B + 4I)\mathbf{x} = \begin{bmatrix} 7 & 0 & 0 & 0 \\ -2 & 0 & 0 & 0 \\ -3 & 1 & 4 & 0 \\ 0 & 0 & -3 & 5 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

requires $1 + 4x_3 = 0$, that is $x_3 = -\frac{1}{4}$. Then the last line of the system requires $\frac{3}{4} + 5x_4 = 0$, that is $x_4 = -\frac{3}{20}$. Hence eigenvectors are proportional to $(0, 1, -\frac{1}{4}, -\frac{3}{20})$. That is, the eigenspace $\mathbb{E}_{-4} = \text{span}\{(0, 1, -\frac{1}{4}, -\frac{3}{20})\}$.

- For $\lambda = 3$, seek an eigenvector $(1, x_2, x_3, x_4)$ then the second line of the system

$$(B - 3I)\mathbf{x} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ -2 & -7 & 0 & 0 \\ -3 & 1 & -3 & 0 \\ 0 & 0 & -3 & -2 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

requires $-2 - 7x_2 = 0$, that is $x_2 = -\frac{2}{7}$. Then the third line of the system requires $-3 - \frac{2}{7} - 3x_3 = 0$, that is $x_3 = -\frac{23}{21}$. Lastly, the last line requires $\frac{23}{7} - 2x_4 = 0$, that is $x_4 = \frac{23}{14}$. Hence eigenvectors are proportional to $(1, -\frac{2}{7}, -\frac{23}{21}, \frac{23}{14})$. That is, the eigenspace $\mathbb{E}_3 = \text{span}\{(1, -\frac{2}{7}, -\frac{23}{21}, \frac{23}{14})\}$.

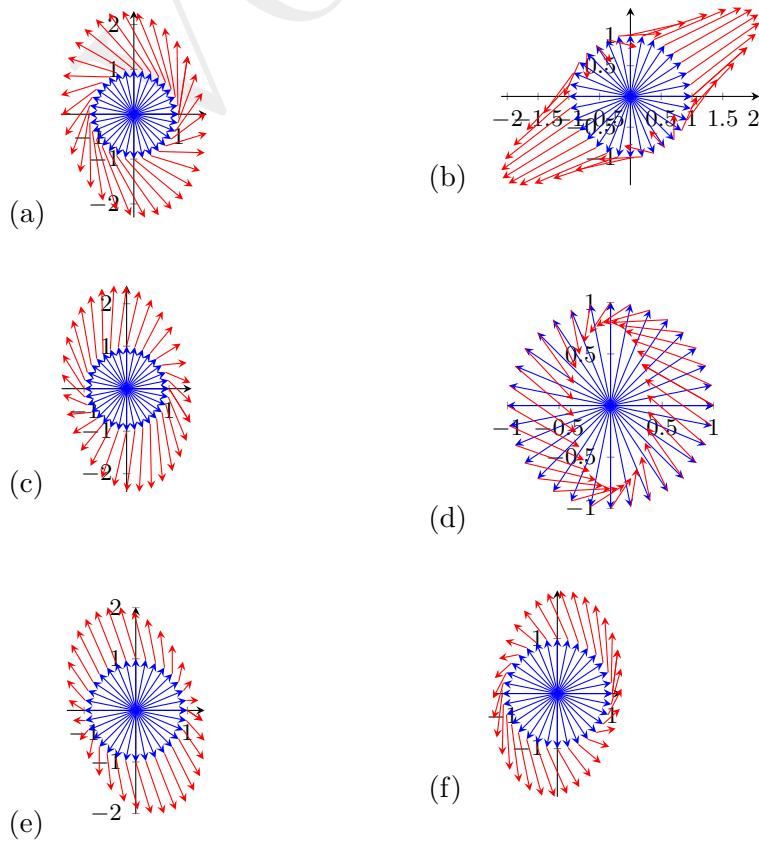
$$(c) \quad C = \begin{bmatrix} -1 & 1 & -8 & -5 & 5 \\ -3 & 6 & 4 & -3 & 0 \\ 1 & -3 & 1 & 0 & 0 \\ -7 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

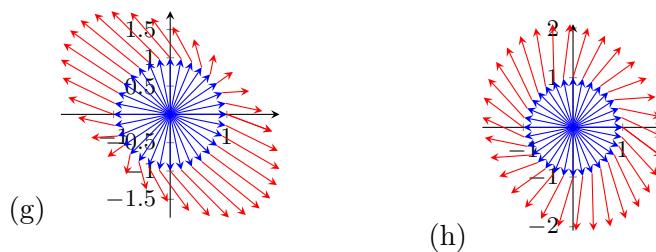
Solution: Matrix C is not a triangular matrix (Definition 6.2.13), so Theorem 7.0.2 does not apply. Row or column swaps could transform it to be triangular, but we have not investigated the effect of such swaps on eigenvalues and eigenvectors.

■

One consequence of the second part of the proof of Theorem 7.0.2 is that, when counted according to multiplicity, there are precisely n eigenvalues of an $n \times n$ diagonal matrix. Correspondingly, the next section establishes there are precisely n eigenvalues of general $n \times n$ matrices, provided we count the eigenvalues according to multiplicity and allow complex valued eigenvalues.

Exercise 7.0.1. Each of the following pictures applies to some specific real matrix, say called A . The pictures plot $A\mathbf{x}$ adjoined to the end of unit vectors \mathbf{x} . By inspection decide whether the matrix, in each case, has real eigenvalues or complex eigenvalues.





Exercise 7.0.2. For each of the following triangular matrices, write down all eigenvalues and then find the corresponding eigenspaces. Show your working.

$$(a) \begin{bmatrix} 2 & 0 \\ -1 & 4 \end{bmatrix}$$

(b) $\begin{bmatrix} 2 & 0 \\ -3 & 2 \end{bmatrix}$

$$(c) \begin{bmatrix} -1 & 0 & 3 \\ 0 & -1 & 2 \\ 0 & 0 & -5 \end{bmatrix}$$

$$(d) \begin{bmatrix} 0 & 0 & 0 \\ -3 & -4 & 0 \\ 1 & 5 & -2 \end{bmatrix}$$

$$(e) \begin{bmatrix} 1 & -5 & 0 \\ 0 & 0 & -4 \\ 0 & 0 & 0 \end{bmatrix}$$

$$(f) \begin{bmatrix} 0 & 0 & -2 \\ 0 & 4 & 1 \\ -3 & -3 & -1 \end{bmatrix}$$

$$(g) \begin{bmatrix} -1 & 0 & 0 \\ -2 & 2 & 0 \\ -2 & -1 & -1 \end{bmatrix}$$

$$(h) \begin{bmatrix} -2 & 4 & -2 & -2 \\ 0 & -2 & 1 & 7 \\ 0 & 0 & -3 & 1 \\ 0 & 0 & 0 & 2 \end{bmatrix}$$

$$(i) \begin{bmatrix} 8 & -2 & 3 & 2 \\ 0 & -6 & 1 & -2 \\ 0 & 0 & 3 & -2 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

$$(j) \begin{bmatrix} 0 & -2 & -5 & 2 \\ 0 & 7 & -1 & 2 \\ 0 & 0 & 3 & -4 \\ 0 & 0 & 0 & 3 \end{bmatrix}$$

7.1 Find eigenvalues and eigenvectors of matrices

Section Contents

7.1.1	A characteristic equation gives eigenvalues . .	559
7.1.2	Repeated eigenvalues are sensitive	574
7.1.3	Application: discrete dynamics of populations	577
7.1.4	Extension: Connect SVDs to eigen-problems	592
7.1.5	Exercises	595

Given the additional determinant methods of Chapter 6, this section begins exploring the properties and some applications of the eigenproblem $A\mathbf{x} = \lambda\mathbf{x}$ for general matrices A . We establish that there are generally n eigenvalues of an $n \times n$ matrix, albeit possibly complex valued, and that repeated eigenvalues are sensitive to errors. Applications include population modelling, connecting to the computation of SVDs, and fitting exponentials to real data.

7.1.1 A characteristic equation gives eigenvalues

The Fundamental Theorem of Algebra asserts that every polynomial equation over the complex field has a root. It is almost beneath the dignity of such a majestic theorem to mention that in fact it has precisely n roots.

J. H. Wilkinson, 1984 (Higham 1996, p.103)

Recall that eigenvalues λ and non-zero eigenvectors \mathbf{x} of a square matrix A must satisfy $(A - \lambda I)\mathbf{x} = \mathbf{0}$. Theorem 6.1.23 then implies the eigenvalues of a square matrix are precisely the solutions of the **characteristic equation** $\det(A - \lambda I) = 0$.

Theorem 7.1.1. *We call $\det(A - \lambda I)$ the **characteristic polynomial** of square matrix A .² The characteristic polynomial of an $n \times n$ matrix A is a polynomial of n th degree in λ . Consequently, there are at most n distinct eigenvalues of an $n \times n$ matrix A .*

Proof. Use induction on the size of the matrix. First, for 1×1 matrix $A = [a_{11}]$, the determinant $\det(A - \lambda I) = \det[a_{11} - \lambda] = a_{11} - \lambda$ which is of degree one in λ . Second, assume that for all $(n - 1) \times (n - 1)$ matrices A , $\det(A - \lambda I)$ is a polynomial of

² Alternatively, many call $\det(\lambda I - A)$ the characteristic polynomial, as does Matlab/Octave. The distinction is immaterial as, for an $n \times n$ matrix A and by Theorem 6.1.6c with multiplicative factor $k = -1$, the only difference in the determinant is a factor of $(-1)^n$. In Matlab/Octave, `poly(A)` computes the characteristic polynomial of the matrix, $\det(\lambda I - A)$, which might be useful for exercises, but is rarely useful in practice due to poor conditioning.

degree $n - 1$ in λ . Then use the Laplace Expansion Theorem 6.2.20 to give the first row expansion, in terms of minors of A and I ,

$$\begin{aligned}\det(A - \lambda I) &= (a_{11} - \lambda) \det(A_{11} - \lambda I_{11}) \\ &\quad - a_{12} \det(A_{12} - \lambda I_{12}) \\ &\quad + \cdots - (-1)^n a_{1n} \det(A_{1n} - \lambda I_{1n}).\end{aligned}$$

Now the minor I_{11} is precisely the $(n - 1) \times (n - 1)$ identity, and hence by assumption $\det(A_{11} - \lambda I_{11})$ is a polynomial of degree $(n - 1)$ in λ . But differently, the minors I_{12}, \dots, I_{1n} have two of the ones removed from the $n \times n$ identity and hence $\lambda I_{12}, \dots, \lambda I_{1n}$ each have only $(n - 2)$ λ s: since each term in a determinant is a product of distinct entries of the matrix (Theorem 6.2.23) it follows that for $j \geq 2$ the determinant $\det(A_{1j} - \lambda I_{1j})$ is a polynomial in λ of degree $\leq n - 2$. Consequently, the first row expansion of

$$\begin{aligned}\det(A - \lambda I) &= (a_{11} - \lambda)(\text{poly degree } n - 1) \\ &\quad - a_{12}(\text{poly degree } \leq n - 2) \\ &\quad + \cdots - (-1)^n a_{1n}(\text{poly degree } \leq n - 2) \\ &= (a_{11} - \lambda)(\text{poly degree } n - 1) \\ &\quad + (\text{poly degree } \leq n - 2) \\ &= (\text{poly degree } n)\end{aligned}$$

as the highest power of λ cannot be cancelled by any term. Induction thus implies the characteristic polynomial of an $n \times n$ matrix A is a polynomial of n th degree in λ .

Lastly, because the characteristic polynomial of A is of n th degree in λ , the Fundamental Theorem of Algebra asserts that the polynomial has at most n roots (possibly complex). Hence there are at most n distinct eigenvalues. \square

Example 7.1.2. Find the characteristic polynomial of each of the following matrices. Where in the coefficients of the polynomial can you see the determinant? and the sum of the diagonal elements?

$$(a) A = \begin{bmatrix} 1 & -1 \\ -2 & 4 \end{bmatrix}$$

Solution: The characteristic polynomial is $\det(A - \lambda I) = \det \begin{bmatrix} 1 - \lambda & -1 \\ -2 & 4 - \lambda \end{bmatrix} = (1 - \lambda)(4 - \lambda) - 2 = \lambda^2 - 5\lambda + 2$. Now $\det A = 4 - 2 = 2$ which is the coefficient of the constant term in this polynomial. Whereas the sum of the diagonal of A is $1 + 4 = 5$ which is the negative of the λ coefficient in the polynomial.

$$(b) B = \begin{bmatrix} 4 & -2 & 1 \\ 1 & -2 & 0 \\ 8 & 2 & 6 \end{bmatrix}$$

Solution: The characteristic polynomial is

$$\begin{aligned}\det(B - \lambda I) &= \det \begin{bmatrix} 4 - \lambda & -2 & 1 \\ 1 & -2 - \lambda & 0 \\ 8 & 2 & 6 - \lambda \end{bmatrix} \\ &= (4 - \lambda)(-2 - \lambda)(6 - \lambda) + 0 + 2 \\ &\quad - (-2 - \lambda)8 - 0 - (-2)(6 - \lambda) \\ &= -48 + 4\lambda + 8\lambda^2 - \lambda^3 + 2 \\ &\quad + 16 + 8\lambda + 12 + 2\lambda \\ &= -\lambda^3 + 8\lambda^2 + 2\lambda - 18.\end{aligned}$$

Now $\det B = 4(-2)6 + 0 + 2 - (-2)8 - 0 - (-2)6 = -48 + 2 + 16 + 12 = -18$ which is the coefficient of the constant term in the polynomial. Whereas the sum of the diagonal of B is $4 - 2 + 6 = 8$ which is the λ^2 coefficient in the polynomial.

■

These observations about the coefficients in the characteristic polynomials leads to the next theorem.

Theorem 7.1.3. *For an $n \times n$ matrix A , the product of the eigenvalues equals $\det A$ and equals the constant term in the characteristic polynomial. The sum of the eigenvalues equals $(-1)^{n-1}$ times the coefficient of λ^{n-1} in the characteristic polynomial and equals the trace of the matrix, $a_{11} + a_{22} + \dots + a_{nn}$.*

This optional theorem helps establish the nature of a characteristic polynomial.

Proof. Theorem 7.1.1, and its proof, establishes that the characteristic polynomial has the form

$$\begin{aligned}\det(A - \lambda I) &= c_0 + c_1\lambda + \dots + c_{n-1}\lambda^{n-1} + (-1)^n\lambda^n \\ &= (\lambda_1 - \lambda)(\lambda_2 - \lambda) \cdots (\lambda_n - \lambda),\end{aligned}\tag{7.1}$$

where the second equality follows from the Fundamental Theorem of Algebra that an n th degree polynomial factors into n linear factors (albeit possibly complex). First, substitute $\lambda = 0$ and this equation (7.1) gives $\det A = c_0 = \lambda_1\lambda_2 \cdots \lambda_n$, as required.

Second, consider the term $c_{n-1}\lambda^{n-1}$ in equation (7.1). From the factorisation on the right-hand side, the λ^{n-1} term arises as

$$\begin{aligned}c_{n-1}\lambda^{n-1} &= \lambda_1(-\lambda)^{n-1} + (-\lambda)\lambda_2(-\lambda)^{n-2} + \cdots + (-\lambda)^{n-1}\lambda_n \\ &= (-1)^{n-1}(\lambda_1 + \lambda_2 + \cdots + \lambda_n)\lambda^{n-1},\end{aligned}$$

and hence the coefficient $c_{n-1} = (-1)^{n-1}(\lambda_1 + \lambda_2 + \cdots + \lambda_n)$, as required. Now use induction to prove the trace formula. Recall that the proof of Theorem 7.1.1 establishes that

$$\det(A - \lambda I) = (a_{11} - \lambda) \det(A_{11} - \lambda I_{11})$$

$$+ (\text{poly degree } \leq n-2).$$

Assume the trace formula holds for $(n-1) \times (n-1)$ matrices such as the minor A_{11} . Then the previous identity gives

$$\begin{aligned} \det(A - \lambda I) &= (a_{11} - \lambda)[(-1)^{n-1}\lambda^{n-1} \\ &\quad + (-1)^{n-2}(a_{22} + \cdots + a_{nn})\lambda^{n-2} \\ &\quad + (\text{poly degree } \leq n-3)] \\ &\quad + (\text{poly degree } \leq n-2) \\ &= (-1)^n\lambda^n + (-1)^{n-1}(a_{11} + a_{22} + \cdots + a_{nn})\lambda^{n-1} \\ &\quad + (\text{poly degree } \leq n-2). \end{aligned}$$

Hence the coefficient $c_{n-1} = (-1)^{n-1}(a_{11} + a_{22} + \cdots + a_{nn})$. Since the formula holds for the basic case $n = 1$, namely $c_0 = +a_{11}$, induction implies the sum of the eigenvalues $\lambda_1 + \lambda_2 + \cdots + \lambda_n = a_{11} + a_{22} + \cdots + a_{nn}$, the trace of the matrix A . \square

Example 7.1.4. (a) What are the two highest order terms and the constant term in the characteristic polynomial of the matrix

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

Solution: First compute the determinant using the Laplace expansion (Theorem 6.2.20). The two zeros in the last row suggest a last row expansion:

$$\begin{aligned} \det A &= (-1)^6 1 \det \begin{bmatrix} -2 & 3 & -2 \\ -1 & -2 & 2 \\ 2 & 0 & 1 \end{bmatrix} \\ &\quad + (-1)^8 (-3) \det \begin{bmatrix} -2 & -1 & 3 \\ -1 & 3 & -2 \\ 2 & -3 & 0 \end{bmatrix} \\ &= (4 + 12 + 0 - 8 - 0 + 3) \\ &\quad - 3(0 + 4 + 9 - 18 + 12 - 0) = -10. \end{aligned}$$

This is the constant term in the characteristic polynomial. Second, the trace of A is $-2 + 3 + 0 - 3 = -2$ so the cubic coefficient in the characteristic polynomial is $(-1)^3(-2) = 2$. That is, the characteristic polynomial of A is of the form $\lambda^4 + 2\lambda^3 + \cdots - 10$.

(b) After laborious calculation you find the characteristic polynomial of the matrix

$$B = \begin{bmatrix} -2 & 5 & -3 & -1 & 2 \\ -2 & -5 & -1 & -1 & 3 \\ 1 & 4 & -2 & 1 & -7 \\ 1 & -5 & 1 & 4 & -5 \\ -1 & 0 & 3 & -3 & 1 \end{bmatrix}$$

is $-\lambda^5 + 2\lambda^4 - 3\lambda^3 + 234\lambda^2 + 884\lambda + 1564$. Could this polynomial be correct?

Solution: No, because the trace of B is $-2 - 5 - 2 + 4 + 1 = -4$ so the coefficient of the λ^4 term must be $(-1)^4(-4) = -4$ instead of the calculated 2.

- (c) After much calculation you find the characteristic polynomial of the matrix

$$C = \begin{bmatrix} 0 & 0 & 3 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & -4 & 0 & 3 & 0 \\ -5 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 2 \\ 0 & 0 & 0 & -6 & 0 & 0 \end{bmatrix}$$

is $\lambda^6 + 4\lambda^5 + 5\lambda^4 + 20\lambda^3 + 108\lambda^2 - 540\lambda + 668$. Could this polynomial be correct?

Solution: No. By the column of zeros in C , $\det C$ must be zero instead of the calculated 668.

- (d) What are the two highest order terms and the constant term in the characteristic polynomial of the matrix

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

Solution: First compute the determinant using the Laplace expansion (Theorem 6.2.20). The nearly all zeros in the last row suggests starting with a last row expansion (although others are just as good):

$$\det D = (-1)^8(-3) \det \begin{bmatrix} 0 & 0 & 0 & 3 & 0 \\ -2 & 0 & 1 & 0 & -2 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & -5 & 0 & -4 & 3 \\ 0 & -3 & 0 & -4 & 0 \end{bmatrix}$$

(using a 3rd row expansion)

$$= -3(-1)^6(-1) \det \begin{bmatrix} 0 & 0 & 3 & 0 \\ -2 & 0 & 0 & -2 \\ 0 & -5 & -4 & 3 \\ 0 & -3 & -4 & 0 \end{bmatrix}$$

(using a 1st column expansion)

$$= 3(-1)^3(-2) \det \begin{bmatrix} 0 & 3 & 0 \\ -5 & -4 & 3 \\ -3 & -4 & 0 \end{bmatrix}$$

$$\begin{aligned}
 & \text{(using a 3rd column expansion)} \\
 &= 6(-1)^5 3 \det \begin{bmatrix} 0 & 3 \\ -3 & -4 \end{bmatrix} \\
 &= -18(0 + 9) = -162.
 \end{aligned}$$

This is the constant term in the characteristic polynomial. Second, the trace of D is $0+0+0+0-4+0 = -4$ so the quintic coefficient in the characteristic polynomial is $(-1)^5(-4) = 4$. That is, the characteristic polynomial of D is of the form $\lambda^6 + 4\lambda^5 + \dots - 162$.

■

Definition 7.1.5. An eigenvalue λ_0 of a matrix A is said to have **multiplicity** m if the characteristic polynomial factorises to $\det(A - \lambda I) = (\lambda - \lambda_0)^m g(\lambda)$ where $g(\lambda_0) \neq 0$, and $g(\lambda)$ is a polynomial of degree $n - m$. Any eigenvalue of multiplicity $m \geq 2$ is also called a **repeated eigenvalue**.

Example 7.1.6. Use the characteristic polynomial to find all eigenvalues and their multiplicity of the following matrices.

$$(a) A = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix}$$

Solution: The characteristic equation is $\det(A - \lambda I) = \det \begin{bmatrix} 3 - \lambda & 1 \\ 0 & 3 - \lambda \end{bmatrix} = (3 - \lambda)^2 - 0 \cdot 1 = (\lambda - 3)^2 = 0$. The only eigenvalue is $\lambda = 3$ with multiplicity two. (Since this is an upper triangular matrix, the eigenvalues are the diagonal elements, namely 3 twice.)

$$(b) B = \begin{bmatrix} -1 & 1 & -2 \\ -1 & 0 & -1 \\ 0 & -3 & 1 \end{bmatrix}$$

Solution: The characteristic equation is

$$\begin{aligned}
 \det(B - \lambda I) &= \det \begin{bmatrix} -1 - \lambda & 1 & -2 \\ -1 & -\lambda & -1 \\ 0 & -3 & 1 - \lambda \end{bmatrix} \\
 &= (1 + \lambda)\lambda(1 - \lambda) + 0 - 6 \\
 &\quad - 0 + 3(1 + \lambda) + (1 - \lambda) \\
 &= -\lambda^3 + 3\lambda - 2 \\
 &= -(\lambda - 1)^2(\lambda + 2) = 0.
 \end{aligned}$$

Eigenvalues are $\lambda = 1$ with multiplicity two, and $\lambda = -2$ with multiplicity one.

$$(c) C = \begin{bmatrix} -1 & 0 & -2 \\ 0 & -3 & 2 \\ 0 & -2 & 1 \end{bmatrix}$$

Solution: The characteristic equation is

$$\begin{aligned}\det(C - \lambda I) &= \det \begin{bmatrix} -1 - \lambda & 0 & -2 \\ 0 & -3 - \lambda & 2 \\ 0 & -2 & 1 - \lambda \end{bmatrix} \\ &= (-1 - \lambda) \det \begin{bmatrix} -3 - \lambda & 2 \\ -2 & 1 - \lambda \end{bmatrix} \\ &= -(1 + \lambda)[(-3 - \lambda)(1 - \lambda) + 4] \\ &= -(\lambda + 1)[\lambda^2 + 2\lambda + 1] \\ &= -(\lambda + 1)^3 = 0.\end{aligned}$$

The only eigenvalue is $\lambda = -1$ with multiplicity three.

$$(d) D = \begin{bmatrix} 2 & 0 & -1 \\ -5 & 3 & -5 \\ 5 & -2 & -2 \end{bmatrix}$$

Solution: The characteristic equation is

$$\begin{aligned}\det(D - \lambda I) &= \det \begin{bmatrix} 2 - \lambda & 0 & -1 \\ -5 & 3 - \lambda & -5 \\ 5 & -2 & -2 - \lambda \end{bmatrix} \\ &= (2 - \lambda)(3 - \lambda)(-2 - \lambda) + 0 - 10 \\ &\quad + 5(3 - \lambda) - 10(2 - \lambda) - 0 \\ &= -\lambda^3 + 3\lambda^2 + 9\lambda - 27 \\ &= -(\lambda - 3)^2(\lambda + 3) = 0.\end{aligned}$$

Eigenvalues are $\lambda = 3$ with multiplicity two, and $\lambda = -3$ with multiplicity one.

$$(e) E = \begin{bmatrix} 0 & 1 \\ -1 & 1 \end{bmatrix}$$

Solution: The characteristic equation is

$$\begin{aligned}\det(E - \lambda I) &= \det \begin{bmatrix} -\lambda & 1 \\ -1 & 1 - \lambda \end{bmatrix} \\ &= -\lambda(1 - \lambda) + 1 \\ &= \lambda^2 - \lambda + 1 = 0.\end{aligned}$$

This quadratic equation does not factor easily so use the formula

$$\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{1 \pm \sqrt{1 - 4}}{2} = \frac{1}{2} \pm \frac{\sqrt{3}}{2}i,$$

for $i = \sqrt{-1}$, are two eigenvalues (complex valued) each of multiplicity one.



Example 7.1.7. Use `eig()` in Matlab/Octave to find the eigenvalues and their multiplicity for the following matrices. Recall (Table 4.1) that executing just `eig(A)` gives a column vector of eigenvalues of A , repeated according to their multiplicity.

$$(a) \begin{bmatrix} 2 & 2 & -1 \\ 0 & 1 & -2 \\ 0 & -1 & 0 \end{bmatrix}$$

Solution: Execute `eig([2 2 -1;0 1 -2;0 -1 0])` to get

```
ans =
    2
    2
   -1
```

So the eigenvalue $\lambda = 2$ has multiplicity two, and the eigenvalue $\lambda = -1$ has multiplicity one.

$$(b) \begin{bmatrix} -2 & -2 & -5 & 0 \\ 0 & -2 & 2 & 1 \\ -1 & 1 & 0 & -1 \\ -2 & 1 & 4 & 0 \end{bmatrix}$$

Solution: In Matlab/Octave execute

```
eig([-2 -2 -5 0
      0 -2 2 1
     -1 1 0 -1
     -2 1 4 0])
```

to get

```
ans =
-3.0000 + 0.0000i
-3.0000 + 0.0000i
1.0000 + 1.4142i
1.0000 - 1.4142i
```

There are two complex-valued eigenvalues, evidently $1 \pm \sqrt{2}i$, each of multiplicity one, and also the (real) eigenvalue $\lambda = -3$ which has multiplicity two.

$$(c) \begin{bmatrix} 3 & -1 & -2 & 1 & -2 \\ 0 & 0 & -2 & -2 & 0 \\ 2 & 1 & 1 & 1 & -1 \\ -1 & -3 & 0 & 1 & 2 \\ 2 & -2 & 1 & 0 & 3 \end{bmatrix}$$

Solution: In Matlab/Octave execute

```
eig([3 -1 -2 1 -2
      0 0 -2 -2 0
      2 1 1 1 -1
     -1 -3 0 1 2
      2 -2 1 0 3])
```



to get

```
ans =
2.0000 + 2.8284i
2.0000 - 2.8284i
4.0000 + 0.0000i
-0.0000 + 0.0000i
-0.0000 - 0.0000i
```

There are three eigenvalues of multiplicity one, namely 4 and $2 \pm \sqrt{8}i$. The last two rows appear to be the eigenvalue $\lambda = 0$ with multiplicity two.

$$(d) \begin{bmatrix} -1 & 0 & 0 & 0 \\ -1 & 2 & -3 & 3 \\ 3 & 1 & -1 & 0 \\ 0 & 3 & -2 & 1 \end{bmatrix}$$

Solution: In Matlab/Octave execute

```
eig([-1 0 0 0
      -1 2 -3 3
      3 1 -1 0
      0 3 -2 1])
```

to get

```
ans =
4.0000 + 0.0000i
-1.0000 + 0.0000i
-1.0000 - 0.0000i
-1.0000 + 0.0000i
```

There is one eigenvalue of multiplicity one, $\lambda = 4$. The last three rows appear to be the eigenvalue $\lambda = -1$ with multiplicity three.

■

To find eigenvalues and eigenvectors, the following restates Procedure 4.1.18 with a little more information, and now empowered to address larger matrices with the determinant tools from Chapter 6.

Procedure 7.1.8 (eigenvalues and eigenvectors). *To find by hand eigenvalues and eigenvectors of a (small) square matrix A:*

1. *find all eigenvalues (possibly complex) by solving the **characteristic equation** of A, $\det(A - \lambda I) = 0$;*
2. *for each eigenvalue λ , solve $(A - \lambda I)\mathbf{x} = \mathbf{0}$ to find the eigenspace \mathbb{E}_λ of all eigenvectors (together with $\mathbf{0}$);*
3. *write each eigenspace as the span of a few chosen eigenvectors (Definition 7.2.15 calls such a set a basis).*

Since, for an $n \times n$ matrix, the characteristic polynomial is of n th degree in λ (Theorem 7.1.1), there are n eigenvalues (when counted according to multiplicity and allowing complex eigenvalues).

Correspondingly, the following restates the computational procedure of section 4.1.1, but slightly more generally: the extra generality caters for non-symmetric matrices.

Compute in Matlab/Octave. For a given square matrix A , execute $[V, D] = \text{eig}(A)$, then the diagonal entries of D , $\text{diag}(D)$, are the eigenvalues of A . Corresponding to the eigenvalue $D(j, j)$ is an eigenvector $v_j = V(:, j)$, the j th column of V .³ If an eigenvalue is repeated (multiplicity more than one), then the corresponding columns of V span the eigenspace (and, as Section 7.2 discusses, the column vectors are a so-called basis for the eigenspace when they have a property called linear independence).

Example 7.1.9. Find the eigenspaces corresponding to the eigenvalues found for the first three matrices of Example 7.1.6.

7.1.6a. $A = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix}$

Solution: The only eigenvalue is $\lambda = 3$ with multiplicity two. Its eigenvectors x satisfy

$$(A - 3I)x = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}x = 0.$$

The second component of this equation is the trivial $0 = 0$. The first component of the equation is $0x_1 + 1x_2 = 0$, hence $x_2 = 0$. All eigenvectors are of the form $x = (1, 0)x_1$. That is the eigenspace $E_3 = \text{span}\{(1, 0)\}$.

In Matlab/Octave, executing $[V, D] = \text{eig}([3 \ 1; 0 \ 3])$ gives us

```
V =
 1.0000 -1.0000
 0.0000  0.0000
D =
 3   0
 0   3
```

Diagonal matrix D confirms the only eigenvalue is three, whereas the two columns of V confirm the eigenspace $E_1 = \text{span}\{(1, 0), (-1, 0)\} = \text{span}\{(1, 0)\}$.

³ Be aware that Matlab/Octave does not use the determinant to find the eigenvalues, nor does it solve the linear equations to find eigenvectors. For any but the smallest matrices such a ‘by hand’ approach is far too expensive. Instead, to find eigenvalues and eigenvectors, just as for the SVD, Matlab/Octave uses an intriguing iteration based upon what is called the QR factorisation.

7.1.6b. $B = \begin{bmatrix} -1 & 1 & -2 \\ -1 & 0 & -1 \\ 0 & -3 & 1 \end{bmatrix}$

Solution: The eigenvalues are $\lambda = -2$ (multiplicity one) and $\lambda = 1$ (multiplicity two).

– For $\lambda = -2$ solve

$$(B + 2I)\mathbf{x} = \begin{bmatrix} 1 & 1 & -2 \\ -1 & 2 & -1 \\ 0 & -3 & 3 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The third component of this equation requires $-3x_2 + 3x_3 = 0$, that is, $x_2 = x_3$. The second component requires $-x_1 + 2x_2 - x_3 = 0$, that is, $x_1 = 2x_2 - x_3 = 2x_3 - x_3 = x_3$. The first component requires $x_1 + x_2 - 2x_3 = 0$ which is also satisfied by $x_1 = x_2 = x_3$. All eigenvectors are of the form $x_3(1, 1, 1)$. That is, the eigenspace $\mathbb{E}_{-2} = \text{span}\{(1, 1, 1)\}$.

– For $\lambda = 1$ solve

$$(B - 1I)\mathbf{x} = \begin{bmatrix} -2 & 1 & -2 \\ -1 & -1 & -1 \\ 0 & -3 & 0 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The third component of this equation requires $-3x_2 = 0$, that is, $x_2 = 0$. The second component requires $-x_1 - x_2 - x_3 = 0$, that is, $x_1 = -x_2 - x_3 = 0 - x_3 = -x_3$. The first component requires $-2x_1 + x_2 - 2x_3 = 0$ which is also satisfied by $x_1 = -x_3$ and $x_2 = 0$. All eigenvectors are of the form $x_3(-1, 0, 1)$. That is, the eigenspace $\mathbb{E}_1 = \text{span}\{(-1, 0, 1)\}$.

Alternatively, in Matlab/Octave, executing

```
B=[-1 1 -2
    -1 0 -1
    0 -3 1]
[V,D]=eig(B)
```

gives us

```
V =
    -0.5774      0.7071     -0.7071
    -0.5774      0.0000      0.0000
    -0.5774     -0.7071      0.7071
D =
    -2            0            0
     0            1            0
     0            0            1
```

Diagonal matrix D confirms the eigenvalues. The first column of V confirms the eigenspace

$$\begin{aligned}\mathbb{E}_{-2} &= \text{span}\{(-0.5774, -0.5774, -0.5774)\} \\ &= \text{span}\{(1, 1, 1)\}.\end{aligned}$$

Whereas the last two columns of V confirm the eigenspace

$$\begin{aligned}\mathbb{E}_1 &= \text{span}\{(0.7071, 0, -0.7071), (-0.7071, 0, 0.7071)\} \\ &= \text{span}\{(-1, 0, 1)\}.\end{aligned}$$

7.1.6c. $C = \begin{bmatrix} -1 & 0 & -2 \\ 0 & -3 & 2 \\ 0 & -2 & 1 \end{bmatrix}$

Solution: The only eigenvalue is $\lambda = -1$ with multiplicity three. Its eigenvectors \mathbf{x} satisfy

$$(C + 1I)\mathbf{x} = \begin{bmatrix} 0 & 0 & -2 \\ 0 & -2 & 2 \\ 0 & -2 & 2 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The first component of this equation requires $x_3 = 0$. The second and third components both require $-2x_2 + 2x_3 = 0$, hence $x_2 = x_3 = 0$. Since x_1 is unconstrained, all eigenvectors are of the form $\mathbf{x} = x_1(1, 0, 0)$. That is the eigenspace $\mathbb{E}_{-1} = \text{span}\{(1, 0, 0)\}$.

Alternatively, in Matlab/Octave, executing

```
C=[-1 0 -2
    0 -3 2
    0 -2 1]
[V,D]=eig(C)
```

gives us

$V =$

$$\begin{matrix} 1 & -1 & -1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

$D =$

$$\begin{matrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{matrix}$$

Diagonal matrix D confirms the only eigenvalue is one with multiplicity three. The three columns of V confirm the eigenspace $\mathbb{E}_{-1} = \text{span}\{(1, 0, 0)\}$.

The matrices in Example 7.1.9 all have repeated eigenvalues. For these repeated eigenvalues the corresponding eigenspaces were all one dimensional. This contrasts with the case of symmetric matrices where the eigenspaces always have the same dimensionality as the multiplicity of the eigenvalue, as illustrated by Examples 4.1.10 and 4.1.15. Subsequent sections work towards Theorem 7.3.12 which establishes that for non-symmetric matrices an eigenspace has dimensionality between one and the multiplicity of the corresponding eigenvalue.

Example 7.1.10. By hand, find the eigenvalues and eigenspaces of the matrix

$$A = \begin{bmatrix} 0 & 3 & 0 & 0 & 0 \\ 1 & 0 & 3 & 0 & 0 \\ 0 & 1 & 0 & 3 & 0 \\ 0 & 0 & 1 & 0 & 3 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

Confirm your answer with `eig()` in Matlab/Octave.

Solution: Adopt Procedure 7.1.8.

(a) The characteristic equation is

$$\det(A - \lambda I) = \begin{vmatrix} -\lambda & 3 & 0 & 0 & 0 \\ 1 & -\lambda & 3 & 0 & 0 \\ 0 & 1 & -\lambda & 3 & 0 \\ 0 & 0 & 1 & -\lambda & 3 \\ 0 & 0 & 0 & 1 & -\lambda \end{vmatrix}$$

(by first row expansion (6.4))

$$= (-\lambda) \begin{vmatrix} -\lambda & 3 & 0 & 0 \\ 1 & -\lambda & 3 & 0 \\ 0 & 1 & -\lambda & 3 \\ 0 & 0 & 1 & -\lambda \end{vmatrix} - 3 \begin{vmatrix} 1 & 3 & 0 & 0 \\ 0 & -\lambda & 3 & 0 \\ 0 & 1 & -\lambda & 3 \\ 0 & 0 & 1 & -\lambda \end{vmatrix}$$

(by first row and first column expansion, respectively)

$$\begin{aligned} &= (-\lambda)^2 \begin{vmatrix} -\lambda & 3 & 0 \\ 1 & -\lambda & 3 \\ 0 & 1 & -\lambda \end{vmatrix} - (-\lambda)3 \begin{vmatrix} 1 & 3 & 0 \\ 0 & -\lambda & 3 \\ 0 & 1 & -\lambda \end{vmatrix} \\ &\quad - 3 \begin{vmatrix} -\lambda & 3 & 0 \\ 1 & -\lambda & 3 \\ 0 & 1 & -\lambda \end{vmatrix} \end{aligned}$$

(by common factor, and first column expansion)

$$= (\lambda^2 - 3) \begin{vmatrix} -\lambda & 3 & 0 \\ 1 & -\lambda & 3 \\ 0 & 1 & -\lambda \end{vmatrix} + 3\lambda \begin{vmatrix} -\lambda & 3 \\ 1 & -\lambda \end{vmatrix}$$

(using (6.1))

$$\begin{aligned} &= (\lambda^2 - 3)[(-\lambda)^3 + 0 + 0 - 0 + 3\lambda + 3\lambda] + 3\lambda(\lambda^2 - 3) \\ &= (\lambda^2 - 3)(-\lambda^3 + 9\lambda) \end{aligned}$$

$$= -\lambda(\lambda^2 - 3)(\lambda^2 - 9) = 0.$$

The five eigenvalues are thus $\lambda = 0, \pm\sqrt{3}, \pm 3$, all of multiplicity one.

(b) Consider each eigenvalue in turn.

$\lambda = 0$ Solve

$$(A - 0I)\mathbf{x} = \begin{bmatrix} 0 & 3 & 0 & 0 & 0 \\ 1 & 0 & 3 & 0 & 0 \\ 0 & 1 & 0 & 3 & 0 \\ 0 & 0 & 1 & 0 & 3 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The last row requires $x_4 = 0$. The fourth row requires $x_3 + 3x_5 = 0$, that is, $x_3 = -3x_5$. The third row requires $x_2 + 3x_4 = 0$, that is, $x_2 = -3x_4 = -3 \cdot 0 = 0$. The second row requires $x_1 + 3x_3 = 0$, that is, $x_1 = -3x_3 = 9x_5$. The first row requires $3x_2 = 0$, which is satisfied as $x_2 = 0$. Since all eigenvectors are of the form $(9x_5, 0, -3x_5, 0, x_5)$, the eigenspace $\mathbb{E}_0 = \text{span}\{(9, 0, -3, 0, 1)\}$.

$\lambda = \pm\sqrt{3}$ Being careful with the upper and lower signs, solve $(A \mp \sqrt{3}I)\mathbf{x} = \mathbf{0}$, that is,

$$\begin{bmatrix} \mp\sqrt{3} & 3 & 0 & 0 & 0 \\ 1 & \mp\sqrt{3} & 3 & 0 & 0 \\ 0 & 1 & \mp\sqrt{3} & 3 & 0 \\ 0 & 0 & 1 & \mp\sqrt{3} & 3 \\ 0 & 0 & 0 & 1 & \mp\sqrt{3} \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The last row requires $x_4 \mp \sqrt{3}x_5 = 0$, that is, $x_4 = \pm\sqrt{3}x_5$. The fourth row requires $x_3 \mp \sqrt{3}x_4 + 3x_5 = 0$, that is, $x_3 = \pm\sqrt{3}x_4 - 3x_5 = 3x_5 - 3x_5 = 0$. The third row requires $x_2 \mp \sqrt{3}x_3 + 3x_4 = 0$, that is, $x_2 = \pm\sqrt{3}x_3 - 3x_4 = \mp 3\sqrt{3}x_5$. The second row requires $x_1 \mp \sqrt{3}x_2 + 3x_3 = 0$, that is, $x_1 = \pm\sqrt{3}x_2 - 3x_3 = -9x_5$. The first row requires $\mp\sqrt{3}x_1 + 3x_2 = 0$, which is satisfied as $\mp\sqrt{3}(-9x_5) + 3(\mp 3\sqrt{3}x_5) = 0$. Since all eigenvectors are of the form $(-9x_5, \mp 3\sqrt{3}x_5, 0, \pm\sqrt{3}x_5, x_5)$, the eigenspaces $\mathbb{E}_{\pm\sqrt{3}} = \text{span}\{(-9, \mp 3\sqrt{3}, 0, \pm\sqrt{3}, 1)\}$.

$\lambda = \pm 3$ Being careful with the upper and lower signs, solve $(A \mp 3I)\mathbf{x} = \mathbf{0}$, that is,

$$\begin{bmatrix} \mp 3 & 3 & 0 & 0 & 0 \\ 1 & \mp 3 & 3 & 0 & 0 \\ 0 & 1 & \mp 3 & 3 & 0 \\ 0 & 0 & 1 & \mp 3 & 3 \\ 0 & 0 & 0 & 1 & \mp 3 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

The last row requires $x_4 \mp 3x_5 = 0$, that is, $x_4 = \pm 3x_5$. The fourth row requires $x_3 \mp 3x_4 + 3x_5 = 0$, that is, $x_3 = \pm 3x_4 - 3x_5 = 9x_5 - 3x_5 = 6x_5$. The third row requires $x_2 \mp 3x_3 + 3x_4 = 0$, that is, $x_2 = \pm 3x_3 - 3x_4 = \pm 9x_5$. The second row requires $x_1 \mp 3x_2 + 3x_3 = 0$, that is, $x_1 = \pm 3x_2 - 3x_3 = 9x_5$. The first row requires $\mp 3x_1 + 3x_2 = 0$, which is satisfied as $\mp 3(9x_5) + 3(\pm 9x_5) = 0$. Since all eigenvectors are of the form $(9x_5, \pm 9x_5, 6x_5, \pm 3x_5, x_5)$, the eigenspaces $\mathbb{E}_{\pm 3} = \text{span}\{(9, \pm 9, 6, \pm 3, 1)\}$.

■

Example 7.1.11. Use Matlab/Octave to confirm the eigenvalues and eigenvectors found for the matrix of Example 7.1.10.

Solution: In Matlab/Octave execute

```
A=[0 3 0 0 0
   1 0 3 0 0
   0 1 0 3 0
   0 0 1 0 3
   0 0 0 1 0]
[V,D]=eig(A)
```

to obtain the report (2 d.p.)

```
V =
  0.62  -0.62   0.94  -0.85  -0.85
  0.62   0.62  -0.00   0.49  -0.49
  0.42  -0.42  -0.31  -0.00   0.00
  0.21   0.21  -0.00  -0.16   0.16
  0.07  -0.07   0.10   0.09   0.09

D =
  3.00      0      0      0      0
      0  -3.00      0      0      0
      0      0  -0.00      0      0
      0      0      0  -1.73      0
      0      0      0      0  1.73
```

The columns of eigenvectors are more easily seen to confirm the hand calculation of Example 7.1.10 when we divide each column by its last element via `V*diag(1./V(5,:))` which gives the more appealing (2 d.p.)

```
ans =
  9.00   9.00   9.00  -9.00  -9.00
  9.00  -9.00   0.00   5.20  -5.20
  6.00   6.00  -3.00   0.00   0.00
  3.00  -3.00   0.00  -1.73   1.73
  1.00   1.00   1.00   1.00   1.00
```

7.1.2 Repeated eigenvalues are sensitive

This optional subsection does not prove the sensitivity: it just uses examples to introduce and illustrate.

Albeit hidden in Example 7.1.7, repeated eigenvalues are exquisitely sensitive to errors in either the matrix or the computation. If the matrix or the computation has an error e , then expect a repeated eigenvalue of multiplicity m to appear as m eigenvalues all within about $e^{1/m}$ of each other. Thus when we find or compute m eigenvalues all within about $e^{1/m}$, then suspect it to actually be one eigenvalue of multiplicity m .

For example, since computers work to a relative error of about 10^{-15} , then expect a repeated eigenvalue of multiplicity m to appear as m eigenvalues within about $10^{-15/m}$ of each other. Repeat some of the previous Examples 7.1.7, preceded by the Matlab/Octave command `format long`, to see that the repeated eigenvalues are sensitive to the computational errors.

Similarly, repeated eigenvalues are sensitive to errors in the matrix. The following examples show the sensitivity when we are uncertain about the components of a matrix.

Example 7.1.12. Compute eigenvalues of the following matrices and comment on the effect on repeated eigenvalues of errors in the matrix.

$$(a) A = \begin{bmatrix} a & 1 \\ 0.0001 & a \end{bmatrix} \text{ for any parameter } a.$$

Solution: By hand, the characteristic equation is

$$\det \begin{bmatrix} a - \lambda & 1 \\ 0.0001 & a - \lambda \end{bmatrix} = (a - \lambda)^2 - 0.0001 = 0.$$

Rearranging gives $(\lambda - a)^2 = 0.0001$. Taking square roots, $\lambda - a = \pm 0.01$; that is, the two eigenvalues are $\lambda = a \pm 0.01$. If we consider that the entry 0.0001 in the matrix is an error, that the entry should really be zero, then the eigenvalues should really be one repeated eigenvalue $\lambda = a$ of multiplicity two. However, the ‘error’ 0.0001 splits the repeated eigenvalue into two by an amount $\sqrt{0.0001} = 0.01$.

$$(b) B = \begin{bmatrix} -5 & -3 & 1 & 0 & 1 \\ -2 & -2 & 2 & -1 & -3 \\ 1 & 1 & -1 & -4 & -5 \\ -2 & -1 & 0 & -4 & 1 \\ 0 & 0 & 0 & 2 & -3 \end{bmatrix}$$

Solution: In Matlab/Octave execute

```
eig([-5 -3 1 0 1
      -2 -2 2 -1 -3
      1 1 -1 -4 -5])
```



```
-2 -1 0 -4 1
0 0 0 2 -3])
```

to get

```
ans =
2.7961e-08
-2.7961e-08
-6.4142
-5.0000
-3.5858
```

There are three eigenvalues of multiplicity one, namely -5 and $-5 \pm \sqrt{2}$. The two values $\pm 2.7961e-08$ at first sight appear to be two eigenvalues, $\pm 2.7961 \cdot 10^{-8}$, each of multiplicity one. However, computers usually work to about 15 digits accuracy, that is, the typical error is about 10^{-15} , so an eigenvalue of multiplicity two is generally computed as two eigenvalues separated by about $\sqrt{10^{-15}} \approx 3 \cdot 10^{-8}$. Thus we suspect that the two values $\pm 2.7961e-08$ actually represent one eigenvalue $\lambda = 0$ with multiplicity two.

- (c) Suppose the previous matrix B is obtained from some experiment where there are experimental errors in the entries with error about 0.0001. Randomly perturb the entries in matrix B to see the effects of such errors on the eigenvalues (use `randn()`, Table 3.1).

Solution: In Matlab/Octave execute

```
B=[-5 -3 1 0 1
-2 -2 2 -1 -3
1 1 -1 -4 -5
-2 -1 0 -4 1
0 0 0 2 -3]
eig(B+0.0001*randn(5))
```

to get something like

```
ans =
0.0307
-0.0308
-6.4143
-4.9996
-3.5862
```

Observe the repeated eigenvalue $\lambda = 0$ splits into two eigenvalues, $\lambda = \pm 0.031$, of size roughly $\sqrt{0.0001} = 0.01$. The other eigenvalues are also perturbed by the errors but only by amounts of size roughly 0.0001.

Depending upon the random numbers, another possible answer is



```
ans =
-0.0001 + 0.0112i
-0.0001 - 0.0112i
-6.4143 + 0.0000i
-4.9999 + 0.0000i
-3.5856 + 0.0000i
```

where the repeated eigenvalue of zero splits to be a pair of complex valued eigenvalues of roughly $\pm i\sqrt{0.0001} = \pm i0.01$.

$$(d) C = \begin{bmatrix} -1 & 0 & 0 & 0 \\ -1 & 2 & -3 & 3 \\ 3 & 1 & -1 & 0 \\ 0 & 3 & -2 & 1 \end{bmatrix} \text{ perturbed by errors of size } 10^{-6}$$

Solution: In Matlab/Octave execute

```
C=[-1 0 0 0
    -1 2 -3 3
    3 1 -1 0
    0 3 -2 1]
eig(C+1e-6*randn(4))
```

to get

```
ans =
4.0000 + 0.0000i
-1.0156 + 0.0000i
-0.9922 + 0.0139i
-0.9922 - 0.0139i
```

The eigenvalue 4 of multiplicity one is not noticeably affected by the errors about 10^{-6} . However, the repeated eigenvalue of $\lambda = -1$ with multiplicity three is split into three close eigenvalues (two complex-valued) all differing by roughly 0.01 which is indeed the cube-root of the perturbation 10^{-6} .



But symmetric matrices are OK The eigenvalues of a symmetric matrix are not so sensitive. This is fortunate as many applications give rise to symmetric matrices (Chapter 4). Such symmetry often reflects some symmetry in the natural world such as Newton's law of every action having an equal and opposite reaction. For symmetric matrices, the eigenvalues and eigenvectors are robust to computational perturbations and experimental errors.

Example 7.1.13. Compute the eigenvalues of the symmetric matrix

$$A = \begin{bmatrix} 1 & 1 & 0 & 2 \\ 1 & 0 & 2 & -1 \\ 0 & 2 & 1 & 4 \\ 2 & -1 & 4 & 1 \end{bmatrix}$$

and see matrix A has an eigenvalue of multiplicity two. Explore the effects on the eigenvalues of errors in the matrix by perturbing the entries by random amounts of size 0.0001.

Solution: In Matlab/Octave execute

```
A=[1 1 0 2
  1 0 2 -1
  0 2 1 4
  2 -1 4 1]
eig(A)
eig(A+0.0001*randn(4))
```

to get

```
ans =
-4.6235
1.0000
1.0000
5.6235
```

showing the eigenvalue $\lambda = 1$ has multiplicity two. Whereas the eigenvalues of the perturbed matrix depend upon the random numbers and so might be

```
ans =
-4.6236
5.6235
0.9998
0.9999
```

or perhaps

```
ans =
-4.6236 + 0.0000i
5.6234 + 0.0000i
1.0001 + 0.0000i
1.0001 - 0.0000i
```

In either case, the perturbation by amounts 0.0001 only change the eigenvalues, whether repeated or not, by an amount of about the same size.



7.1.3 Application: discrete dynamics of populations

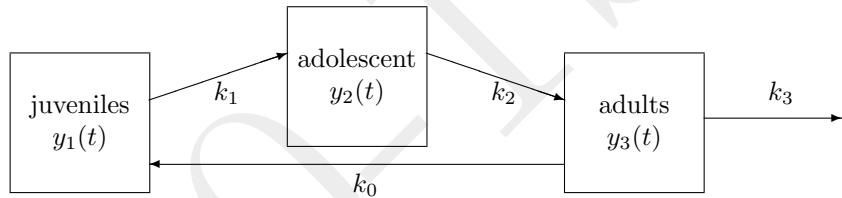
Age structured populations are one case where matrix properties and methods are crucial. The approach of this section is also akin to much mathematical modelling of diseases and epidemics. This section aims to show how to derive and use a matrix-vector model for the change in time t of interesting properties \mathbf{y} of a

population. Specifically, this subsection derives and analyses the model $\mathbf{y}(t+1) = A\mathbf{y}(t)$.

For a given species, let's define

- $y_1(t)$ to be the number of juveniles (including infants),
- $y_2(t)$ the number of adolescents, and
- $y_3(t)$ the number of adults.

Mostly, biologists only count females as females are the determining sex for reproduction. (Some bacteria/algae have seven sexes!) How do these numbers of females evolve over time? from generation to generation? First we need to choose a basic time interval (the unit of time): it could be one year, one month, one day, or maybe six months. Whatever we choose as convenient, we then quantify the number of events that happen to the females in each time interval as shown schematically in the diagram below:



Over any one time interval:

- a fraction k_1 of the juveniles become adolescents;
- a fraction k_2 of the adolescents become adults;
- a fraction k_3 of the adults die;
- but adults also give birth to juveniles at rate k_0 per adult.

Model this scenario with a system of discrete dynamical equations which are of the form that the numbers at the next time, $t + 1$, depend upon the numbers at the time t :

$$\begin{aligned} y_1(t+1) &= \dots, \\ y_2(t+1) &= \dots, \\ y_3(t+1) &= \dots. \end{aligned}$$

Let's fill in the right-hand sides from the given information about the rate of particular events per time interval.

- A fraction k_1 of the juveniles $y_1(t)$ becoming adolescents also means a fraction $(1 - k_1)$ of the juveniles remain juveniles, hence

$$\begin{aligned} y_1(t+1) &= (1 - k_1)y_1(t) + \dots, \\ y_2(t+1) &= +k_1y_1(t) + \dots, \\ y_3(t+1) &= \dots. \end{aligned}$$

- A fraction k_2 of the adolescents $y_2(t)$ becoming adults also means a fraction $(1 - k_2)$ of the adolescents remain adolescents, hence additionally

$$\begin{aligned}y_1(t+1) &= (1 - k_1)y_1(t) + \dots, \\y_2(t+1) &= +k_1y_1(t) + (1 - k_2)y_2(t), \\y_3(t+1) &= +k_2y_2(t) + \dots.\end{aligned}$$

- A fraction k_3 of the adults die mean that a fraction $(1 - k_3)$ of the adults remain adults, hence

$$\begin{aligned}y_1(t+1) &= (1 - k_1)y_1(t) + \dots, \\y_2(t+1) &= +k_1y_1(t) + (1 - k_2)y_2(t), \\y_3(t+1) &= +k_2y_2(t) + (1 - k_3)y_3(t).\end{aligned}$$

- But adults also give birth to juveniles at rate k_0 per adult so the number of juveniles increases by k_0y_3 from births:

$$\begin{aligned}y_1(t+1) &= (1 - k_1)y_1(t) + k_0y_3(t), \\y_2(t+1) &= +k_1y_1(t) + (1 - k_2)y_2(t), \\y_3(t+1) &= +k_2y_2(t) + (1 - k_3)y_3(t).\end{aligned}$$

This is our mathematical model of the age structure of the population.

Finally, write the mathematical model as the matrix-vector system

$$\begin{bmatrix} y_1(t+1) \\ y_2(t+1) \\ y_3(t+1) \end{bmatrix} = \begin{bmatrix} 1 - k_1 & 0 & k_0 \\ k_1 & 1 - k_2 & 0 \\ 0 & k_2 & 1 - k_3 \end{bmatrix} \begin{bmatrix} y_1(t) \\ y_2(t) \\ y_3(t) \end{bmatrix},$$

that is, $\mathbf{y}(t+1) = A\mathbf{y}(t)$. Such a model empowers predictions.

Example 7.1.14 (orangutans). From the following extract of the Wikipedia entry on orangutans [20 Mar 2014] derive a mathematical model for the age structure of the orangutans from one year to the next.



Gestation lasts for 9 months, with females giving birth to their first offspring between the ages of 14 and 15 years. Female orangutans have [seven to] eight-year intervals between births, the longest interbirth intervals among the great apes. ... Infant orangutans are completely dependent on their mothers for the first two years of their lives. The mother will carry the infant during travelling, as well as feed it and sleep with it in the same night nest. For the first four months, the infant is carried on its belly and never relieves physical contact. In the following months, the time an infant spends with its mother decreases. When an orangutan reaches the age of two, its climbing skills improve and it will travel through

the canopy holding hands with other orangutans, a behaviour known as “buddy travel”. Orangutans are juveniles from about two to five years of age and will start to temporarily move away from their mothers. Juveniles are usually weaned at about four years of age. Adolescent orangutans will socialize with their peers while still having contact with their mothers. Typically, orangutans live over 30 years in both the wild and captivity.

Suppose the initial population of orangutans in some area at year zero of a study is that of 30 adolescent females and 15 adult females. Use the mathematical model to predict the population for the next five years.

Solution: Choose level: we choose a time unit of one year, and choose to model three age categories.⁴

Define:

- $y_1(t)$ is the number of juvenile females (including infant) at time t (years), say age ≤ 5 years;
- $y_2(t)$ is the number of adolescent females, say $6 \leq$ age ≤ 14 years;
- $y_3(t)$ is the number of adult females, say age ≥ 15 years.

Quantify changes in a year: from Wikipedia information (with numbers slightly modified here to make the results numerically simpler):

- Orangutans are juvenile for say 5 years, so in any one year a fraction $1/5$ of them grow to be adolescent and a fraction $4/5$ remain as juveniles. Say an adult female gives birth every 7–8 years, so on average it gives birth to a juvenile female every 15 years. Thus a fraction $1/15$ of adults give birth to a juvenile female in any year. Consequently, model as $y_1(t+1) = \frac{4}{5}y_1(t) + \frac{1}{15}y_3(t)$.
- Adolescent orangutans become breeding adults after another 9–10 years, so in any one year a fraction $1/10$ of them grow to adults and $9/10$ ths remain adolescents. Consequently, $y_2(t+1) = +\frac{1}{5}y_1(t) + \frac{9}{10}y_2(t)$.
- Orangutans live to 30 years, about 15 years of adulthood so in any one year a fraction $14/15$ of adult females live to the next year. Consequently, $y_3(t+1) = \frac{1}{10}y_2(t) + \frac{9}{10}y_3(t)$.

The mathematical model of the age structure is to let vector $\mathbf{y} =$

⁴ A coarse model considers just the total number in a species; a fine model might count the number of each age in years (here 30 years); a ‘micro’ model might simulate each and every orangutan as individuals (thousands).

(y_1, y_2, y_3) then

$$\mathbf{y}(t+1) = A\mathbf{y}(t) = \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix} \mathbf{y}(t)$$

To predict the population we need to know the current population. The given information is that there are initially 30 adolescent females and 15 adult females. This information specifies that at time zero the population vector $\mathbf{y}(0) = (0, 30, 15)$ in the study area.

(a) Then the rule $\mathbf{y}(t+1) = A\mathbf{y}(t)$ with time $t = 0$ gives

$$\begin{aligned} \mathbf{y}(1) &= A\mathbf{y}(0) \\ &= \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix} \begin{bmatrix} 0 \\ 30 \\ 15 \end{bmatrix} \\ &= \begin{bmatrix} 1 \\ 27 \\ 17 \end{bmatrix}. \end{aligned}$$

That is, during the first year there is one birth of a female juvenile, three adolescents matured to adults, and one adult died.

(b) Then the rule $\mathbf{y}(t+1) = A\mathbf{y}(t)$ with time $t = 1$ year gives

$$\begin{aligned} \mathbf{y}(2) &= A\mathbf{y}(1) \\ &= \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix} \begin{bmatrix} 1 \\ 27 \\ 17 \end{bmatrix} \\ &= \begin{bmatrix} \frac{29}{15} \\ \frac{49}{2} \\ \frac{557}{30} \end{bmatrix} = \begin{bmatrix} 1.93 \\ 24.50 \\ 18.57 \end{bmatrix} \text{ (2 d.p.)}. \end{aligned}$$

In real life we cannot have such fractions of an orangutan. These predictions are averages, or expectations on average, and must be interpreted as such. Thus after two years, we expect on average nearly two juveniles, 24 or 25 adolescents, and 18 or 19 adults.

(c) Continuing on with the aid of Matlab/Octave or calculator, the rule $\mathbf{y}(t+1) = A\mathbf{y}(t)$ with time $t = 2$ years gives

$$\begin{aligned} \mathbf{y}(3) &= A\mathbf{y}(2) \\ &= \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix} \begin{bmatrix} 1.93 \\ 24.50 \\ 18.57 \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} 2.78 \\ 22.44 \\ 19.78 \end{bmatrix} \text{ (2 d.p.)}.$$

In Matlab/Octave do this calculation via

```
A=[4/5 0 1/15;1/5 9/10 0;0 1/10 14/15]
y0=[0;30;15]
y1=A*y0
y2=A*y1
y3=A*y2
y4=A*y3
y5=A*y4
```



- (d) Consequently, the rule $\mathbf{y}(t+1) = A\mathbf{y}(t)$ with time $t = 3$ years gives

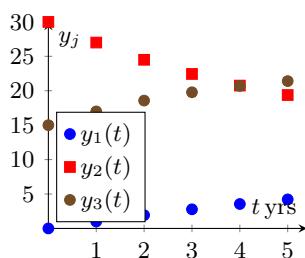
$$\begin{aligned} \mathbf{y}(4) &= A\mathbf{y}(3) \\ &= \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix} \begin{bmatrix} 2.78 \\ 22.44 \\ 19.78 \end{bmatrix} \\ &= \begin{bmatrix} 3.55 \\ 20.75 \\ 20.70 \end{bmatrix} \text{ (2 d.p.)}. \end{aligned}$$

- (e) Lastly, the rule $\mathbf{y}(t+1) = A\mathbf{y}(t)$ with time $t = 4$ years gives

$$\begin{aligned} \mathbf{y}(5) &= A\mathbf{y}(4) \\ &= \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix} \begin{bmatrix} 3.55 \\ 20.75 \\ 20.70 \end{bmatrix} \\ &= \begin{bmatrix} 4.22 \\ 19.38 \\ 21.40 \end{bmatrix} \text{ (2 d.p.)}. \end{aligned}$$

Thus after five years the mathematical model predicts about 4 juveniles, 19 adolescents, and 21 adults (on average).

Notice that the five-year population of 44 females (45 if you add all the fractions) is the same as the starting population. This nearly static total population is no accident, as we next see, and contributes to why orangutans are critically endangered. ■



The mathematical model $\mathbf{y}(t+1) = A\mathbf{y}(t)$ does forecast the future populations. However, to make predictions for many years and for general initial populations we prefer the formula solution given by the next Theorem 7.1.16 and introduced in the next example.

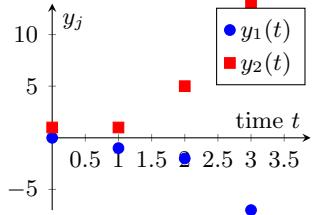
Example 7.1.15. A vector $\mathbf{y}(t) \in \mathbb{R}^2$ changes with time t according to

$$\mathbf{y}(t+1) = A\mathbf{y}(t) = \begin{bmatrix} 1 & -1 \\ -4 & 1 \end{bmatrix} \mathbf{y}(t).$$

First, what is $\mathbf{y}(3)$ if the initial $\mathbf{y}(0) = (0, 1)$? Second, find a general formula for $\mathbf{y}(t)$ from any $\mathbf{y}(0)$.

Solution: First, given $\mathbf{y}(0) = (0, 1)$ just compute (as drawn in the margin)

$$\begin{aligned}\mathbf{y}(1) &= A\mathbf{y}(0) = \begin{bmatrix} 1 & -1 \\ -4 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \\ \mathbf{y}(2) &= A\mathbf{y}(1) = \begin{bmatrix} 1 & -1 \\ -4 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -2 \\ 5 \end{bmatrix}, \\ \mathbf{y}(3) &= A\mathbf{y}(2) = \begin{bmatrix} 1 & -1 \\ -4 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 5 \end{bmatrix} = \begin{bmatrix} -7 \\ 13 \end{bmatrix}.\end{aligned}$$



Second, let's suppose there may be solutions in the form $\mathbf{y} = \lambda^t \mathbf{v}$ for some non-zero scalar λ and some vector $\mathbf{v} \in \mathbb{R}^2$. Substitute into $\mathbf{y}(t+1) = A\mathbf{y}(t)$ to find $\lambda^{t+1}\mathbf{v} = A\lambda^t\mathbf{v}$. Divide by (non-zero) λ^t to find we require $A\mathbf{v} = \lambda\mathbf{v}$: this is an eigen-problem. That is, for every eigenvalue λ with corresponding eigenvector \mathbf{v} there is a solution $\mathbf{y} = \lambda^t \mathbf{v}$.

To find the eigenvalues here solve the characteristic equation

$$\det(A - \lambda I) = \begin{vmatrix} 1 - \lambda & -1 \\ -4 & 1 - \lambda \end{vmatrix} = (1 - \lambda)^2 - 4 = 0.$$

Rearrange to $(\lambda - 1)^2 = 4$ and take square roots to give $\lambda - 1 = \pm 2$, that is, eigenvalues $\lambda = 1 \pm 2 = -1, 3$.

- For eigenvalue $\lambda_1 = -1$ the corresponding eigenvector has to satisfy

$$(A + I)\mathbf{v}_1 = \begin{bmatrix} 2 & -1 \\ -4 & 2 \end{bmatrix} \mathbf{v}_1 = \mathbf{0}.$$

That is, $\mathbf{v}_1 \propto (1, 2)$; let's take $\mathbf{v}_1 = (1, 2)$.

- For eigenvalue $\lambda_2 = 3$ the corresponding eigenvector has to satisfy

$$(A - 3I)\mathbf{v}_2 = \begin{bmatrix} -2 & -1 \\ -4 & -2 \end{bmatrix} \mathbf{v}_2 = \mathbf{0}.$$

That is, $\mathbf{v}_2 \propto (1, -2)$; let's take $\mathbf{v}_2 = (-1, 2)$.

Summarising so far: $\mathbf{y}'(t) = (-1)^t (1, 2)$ and $\mathbf{y}''(t) = 3^t (-1, 2)$ are both solutions of $\mathbf{y}(t+1) = A\mathbf{y}(t)$.

Further, the equation $\mathbf{y}(t+1) = A\mathbf{y}(t)$ is linear in \mathbf{y} , so any linear combination of solutions is also a solution. Let's try $\mathbf{y}(t) = c_1 \mathbf{y}'(t) + c_2 \mathbf{y}''(t)$. Substituting

$$A\mathbf{y}(t) = A(c_1 \mathbf{y}'(t) + c_2 \mathbf{y}''(t))$$

$$\begin{aligned} &= c_1 A \mathbf{y}'(t) + c_2 A \mathbf{y}''(t) \\ &= c_1 \mathbf{y}'(t+1) + c_2 \mathbf{y}''(t+1) = \mathbf{y}(t+1), \end{aligned}$$

as required. So, $\mathbf{y}(t) = c_1(-1)^t(1, 2) + c_23^t(-1, 2)$ are solutions for all constants c_1 and c_2 , and over all time t .

Now, what values of constants c_1 and c_2 should be chosen for any given initial $\mathbf{y}(0)$? Substitute time $t = 0$ into $\mathbf{y}(t) = c_1(-1)^t(1, 2) + c_23^t(-1, 2)$ to require $\mathbf{y}(0) = c_1(-1)^0(1, 2) + c_23^0(-1, 2)$. Since the zero powers $(-1)^0 = 3^0 = 1$, this requires $\mathbf{y}(0) = c_1(1, 2) + c_2(-1, 2)$. Write as the matrix-vector system

$$\begin{bmatrix} 1 & -1 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \mathbf{y}(0) \iff \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} \frac{1}{2} & \frac{1}{4} \\ -\frac{1}{2} & \frac{1}{4} \end{bmatrix} \mathbf{y}(0)$$

upon invoking the inverse (Theorem 3.2.6) of the matrix of eigenvectors, $P = [\mathbf{v}_1 \ \mathbf{v}_2]$. That is, because the matrix of eigenvectors is invertible, we find constants c_1 and c_2 to suit any initial $\mathbf{y}(0)$.

For example, with initial $\mathbf{y}(0) = (0, 1)$ the above formula gives $(c_1, c_2) = (\frac{1}{4}, \frac{1}{4})$ and so the corresponding formula solution is $\mathbf{y}(t) = \frac{1}{4}(-1)^t(1, 2) + \frac{1}{4}3^t(-1, 2)$. To check, evaluate at say $t = 3$ to find $\mathbf{y}(3) = -\frac{1}{4}(1, 2) + \frac{27}{4}(-1, 2) = (-7, 13)$, as before. In general, as here, as time t increases, the solution $\mathbf{y}(t)$ grows like 3^t with a little oscillation from the $(-1)^t$ term.

■

Theorem 7.1.16. Suppose the $n \times n$ square matrix A governs the dynamics of $\mathbf{y}(t) \in \mathbb{R}^n$ according to $\mathbf{y}(t+1) = A\mathbf{y}(t)$.

- (a) Let $\lambda_1, \lambda_2, \dots, \lambda_m$ be eigenvalues of A and $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$ be corresponding eigenvectors, then a solution of $\mathbf{y}(t+1) = A\mathbf{y}(t)$ is the linear combination

$$\mathbf{y}(t) = c_1 \lambda_1^t \mathbf{v}_1 + c_2 \lambda_2^t \mathbf{v}_2 + \cdots + c_m \lambda_m^t \mathbf{v}_m \quad (7.2)$$

for all constants c_1, c_2, \dots, c_m .

- (b) Further, if the number of eigenvectors $m = n$, the size of the matrix, and the matrix of eigenvectors $P = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$ is invertible, then the general linear combination (7.2) is a **general solution** in that constants c_1, c_2, \dots, c_n can be found for every given initial $\mathbf{y}(0)$.

Proof.

7.1.16a Just premultiply (7.2) by matrix A to find that

$$\begin{aligned} A\mathbf{y}(t) &= A(c_1 \lambda_1^t \mathbf{v}_1 + c_2 \lambda_2^t \mathbf{v}_2 + \cdots + c_m \lambda_m^t \mathbf{v}_m) \\ &\quad (\text{using distributivity Thm. 3.1.16}) \end{aligned}$$

$$\begin{aligned}
&= c_1 \lambda_1^t A \mathbf{v}_1 + c_2 \lambda_2^t A \mathbf{v}_2 + \cdots + c_m \lambda_m^t A \mathbf{v}_m \\
&\quad (\text{as eigenvectors } A \mathbf{v}_j = \lambda_j \mathbf{v}_j) \\
&= c_1 \lambda_1^t \lambda_1 \mathbf{v}_1 + c_2 \lambda_2^t \lambda_2 \mathbf{v}_2 + \cdots + c_m \lambda_m^t \lambda_m \mathbf{v}_m \\
&= c_1 \lambda_1^{t+1} \mathbf{v}_1 + c_2 \lambda_2^{t+1} \mathbf{v}_2 + \cdots + c_m \lambda_m^{t+1} \mathbf{v}_m
\end{aligned}$$

which is the given formula (7.2) for $\mathbf{y}(t+1)$. Hence (7.2) is a solution of $\mathbf{y}(t+1) = A\mathbf{y}(t)$.

7.1.16b For any given initial value $\mathbf{y}(0)$, the solution (7.2) will hold if we can find constants c_1, c_2, \dots, c_m such that the solution (7.2) evaluates to $\mathbf{y}(0)$ at time $t = 0$. Let's do this given the preconditions that the number of eigenvectors $m = n$ and the matrix $P = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$ is invertible. Evaluating the solution (7.2) at $t = 0$ we need to solve, since the zeroth power $\lambda_j^0 = 1$,

$$\mathbf{y}(0) = c_1 \mathbf{v}_1 + c_2 \mathbf{v}_2 + \cdots + c_n \mathbf{v}_n.$$

Writing as a matrix-vector system this equation requires $P\mathbf{c} = \mathbf{y}(0)$ for constant vector $\mathbf{c} = (c_1, c_2, \dots, c_n)$. Since P is invertible, $P\mathbf{c} = \mathbf{y}(0)$ always has the unique solution $\mathbf{c} = P^{-1}\mathbf{y}(0)$ (Theorem 3.4.35) which are the requisite constants.

□

Example 7.1.17. Consider the dynamics of $\mathbf{y}(t+1) = A\mathbf{y}(t)$ for matrix $A = \begin{bmatrix} 1 & 3 \\ -1 & 1 \end{bmatrix}$. First, what is $\mathbf{y}(3)$ when the initial $\mathbf{y}(0) = (1, 0)$? Second, find a general solution.

Solution: First, just compute (as illustrated)

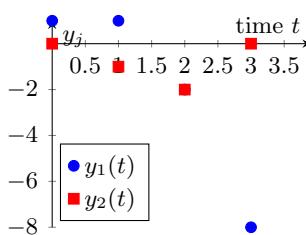
$$\begin{aligned}
\mathbf{y}(1) &= A\mathbf{y}(0) = \begin{bmatrix} 1 & 3 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \\
\mathbf{y}(2) &= A\mathbf{y}(1) = \begin{bmatrix} 1 & 3 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -2 \\ -2 \end{bmatrix}, \\
\mathbf{y}(3) &= A\mathbf{y}(2) = \begin{bmatrix} 1 & 3 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ -2 \end{bmatrix} = \begin{bmatrix} -8 \\ 0 \end{bmatrix}.
\end{aligned}$$

Interestingly, after three steps in time $\mathbf{y}(3)$ is (-8) times the initial $\mathbf{y}(0)$. This suggest after six steps in time $\mathbf{y}(6)$ will be $(-8)^2 = 64$ times the initial $\mathbf{y}(0)$, and so on. Perhaps the solution grows in size roughly like 2^t but in some irregular manner.

Second, find a general solution via the eigenvalues and eigenvectors of the matrix A . Its characteristic equation is

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

That is, $(\lambda - 1)^2 = -3$ which upon taking square roots gives the complex conjugate pair of eigenvalues $\lambda = 1 \pm i\sqrt{3}$. Theorem 7.1.16 applies for complex eigenvalues and eigenvectors so we proceed.



- For eigenvalue $\lambda_1 = 1 + i\sqrt{3}$ the corresponding eigenvectors \mathbf{v}_1 satisfy

$$(A - (1 + i\sqrt{3})I)\mathbf{v}_1 = \begin{bmatrix} -i\sqrt{3} & 3 \\ -1 & -i\sqrt{3} \end{bmatrix} \mathbf{v}_1 = \mathbf{0}.$$

Solutions are proportional to $\mathbf{v}_1 = (-i\sqrt{3}, 1)$.

- For eigenvalue $\lambda_2 = 1 - i\sqrt{3}$ the corresponding eigenvectors \mathbf{v}_2 satisfy

$$(A - (1 - i\sqrt{3})I)\mathbf{v}_2 = \begin{bmatrix} i\sqrt{3} & 3 \\ -1 & i\sqrt{3} \end{bmatrix} \mathbf{v}_2 = \mathbf{0}.$$

Solutions are proportional to $\mathbf{v}_2 = (+i\sqrt{3}, 1)$.

Theorem 7.1.16 then establishes that a solution to $\mathbf{y}(t+1) = A\mathbf{y}(t)$ is

$$\mathbf{y}(t) = c_1(1 + i\sqrt{3})^t \begin{bmatrix} -i\sqrt{3} \\ 1 \end{bmatrix} + c_2(1 - i\sqrt{3})^t \begin{bmatrix} +i\sqrt{3} \\ 1 \end{bmatrix}.$$

This is a general solution since the matrix of the two eigenvectors (albeit complex valued) is invertible:

$$P = [\mathbf{v}_1 \ \mathbf{v}_2] = \begin{bmatrix} -i\sqrt{3} & i\sqrt{3} \\ 1 & 1 \end{bmatrix} \quad \text{has } P^{-1} = \begin{bmatrix} \frac{i}{2\sqrt{3}} & \frac{1}{2} \\ \frac{-i}{2\sqrt{3}} & \frac{1}{2} \end{bmatrix}$$

as its inverse (Theorem 3.2.6).

For example, if $\mathbf{y}(0) = (1, 0)$, then the coefficient constants are $(c_1, c_2) = P^{-1}(1, 0) = (i, -i)/(2\sqrt{3})$. Then the solution becomes

$$\begin{aligned} \mathbf{y}(t) &= \frac{i}{2\sqrt{3}}(1 + i\sqrt{3})^t \begin{bmatrix} -i\sqrt{3} \\ 1 \end{bmatrix} - \frac{i}{2\sqrt{3}}(1 - i\sqrt{3})^t \begin{bmatrix} +i\sqrt{3} \\ 1 \end{bmatrix} \\ &= \frac{1}{2}(1 + i\sqrt{3})^t \begin{bmatrix} 1 \\ \frac{i}{\sqrt{3}} \end{bmatrix} + \frac{1}{2}(1 - i\sqrt{3})^t \begin{bmatrix} 1 \\ -\frac{i}{\sqrt{3}} \end{bmatrix}. \end{aligned}$$

Through the magic of the complex conjugate form of the two terms in this expression, the complex parts cancel to always give a real result. For example, this complex formula predicts at time step $t = 1$

$$\begin{aligned} \mathbf{y}(1) &= \frac{1}{2}(1 + i\sqrt{3}) \begin{bmatrix} 1 \\ \frac{i}{\sqrt{3}} \end{bmatrix} + \frac{1}{2}(1 - i\sqrt{3}) \begin{bmatrix} 1 \\ -\frac{i}{\sqrt{3}} \end{bmatrix} \\ &= \frac{1}{2} \begin{bmatrix} 1 + i\sqrt{3} + 1 - i\sqrt{3} \\ \frac{i}{\sqrt{3}} - 1 - \frac{i}{\sqrt{3}} - 1 \end{bmatrix} \\ &= \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \end{aligned}$$

as computed directly at the start of this example.

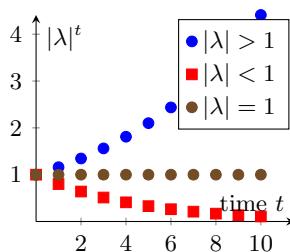
One crucial qualitative aspect we need to know is whether components in the solution (7.2) grow, decay, or stay the same size as time increases. This is determined by the eigenvalues as λ_j^t is the only place that the time appears in the formula (7.2).

- For example, in the general solution for Example 7.1.15, $\mathbf{y}(t) = c_1(-1)^t(1, 2) + c_23^t(-1, 2)$, the 3^t factor grows in time since $3^1 = 3$, $3^2 = 9$, $3^3 = 27$, and so on. Whereas the $(-1)^t$ factor just oscillates in time since $(-1)^1 = -1$, $(-1)^2 = 1$, $(-1)^3 = -1$, and so on. Thus for long times we know that the term involving the factor 3^t will dominate the solution as it grows.
- In Example 7.1.17 with complex conjugate eigenvalues the situation is more complicated. Let's write any given complex eigenvalue in polar form $\lambda = r(\cos \theta + i \sin \theta)$ where magnitude $r = |\lambda|$ and angle θ such that $\tan \theta = (\Im \lambda) / (\Re \lambda)$. For example, $1 - i$ has magnitude $r = |1 - i| = \sqrt{1^2 + (-1)^2} = \sqrt{2}$ and angle $\theta = -\frac{\pi}{4}$ since $\tan(-\frac{\pi}{4}) = -1/1$. Question: how does this help understand the solution which has λ_j^t in it? Answer: De Moivre's theorem says that if $\lambda = r[\cos \theta + i \sin \theta]$, then $\lambda^t = r^t[\cos(\theta t) + i \sin(\theta t)]$. Since the magnitude $|\cos(\theta t) + i \sin(\theta t)| = \sqrt{\cos^2(\theta t) + \sin^2(\theta t)} = \sqrt{1} = 1$, the magnitude $|\lambda^t| = r^t$. For example, the magnitude $|(1 - i)^2| = (\sqrt{2})^2 = 2$ which we check by computing $(1 - i)^2 = 1^2 - 2i + i^2 = -2i$ and $|-2i| = 2$.

In Example 7.1.17, the eigenvalue $\lambda_1 = 1 + i\sqrt{3}$ so its magnitude is $r_1 = |\lambda_1| = |1 + i\sqrt{3}| = \sqrt{1+3} = 2$. Hence the magnitude $|\lambda_1^t| = 2^t$ at any time step t . Similarly, the magnitude $|\lambda_2^t| = 2^t$ at any time step t . Consequently, the general solution

$$\mathbf{y}(t) = c_1 \lambda_1^t \begin{bmatrix} -i\sqrt{3} \\ 1 \end{bmatrix} + c_2 \lambda_2^t \begin{bmatrix} +i\sqrt{3} \\ 1 \end{bmatrix}$$

will grow in magnitude roughly like 2^t as both components grow like 2^t . It is a ‘rough’ growth because the components $\cos(\theta t)$ and $\sin(\theta t)$ cause ‘oscillations’ in time t . Nonetheless the overall growth like $|\lambda_1|^t = |\lambda_2|^t = 2^t$ is inexorable—and seen previously in the particular solution where we observe $\mathbf{y}(3)$ is eight times the magnitude of $\mathbf{y}(0)$.



In general, for both real or complex eigenvalues λ , a term involving the factor λ^t will, as time t increases,

- grow to infinity if $|\lambda| > 1$,
- decay to zero if $|\lambda| < 1$, and
- remain the same magnitude if $|\lambda| = 1$.

Example 7.1.18 (orangutans over many years). Extend orangutan analysis of Example 7.1.14. Use Theorem 7.1.16 to predict the population over many years: from an initial population of 30 adolescent females and 15 adult females; and from a general initial population.

Solution: Example 7.1.14 derived that the age structure population $\mathbf{y} = (y_1, y_2, y_3)$ satisfies $\mathbf{y}(t+1) = A\mathbf{y}(t)$ for matrix

$$A = \begin{bmatrix} \frac{4}{5} & 0 & \frac{1}{15} \\ \frac{1}{5} & \frac{9}{10} & 0 \\ 0 & \frac{1}{10} & \frac{14}{15} \end{bmatrix}.$$

Let's find the eigenvalues and eigenvectors of the matrix A using Matlab/Octave via



```
A=[4/5 0 1/15;1/5 9/10 0;0 1/10 14/15]
[V,D]=eig(A)
```

to find

```
V =
-0.3077+0.2952i -0.3077-0.2952i 0.2673+0.0000i
 0.7385+0.0000i 0.7385+0.0000i 0.5345+0.0000i
-0.4308-0.2952i -0.4308+0.2952i 0.8018+0.0000i
D =
 0.8167+0.0799i 0.0000+0.0000i 0.0000+0.0000i
 0.0000+0.0000i 0.8167-0.0799i 0.0000+0.0000i
 0.0000+0.0000i 0.0000+0.0000i 1.0000+0.0000i
```

Evidently there is one real eigenvalue of $\lambda_3 = 1$ and two complex conjugate eigenvalues $\lambda_{1,2} = 0.8167 \pm i0.0799$. Corresponding eigenvectors are the columns \mathbf{v}_j of V . Thus a solution for the orangutan population is

$$\mathbf{y}(t) = c_1 \lambda_1^t \mathbf{v}_1 + c_2 \lambda_2^t \mathbf{v}_2 + c_3 \lambda_3^t \mathbf{v}_3.$$

- For the initial population $\mathbf{y}(0) = (0, 30, 15)$ we need to find constants $\mathbf{c} = (c_1, c_2, c_3)$ such that $V\mathbf{c} = \mathbf{y}(0)$. Solve this linear equation in Matlab/Octave with

```
y0=[0;30;15]
rcond(V)
c=V\y0
```

which gives the answer

```
ans =
 0.1963
c =
 10.1550+2.1175i
 10.1550-2.1175i
 28.0624+0.0000i
```

The `rcond` value of 0.1963 indicates that matrix V is invertible. Then the backslash operator computes the above coefficients \mathbf{c} . Via the magic of complex conjugates cancelling, the real population of orangutans is for all times predicted to be (2 d.p.)

$$\begin{aligned}\mathbf{y}(t) &= (10.16 + 2.12i)(0.82 + 0.08i)^t \begin{bmatrix} -0.31 + 0.30i \\ 0.74 \\ -0.43 - 0.30i \end{bmatrix} \\ &\quad + (10.16 - 2.12i)(0.82 + 0.08i)^t \begin{bmatrix} -0.31 - 0.30i \\ 0.74 \\ -0.43 + 0.30i \end{bmatrix} \\ &\quad + 28.06 \begin{bmatrix} 0.27 \\ 0.53 \\ 0.80 \end{bmatrix}\end{aligned}$$

since $\lambda_3^t = 1^t = 1$.

Since the magnitude $|\lambda_1| = |\lambda_2| = 0.82$ (2 d.p.), the first two terms in this expression decay to zero as time t increases. For example, $|\lambda_1^{12}| = |\lambda_2^{12}| = 0.09$. Hence the model predicts that over long times the population

$$\mathbf{y}(t) \approx 28.06 \begin{bmatrix} 0.27 \\ 0.53 \\ 0.80 \end{bmatrix} = \begin{bmatrix} 7.5 \\ 15.0 \\ 22.5 \end{bmatrix}$$

Such a static population means that the orangutans are highly sensitive to disease, or deforestation, or chance events, and so on.

- Such unfortunate sensitivity is typical for orangutans. It is not a quirk of the initial population. Recall the general prediction for the orangutans is

$$\mathbf{y}(t) = c_1 \lambda_1^t \mathbf{v}_1 + c_2 \lambda_2^t \mathbf{v}_2 + c_3 \lambda_3^t \mathbf{v}_3.$$

The initial population determines the constants \mathbf{c} . However, the long term population is always predicted to be static. The reason is that the magnitude of the eigenvalues $|\lambda_1| = |\lambda_2| = 0.82$ and so the first two terms in this general solution will in time always decay to zero. Further, the remaining third eigenvalue has magnitude $|\lambda_3| = |1| = 1$ and so the third term in the population prediction is always constant in time t . That is, over long times the population is always

$$\mathbf{y}(t) \approx c_3 \mathbf{v}_3.$$

Such a static population means that the orangutans are always sensitive to disease, or deforestation, or chance events, and so on.

Bliss et al. (2016) discuss mathematical modelling. On page 23 they comment “*Modelling (like real life) is open-ended and messy*”: in our two examples here you have to extract the important factors from many unneeded details, and use them in the context of an imperfect model. Bliss et al. (2016) [p.23] also comment that modellers “*must be making genuine choices*”: in these problems, as in all modelling, there are choices that lead to different models—we have to operate and sensibly predict with such uncertainty. Lastly, Bliss et al. (2016) recommend to “*focus on the process, not the product*”: depending upon your choices and interpretations you will develop alternative plausible models in these scenarios—it is the process of forming plausible models and interpreting the results that are important here.

Example 7.1.19 (servals grow). The serval is a member of the cat family that lives in Africa. Given next is an extract from Wikipedia of a serval’s Reproduction and Life History.



Kittens are born shortly before the peak breeding period of local rodent populations. A serval is able to give birth to multiple litters throughout the year, but commonly does so only if the earlier litters die shortly after birth. Gestation lasts from 66 to 77 days and commonly results in the birth of two kittens, although sometimes as few as one or as many as four have been recorded.

The kittens are born in dense vegetation or sheltered locations such as abandoned aardvark burrows. If such an ideal location is not available, a place beneath a shrub may be sufficient. The kittens weigh around 250 gm at birth, and are initially blind and helpless, with a coat of greyish woolly hair. They open their eyes at 9 to 13 days of age, and begin to take solid food after around a month. At around six months, they acquire their permanent canine teeth and begin to hunt for themselves; they leave their mother at about 12 months of age. They may reach sexual maturity from 12 to 25 months of age.

Life expectancy is about 10 years in the wild.

From the information in this extract, create a plausible, age structured, population model of the serval: give reasons for estimates of the coefficients in the model. Choose three age categories of kittens, juveniles, sexually mature adults. What does the model predict over long times?

Solution: Recall we only model the number of *female* servals as females are the limiting breeders. Define

- $y_1(t)$ is the number of female kittens, less than 0.5 years old from when they “begin to hunt for themselves”;
- $y_2(t)$ is the number of female juveniles, between 0.5 years and 1.5 years which is when they “reach sexual maturity” on average;
- $y_3(t)$ is the number of female breeding adults, older than 1.5 years, and dying at about the “life expectancy” of 10 years;
- since servals transition from one age category to another in multiples of six months (0.5 years), let the unit of time be six months, equivalently a half-year. Consequently, time $t + 1$ is the time a half-year later than time t .

Modelling of the servals leads to the following equations.

- Kittens are commonly born once a year to each female, and the common litter size is two, so *on average* one female kitten is born per year per adult, that is, on average $\frac{1}{2}$ female kitten is born per half-year per adult female: Also all kittens age to juveniles after 0.5 years, so none remain as kittens. Hence the kitten equation is $y_1(t + 1) = 0y_1(t) + \frac{1}{2}y_3(t)$.
- Juveniles mature from the kittens, and age to an adult after about one year: that is, on average half of them become adults every half-year, and half remain juveniles. So the juvenile equation is $y_2(t + 1) = y_1(t) + \frac{1}{2}y_2(t)$.
- Adults mature from the juveniles, and die after about 8.5 years which is about a rate $1/8.5$ per year: that is, a rate of $\frac{1}{17}$ per half-year leaving $\frac{16}{17}$ of them to live into the next half-year. So the adult equation is $y_3(t + 1) = \frac{1}{2}y_2(t) + \frac{16}{17}y_3(t)$.

Bring these equations together, the age structure population $\mathbf{y} = (y_1, y_2, y_3)$ satisfies $\mathbf{y}(t + 1) = A\mathbf{y}(t)$ for matrix

$$A = \begin{bmatrix} 0 & 0 & \frac{1}{2} \\ 1 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & \frac{16}{17} \end{bmatrix}.$$

Find the eigenvalues and eigenvectors of the matrix A using Matlab/Octave via

```
A=[0 0 1/2;1 1/2 0;0 1/2 16/17]
[V,D]=eig(A)
```

to find

```
V =
-0.3066+0.3439i -0.3066-0.3439i 0.3352+0.0000i
 0.7838+0.0000i 0.7838+0.0000i 0.4633+0.0000i
-0.3684-0.1942i -0.3684+0.1942i 0.8203+0.0000i
D =
```



$$\begin{array}{ccc} 0.1088+0.4387i & 0.0000+0.0000i & 0.0000+0.0000i \\ 0.0000+0.0000i & 0.1088-0.4387i & 0.0000+0.0000i \\ 0.0000+0.0000i & 0.0000+0.0000i & 1.2236+0.0000i \end{array}$$

Evidently there is one real eigenvalue of $\lambda_3 = 1.2236$ and two complex conjugate eigenvalues $\lambda_{1,2} = 0.1088 \pm i0.4387$. Corresponding eigenvectors are the columns \mathbf{v}_j of V . Thus a general solution for the serval population is (Theorem 7.1.16)

$$\mathbf{y}(t) = c_1 \lambda_1^t \mathbf{v}_1 + c_2 \lambda_2^t \mathbf{v}_2 + c_3 \lambda_3^t \mathbf{v}_3.$$

In this general solution, the first two terms will decay in time to zero. The reason is that the magnitudes $|\lambda_1| = |\lambda_2| = |0.1088 \pm i0.4387| = \sqrt{0.1088^2 + 0.4387^2} = 0.4520$, and since this magnitude is less than one, then λ_1^t and λ_2^t will decay to zero with increasing time t . However, the third term increases in time as $\lambda_3 = 1.2236 > 1$. The model predicts the serval population increases by about 22% per half-year (about 50% per year).

Predation, disease, and food shortages are just some processes not included in this model which act to limit the serval's population in ways not included in this model. ■

Crucial in this section—so that we find a solution for all initial states—is that the matrix of eigenvectors is invertible. The next Section 7.2 relates the invertibility of a matrix of eigenvectors to the concept of ‘linear independence’ defined in that section (Theorem 7.2.31).

7.1.4 Extension: Connect SVDs to eigen-problems

This optional section connects the SVD of a general matrix to a symmetric eigen-problem, in principle.

Recall that Chapter 4 starts by illustrating the close connection between the SVD of a symmetric matrix and the eigenvalues and eigenvectors of that symmetric matrix. This subsection establishes that an SVD of a general matrix is closely connected to the eigenvalues and eigenvectors of a specific matrix of double the size. The connection depends upon determinants and solving linear systems and so, in principle, is an approach to compute an SVD distinct from the inductive maximisation.

Example 7.1.20. Compute the eigenvalues and eigenvectors of the (symmetric) matrix

$$B = \begin{bmatrix} 0 & 0 & 10 & 2 \\ 0 & 0 & 5 & 11 \\ 10 & 5 & 0 & 0 \\ 2 & 11 & 0 & 0 \end{bmatrix}.$$

Solution: In Matlab/Octave execute



```
B=[0 0 10 2
   0 0 5 11
   10 5 0 0
   2 11 0 0]
[V,D]=eig(B)
```

and obtain (2 d.p.)

```
V =
  0.42   0.57   0.57  -0.42
  0.57  -0.42  -0.42  -0.57
 -0.50  -0.50   0.50  -0.50
 -0.50   0.50  -0.50  -0.50
D =
 -14.14      0      0      0
      0   -7.07      0      0
      0      0    7.07      0
      0      0      0  14.14
```

The eigenvalues are the pairs ± 7.07 and ± 14.14 , with corresponding eigenvector pairs $(0.57, -0.42, \pm 0.50, \mp 0.50)$ and $(\mp 0.42, \mp 0.57, -0.50, -0.50)$.

These eigenvalues/vectors occur in \pm pairs because this matrix has the form

$$B = \begin{bmatrix} O_2 & A \\ A^T & O_2 \end{bmatrix}, \quad \text{here for matrix } A = \begin{bmatrix} 10 & 2 \\ 5 & 11 \end{bmatrix}$$

from Example 3.3.2. Observe that not only are the eigenvectors orthogonal, because B is symmetric, but also the two parts of the eigenvectors are orthogonal: $(0.57, -0.42)$ from the first pair is orthogonal to $(-0.42, -0.57)$ from the second pair; and $(0.50, -0.50)$ from the first pair is orthogonal to $(-0.50, -0.50)$ from the second pair. The next Theorem 7.1.21 establishes how these properties relate to an SVD for the matrix A . ■

Procedure 7.1.8 computes eigenvalues and eigenvectors by hand (in principle no matter how large the matrix). The procedure is independent of the svd. We now invoke this procedure to find another method to find an SVD distinct from the inductive maximisation of the proof in Section 3.3.3. The following theorem is a step towards an efficient numerical computation of an SVD (Trefethen & Bau 1997, p.234).

Theorem 7.1.21 (svd as an eigenproblem). *Given an $m \times n$ matrix A , the singular values of A are the non-negative eigenvalues of the $(m+n) \times (m+n)$ matrix $B = \begin{bmatrix} O_m & A \\ A^T & O_n \end{bmatrix}$. A corresponding eigenvector $w \in \mathbb{R}^{m+n}$ of B gives corresponding singular vectors of A , namely $w = (u, v)$ for singular vectors $u \in \mathbb{R}^m$ and $v \in \mathbb{R}^n$.*

Proof. First prove the SVD of A gives eigenvalues/vectors of B , and second the converse. For simplicity this proof considers only the case $m = n$; the case $m \neq n$ is similar but the more intricate details are of little interest.

First, let $n \times n$ matrix $A = USV^T$ be an SVD (Theorem 3.3.5) for $n \times n$ orthogonal $U = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_n]$, orthogonal $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$, and diagonal $S = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n)$. Post-multiply the SVD by orthogonal V gives $AV = US$. Also, transpose the SVD to give $A^T = (USV^T)^T = VS^T U^T = VSU^T$ and post-multiply by orthogonal U to give $A^T U = VS$. Now consider each of the \pm cases of

$$B \begin{bmatrix} U \\ \pm V \end{bmatrix} = \begin{bmatrix} O_n & A \\ A^T & O_n \end{bmatrix} \begin{bmatrix} U \\ \pm V \end{bmatrix} = \begin{bmatrix} \pm AV \\ A^T U \end{bmatrix} = \begin{bmatrix} \pm US \\ VS \end{bmatrix} = \begin{bmatrix} U \\ \pm V \end{bmatrix} (\pm S).$$

Letting $\mathbf{w} = (\mathbf{u}_j, \pm \mathbf{v}_j) \neq \mathbf{0}$, the j th column of the above equation is $B\mathbf{w} = \pm \sigma_j \mathbf{w}$ and hence $(\mathbf{u}_j, \pm \mathbf{v}_j)$ is an eigenvector of B corresponding to eigenvalue $\pm \sigma_j$, respectively, for $j = 1, \dots, n$ and each of the \pm cases.

Second, let $\mathbf{w} \in \mathbb{R}^{2n}$ be an eigenvector of B corresponding to an eigenvalue λ (real as B is symmetric, Theorem 4.2.8) and normalised so that $|\mathbf{w}| = 1$. Partition $\mathbf{w} = (\mathbf{u}, \mathbf{v})$ for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$. Then the fundamental eigen-problem $B\mathbf{w} = \lambda\mathbf{w}$ (Definition 4.1.1) partitions into

$$\begin{bmatrix} O & A \\ A^T & O \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} \iff A\mathbf{v} = \lambda\mathbf{u} \text{ and } A^T\mathbf{u} = \lambda\mathbf{v}. \quad (7.3)$$

Now eigenvalues and eigenvectors of B come in pairs as the eigenvalue $\lambda' = -\lambda$ corresponds to eigenvector $\mathbf{w}' = (\mathbf{u}, -\mathbf{v})$: substitute into (7.3) to check, $A(-\mathbf{v}) = -A\mathbf{v} = -\lambda\mathbf{u} = \lambda'\mathbf{u}$ and $A^T\mathbf{u} = \lambda\mathbf{v} = (-\lambda)(-\mathbf{v}) = \lambda'(-\mathbf{v})$. If the eigenvalue $\lambda \neq 0$, then corresponding to the distinct eigenvalues $\pm\lambda$ the eigenvectors $(\mathbf{u}, \pm \mathbf{v})$ of symmetric B are orthogonal (Theorem 4.2.10) and so $0 = (\mathbf{u}, \mathbf{v}) \cdot (\mathbf{u}, -\mathbf{v}) = \mathbf{u} \cdot \mathbf{u} - \mathbf{v} \cdot \mathbf{v} = |\mathbf{u}|^2 - |\mathbf{v}|^2$ and hence $|\mathbf{u}| = |\mathbf{v}| = \frac{1}{\sqrt{2}}$ as $|\mathbf{w}|^2 = |\mathbf{u}|^2 + |\mathbf{v}|^2 = 1$ (see Example 7.1.20). Further, for any two distinct eigenvalues $|\lambda_i| \neq |\lambda_j| \neq 0$, the eigenvectors $(\mathbf{u}_i, \mathbf{v}_i)$ and $(\mathbf{u}_j, \pm \mathbf{u}_j)$ are also orthogonal, hence $0 = (\mathbf{u}_i, \mathbf{v}_i) \cdot (\mathbf{u}_j, \pm \mathbf{u}_j) = \mathbf{u}_i \cdot \mathbf{u}_j \pm \mathbf{v}_i \cdot \mathbf{v}_j$. Taking the sum and the difference of the \pm cases of this equation gives $\mathbf{u}_i \cdot \mathbf{u}_j = \mathbf{v}_i \cdot \mathbf{v}_j = 0$; that is, (Definition 3.2.31) $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ and $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ are both orthogonal sets. In the cases when matrix B has $2n$ distinct non-zero eigenvalues, choose $(\mathbf{u}_j, \mathbf{v}_j)$ to be a normalised eigenvector corresponding to positive eigenvalue λ_j , $j = 1, \dots, n$. Upon setting $U = \sqrt{2} [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_n]$, $V = \sqrt{2} [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_n]$, and $S = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, equation (7.3) gives the columns of $AV = US$ which since the columns of U and V are orthonormal gives the SVD $A = USV^T$ for singular vectors $\mathbf{u}_j, \mathbf{v}_j \in \mathbb{R}^n$ and singular values $\lambda_j (> 0)$.

Extensions of this proof cater for the case when zero is an eigenvalue and/or eigenvalues are repeated and/or the dimensions $m \neq n$. \square

7.1.5 Exercises

Exercise 7.1.1. For each of the following list of numbers, could the numbers be all the eigenvalues of a 4×4 matrix? Justify your answer.

(a) $-1.2, -0.6, 0.2, -1.4$ (b) $\pm 2, -3$

(c) $0, 3, \pm 5, 8$ (d) $0, 3 \pm 5i, 8$

(e) $-1.4 \pm \sqrt{7}i, -4, 3 \pm 2i$

Exercise 7.1.2. For each of the following matrices, determine the two highest order terms and the constant term in the characteristic polynomial of the matrix.

(a) $\begin{bmatrix} -7 & 1 \\ -2 & 2 \end{bmatrix}$

(b) $\begin{bmatrix} 3 & -3 \\ 6 & 2 \end{bmatrix}$

(c) $\begin{bmatrix} 0 & 0 & 2 \\ 1 & 0 & 2 \\ 4 & 2 & 0 \end{bmatrix}$

(d) $\begin{bmatrix} 3 & -3 & 6 \\ -1 & 4 & 0 \\ 0 & 0 & -4 \end{bmatrix}$

(e) $\begin{bmatrix} -1 & 0 & -6 \\ -7 & 0 & -4 \\ 0 & -6 & 0 \end{bmatrix}$

(f) $\begin{bmatrix} -1 & 0 & -2 & 0 \\ 0 & 0 & 0 & 0 \\ -3 & -1 & 0 & -2 \\ -4 & 0 & -5 & 0 \end{bmatrix}$

(g) $\begin{bmatrix} -3 & 0 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \\ -4 & -4 & 0 & 4 \end{bmatrix}$

(h) $\begin{bmatrix} 0 & 4 & 0 & 1 \\ 4 & -5 & 0 & -3 \\ 0 & -4 & 0 & 0 \\ 0 & 0 & -1 & 0 \end{bmatrix}$

Exercise 7.1.3. For each of the following characteristic polynomials, write down the size of the corresponding matrix, and the matrix's trace and determinant.

(a) $\lambda^2 + 5\lambda - 6$

(b) $\lambda^2 + 2\lambda - 10$

(c) $\lambda^2 - \lambda$

(d) $-\lambda^3 + 5\lambda^2 + 7\lambda - 20$

(e) $-\lambda^3 + 28\lambda - 199$

(f) $-\lambda^3 + 8\lambda^2 - 5\lambda$

(g) $\lambda^4 + 3\lambda^3 + 143\lambda - 56$

(h) $\lambda^4 + 5\lambda^2 - 41\lambda - 5$

Exercise 7.1.4. For each the following matrices, determine the characteristic polynomial by hand, and hence find all eigenvalues of the matrix and their multiplicity. Show your working.

$$(a) \begin{bmatrix} 0 & -3 \\ -1 & -2 \end{bmatrix}$$

$$(b) \begin{bmatrix} 0 & 5 \\ -2 & 2 \end{bmatrix}$$

$$(c) \begin{bmatrix} 0 & -3 \\ 3 & 6 \end{bmatrix}$$

$$(d) \begin{bmatrix} 3 & -4 \\ 2 & -3 \end{bmatrix}$$

$$(e) \begin{bmatrix} 4.5 & 16 \\ -1 & -3.5 \end{bmatrix}$$

$$(f) \begin{bmatrix} -1 & 1 & -1 \\ -6 & -6 & 2 \\ -5 & -3 & -3 \end{bmatrix}$$

$$(g) \begin{bmatrix} -2 & -5 & -1 \\ 0 & 3 & 1 \\ 0 & -6 & -2 \end{bmatrix}$$

$$(h) \begin{bmatrix} 9 & 3 & 0 \\ -12 & -3 & 0 \\ 2 & -4 & -2 \end{bmatrix}$$

$$(i) \begin{bmatrix} -14 & 24 & 52 \\ -4 & 8 & 18 \\ -2 & 3 & 6 \end{bmatrix}$$

$$(j) \begin{bmatrix} -1 & -2 & 2 \\ 7 & 18 & -12 \\ 7 & 17 & -11 \end{bmatrix}$$

$$(k) \begin{bmatrix} -10 & -10 & -16 \\ 4 & 4 & 6 \\ 3 & 3 & 5 \end{bmatrix}$$

$$(l) \begin{bmatrix} 1 & -15 & 7 \\ -1 & -1 & -1 \\ -5 & -15 & -1 \end{bmatrix}$$



Exercise 7.1.5. For each the following matrices, use Matlab/Octave to find all eigenvalues of the matrix and their multiplicity.

$$(a) \begin{bmatrix} -2.7 & 0 & 1.6 \\ 6.3 & -1 & -27.8 \\ -0.1 & 0 & -1.9 \end{bmatrix}$$

$$(b) \begin{bmatrix} 4.9 & -8.1 & 5.4 \\ 8 & -11.2 & 8 \\ 3.9 & -3.9 & 3.4 \end{bmatrix}$$

$$(c) \begin{bmatrix} -6.7 & -0.6 & -6.6 & 3.6 \\ 3 & 0.1 & 3 & -2 \\ 2.8 & 0.6 & 2.7 & -1.6 \\ -6 & 0 & -6 & 3.1 \end{bmatrix}$$

$$(d) \begin{bmatrix} 11 & 17.9 & -33.4 & 46.4 \\ 1.2 & 0.9 & 2.8 & 2.2 \\ -12.8 & -21 & 37.2 & -54.8 \\ -12.3 & -19.7 & 33.7 & -51.3 \end{bmatrix}$$

$$(e) \begin{bmatrix} 0.6 & 0 & 0 & 0 \\ 9.6 & -1 & 33.6 & 17.6 \\ -9.6 & 1.6 & -33 & -17.6 \\ 19.2 & -3.2 & 67.2 & 35.8 \end{bmatrix}$$





$$(f) \begin{bmatrix} 52.8 & -93.3 & 73.4 & 57.1 & 104 \\ 18.3 & -30.7 & 28.5 & 20.6 & 36.6 \\ -20 & 34.9 & -29.4 & -22.3 & -40 \\ 18.2 & -31.8 & 27.4 & 20.8 & 36.4 \\ -6.8 & 13.6 & -6.8 & -6.8 & -12.8 \end{bmatrix}$$



$$(g) \begin{bmatrix} -356.4 & 264.6 & -880.8 & -689.9 & 455.5 \\ -58.5 & 41.7 & -142.8 & -111.4 & 77 \\ 157 & -115.6 & 392.9 & 311.4 & -201 \\ -78.5 & 57.8 & -198.4 & -159.6 & 100.5 \\ -58.5 & 45.6 & -142.8 & -111.4 & 73.1 \end{bmatrix}$$



$$(h) \begin{bmatrix} 309.4 & -29.7 & 451.3 & 337.3 & 305.9 \\ -217.6 & 20.3 & -313.5 & -236.1 & -215.9 \\ 3 & 0 & 0.7 & 3 & 3 \\ -232.6 & 24.1 & -336 & -254.9 & -230.9 \\ -83.6 & 5.6 & -119.8 & -89.2 & -81.8 \end{bmatrix}$$

Exercise 7.1.6. For each of the following matrices, find by hand the eigenspace of the nominated eigenvalue. Confirm your answer with Matlab/Octave. Show your working.

$$(a) \begin{bmatrix} -12 & 10 \\ -15 & 13 \end{bmatrix}, \lambda = 3$$

$$(b) \begin{bmatrix} -1 & 9 \\ -1 & 5 \end{bmatrix}, \lambda = 2$$

$$(c) \begin{bmatrix} -1 & 0 \\ -2 & 1 \end{bmatrix}, \lambda = -1$$

$$(d) \begin{bmatrix} 11 & -4 & -12 \\ -27 & 10 & 27 \\ 19 & -7 & -20 \end{bmatrix}, \lambda = 1$$

$$(e) \begin{bmatrix} -1 & -7 & -2 \\ 8 & 14 & 2 \\ 0 & 0 & 7 \end{bmatrix}, \lambda = 7$$

$$(f) \begin{bmatrix} -12 & -82 & -17 \\ 3 & 18 & 3 \\ -6 & -26 & -1 \end{bmatrix}, \lambda = 0$$

$$(g) \begin{bmatrix} -4 & 0 & -4 \\ -2 & -4 & 0 \\ 8 & 4 & 6 \end{bmatrix}, \lambda = -2$$

$$(h) \begin{bmatrix} -3 & 2 & -6 \\ -4 & 3 & -6 \\ 2 & -1 & 4 \end{bmatrix}, \lambda = 1$$

Exercise 7.1.7. For each of the following matrices, Use Matlab/Octave to find their eigenvalues, with multiplicity, and to find eigenvectors corresponding to each eigenvalue (2 d.p.). (The next Section 7.2

discusses that for repeated eigenvalues we generally want to record the so-called ‘linearly independent’ eigenvectors.)

(a)
$$\begin{bmatrix} 14 & 3 & -6 \\ 14 & -2 & -2 \\ 31 & 4 & -11 \end{bmatrix}$$

(b)
$$\begin{bmatrix} -1 & 0 & 0 \\ 5 & 5 & 6 \\ -5 & -6 & -7 \end{bmatrix}$$

(c)
$$\begin{bmatrix} -144 & -374 & 316 & 18 \\ 21 & 45 & -42 & 0 \\ -49 & -138 & 112 & 9 \\ 134 & 336 & -286 & -13 \end{bmatrix}$$

(d)
$$\begin{bmatrix} -50 & 30 & 46 & -62 \\ 0 & 2 & 0 & 0 \\ -104 & 62 & 104 & -142 \\ -39 & 24 & 42 & -58 \end{bmatrix}$$

(e)
$$\begin{bmatrix} 4 & -11 & -12 & 9 & -19 \\ 0 & 0 & 0 & 3 & -2 \\ 4 & 16 & 19 & -12 & 23 \\ 6 & 16 & 20 & -7 & 31 \\ -1 & -3 & -3 & 3 & 1 \end{bmatrix}$$

(f)
$$\begin{bmatrix} 75 & 7 & -13 & -51 & -129 \\ 120 & 12 & -24 & -84 & -208 \\ 62 & 6 & -12 & -42 & -106 \\ -48 & -5 & 9 & 32 & 83 \\ 62 & 6 & -11 & -42 & -107 \end{bmatrix}$$

Exercise 7.1.8. Consider the following three matrices which are perturbed versions of the matrix $\begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix}$. Which perturbations show the high sensitivity of the repeated eigenvalue? Give reasons.

(a) $A = \begin{bmatrix} 2 & 1 \\ -0.0001 & 2 \end{bmatrix}$ (b) $B = \begin{bmatrix} 2 & 1.0001 \\ 0 & 2 \end{bmatrix}$

(c) $C = \begin{bmatrix} 2.0001 & 1 \\ 0 & 2 \end{bmatrix}$

Exercise 7.1.9. For each of the following matrices, use Matlab/Octave to compute the eigenvalues, and to compute the eigenvalues of matrices obtained by adding random perturbations of size 0.0001 (use `randn`). Give reasons for which eigenvalues appear sensitive and which appear to be not sensitive.

(a) $A = \begin{bmatrix} 0 & -3 \\ -1 & -2 \end{bmatrix}$

$$(b) \quad B = \begin{bmatrix} 0 & -3 \\ 3 & 6 \end{bmatrix}$$

$$(c) \quad C = \begin{bmatrix} -10 & -10 & -16 \\ 4 & 4 & 6 \\ 3 & 3 & 5 \end{bmatrix}$$

$$(d) \quad D = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}$$

$$(e) \quad E = \begin{bmatrix} -1 & 1 & -1 \\ -6 & -6 & 2 \\ -5 & -3 & -3 \end{bmatrix}$$

$$(f) \quad F = \begin{bmatrix} -6.7 & -0.6 & -6.6 & 3.6 \\ 3 & 0.1 & 3 & -2 \\ 2.8 & 0.6 & 2.7 & -1.6 \\ -6 & 0 & -6 & 3.1 \end{bmatrix}$$

$$(g) \quad G = \begin{bmatrix} 1.4 & -7.1 & -0.7 & 6.2 \\ -7.1 & -1.0 & -2.2 & -2.5 \\ -0.7 & -2.2 & -3.4 & -4.1 \\ 6.2 & -2.5 & -4.1 & -1.0 \end{bmatrix}$$

$$(h) \quad H = \begin{bmatrix} 0.6 & 0 & 0 & 0 \\ 9.6 & -1 & 33.6 & 17.6 \\ -9.6 & 1.6 & -33 & -17.6 \\ 19.2 & -3.2 & 67.2 & 35.8 \end{bmatrix}$$



Exercise 7.1.10. For each of the following matrices, predict $\mathbf{y}(1)$, $\mathbf{y}(2)$ and $\mathbf{y}(3)$, for the given initial $\mathbf{y}(0)$, and given that $\mathbf{y}(t+1) = A\mathbf{y}(t)$.

$$(a) \quad \begin{bmatrix} 3 & 0 \\ 3 & 2 \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 6 \\ -1 \end{bmatrix}$$

$$(b) \quad \begin{bmatrix} 0 & -1 \\ -4 & 2 \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 3 \\ 0 \end{bmatrix}$$

$$(c) \quad \begin{bmatrix} 26 & 21 \\ -28 & -23 \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 3 \\ -4 \end{bmatrix}$$

$$(d) \quad \begin{bmatrix} -2 & 5 \\ -2 & 4 \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$(e) \quad \begin{bmatrix} 11 & -14 & -4 \\ 7 & -10 & -2 \\ 4 & -4 & -2 \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} 0 \\ 2 \\ -2 \end{bmatrix}$$

$$(f) \quad \begin{bmatrix} 9 & 7 & -5 \\ -16 & -8 & 8 \\ 10 & 10 & -6 \end{bmatrix}, \quad \mathbf{y}(0) = \begin{bmatrix} -1 \\ 1 \\ 1 \end{bmatrix}$$

$$(g) \begin{bmatrix} 2 & -2 & 0 \\ 4 & -2 & 0 \\ 0 & -5 & -1 \end{bmatrix}, \mathbf{y}(0) = \begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix}$$

$$(h) \begin{bmatrix} -4 & 14 & 5 \\ 2 & -1 & 0 \\ -13 & 26 & 8 \end{bmatrix}, \mathbf{y}(0) = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix}$$

Exercise 7.1.11. For each of the matrices of the previous Exercise 7.1.10, find a general solution of $\mathbf{y}(t+1) = A\mathbf{y}(t)$, if possible. Then use the corresponding given initial $\mathbf{y}(0)$ to find a formula for the specific $\mathbf{y}(t)$. Finally, check that the formula reproduces the values of $\mathbf{y}(1)$, $\mathbf{y}(2)$ and $\mathbf{y}(3)$ found in Exercise 7.1.10. Show your working.

Exercise 7.1.12. From the following partial description of the Tasmanian Devil, derive a mathematical model in the form $\mathbf{y}(t+1) = A\mathbf{y}(t)$ for the age structure of the Tasmanian Devil. By finding eigenvalues and an eigenvector, predict the long-term growth of the population, and predict the long-term relative numbers of Devils of various ages.



https://en.wikipedia.org/wiki/Tasmanian_devil

Devils are not monogamous, and their reproductive process is very robust and competitive. Males fight one another for the females, and then guard their partners to prevent female infidelity. Females can ovulate three times in as many weeks during the mating season, and 80% of two-year-old females are seen to be pregnant during the annual mating season. Females average four breeding seasons in their life and give birth to 20–30 live young after three weeks' gestation. The newborn are pink, lack fur, have indistinct facial features and weigh around 0.20 g (0.0071 oz) at birth. As there are only four nipples in the pouch, competition is fierce and few newborns survive. The young grow rapidly and are ejected from the pouch after around 100 days, weighing roughly 200 g (7.1 oz). The young become independent after around nine months, so the female spends most of her year in activities related to birth and rearing.

Wikipedia, 2016

Exercise 7.1.13. From the following partial description of the elephant, derive a mathematical model in the form $\mathbf{y}(t+1) = A\mathbf{y}(t)$ for the age structure of the elephant. By finding eigenvalues and an eigenvector, predict the long-term growth of the population, and predict the long-term relative numbers of elephants of various ages.



<https://en.wikipedia.org/wiki/Elephant>

Gestation in elephants typically lasts around two years with interbirth intervals usually lasting four to five years. Births tend to take place during the wet season. Calves are born 85 cm (33 in) tall and weigh around 120 kg

(260 lb). Typically, only a single young is born, but twins sometimes occur. The relatively long pregnancy is maintained by five corpus luteums (as opposed to one in most mammals) and gives the foetus more time to develop, particularly the brain and trunk. As such, newborn elephants are precocial and quickly stand and walk to follow their mother and family herd. A new calf is usually the centre of attention for herd members. Adults and most of the other young will gather around the newborn, touching and caressing it with their trunks. For the first few days, the mother is intolerant of other herd members near her young. Alloparenting—where a calf is cared for by someone other than its mother—takes place in some family groups. Allomothers are typically two to twelve years old. When a predator is near, the family group gathers together with the calves in the centre.

For the first few days, the newborn is unsteady on its feet, and needs the support of its mother. It relies on touch, smell and hearing, as its eyesight is poor. It has little precise control over its trunk, which wiggles around and may cause it to trip. By its second week of life, the calf can walk more firmly and has more control over its trunk. After its first month, a calf can pick up, hold and put objects in its mouth, but cannot suck water through the trunk and must drink directly through the mouth. It is still dependent on its mother and keeps close to her.

For its first three months, a calf relies entirely on milk from its mother for nutrition after which it begins to forage for vegetation and can use its trunk to collect water. At the same time, improvements in lip and leg coordination occur. Calves continue to suckle at the same rate as before until their sixth month, after which they become more independent when feeding. By nine months, mouth, trunk and foot coordination is perfected. After a year, a calf's abilities to groom, drink, and feed itself are fully developed. It still needs its mother for nutrition and protection from predators for at least another year. Suckling bouts tend to last 2–4 min/hr for a calf younger than a year and it continues to suckle until it reaches three years of age or older. Suckling after two years may serve to maintain growth rate, body condition and reproductive ability. Play behaviour in calves differs between the sexes; females run or chase each other, while males play-fight. The former are sexually mature by the age of nine years while the latter become mature around 14–15 years. Adulthood starts

at about 18 years of age in both sexes. Elephants have long lifespans, reaching 60–70 years of age.

Wikipedia, 2016

Exercise 7.1.14.

From the following partial description of the giant mouse lemur, derive a mathematical model in the form $\mathbf{y}(t+1) = A\mathbf{y}(t)$ for the age structure of the giant mouse lemur. By finding eigenvalues and an eigenvector, predict the long-term growth of the population, and predict the long-term relative numbers of giant mouse lemurs of various ages.

Reproduction starts in November for Coquerel's giant mouse lemur at Kirindy Forest; the estrous cycle runs approximately 22 days, while estrus lasts only a day or less. . . .

One to three offspring (typically two) are born after 90 days of gestation, weighing approximately 12 g (0.42 oz). Because they are poorly developed, they initially remain in their mother's nest for up to three weeks, being transported by mouth between nests. Once they have grown sufficiently, typically after three weeks, the mother will park her offspring in vegetation while she forages nearby. After a month, the young begin to participate in social play and grooming with their mother, and between the first and second month, young males begin to exhibit early sexual behaviors (including mounting, neck biting, and pelvic thrusting). By the third month, the young forage independently, though they maintain vocal contact with their mother and use a small part of her range.

Females start reproducing after ten months, while males develop functional testicles by their second mating season. Testicle size in the northern giant mouse lemur does not appear to fluctuate by season, and is so large relative to the animal's body mass that it is the highest among all primates. This emphasis on sperm production in males, as well as the use of copulatory plugs, suggests a mating system best described as polygynandrous where males use scramble competition (roaming widely to find many females). In contrast, male Coquerel's giant mouse lemurs appear to fight for access to females (contest competition) during their breeding season. Males disperse from their natal range, and the age at which they leave varies from two years to several. Females reproduce every year, although postpartum estrus has been observed in captivity. In the wild, the lifespan of giant mouse lemurs is thought to rarely ex-



https://en.wikipedia.org/wiki/Giant_mouse_lemur

ceed five or six years

Wikipedia, 2016

Exercise 7.1.15. From the following partial description of the dolphin (Indo-Pacific bottlenose dolphin), derive a mathematical model in the form $\mathbf{y}(t+1) = A\mathbf{y}(t)$ for the age structure of the dolphin. (Assume only one calf is born at a time.) By finding eigenvalues and an eigenvector, predict the long-term growth of the population, and predict the long-term relative numbers of dolphins of various ages.



https://en.wikipedia.org/wiki/Indo-Pacific_bottlenose_dolphin

Indo-Pacific bottlenose dolphins live in groups that can number in the hundreds, but groups of five to 15 dolphins are most common. In some parts of their range, they associate with the common bottlenose dolphin and other dolphin species, such as the humpback dolphin.

The peak mating and calving seasons are in the spring and summer, although mating and calving occur throughout the year in some regions. Gestation period is about 12 months. Calves are between 0.84 and 1.5 metres (2.8 and 4.9 ft) long, and weigh between 9 and 21 kilograms (20 and 46 lb). The calves are weaned between 1.5 and two years, but can remain with their mothers for up to five years. The interbirth interval for females is typically four to six years.

In some parts of its range, this dolphin is subject to predation by sharks; its life span is more than 40 years.

Wikipedia, 2016

Exercise 7.1.16. You are given that a mathematical model of the age structure of some animal population is

$$\begin{aligned} y_1(t+1) &= 0.5y_1(t) + y_3(t), \\ y_2(t+1) &= 0.5y_1(t) + 0.7y_2(t), \\ y_3(t+1) &= 0.3y_2(t) + 0.9y_3(t). \end{aligned}$$

Invent an animal species, and time scale, and create details of a plausible scenario for the breeding and life cycle of the species that could lead to this mathematical model. Write a coherent paragraph about the breeding and life cycle of the species with enough information that someone could deduce this mathematical model from your description. Be creative.

Exercise 7.1.17. For each of the following matrices, say A for instance, find by hand calculation the eigenvalues and eigenvectors of the larger matrix $\begin{bmatrix} O & A \\ A^T & O \end{bmatrix}$. Show your working. Relate these to an SVD of the matrix A .

(a) $A = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$

(b) $B = [-5 \ 12]$

(c) $C = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}$

(d) $D = \begin{bmatrix} 0 & 1 \\ -4 & 0 \end{bmatrix}$

VO-1b

7.2 Linear independent vectors may form a basis

Section Contents

7.2.1	Linearly (in)dependent sets	606
7.2.2	Form a basis for subspaces	616
	Revisit unique solutions	631
7.2.3	Exercises	631

In Chapter 4 on symmetric matrices, the eigenvectors from distinct eigenvalues are orthogonal. For general matrices the eigenvectors are not orthogonal—as introduced at the start of this Chapter 7. But the orthogonal property is extremely useful. Question: is there a similarly useful analogue of orthogonality? Answer: yes; we now extend “orthogonal” with the more general concept of “linear independence” which for problems with general matrices replaces orthonormality.

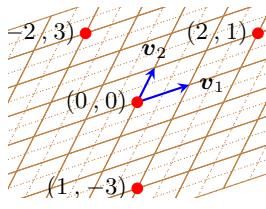
One reason that orthogonal vectors are useful is that they can form an orthonormal basis and hence act as the unit vectors of an orthogonal coordinate systems. The concept of linear independence is closely connected to coordinate systems which are not orthogonal.

Subspace coordinate systems In any given problem we want two things from a general solution:

- firstly, the general solution must encompass every possibility (the solution must span the possibilities); and
- secondly, each possible solution should have a unique algebraic description in the general solution.

For an example of the need for a unique algebraic form, let’s suppose we wanted to find solutions to the differential equation $d^2y/dt^2 - y = 0$. You might find $y = 3e^x + 2e^{-x}$, whereas I find $y = 5 \cosh x + \sinh x$, and a friend finds $y = e^x + 4 \cosh x$. By looking at these disparate algebraic forms it is apparent that we all disagree. Should we all go and search for errors in the solution process? No. The reason is that all these solutions are the same. The apparent differences arise only because you choose exponentials to represent the solution, whereas I choose hyperbolic functions: the solutions are the same, it is only the algebraic representation that appears different. In general, when we cannot immediately distinguish identical solutions, all algebraic manipulation becomes immensely more difficult due to algebraic ambiguity. To avoid such immense difficulties, in both calculus and linear algebra, we need to introduce the concept of linear independence.

Linear independence empowers us, often implicitly, to use a non-orthogonal coordinate system in a subspace. We replace the or-



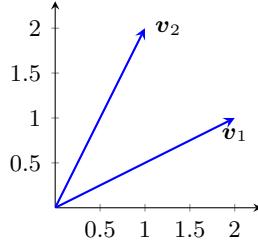
thonormal standard unit vectors by any suitable set of basis vectors. For example, in the plane any two vectors at an angle to each other suffice to be able to describe uniquely every vector (point) in the plane. As illustrated in the margin, every point in the plane (end point of a vector) is a unique linear combination of the two basis vectors v_1 and v_2 . Such a pair of basis vectors, termed linearly independent, avoid the immense difficulties of algebraic ambiguity.

7.2.1 Linearly (in)dependent sets

This section defines linear (in)dependence, and then relates the concept to homogeneous linear equations, orthogonality, and sets of eigenvectors.

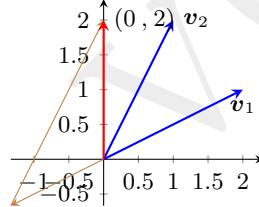
Example 7.2.1 (2D non-orthogonal coordinates). Show that every vector in the plane \mathbb{R}^2 can be written uniquely as a linear combination of the two vectors $v_1 = (2, 1)$ and $v_2 = (1, 2)$ that are shown in the margin.

Solution: Let's start with some specific example vectors.



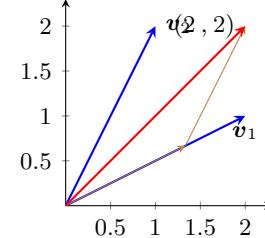
(a) The vector $(0, 2)$ may be written as the linear combination

$$(0, 2) = -\frac{2}{3}v_1 + \frac{4}{3}v_2 \text{ as shown.}$$



(b) The vector $(2, 2)$ may be written as the linear combination

$$(2, 2) = \frac{2}{3}v_1 + \frac{2}{3}v_2 \text{ as shown.}$$



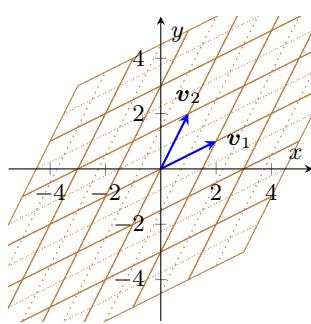
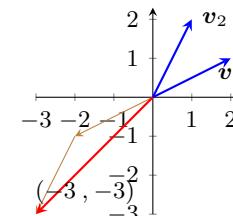
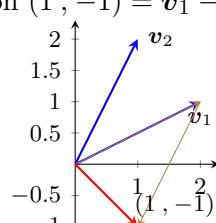
(c) The vector $(1, -1)$ may be written as the linear combination

$$(1, -1) = v_1 - v_2$$

as shown.

(d) The vector $(-3, -3)$ may be written as the linear combination

$$(-3, -3) = -v_1 - v_2.$$



Now proceed to consider a general vector (x, y) and seek it as a linear combination of v_1 and v_2 , namely $(x, y) = c_1v_1 + c_2v_2$. That is, let's write each and every point in the plane as a linear combination of v_1 and v_2 as illustrated in the margin. Rewrite the

equation in matrix-vector form as

$$[\mathbf{v}_1 \ \mathbf{v}_2] \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} = \begin{bmatrix} x \\ y \end{bmatrix}, \quad \text{that is, } V\mathbf{c} = \begin{bmatrix} x \\ y \end{bmatrix} \quad \text{for } V = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}.$$

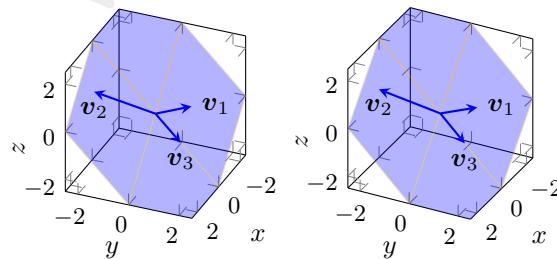
For any given (x, y) , $V\mathbf{c} = (x, y)$ is a system of linear equations for the coefficients \mathbf{c} . Theorem 3.4.35 asserts the system has a unique solution \mathbf{c} if and only iff the matrix V is invertible. Here the unique solution is then that the vector of coefficients

$$\mathbf{c} = V^{-1} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}.$$

Equivalently, Theorem 3.4.35c asserts the system has a unique solution \mathbf{c} —unique coefficients \mathbf{c} —if and only iff the homogeneous system $V\mathbf{c} = \mathbf{0}$ has only the zero solution $\mathbf{c} = \mathbf{0}$. It is this last statement that leads to the upcoming Definition 7.2.3 of linear independence. ■

Example 7.2.2 (3D failure). Show that vectors in \mathbb{R}^3 are not written uniquely as a linear combination of $\mathbf{v}_1 = (-1, 1, 0)$, $\mathbf{v}_2 = (1, -2, 1)$ and $\mathbf{v}_3 = (0, 1, -1)$.

One reason for the failure is that these three vectors only span a plane, as shown below in stereo. The solution here looks at the different issue of unique representation.



Solution: As one example, consider the vector $(1, 0, -1)$:

$$\begin{aligned} (1, 0, -1) &= -1\mathbf{v}_1 + 0\mathbf{v}_2 + 1\mathbf{v}_3; \\ (1, 0, -1) &= 1\mathbf{v}_1 + 2\mathbf{v}_2 + 3\mathbf{v}_3; \\ (1, 0, -1) &= -2\mathbf{v}_1 - 1\mathbf{v}_2 + 0\mathbf{v}_3; \\ (1, 0, -1) &= (-1+t)\mathbf{v}_1 + t\mathbf{v}_2 + (1+t)\mathbf{v}_3, \quad \text{for any } t. \end{aligned}$$

This last combination shows there are an infinite number of ways to write $(1, 0, -1)$ as a linear combination of \mathbf{v}_1 , \mathbf{v}_2 and \mathbf{v}_3 . Such an infinity of linear combinations means that \mathbf{v}_1 , \mathbf{v}_2 and \mathbf{v}_3 do not form a useful ‘coordinate system’ because we cannot easily distinguish between the different combinations all giving the same vector. The reason for the infinity of combinations is that there

is a nontrivial linear combination of \mathbf{v}_1 , \mathbf{v}_2 and \mathbf{v}_3 which is zero, namely $\mathbf{v}_1 + \mathbf{v}_2 + \mathbf{v}_3 = \mathbf{0}$. It is this last statement that leads to the Definition 7.2.3 of linear dependence.

■

Definition 7.2.3. A set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is **linearly dependent** if there are scalars c_1, c_2, \dots, c_k , at least one of which is nonzero, such that $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k = \mathbf{0}$. A set of vectors that is not linearly dependent is called **linearly independent** (characterised by only the linear combination with $c_1 = c_2 = \dots = c_k = 0$ gives the zero vector).

When reading the terms “linearly in/dependent” be very careful: it is all too easy to misread the presence or absence of the crucial “in” syllable. The presence or absence of the “in” syllable makes all the difference between the property and its opposite.

Example 7.2.4. Are the following sets of vectors linearly dependent or linearly independent. Give reasons.

(a) $\{(-1, 1, 0), (1, -2, 1), (0, 1, -1)\}$

Solution: The set is linearly dependent as the linear combination $(-1, 1, 0) + (1, -2, 1) + (0, 1, -1) = (0, 0, 0)$.

(b) $\{(2, 1), (1, 2)\}$

Solution: The set is linearly independent because the linear combination equation $c_1(2, 1) + c_2(1, 2) = (0, 0)$ is equivalent to the homogeneous matrix-vector system $\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \mathbf{c} = \mathbf{0}$ which has *only* the zero solution $\mathbf{c} = \mathbf{0}$.

(c) $\{(-2, 4, 1, -1, 0)\}$

Solution: This set of one vector in \mathbb{R}^5 is linearly independent as $c_1(-2, 4, 1, -1, 0) = \mathbf{0}$ can only be satisfied with $c_1 = 0$.

Indeed, any one non-zero vector \mathbf{v} in \mathbb{R}^n forms a linearly independent set, $\{\mathbf{v}\}$, for the same reason.

(d) $\{(2, 1), (0, 0)\}$

Solution: The set is linearly dependent because the linear combination $0(2, 1) + c_2(0, 0) = (0, 0)$ for any non-zero c_2 .

(e) $\{\mathbf{0}, \mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_k\}$

Solution: Any set that includes the zero vector is linearly dependent as $c_1\mathbf{0} + 0\mathbf{v}_2 + \dots + 0\mathbf{v}_k = \mathbf{0}$ for any non-zero c_1 .

(f) $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$, the set of standard unit vectors in \mathbb{R}^3 .

Solution: This set is linearly independent as $c_1\mathbf{e}_1 + c_2\mathbf{e}_2 + c_3\mathbf{e}_3 = (c_1, c_2, c_3) = \mathbf{0}$ only when all three components are zero, $c_1 = c_2 = c_3 = 0$.

$$(g) \left\{ \left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3} \right), \left(\frac{2}{3}, \frac{1}{3}, -\frac{2}{3} \right) \right\}$$

Solution: This set is linearly independent. Seek some linear combination $c_1\left(\frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right) + c_2\left(\frac{2}{3}, \frac{1}{3}, -\frac{2}{3}\right) = \mathbf{0}$. Take the dot product of both sides of this equation with $(1, 2, 2)$:

$$\begin{aligned} c_1 \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \\ \frac{2}{3} \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} + c_2 \begin{bmatrix} \frac{2}{3} \\ \frac{1}{3} \\ -\frac{2}{3} \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} &= \mathbf{0} \cdot \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \\ \implies c_1 3 + c_2 0 &= 0 \\ \implies c_1 &= 0. \end{aligned}$$

Similarly, take the dot product with $(2, 1, -2)$:

$$\begin{aligned} c_1 \begin{bmatrix} \frac{1}{3} \\ \frac{2}{3} \\ \frac{2}{3} \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix} + c_2 \begin{bmatrix} \frac{2}{3} \\ \frac{1}{3} \\ -\frac{2}{3} \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix} &= \mathbf{0} \cdot \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix} \\ \implies c_1 0 + c_2 3 &= 0 \\ \implies c_2 &= 0. \end{aligned}$$

Hence $c_1 = c_2 = 0$ and so the vectors are linearly independent.

These last two cases generalise to the next Theorem 7.2.6 about the linear independence of any orthonormal set of vectors. ■

Example 7.2.5 (calculus extension). In calculus the notion of a function corresponds precisely to the notion of a vector in our linear algebra. For the purposes of this example, consider ‘vector’ and ‘function’ to be synonymous, and that ‘all components’ and ‘all x ’ are synonymous. Show that the set $\{e^x, e^{-x}, \cosh x, \sinh x\}$ is linearly dependent. What is a subset that is linearly independent?

Solution: The definition of the hyperbolic functions, namely that $\cosh x = (e^x + e^{-x})/2$ and $\sinh x = (e^x - e^{-x})/2$, immediately give two nontrivial linear combinations that are zero for all x , namely $2\cosh x - e^x - e^{-x} = 0$ and $2\sinh x - e^x + e^{-x} = 0$ for all x . Either one of these implies the set $\{e^x, e^{-x}, \cosh x, \sinh x\}$ is linearly dependent.

Because e^x and e^{-x} are not proportional to each other, there is no linear combination which is zero for all x , and hence the set $\{e^x, e^{-x}\}$ is linearly independent (as are any other pairs of the four functions).

Theorem 7.2.6. Any orthonormal set (Definition 3.2.31) of vectors is linearly independent.

Proof. Let $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ be an orthonormal set of vectors in \mathbb{R}^n . Let's find all possible scalars c_1, c_2, \dots, c_k such that $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k = \mathbf{0}$. Taking the dot product of this equation with \mathbf{v}_1 implies $c_1\mathbf{v}_1 \cdot \mathbf{v}_1 + c_2\mathbf{v}_2 \cdot \mathbf{v}_1 + \dots + c_k\mathbf{v}_k \cdot \mathbf{v}_1 = \mathbf{0} \cdot \mathbf{v}_1$; by orthonormality this identity becomes $c_11 + c_20 + \dots + c_k0 = 0$; that is, $c_1 = 0$. Similarly, taking the dot product with \mathbf{v}_2 implies $c_1\mathbf{v}_1 \cdot \mathbf{v}_2 + c_2\mathbf{v}_2 \cdot \mathbf{v}_2 + \dots + c_k\mathbf{v}_k \cdot \mathbf{v}_2 = \mathbf{0} \cdot \mathbf{v}_2$; by orthonormality this identity becomes $c_10 + c_21 + \dots + c_k0 = 0$; that is, $c_2 = 0$. And so on for all vectors in the set, implying the coefficients $c_1 = c_2 = \dots = c_k = 0$ is the only possibility. By Definition 7.2.3, the orthonormal set must be linearly independent. \square

In contrast to orthonormal vectors which are always linearly independent, a set of two vectors proportional to each other is always linearly dependent as seen in the following examples. This linear dependence of proportional vectors then generalises in the next Theorem 7.2.8.

Example 7.2.7. Show the following sets are linearly dependent.

(a) $\{(1, 2), (3, 6)\}$

Solution: Since $(3, 6) = 3(1, 2)$ then the linear combination $1(3, 6) - 3(1, 2) = \mathbf{0}$ and the set is linearly dependent.

(b) $\{(2.2, -2.1, 0, 1.5), (-8.8, 8.4, 0, -6)\}$

Solution: Since $(-8.8, 8.4, 0, -6) = -4(2.2, -2.1, 0, 1.5)$ then the linear combination

$$(-8.8, 8.4, 0, -6) + 4(2.2, -2.1, 0, 1.5) = \mathbf{0},$$

and so the set is linearly dependent. \blacksquare

Theorem 7.2.8. A set of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly dependent if and only if at least one of the vectors can be expressed as a linear combination of the other vectors. In particular, a set of two vectors $\{\mathbf{v}_1, \mathbf{v}_2\}$ is linearly dependent if and only if one of the vectors is a multiple of the other.

Proof. Exercise 7.2.3 establishes the particular case of a set of two vectors.

In the general case of m vectors, first establish that if one of the vectors can be expressed as a linear combination of the others, then the set is linearly dependent. Suppose we have labelled the set

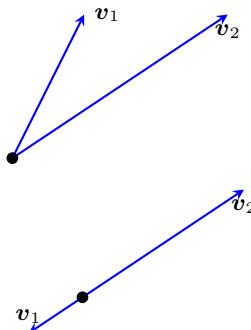
of vectors so that it is vector \mathbf{v}_1 which is a linear combination of the others; that is, $\mathbf{v}_1 = c_2\mathbf{v}_2 + c_3\mathbf{v}_3 + \cdots + c_m\mathbf{v}_m$. Rearranging, $(-1)\mathbf{v}_1 + c_2\mathbf{v}_2 + c_3\mathbf{v}_3 + \cdots + c_m\mathbf{v}_m = \mathbf{0}$; that is, there is a non-trivial (as at least $c_1 = -1 \neq 0$) linear combination of the set of vectors which is zero. Hence the set is linearly dependent.

Second, establish the converse. Given the set is linearly dependent, there exist coefficients, not all zero, such that $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \cdots + c_m\mathbf{v}_m = \mathbf{0}$. Suppose that we have labelled the vectors so that $c_1 \neq 0$. Then rearranging the equation gives $c_1\mathbf{v}_1 = -c_2\mathbf{v}_2 - c_3\mathbf{v}_3 - \cdots - c_m\mathbf{v}_m$. Divide by the non-zero c_1 to deduce $\mathbf{v}_1 = -(c_2/c_1)\mathbf{v}_2 - (c_3/c_1)\mathbf{v}_3 - \cdots - (c_m/c_1)\mathbf{v}_m$; that is, \mathbf{v}_1 is a linear combination of the other vectors. \square

Example 7.2.9. Invoke Theorem 7.2.8 to deduce whether the following sets are linearly independent or linearly dependent.

- (a) $\{(-1, 1, 0), (1, -2, 1), (0, 1, -1)\}$

Solution: Since $(1, -2, 1) = -(-1, 1, 0) - (0, 1, -1)$ the set must be linearly dependent.



- (b) The set of two vectors shown in the margin.

Solution: Since they are not proportional to each other, we cannot write either as a multiple of the other, and so the pair are linearly independent.

- (c) The set of two vectors shown in the margin.

Solution: Since they appear proportional to each other, $\mathbf{v}_2 \approx (-3)\mathbf{v}_1$, so the pair appear linearly dependent.

- (d) $\{(1, 3, 0, -1), (1, 0, -4, 2), (-2, 3, 0, -3), (0, 6, -4, -2)\}$

Solution: Notice that the last vector is the sum of the first three, $(0, 6, -4, -2) = (1, 3, 0, -1) + (1, 0, -4, 2) + (-2, 3, 0, -3)$, and so the set is linearly dependent. \blacksquare

Recall that Theorem 4.2.10 established that for every two distinct eigenvalues of a symmetric matrix A , any corresponding two eigenvectors are orthogonal. Consequently, for a symmetric A , a set of eigenvectors from distinct eigenvalues forms an orthogonal set. The following Theorem 7.2.10 generalises this property to non-symmetric matrices using the concept of linear independence.

Theorem 7.2.10. For an $n \times n$ matrix A , let $\lambda_1, \lambda_2, \dots, \lambda_m$ be distinct eigenvalues of A with corresponding eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$. The set $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly independent.

Proof. Use contradiction. Assume the set $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly dependent. Choose $k < m$ such that the set $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is linearly independent, whereas the set $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k+1}\}$ is linearly dependent. Hence there exists non-trivial coefficients such that

$$c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \cdots + c_k\mathbf{v}_k + c_{k+1}\mathbf{v}_{k+1} = \mathbf{0};$$

further, $c_{k+1} \neq 0$ as $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is linearly independent. Pre-multiply the linear combination by matrix A :

$$\begin{aligned} c_1A\mathbf{v}_1 + c_2A\mathbf{v}_2 + \cdots + c_kA\mathbf{v}_k + c_{k+1}A\mathbf{v}_{k+1} &= A\mathbf{0} \\ \implies c_1\lambda_1\mathbf{v}_1 + c_2\lambda_2\mathbf{v}_2 + \cdots + c_k\lambda_k\mathbf{v}_k + c_{k+1}\lambda_{k+1}\mathbf{v}_{k+1} &= \mathbf{0}. \end{aligned}$$

Now subtract $\lambda_{k+1} \times$ the original linear combination:

$$\begin{aligned} c_1\lambda_1\mathbf{v}_1 + c_2\lambda_2\mathbf{v}_2 + \cdots + c_k\lambda_k\mathbf{v}_k + c_{k+1}\lambda_{k+1}\mathbf{v}_{k+1} \\ - (c_1\lambda_{k+1}\mathbf{v}_1 + c_2\lambda_{k+1}\mathbf{v}_2 + \cdots + c_k\lambda_{k+1}\mathbf{v}_k + c_{k+1}\lambda_{k+1}\mathbf{v}_{k+1}) &= \mathbf{0} \\ \implies c_1(\lambda_1 - \lambda_{k+1})\mathbf{v}_1 + c_2(\lambda_2 - \lambda_{k+1})\mathbf{v}_2 + \cdots + c_k(\lambda_k - \lambda_{k+1})\mathbf{v}_k &= \mathbf{0} \\ \implies c'_1\mathbf{v}_1 + c'_2\mathbf{v}_2 + \cdots + c'_k\mathbf{v}_k &= \mathbf{0} \end{aligned}$$

for coefficients $c'_j = c_j(\lambda_j - \lambda_{k+1})$. Since all the eigenvalues are distinct, $\lambda_j - \lambda_{k+1} \neq 0$, and since the coefficients c_j are not all zero, hence c'_j are not all zero. Thus we have created a non-trivial linear combination of $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ which is zero, and so the set $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is linearly dependent. This contradiction of the choice of k proves the assumption must be wrong. Hence the set $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly independent, as required. \square

Example 7.2.11. For each of the following matrices, show the eigenvectors from distinct eigenvalues form linearly independent sets.

(a) Recall from Example 7.1.9 the matrix $B = \begin{bmatrix} -1 & 1 & -2 \\ -1 & 0 & -1 \\ 0 & -3 & 1 \end{bmatrix}$

Solution: In Matlab/Octave, executing

```
B=[-1 1 -2
   -1 0 -1
   0 -3 1]
[V,D]=eig(B)
```

gives eigenvectors and corresponding eigenvalues in

```
V =
    -0.5774      0.7071     -0.7071
    -0.5774      0.0000      0.0000
    -0.5774     -0.7071      0.7071

D =
      -2            0            0
      0            1            0
      0            0            1
```



Recognising $0.7071 = 1/\sqrt{2}$, the last two eigenvectors, $(1/\sqrt{2}, 0, -1/\sqrt{2})$ and $(-1/\sqrt{2}, 0, 1/\sqrt{2})$, form a linearly dependent set because they are proportional to each other. This linear dependence does not confound Theorem 7.2.12 because the corresponding eigenvalues are the same, not distinct, namely $\lambda = 1$. The theorem only applies to eigenvectors of distinct eigenvalues.

Here the two distinct eigenvalues are $\lambda = -2$ and $\lambda = 1$. Recognising $0.5774 = 1/\sqrt{3}$, two corresponding eigenvectors are $(-1/\sqrt{3}, -1/\sqrt{3}, -1/\sqrt{3})$ and $(1/\sqrt{2}, 0, -1/\sqrt{2})$. Because of the zero component in the second, these cannot be proportional to each other, and so the pair form a linearly independent set.

- (b) Example 7.1.10 found the eigenvalues and eigenvectors of matrix

$$A = \begin{bmatrix} 0 & 3 & 0 & 0 & 0 \\ 1 & 0 & 3 & 0 & 0 \\ 0 & 1 & 0 & 3 & 0 \\ 0 & 0 & 1 & 0 & 3 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

In Matlab/Octave execute

```
A=[0 3 0 0 0
   1 0 3 0 0
   0 1 0 3 0
   0 0 1 0 3
   0 0 0 1 0]
[V,D]=eig(A)
```

to obtain the report (2 d.p.)

```
V =
  0.62 -0.62  0.94 -0.85 -0.85
  0.62  0.62 -0.00  0.49 -0.49
  0.42 -0.42 -0.31 -0.00  0.00
  0.21  0.21 -0.00 -0.16  0.16
  0.07 -0.07  0.10  0.09  0.09
D =
  3.00      0      0      0      0
      0 -3.00      0      0      0
      0      0 -0.00      0      0
      0      0      0 -1.73      0
      0      0      0      0  1.73
```

The five eigenvalues are all distinct, so Theorem 7.2.12 asserts a set of corresponding eigenvectors will be linearly independent. The five columns of V , call them $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_5$, are a set of corresponding eigenvectors. To confirm their linear independence let's seek a linear combination being zero, that is, $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_5\mathbf{v}_5 = \mathbf{0}$. Written as a matrix-vector



system we seek $\mathbf{c} = (c_1, c_2, \dots, c_5)$ such that $V\mathbf{c} = \mathbf{0}$. Because the five singular values of square matrix V are all non-zero,⁵ obtained from `svd(V)` as

```
ans =
1.7703
1.1268
0.6542
0.3625
0.1922
```

consequently Theorem 3.4.35 asserts $V\mathbf{c} = \mathbf{0}$ has only the zero solution. Hence, by Definition 7.2.3 the set of eigenvectors in the columns of V are linearly independent.

■

This last case of Example 7.2.11b connects the concept of linear (in)dependence to the existence or otherwise of non-zero solutions to a homogeneous system of linear equations, $V\mathbf{c} = \mathbf{0}$. So does Example 7.2.4b. The great utility of this connection is that we understand a lot about homogeneous systems of linear equations. The next Theorem 7.2.12 establishes this connection in general.

Theorem 7.2.12. *Let $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m$ be vectors in \mathbb{R}^n . Let $n \times m$ matrix $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_m]$. Then $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly dependent if and only if $V\mathbf{c} = \mathbf{0}$ has a nonzero solution \mathbf{c} .*

Proof. Now $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly dependent if and only if there are scalars, not all zero, such that the equation $c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_m\mathbf{v}_m = \mathbf{0}$ holds (Definition 7.2.3). Let the vector $\mathbf{c} = (c_1, c_2, \dots, c_m)$, then this equation is equivalent to the statement $V\mathbf{c} = \mathbf{0}$. That is, if and only if $V\mathbf{c} = \mathbf{0}$ has a nonzero solution. □

Recall Theorem 1.3.20 that in \mathbb{R}^n there can be no more than n vectors in an orthogonal set of vectors. The following theorem is the generalisation: in \mathbb{R}^n there can be no more than n vectors in a linearly independent set of vectors.

Theorem 7.2.13. *Any set of m vectors in \mathbb{R}^n is linearly dependent when the number of vectors $m > n$.*

Proof. Form the m vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m \in \mathbb{R}^n$ into the $n \times m$ matrix $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_m]$. Consider the homogeneous system $V\mathbf{c} = \mathbf{0}$: as $m > n$, Theorem 2.2.25 (with the meaning of m and n swapped) asserts $V\mathbf{c} = \mathbf{0}$ has infinitely many solutions. Thus

⁵ One could alternatively compute the determinant `det(V)` = 0.09090 and because it is non-zero Theorem 7.2.31 asserts that the equation has only the solution $\mathbf{c} = \mathbf{0}$.

$V\mathbf{c} = \mathbf{0}$ has nonzero solutions, so Theorem 7.2.12 implies the set of eigenvectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ is linearly dependent. \square

Example 7.2.14. Determine if the following sets of vectors are linearly dependent or independent. Give reasons.

- (a) $\{(-1, -2), (-1, 4), (0, 5), (2, 3)\}$

Solution: As there are four vectors in \mathbb{R}^2 so Theorem 7.2.13 asserts the set is linearly dependent.

- (b) $\{(-6, -4, -1, -2), (2, 0, 1, -2), (2, -1, -1, 1)\}$

Solution: In Matlab/Octave form the matrix with these vectors as columns

```
V=[-6 2 2
   -4 0 -1
   -1 1 -1
   -2 -2 1]
svd(V)
```

and find the three singular values are all non-zero (namely 7.7568, 2.7474, and 2.2988). Hence there are no free variables when solving $V\mathbf{c} = \mathbf{0}$ (Procedure 3.3.13), and consequently there is only the unique solution $\mathbf{c} = \mathbf{0}$. By Theorem 7.2.12, the set of vectors is linearly independent.

- (c) $\{(-1, -2, 2, -1), (1, 3, 1, -1), (-2, -4, 4, -2)\}$

Solution: By inspection, the third vector is twice the first. Hence the linear combination $2(-1, -2, 2, -1) + 0(1, 3, 1, -1) - (-2, -4, 4, -2) = \mathbf{0}$ and so the set of vectors is linearly dependent.

- (d) $\{(3, 3, -1, -1), (0, -3, -1, -7), (1, 2, 0, 2)\}$

Solution: In Matlab/Octave form the matrix with these vectors as columns

```
V=[3 0 1
   3 -3 2
   -1 -1 0
   -1 -7 2]
svd(V)
```

and find the three singular values are 8.1393, 4.6638, and 0.0000. The zero singular value implies there is a free variables when solving $V\mathbf{c} = \mathbf{0}$ (Procedure 3.3.13), and consequently there are infinitely many non-zero \mathbf{c} that solve $V\mathbf{c} = \mathbf{0}$. By Theorem 7.2.12, the set of vectors is linearly dependent.

- (e) $\{(10, 3, 3, 1), (2, -3, 0, -1), (1, -1, 2, -1), (2, -1, -3, 0), (-2, 0, 2, -1)\}$

Solution: As there are five vectors in \mathbb{R}^4 so Theorem 7.2.13 asserts the set is linearly dependent.



- (f) $\{(-0.4, -1.8, -0.2, 0.7, -0.2), (-1.1, 2.8, 2.7, -3.0, -2.6), (-2.3, -2.3, 4.1, 3.4, -2.6), (-5.3, -3.3, -1.3, -4.1), (1.4, 5.2, -6.9, -0.7, 0.6)\}$

Solution: In Matlab/Octave form the matrix V with these vectors as columns

```
V=[-0.4 -1.1 -2.3 -2.6 1.4
-1.8 2.8 -2.3 -5.3 5.2
-0.2 2.7 4.1 -3.3 -6.9
0.7 -3.0 3.4 -1.3 -0.7
-0.2 -2.6 -1.6 -4.1 0.6]
svd(V)
```



and find the five singular values are 10.6978, 8.0250, 5.5920, 3.0277 and 0.0024. As the singular values are all non-zero, the homogeneous system $V\mathbf{c} = \mathbf{0}$ has the unique solution $\mathbf{c} = \mathbf{0}$ (Procedure 3.3.13), and hence the set of five vectors are linearly independent.

However, the answer depends upon the context. In the strict mathematical context the vectors are unequivocally linearly independent. But in the context of practical problems, where errors in matrix entries are likely, there are ‘shades of grey’. Here, one of the singular values is quite small, namely 0.0024. If the context informs us that the entries in the matrix had errors of say 0.01, then this singular value is effectively zero (Section 5.2). In the context of such errors, this set of five vectors would be effectively linearly dependent.

7.2.2 Form a basis for subspaces

Recall from Sections 2.3 and 3.4 the definition of subspaces and the span: namely that a subspace is a set of vectors closed under addition and scalar multiplication; and a span is a set of vectors whose linear combinations give all vectors in a given subspace. Also, Definition 3.4.14 defined an “orthonormal basis” for a subspace to be a set of orthonormal vectors that span a subspace. This section generalises orthonormal basis by relaxing the requirement of orthonormality.

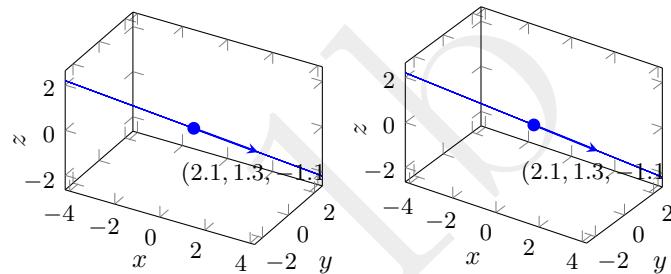
Definition 7.2.15. A **basis** for a subspace \mathbb{W} of \mathbb{R}^n is a set of vectors that both span \mathbb{W} and is linearly independent.

Example 7.2.16. (a) Recall Examples 7.2.4b and 7.2.1 showed that the two vectors $(2, 1)$ and $(1, 2)$ are linearly independent and span \mathbb{R}^2 . Hence the set $\{(2, 1), (1, 2)\}$ is a basis of \mathbb{R}^2 .

(b) Recall that Example 7.2.4a showed the set $\{(-1, 1, 0), (1, -2, 1), (0, 1, -1)\}$ is linearly dependent so it cannot be a basis.

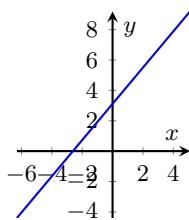
However, remove one vector, such as the middle one, and consider the set $\{(-1, 1, 0), (0, 1, -1)\}$. As the two vectors are not proportional to each other, this set is linearly independent (Theorem 7.2.8). Also, the plane $x + y + z = 0$ is a subspace, say \mathbb{W} . It is characterised by $y = -x - z$. So every vector in \mathbb{W} can be written as $(x, -x - z, z) = (x, -x, 0) + (0, -z, z) = (-x)(-1, 1, 0) + (-z)(0, 1, -1)$. That is, $\text{span}\{(-1, 1, 0), (0, 1, -1)\} = \mathbb{W}$. Hence $\{(-1, 1, 0), (0, 1, -1)\}$ is a basis for the plane \mathbb{W} .

- (c) Find a basis for the line given parametrically as $x = 2.1t$, $y = 1.3t$ and $z = -1.1t$ (shown below in stereo).



Solution: The vectors in the line may be written as $\mathbf{x} = (x, y, z) = (2.1t, 1.3t, -1.1t) = (2.1, 1.3, -1.1)t$. Since the parameter t may vary over all values, vectors in the line form $\text{span}\{(2.1, 1.3, -1.1)\}$. Since $\{(2.1, 1.3, -1.1)\}$ is a linearly independent set of vectors (Example 7.2.4c), it thus forms a basis for the vectors in the given line.

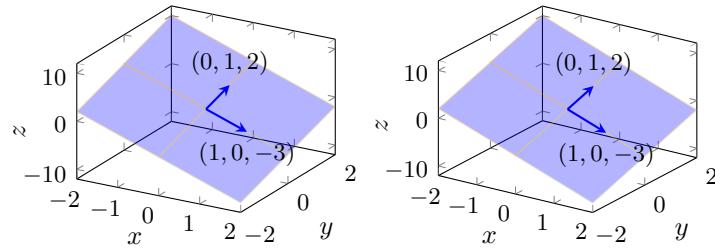
- (d) Find a basis for the line given parametrically as $x = 5.7t - 0.6$ and $y = 6.8t + 2.4$.



Solution: The vectors in the line may be written as $\mathbf{x} = (5.7t - 0.6, 6.8t + 2.4)$. But this does not form a subspace as it does not include the zero vector $\mathbf{0}$ (as illustrated in the margin): the x -component is zero for positive t whereas the y -component is zero for negative t so they are never zero for the same value of parameter t . Since this line is not a subspace, it cannot have a basis.

- (e) Find a basis for the plane $3x - 2y + z = 0$.

Solution: Writing the equation of the plane as $z = -3x + 2y$ we then write the plane parametrically (section 1.3.4) as the vectors $\mathbf{x} = (x, y, -3x + 2y) = (x, 0, -3x) + (0, y, 2y) = x(1, 0, -3) + y(0, 1, 2)$. Since x and y may vary over all values, the plane is the subspace $\text{span}\{(1, 0, -3), (0, 1, 2)\}$ (as illustrated below in stereo). Since $(1, 0, -3)$ and $(0, 1, 2)$ are not proportional to each other, they form a linearly independent set. Hence $\{(1, 0, -3), (0, 1, 2)\}$ is a basis for the plane.



- (f) Prove that every orthonormal basis of a subspace \mathbb{W} is also a basis of \mathbb{W} .

Solution: Theorem 7.2.6 establishes that any orthonormal basis is linearly independent. By Definition 3.4.14, an orthonormal basis of \mathbb{W} spans \mathbb{W} . Hence an orthonormal basis of \mathbb{W} is also a basis of \mathbb{W} .

■

Recall that Theorem 3.4.22 establishes that an orthonormal basis of a given subspace always has the same number of vectors. The following theorem establishes the same is true for general bases. The proof is direct generalisation of that for Theorem 3.4.22.

Theorem 7.2.17. *Any two bases for a given subspace have the same number of vectors.*

Proof. Let $\mathcal{U} = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ and $\mathcal{V} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s\}$ be any two bases for a subspace in \mathbb{R}^n . Prove the number of vectors $r = s$ by contradiction. In the first case, assume $r < s$. Since \mathcal{U} is a basis for the subspace every vector in the set \mathcal{V} can be written as a linear combination of vectors in \mathcal{U} with some coefficients a_{ij} :

$$\begin{aligned}\mathbf{v}_1 &= \mathbf{u}_1 a_{11} + \mathbf{u}_2 a_{21} + \cdots + \mathbf{u}_r a_{r1}, \\ \mathbf{v}_2 &= \mathbf{u}_1 a_{12} + \mathbf{u}_2 a_{22} + \cdots + \mathbf{u}_r a_{r2}, \\ &\vdots \\ \mathbf{v}_s &= \mathbf{u}_1 a_{1s} + \mathbf{u}_2 a_{2s} + \cdots + \mathbf{u}_r a_{rs}.\end{aligned}$$

Write each of these, such as the first one, in the form

$$\mathbf{v}_1 = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_r] \begin{bmatrix} a_{11} \\ a_{21} \\ \vdots \\ a_{r1} \end{bmatrix} = U \mathbf{a}_1,$$

where $n \times r$ matrix $U = [\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_r]$. Similarly for the other equations $\mathbf{v}_2 = \cdots = U \mathbf{a}_2$ through to $\mathbf{v}_s = \cdots = U \mathbf{a}_s$. Then the $n \times s$ matrix

$$V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_s]$$

$$\begin{aligned}
 &= [U\mathbf{a}_1 \ U\mathbf{a}_2 \ \cdots \ U\mathbf{a}_s] \\
 &= U [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_s] = UA
 \end{aligned}$$

for the $r \times s$ matrix $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_s]$. By assumption $r < s$ and so Theorem 2.2.25 assures us that the homogeneous system $A\mathbf{x} = \mathbf{0}$ has infinitely many solutions, choose any non-trivial solution $\mathbf{x} \neq \mathbf{0}$. Consider

$$\begin{aligned}
 V\mathbf{x} &= UA\mathbf{x} \quad (\text{from above}) \\
 &= U\mathbf{0} \quad (\text{since } A\mathbf{x} = \mathbf{0}) \\
 &= \mathbf{0}.
 \end{aligned}$$

The identity $V\mathbf{x} = \mathbf{0}$ implies there is a linear combination of the columns $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s$ of V which gives zero, hence the set \mathcal{V} is linearly dependent (Theorem 7.2.12). But this is a contradiction, so we cannot have $r < s$.

Second, a corresponding argument establishes we cannot have $s < r$. Hence $s = r$: all bases of a given subspace must have the same number of vectors. \square

Example 7.2.18. Consider the plane $x + y + z = 0$ in \mathbb{R}^3 . Each of the following are a basis for the plane:

- $\{(-1, 1, 0), (1, -2, 1)\}$;
- $\{(1, -2, 1), (0, 1, -1)\}$;
- $\{(0, 1, -1), (-1, 1, 0)\}$.

The reasons are that all three vectors involved are in the plane, and that each pair are linearly independent (as, in each pair, one is not proportional to the other). However, although each of the three vectors in the set $\{(-1, 1, 0), (1, -2, 1), (0, 1, -1)\}$ is in the plane, this set is not a basis because it is not linearly independent (Example 7.2.4a). Also, each individual vector, say $(-1, 1, 0)$, cannot form a basis for the plane because the span of one vector, such as $\text{span}\{(-1, 1, 0)\}$, is a line not the whole plane.

The orthonormal basis $\{(1, 0, -1)/\sqrt{2}, (1, -2, 1)/\sqrt{6}\}$ is another basis for the plane: both vectors satisfy $x + y + z = 0$ and are orthogonal and so linearly independent (Theorem 7.2.6). All these bases possess two vectors. \blacksquare

That all bases for a given subspace, including orthonormal bases, have the same number of vectors leads to the following theorem about the dimensionality.

Theorem 7.2.19. Let \mathbb{W} be a subspace of \mathbb{R}^n . The **dimension** of \mathbb{W} , denoted $\dim \mathbb{W}$, is the number of vectors in any basis for \mathbb{W} .

Proof. Recall Definition 3.4.24 defined $\dim \mathbb{W}$ to be the number of vectors in any orthonormal basis for \mathbb{W} . Theorem 7.2.6 certifies that all orthonormal bases are also bases (Definition 7.2.15), so Theorem 7.2.17 implies every basis of \mathbb{W} has $\dim \mathbb{W}$ vectors. \square

Procedure 7.2.20 (basis for a span). *To find a basis for the subspace $\mathbb{A} = \text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ given $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ is a set of n vectors in \mathbb{R}^m . Recall Theorem 3.4.18 underpins finding an orthonormal basis by the following.*

1. Form matrix $A := [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$.
2. Factorise A into its SVD, $A = USV^T$, and let $r = \text{rank } A$ be the number of nonzero singular values (or effectively nonzero when the matrix has experimental errors, Section 5.2).
3. The set $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r\}$ (where \mathbf{u}_j denotes the columns of U) is a basis, specifically an orthonormal basis, for the rD subspace \mathbb{A} .

Alternatively, if $r = n$, then the set $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ is linearly independent and span the subspace \mathbb{A} , and so is also a basis for the nD subspace \mathbb{A} .

Example 7.2.21. Apply Procedure 7.2.20 to find a basis for the following sets.

- (a) Recall Example 7.2.18 identified that any pair of vectors in the set $\{(-1, 1, 0), (1, -2, 1), (0, 1, -1)\}$ forms a basis for the plane that they span. Find another basis for the plane.

Solution: In Matlab/Octave form the matrix with these vectors as columns:

```
A=[-1 1 0
    1 -2 1
    0 1 -1]
[U,S,V]=svd(A)
```

Then the SVD obtains (2 d.p.)

```
U =
    -0.41   -0.71    0.58
    0.82     0.00    0.58
    -0.41    0.71    0.58
S =
    3.00      0      0
        0    1.00      0
        0      0    0.00
V = ...
```

The two non-zero singular values determine $\text{rank } A = 2$ and hence the first two columns of U form an (orthonormal) basis for $\text{span}\{(-1, 1, 0), (1, -2, 1), (0, 1, -1)\}$. That is, $\{0.41(-1, 2, -1), 0.71(-1, 0, 1)\}$ is an (orthonormal) basis.



(b) The span of the three vectors

$$(-2, 0, -4, 1, 1), (7, 1, 2, -1, -5), (-5, -1, 2, 3, -2).$$

Solution: In Matlab/Octave it is often easiest to enter these vectors as rows, and then transpose with the dash operator to form the matrix with these as columns:

```
A=[ -2 0 -4 1 1
    7 1 2 -1 -5
    -5 -1 2 3 -2]',
```

```
[U,S,V]=svd(A)
```

Then the SVD obtains (2 d.p.)

```
U =
  0.86 -0.32 -0.02  0.40  0.02
  0.12 -0.11 -0.06 -0.40  0.90
  0.22  0.65  0.72  0.07  0.13
 -0.23  0.35 -0.38  0.73  0.37
 -0.38 -0.59  0.58  0.37  0.19

S =
 10.07      0      0
  0   5.87      0
  0      0   3.01
  0      0      0
  0      0      0

V = ...
```

The three non-zero singular values determine $\text{rank } A = 3$ and so the original three vectors are linearly independent. Consequently the original three vectors form a basis for their span.

If you prefer an orthonormal basis, then use the first three columns of U.

(c) The span of the four vectors $(1, 0, 3, -4, 0)$, $(-1, -1, 1, 4, 2)$, $(-3, 2, 2, 2, 1)$, $(3, -3, 2, -2, 1)$.

Solution: In Matlab/Octave, enter these vectors as rows, and then transpose with the dash operator to form the matrix with these as columns:

```
A=[1 0 3 -4 0
   -1 -1 1 4 2
   -3 2 2 2 1
   3 -3 2 -2 1]',
```

```
[U,S,V]=svd(A)
```

Then the SVD obtains (2 d.p.)

```
U =
 -0.52    0.11    0.46   -0.71    0.02
```

$$\begin{array}{cccccc}
 0.28 & -0.44 & -0.55 & -0.61 & 0.24 \\
 -0.19 & 0.66 & -0.60 & -0.17 & -0.37 \\
 0.78 & 0.32 & 0.36 & -0.30 & -0.27 \\
 0.10 & 0.51 & -0.00 & 0.03 & 0.86 \\
 \mathbf{S} = & & & & & \\
 7.64 & 0 & 0 & 0 \\
 0 & 4.59 & 0 & 0 \\
 0 & 0 & 4.30 & 0 \\
 0 & 0 & 0 & 0.00 \\
 0 & 0 & 0 & 0 \\
 \mathbf{V} = \dots & & & & &
 \end{array}$$

The three non-zero singular values determine $\text{rank } A = 3$. Consequently, the first three columns of \mathbf{U} form an orthonormal basis for the span.

Alternatively, you might notice that the sum of the first two vectors is the sum of the last two vectors. Consequently, given the rank is three, we obtain three linearly independent vectors by omitting any one. That is, any three of the given vectors form a basis for the span of the four.

The procedure is different if the subspace of interest is defined by a system of equations instead of the span of some vectors.

Example 7.2.22. Find a basis for the solutions of the system in \mathbb{R}^3 of $3x + y = 0$ and $3x + 2y + 3z = 0$.

Solution: By hand manipulation, the first equation gives $y = -3x$; which when substituted into the second gives $3x - 6x + 3z = 0$, namely $z = x$. That is, all solutions are of the form $(x, -3x, x)$, namely $\text{span}\{(1, -3, 1)\}$. Thus a basis for the subspace of solutions is $\{(1, -3, 1)\}$. (Infinitely many other bases are possible answers.)

Example 7.2.23. Find a basis for the solutions of $-2x - y + 3z = 0$ in \mathbb{R}^3 .

Solution: Rearrange the equation so that one variable is a function of the others, say $y = -2x + 3z$. Then the vector form of solutions are $(x, y, z) = (x, -2x + 3z, z) = (1, -2, 0)x + (0, 3, 1)z$ in terms of free variables x and z . Since $(1, -2, 0)$ and $(0, 3, 1)$ are not proportional to each other, they are linearly independent, and so a basis for the solutions is $\{(1, -2, 0), (0, 3, 1)\}$. (Infinitely many other bases are possible answers.)

Procedure 7.2.24 (basis from equations). Suppose we seek a basis for a subspace \mathbb{W} defined as the solutions of a system of equations.

1. Rewrite the system of equations as the homogeneous system $A\mathbf{x} = \mathbf{0}$, hence the subspace \mathbb{W} is the nullspace of $m \times n$ matrix A .
2. Adapting Procedure 3.3.13 for the specific case of homogeneous systems, first find the SVD factorisation $A = USV^T$ and let $r = \text{rank } A$ be the number of nonzero singular values (or effectively nonzero when the matrix has experimental errors, Section 5.2).
3. Then the general solution of $S\mathbf{y} = \mathbf{z} = \mathbf{0}$ is $\mathbf{y} = (0, \dots, 0, y_{r+1}, \dots, y_n)$. Consequently, all possible solutions $\mathbf{x} = V\mathbf{y}$ are spanned by the last $n - r$ columns of V , which thus form an orthonormal basis for the subspace \mathbb{W} .

Example 7.2.25. Find a basis for all solutions to each of the following systems of equations.

- (a) $3x + y = 0$ and $3x + 2y + 3z = 0$ from Example 7.2.22.

Solution: Form matrix $A = \begin{bmatrix} 3 & 1 & 0 \\ 3 & 2 & 3 \end{bmatrix}$ and compute an SVD with $[U, S, V] = \text{svd}(A)$ to obtain (2 d.p.)

$$\begin{aligned} U &= \dots \\ S &= \\ &\begin{array}{ccc} 5.34 & 0 & 0 \\ 0 & 1.86 & 0 \end{array} \\ V &= \\ &\begin{array}{ccc} 0.77 & 0.56 & 0.30 \\ 0.42 & -0.09 & -0.90 \\ 0.48 & -0.82 & 0.30 \end{array} \end{aligned}$$

The two non-zero singular values determine $\text{rank } A = 2$. Hence the solutions of the system are spanned by the last one column of V . That is, a basis for the solutions is $\{(0.3, -0.9, 0.3)\}$.

- (b) $7x = 6y + z + 3$ and $4x + 9y + 2z + 2 = 0$.

Solution: This system is not homogeneous (due to the constant terms, Definition 2.2.23), therefore $\mathbf{x} = \mathbf{0}$ is not a solution. Consequently, the solutions of the system cannot form a subspace (Definition 3.4.2). Thus the concept of a basis does not apply (Definition 7.2.15).

- (c) $w + x = z$, $3w = x + y + 5z$, $4x + y + 2z = 0$.

Solution: Rearrange to the matrix-vector system $A\mathbf{x} = \mathbf{0}$ for vector $\mathbf{x} = (w, x, y, z) \in \mathbb{R}^4$ and matrix

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



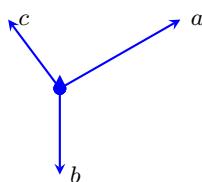
Enter into Matlab/Octave as above and then find an SVD with $[U, S, V] = \text{svd}(A)$ to obtain (2 d.p.)

```

U = ...
S =
    6.77      0      0      0
        0    3.76      0      0
        0      0   0.00      0
V =
   -0.40    0.45    0.09    0.80
    0.41    0.86   -0.19   -0.25
    0.20    0.10    0.97   -0.07
    0.80   -0.24   -0.10    0.54

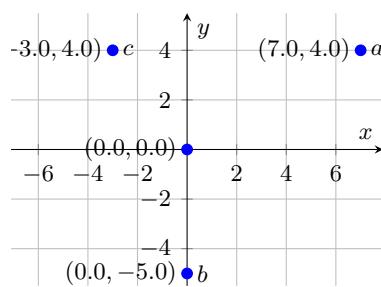
```

There are two non-zero singular values, so $\text{rank } A = 2$. There are thus $4 - 2 = 2$ free variables in solving $Ax = \mathbf{0}$ leading to a 2D subspace with orthonormal basis of the last two columns of V . That is, an orthonormal basis for the subspace of all solutions is $\{(0.09, -0.19, 0.97, -0.10), (0.80, -0.25, -0.07, 0.54)\}$.

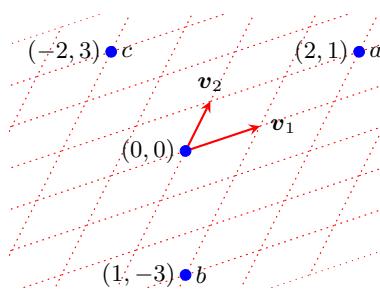


Recall this Section 7.2 started by discussing the need to have a unique representation of solutions to problems. If we do not have uniqueness, then the ambiguity in algebraic representation ruins basic algebra. The following theorem assures us that the linear independence of a basis ensures the unique representation that we need. In essence it says that any basis, whether orthogonal or not, can be used to form a coordinate system.

Example 7.2.26 (a tale of two coordinate systems). In the margin are plotted three vectors and the origin. Take the view that these are fixed physically meaningful vectors: the issue of this example is how do we code such vectors in mathematics.



In the standard orthogonal coordinate system these three vectors and the origin have coordinates as plotted left by their end-points. We write $\mathbf{a} = (7, 4)$, $\mathbf{b} = (0, -5)$ and $\mathbf{c} = (-3, 4)$.



Now use the (red) basis $\mathcal{B} = \{\mathbf{v}_1, \mathbf{v}_2\}$ to form a non-orthogonal coordinate system (represented by the dotted grid). Then in this system the three vectors have coordinates $\mathbf{a} = (2,1)$, $\mathbf{b} = (1, -3)$ and $\mathbf{c} = (-2, 3)$.

But we cannot say both $\mathbf{a} = (7,4)$ and $\mathbf{a} = (2,1)$: it appears nonsense. The reason for the different numbers representing the one vector \mathbf{a} is that the underlying coordinate systems are different. For example, we can say both $\mathbf{a} = 7\mathbf{e}_1 + 4\mathbf{e}_2$ and $\mathbf{a} = 2\mathbf{v}_1 + \mathbf{v}_2$ without any apparent contradiction: these statements recognise the underlying standard unit vectors in the first expression and the underlying non-orthogonal basis vectors in the second.

Consequently we invent a new better notation. We write $[\mathbf{a}]_{\mathcal{B}} = (2,1)$ to represent that the coordinates of vector \mathbf{a} in the basis \mathcal{B} are $(2,1)$. Correspondingly, letting $\mathcal{E} = \{\mathbf{e}_1, \mathbf{e}_2\}$ denote the basis of the standard unit vectors, we write $[\mathbf{a}]_{\mathcal{E}} = (7,4)$ to represent that the coordinates of vector \mathbf{a} in the standard basis \mathcal{E} are $(7,4)$. Similarly, $[\mathbf{b}]_{\mathcal{E}} = (0, -5)$ and $[\mathbf{b}]_{\mathcal{B}} = (1, -3)$; and $[\mathbf{c}]_{\mathcal{E}} = (-3, 4)$ and $[\mathbf{c}]_{\mathcal{B}} = (-2, 3)$.

The endemic practice of just writing $\mathbf{a} = (2,1)$, $\mathbf{b} = (1, -3)$ and $\mathbf{c} = (-2, 3)$ is rationalised in this new notation by the convention that if no basis is specified, then the standard basis \mathcal{E} is assumed.

⁶



Theorem 7.2.27. *Let \mathbb{W} be a subspace of \mathbb{R}^n and let $\mathcal{B} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ be a basis for \mathbb{W} . There is exactly one way to write each and every vector $\mathbf{w} \in \mathbb{W}$ as a linear combination of the basis vectors: $\mathbf{w} = c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k$. Then c_1, c_2, \dots, c_k are called the **coordinates of \mathbf{w} with respect to \mathcal{B}** , and the column vector $[\mathbf{w}]_{\mathcal{B}} = (c_1, c_2, \dots, c_k)$ is called the **coordinate vector of \mathbf{w} with respect to \mathcal{B}** .*

Proof. Consider any vector $\mathbf{w} \in \mathbb{W}$. Since $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is a basis for the subspace \mathbb{W} , \mathbf{w} can be written as a linear combination of the basis vectors. Let two such linear combinations be

$$\begin{aligned}\mathbf{w} &= c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_k\mathbf{v}_k , \\ \mathbf{w} &= d_1\mathbf{v}_1 + d_2\mathbf{v}_2 + \dots + d_k\mathbf{v}_k .\end{aligned}$$

⁶ Given that the numerical representation of a vector changes with the coordinate basis, some of you will wonder whether the same thing happens for matrices. The answer is yes.

Subtract the second of these equations from the first, grouping common vectors:

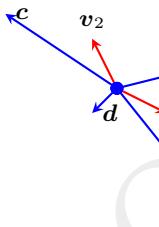
$$\mathbf{0} = (c_1 - d_1)\mathbf{v}_1 + (c_2 - d_2)\mathbf{v}_2 + \cdots + (c_k - d_k)\mathbf{v}_k.$$

Since $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is linearly independent, this equation implies all the coefficients in parentheses are zero:

$$0 = (c_1 - d_1) = (c_2 - d_2) = \cdots = (c_k - d_k).$$

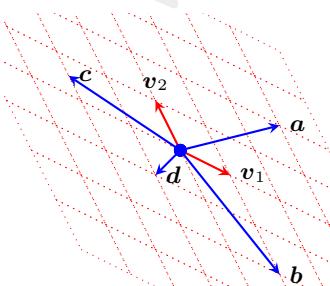
That is, $c_1 = d_1$, $c_2 = d_2$, ..., $c_k = d_k$, and the two linear combinations are identical. That is, there is exactly one way to write a vector $\mathbf{w} \in \mathbb{W}$ as a linear combination of the basis vectors. \square

Example 7.2.28. (a) Consider the diagram of six labelled vectors drawn below.



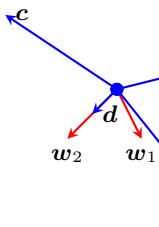
a Estimate the coordinates of the four shown vectors \mathbf{a} , \mathbf{b} , \mathbf{c} and \mathbf{d} in the shown basis $\mathcal{B} = \{\mathbf{v}_1, \mathbf{v}_2\}$.

Solution: Draw in a grid corresponding to multiples of \mathbf{v}_1 and \mathbf{v}_2 in both directions, and parallel to \mathbf{v}_1 and \mathbf{v}_2 , as shown below. Then from the grid, estimate that $\mathbf{a} \approx 3\mathbf{v}_1 + 2\mathbf{v}_2$ hence the coordinates $[\mathbf{a}]_{\mathcal{B}} \approx (3, 2)$.



Similarly, $\mathbf{b} \approx \mathbf{v}_1 - 2\mathbf{v}_2$ hence the coordinates $[\mathbf{b}]_{\mathcal{B}} \approx (1, -2)$. Also, $\mathbf{c} \approx -2\mathbf{v}_1 + 0.5\mathbf{v}_2$ hence the coordinates $[\mathbf{c}]_{\mathcal{B}} \approx (-2, 0.5)$. And lastly, $\mathbf{d} \approx -\mathbf{v}_1 - \mathbf{v}_2$ hence the coordinates $[\mathbf{d}]_{\mathcal{B}} \approx (-1, -1)$.

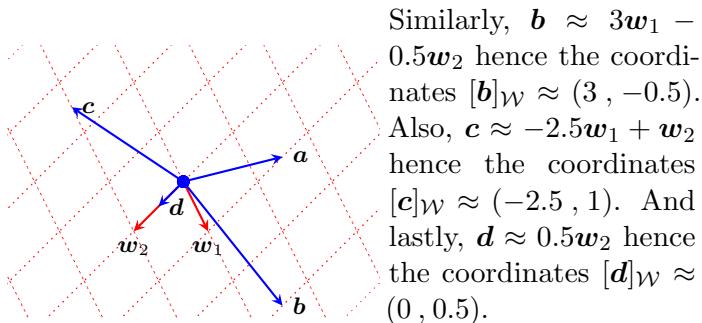
(b) Consider the same four vectors but with a pair of different basis vectors: let's see that although the vectors are the same, the coordinates in the different basis are different.



a Estimate the coordinates of the four shown vectors \mathbf{a} , \mathbf{b} , \mathbf{c} and \mathbf{d} in the shown basis $\mathcal{W} = \{\mathbf{w}_1, \mathbf{w}_2\}$.

Solution: Draw in a grid corresponding to multiples of \mathbf{w}_1 and \mathbf{w}_2 in both directions, and parallel to \mathbf{w}_1 and \mathbf{w}_2 ,

as shown below. Then from the grid, estimate that $\mathbf{a} \approx \mathbf{w}_1 - 1.5\mathbf{w}_2$ hence the coordinates $[\mathbf{a}]_{\mathcal{W}} \approx (1, -1.5)$.



Example 7.2.29. Let the basis $\mathcal{B} = \{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\}$ for the three given vectors $\mathbf{v}_1 = (-1, 1, -1)$, $\mathbf{v}_2 = (1, -2, 0)$ and $\mathbf{v}_3 = (0, 4, 5)$ (specified in the standard basis \mathcal{E} of the standard unit vectors \mathbf{e}_1 , \mathbf{e}_2 and \mathbf{e}_3).

- (a) What is the vector with coordinates $[\mathbf{a}]_{\mathcal{B}} = (3, -2, 1)$?

Solution: $\mathbf{a} = 3\mathbf{v}_1 - 2\mathbf{v}_2 + \mathbf{v}_3$ which has standard coordinates $[\mathbf{a}]_{\mathcal{E}} = 3(-1, 1, -1) - 2(1, -2, 0) + (0, 4, 5) = (-5, 11, 2)$.

- (b) What is the vector with coordinates $[\mathbf{b}]_{\mathcal{B}} = (-1, 1, 1)$?

Solution: $\mathbf{b} = -\mathbf{v}_1 + \mathbf{v}_2 + \mathbf{v}_3$ which has standard coordinates $[\mathbf{b}]_{\mathcal{E}} = -(-1, 1, -1) + (1, -2, 0) + (0, 4, 5) = (2, 1, 6)$.

- (c) What are the coordinates in the basis \mathcal{B} of the vector $\mathbf{c} = (-1, 3, 3)$ in the standard basis \mathcal{E} ?

Solution: We seek coordinate values c_1, c_2, c_3 such that $\mathbf{c} = c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + c_3\mathbf{v}_3$. Expressed in the standard basis this equation is

$$\begin{bmatrix} -1 \\ 3 \\ 3 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} c_1 + \begin{bmatrix} 1 \\ -2 \\ 0 \end{bmatrix} c_2 + \begin{bmatrix} 0 \\ 4 \\ 5 \end{bmatrix} c_3.$$

A small system like this we solve by hand (recall section 2.2.2): write in component form as

$$\begin{cases} -c_1 + c_2 = -1 \\ c_1 - 2c_2 + 4c_3 = 3 \\ -c_1 + 5c_3 = 3 \end{cases}$$

(add 1st row to 2nd and take from 3rd)

$$\begin{cases} -c_1 + c_2 = -1 \\ -c_2 + 4c_3 = 2 \\ -c_2 + 5c_3 = 4 \end{cases}$$

(subtract 2nd row from 3rd)

$$\begin{cases} -c_1 + c_2 = -1 \\ -c_2 + 4c_3 = 2 \\ c_3 = 2 \end{cases}$$

Solving this triangular system gives $c_3 = 2$, $c_2 = 4c_3 - 2 = 6$, and $c_1 = c_2 + 1 = 7$. Thus the coordinates $[\mathbf{c}]_{\mathcal{B}} = (7, 6, 2)$ in the basis \mathcal{B} .

- (d) What are the coordinates in the basis \mathcal{B} of the vector $\mathbf{d} = (-3, 2, 0)$ in the standard basis \mathcal{E} ?

Solution: We seek coordinate values d_1, d_2, d_3 such that $\mathbf{d} = d_1 \mathbf{v}_1 + d_2 \mathbf{v}_2 + d_3 \mathbf{v}_3$. Expressed in the standard basis this equation is

$$\begin{bmatrix} -3 \\ 2 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} d_1 + \begin{bmatrix} 1 \\ -2 \\ 0 \end{bmatrix} d_2 + \begin{bmatrix} 0 \\ 4 \\ 5 \end{bmatrix} d_3.$$

To solve this system with Matlab/Octave (procedure 2.2.4), enter the matrix (easiest by transposing rows of \mathbf{v}_1 , \mathbf{v}_2 and \mathbf{v}_3) and the standard coordinates of \mathbf{d} :

```
A=[-1 1 -1
    1 -2 0
    0 4 5];
d=[-3;2;0]
```



Then compute the coordinates $[\mathbf{d}]_{\mathcal{B}}$ with $\mathbf{dB}=\mathbf{A}\backslash\mathbf{d}$ to determine $[\mathbf{d}]_{\mathcal{B}} = (20, 17, 4)$.

But remember, before using $\mathbf{A}\backslash$ always first check `rcond(A)` which here is the poor 0.0053 (Procedure 2.2.4). Interestingly, this poor small value of `rcond` indicates that although the basis vectors in \mathcal{B} are linearly independent, they are ‘only just’ linearly independent. A small change or error might make them linearly dependent and thus \mathcal{B} be ruined as a basis for \mathbb{R}^3 . The poor `rcond` indicates that \mathcal{B} is a poor basis in practice.

■

Example 7.2.30. You are given a basis $\mathcal{W} = \{\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3\}$ for a 3D subspace \mathbb{W} of \mathbb{R}^5 where the three basis vectors are $\mathbf{w}_1 = (1, 3, -4, -3, 3)$, $\mathbf{w}_2 = (-4, 1, -2, -4, 1)$, and $\mathbf{w}_3 = (-1, 1, 0, 2, -3)$ (in the standard basis \mathcal{E}).

- (a) What are the coordinates in the standard basis of the vector $\mathbf{a} = 2\mathbf{w}_1 + 3\mathbf{w}_2 + \mathbf{w}_3$?

Solution: In the standard basis

$$[\mathbf{a}]_{\mathcal{E}} = 2 \begin{bmatrix} 1 \\ 3 \\ -4 \\ -3 \\ 3 \end{bmatrix} + 3 \begin{bmatrix} -4 \\ 1 \\ -2 \\ -4 \\ 1 \end{bmatrix} + \begin{bmatrix} -1 \\ 1 \\ 0 \\ 2 \\ -3 \end{bmatrix} = \begin{bmatrix} -11 \\ 10 \\ -14 \\ -16 \\ 6 \end{bmatrix}.$$

- (b) What are the coordinates in the basis \mathcal{W} of the vector $\mathbf{b} = (-1, 2, -6, -11, 10)$ (in the standard coordinates \mathcal{E})?

Solution: We need to find coefficients c_1, c_2, c_3 such that $\mathbf{b} = c_1 \mathbf{w}_1 + c_2 \mathbf{w}_2 + c_3 \mathbf{w}_3$. This forms the set of linear equations

$$\begin{bmatrix} -1 \\ 2 \\ -6 \\ -11 \\ 10 \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \\ -4 \\ -3 \\ 3 \end{bmatrix} c_1 + \begin{bmatrix} -4 \\ 1 \\ -2 \\ -4 \\ 1 \end{bmatrix} c_2 + \begin{bmatrix} -1 \\ 1 \\ 0 \\ 2 \\ -3 \end{bmatrix} c_3.$$

Form as the matrix-vector system

$$\begin{bmatrix} 1 & -4 & -1 \\ 3 & 1 & 1 \\ -4 & -2 & 0 \\ -3 & -4 & 2 \\ 3 & 1 & -3 \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix} = \begin{bmatrix} -1 \\ 2 \\ -6 \\ -11 \\ 10 \end{bmatrix},$$

and perhaps solve with Matlab/Octave. Since there are more equations than unknowns, we should use an SVD in order to check the system is consistent, namely, to check that $\mathbf{b} \in \mathbb{W}$.

- i. Code the matrix and the vector:

```
W=[1 -4 -1
    3 1 1
   -4 -2 0
   -3 -4 2
    3 1 -3]
b=[-1;2;-6;-11;10]
```

- ii. Then obtain an SVD with $[\mathbf{U}, \mathbf{S}, \mathbf{V}] = \text{svd}(\mathbf{W})$ (2 d.p.)

```
U =
  0.18  -0.88   0.04  -0.24  -0.37
 -0.32  -0.09   0.65   0.62  -0.28
  0.51   0.11  -0.46   0.56  -0.44
  0.64  -0.18   0.34   0.22   0.63
 -0.44  -0.42  -0.49   0.44   0.45
S =
  8.18      0      0
      0   4.52      0
```



```

          0      0   3.11
          0      0      0
          0      0      0
V =
 -0.74  -0.51   0.43
 -0.62   0.77  -0.15
  0.26   0.38   0.89

```

The three non-zero singular values establish that the three vectors in the basis \mathcal{W} are indeed linearly independent (and since no singular value is small, then the vectors are robustly linearly independent).

iii. Find $\mathbf{z} = U^T \mathbf{b}$ with $\mathbf{z} = \mathbf{U}' * \mathbf{b}$ to get

```

z =
-15.3413
-2.2284
-4.6562
-0.0000
 0.0000

```

The last two values of \mathbf{z} being zero confirm the system of equations is consistent and so vector \mathbf{b} is in the range of \mathcal{W} , that is, \mathbf{b} is in the subspace \mathbb{W} .

iv. Find $y_j = z_j / \sigma_j$ with $\mathbf{y} = \mathbf{z}(1:3) ./ \text{diag}(\mathbf{S})$ to get

```

y =
-1.8761
-0.4929
-1.4958

```

v. Lastly, find the coefficients $[\mathbf{b}]_{\mathcal{W}} = V \mathbf{y}$ with $\mathbf{bw} = \mathbf{V} * \mathbf{y}$ to get

```

bw =
 1.00000
 1.00000
-2.00000

```

That is, $[\mathbf{b}]_{\mathcal{W}} = (1, 1, -2)$.

That $[\mathbf{b}]_{\mathcal{W}}$ has three components and $[\mathbf{b}]_{\mathcal{E}}$ has five components is not a contradiction. The difference in components occurs because the subspace \mathbb{W} is 3D but lies in \mathbb{R}^5 . Using the basis \mathcal{W} implicitly builds in the information that the vector \mathbf{b} is in a lower dimensional space, and so needs fewer components.



Revisit unique solutions

Lastly, with all these extra concepts of determinants, eigenvalues, linear independence and a basis, we now revisit the issue of when there is a unique solution to a set of linear equations.

Theorem 7.2.31 (Unique Solutions: version 3). *Let A be an $n \times n$ square matrix. Extending Theorem 3.4.35, the following statements are equivalent:*

- (a) A is invertible;
- (b) $Ax = b$ has a unique solution for every $b \in \mathbb{R}^n$;
- (c) $Ax = \mathbf{0}$ has only the zero solution;
- (d) all n singular values of A are nonzero;
- (e) $\text{rank } A = n$;
- (f) $\text{nullity } A = 0$;
- (g) the column vectors of A span \mathbb{R}^n ;
- (h) the row vectors of A span \mathbb{R}^n .
- (i) $\det A \neq 0$;
- (j) 0 is not an eigenvalue of A ;
- (k) the n column vectors of A are linearly independent;
- (l) the n row vectors of A are linearly independent.

Proof. Recall Theorems 3.3.21 and 3.4.35 established the equivalence of 7.2.31a–7.2.31h, and Theorem 6.1.23 proved the equivalence of 7.2.31a and 7.2.31i. To establish that Property 7.2.31j is equivalent to 7.2.31i, recall Theorem 7.1.3 proved that $\det A$ equals the product of the eigenvalues of A . Hence $\det A$ is not zero if and only if all the eigenvalues are non-zero.

Lastly, Property 7.2.31g says the n column vectors span \mathbb{R}^n , so they must be a basis, and hence linearly independent. Conversely, if the n columns of A are linearly independent then they must span \mathbb{R}^n . Hence Property 7.2.31k is equivalent to 7.2.31g. Similarly for Property 7.2.31h and the row vectors of 7.2.31l. \square

7.2.3 Exercises

Exercise 7.2.1. By inspection or basic arguments, decide whether the following sets of vectors are linearly dependent, or linearly independent. Give reasons.

- (a) $\{(-2, 3, 3), (-1, 2, -1)\}$
- (b) $\{(0, 2), (2, -2), (0, -1)\}$

- (c) $\{(-3, 0, -3), (3, 2, -2)\}$
 (d) $\{(0, 2, 2)\}$
 (e) $\{(-2, 0, -1), (0, -2, 2), (2, 0, 1)\}$
 (f) $\{(0, 3, -2), (-2, -1, 1), (1, -2, -4)\}$
 (g) $\{(-2, 1), (0, 0)\}$
 (h) $\{(-2, -1, 1), (-2, -2, 2), (2, -1, -1), (2, -2, -2)\}$
 (i) $\{(2, 4), (1, 2)\}$
 (j) $\{(1, -2, 3), (1, 1, 2)\}$

Exercise 7.2.2. Compute an SVD to decide whether the following sets of vectors are linearly dependent, or linearly independent. Give reasons.

(a) $\begin{bmatrix} 5 \\ 0 \\ -1 \end{bmatrix}, \begin{bmatrix} -2 \\ 5 \\ -1 \end{bmatrix}$	(b) $\begin{bmatrix} -2 \\ -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 3 \\ -4 \\ -1 \end{bmatrix}, \begin{bmatrix} 5 \\ -3 \\ -2 \end{bmatrix}$
(c) $\begin{bmatrix} 2 \\ -4 \\ 1 \\ 6 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ 10 \\ -5 \end{bmatrix}, \begin{bmatrix} 1 \\ 6 \\ -2 \\ 4 \end{bmatrix}$	(d) $\begin{bmatrix} 1 \\ 0 \\ -5 \\ 2 \end{bmatrix}, \begin{bmatrix} -1 \\ 0 \\ -2 \\ 4 \end{bmatrix}, \begin{bmatrix} -1 \\ 3 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 3 \\ 4 \\ -4 \\ -4 \end{bmatrix}$
(e) $\begin{bmatrix} 0 \\ -2 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} -1 \\ -4 \\ -6 \\ -1 \end{bmatrix}, \begin{bmatrix} 4 \\ 3 \\ -5 \\ 0 \end{bmatrix}$	(f) $\begin{bmatrix} 0 \\ 3 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} -3 \\ -3 \\ 1 \\ -2 \end{bmatrix}, \begin{bmatrix} -3 \\ 0 \\ 3 \\ -2 \end{bmatrix}$
(g) $\begin{bmatrix} -3 \\ 2 \\ 6 \\ 3 \\ -2 \end{bmatrix}, \begin{bmatrix} -3 \\ -3 \\ 1 \\ 6 \\ -4 \end{bmatrix}, \begin{bmatrix} 3 \\ -2 \\ -2 \\ 4 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ -2 \\ -2 \\ 3 \\ 2 \end{bmatrix}$	(h) $\begin{bmatrix} 3 \\ 0 \\ 4 \\ 1 \\ 2 \end{bmatrix}, \begin{bmatrix} -2 \\ 1 \\ 2 \\ 2 \\ 2 \end{bmatrix}, \begin{bmatrix} 2 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} -9 \\ -1 \\ -2 \\ 1 \\ 0 \end{bmatrix}$

Exercise 7.2.3. Prove the particular case of Theorem 7.2.8, namely that a set of two vectors $\{\mathbf{v}_1, \mathbf{v}_2\}$ is linearly dependent if and only if one of the vectors is a scalar multiple of the other.

Exercise 7.2.4. Prove that every (non-empty) subset of a linearly independent set is also linearly independent. (Perhaps use contradiction.)

Exercise 7.2.5. Let $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$ be a linearly independent set of vectors in \mathbb{R}^n . Given that a vector $\mathbf{u} = c_1\mathbf{v}_1 + c_2\mathbf{v}_2 + \dots + c_m\mathbf{v}_m$ with coefficient $c_1 \neq 0$, prove that the set $\{\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_m, \mathbf{u}\}$ is linearly independent.

Exercise 7.2.6. For each of the following systems of equations find by hand two different bases for their solution set (among the infinitely many bases that are possible). Show your working.

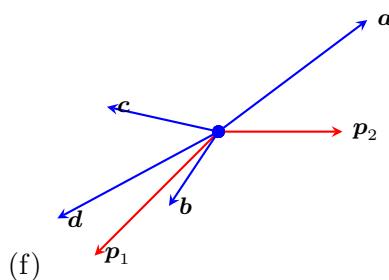
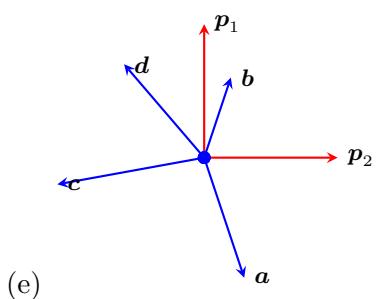
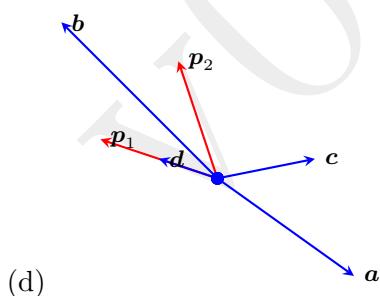
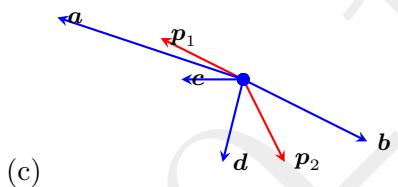
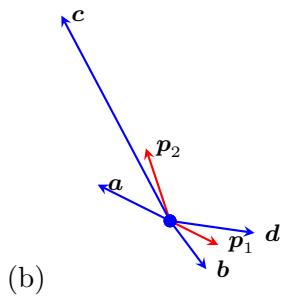
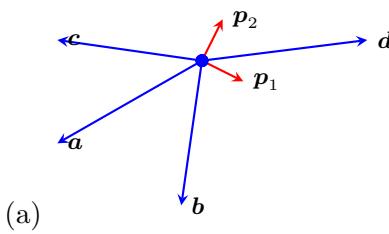
- (a) $-x - 5y = 0$ and $y - 3z = 0$
- (b) $6x + 4y + 2z = 0$ and $-2x - y - 2z = 0$
- (c) $-2y - z + 2 = 0$ and $3x + 4z = 0$
- (d) $-7x + y - z = 0$ and $-3x + 2 - 2z = 0$
- (e) $x + 2y + 2z = 0$
- (f) $2x + 0y - 4z = 0$
- (g) $-2x + 3y - 6z = 0$
- (h) $9x + 4y - 9z = 6$

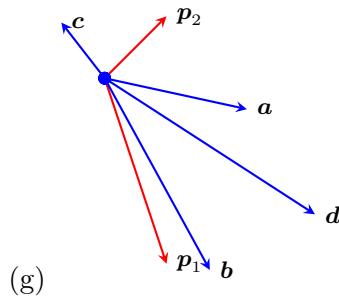
Exercise 7.2.7. Use Procedure 7.2.24 to compute in Matlab/Octave an orthonormal basis for all solutions to each of the following systems of equations.

- (a) $2x - 6y - 9z = 0$, $2x - 2z = 0$, $-x + z = 0$
- (b) $3x - 3y + 8z = 0$, $-2x - 4y + 2z = 0$, $-4x + y - 7z = 0$
- (c) $-2x + 2y - 2z = 0$, $x + 3y - z = 0$, $3x + 3y = 0$
- (d) $2w + x + 2y + z = 0$, $w + 4x - 4y - z = 0$, $3w - 2x + 5y = 0$,
 $2w - x + y - 2z = 0$
- (e) $5w + y - 3z = 0$, $-5x - 5y = 0$, $-3x - y + 4z = 0$, $3x + y - 4z = 0$
- (f) $-w - 2y + 4z = 0$, $2w + 2y + 2z = 0$, $-2w + 3x + y + z = 0$,
 $-w + x - y + 5z = 0$
- (g) $-2w + x + 2y - 6z = 0$, $-2w + 3x + 4y = 0$, $-2w + 2x + 3x - 3z = 0$
- (h) $-w - 2x - 3y + 2z = 0$, $2x + 2y - 2z = 0$, $-w - 3x - 4y + 3z = 0$

Exercise 7.2.8. Recall that Theorem 4.2.14 establishes there are at most n eigenvalues of an $n \times n$ symmetric matrix. Adapt the proof of that theorem, using linear independence, to prove there are at most n eigenvalues of an $n \times n$ non-symmetric matrix. (This is an alternative to the given proof of Theorem 7.1.1.)

Exercise 7.2.9. For each diagram, estimate roughly the components of each of the four vectors \mathbf{a} , \mathbf{b} , \mathbf{c} and \mathbf{d} in the basis $\mathcal{P} = \{\mathbf{p}_1, \mathbf{p}_2\}$.





Exercise 7.2.10. Let the three given vectors $\mathbf{b}_1 = (-1, 1, -1)$, $\mathbf{b}_2 = (1, -2, 0)$ and $\mathbf{b}_3 = (0, 4, 5)$ form a basis $\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3\}$ (where these vectors are specified in the standard basis \mathcal{E} of the standard unit vectors \mathbf{e}_1 , \mathbf{e}_2 and \mathbf{e}_3). For each of the following vectors with specified coordinates in basis \mathcal{B} , what is the vector when written in the standard basis?

- (a) $[\mathbf{p}]_{\mathcal{B}} = (1, -1, 2)$
- (b) $[\mathbf{q}]_{\mathcal{B}} = (0, 2, 3)$
- (c) $[\mathbf{r}]_{\mathcal{B}} = (1, -3, -2)$
- (d) $[\mathbf{s}]_{\mathcal{B}} = (1, 2, 1)$
- (e) $[\mathbf{t}]_{\mathcal{B}} = (1/2, -1/2, 1)$
- (f) $[\mathbf{u}]_{\mathcal{B}} = (-1/2, 1/2, -1/2)$
- (g) $[\mathbf{v}]_{\mathcal{B}} = (0, 1/2, -1/2)$
- (h) $[\mathbf{w}]_{\mathcal{B}} = (-0.7, 0.5, 1.1)$
- (i) $[\mathbf{x}]_{\mathcal{B}} = (0.2, -0.1, 0.9)$
- (j) $[\mathbf{y}]_{\mathcal{B}} = (2.1, -0.2, 0.1)$

Exercise 7.2.11. Repeat Exercise 7.2.10 but with the three basis vectors $\mathbf{b}_1 = (6, 2, 1)$, $\mathbf{b}_2 = (-2, -1, -2)$ and $\mathbf{b}_3 = (-3, -1, 5)$.

Exercise 7.2.12. Let the two given vectors $\mathbf{b}_1 = (1, -2, 2)$ and $\mathbf{b}_2 = (1, -1, -1)$ form a basis $\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2\}$ for the subspace \mathbb{B} of \mathbb{R}^3 (specified in the standard basis \mathcal{E} of the standard unit vectors \mathbf{e}_1 , \mathbf{e}_2 and \mathbf{e}_3). For each of the following vectors, specified in the standard basis \mathcal{E} , what is the vector when written in the basis \mathcal{B} ? if possible.

- (a) $[\mathbf{p}]_{\mathcal{E}} = (0, 1, 3)$
- (b) $[\mathbf{q}]_{\mathcal{E}} = (-4, 9, -11)$
- (c) $[\mathbf{r}]_{\mathcal{E}} = (-4, 7, -5)$
- (d) $[\mathbf{s}]_{\mathcal{E}} = (0, 2, -4)$
- (e) $[\mathbf{t}]_{\mathcal{E}} = (0, -2, 5)$
- (f) $[\mathbf{u}]_{\mathcal{E}} = (-2/3, 0, 8/3)$
- (g) $[\mathbf{v}]_{\mathcal{E}} = (-8/3, 19/3, -19/3)$
- (h) $[\mathbf{w}]_{\mathcal{E}} = (-10/3, 5, -5/3)$
- (i) $[\mathbf{x}]_{\mathcal{E}} = (-0.3, 0.4, 0)$
- (j) $[\mathbf{y}]_{\mathcal{E}} = (0.5, -0.7, 0.1)$

Exercise 7.2.13. Repeat Exercise 7.2.12 but with the two basis vectors $\mathbf{b}_1 = (-2, 3, -1)$ and $\mathbf{b}_2 = (0, -1, 3)$.

Exercise 7.2.14. Let the three vectors $\mathbf{b}_1 = (-1, -2, 5, 3, -2)$, $\mathbf{b}_2 = (-2, -2, 2, -1, -1)$ and $\mathbf{b}_3 = (-4, 6, -4, 2, -1)$ form a basis $\mathcal{B} = \{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3\}$ for the subspace \mathbb{B} in \mathbb{R}^5 (specified in the standard basis \mathcal{E}). For each of the following vectors, use Matlab/Octave to find the requested coordinates, if possible.

- (a) Find $[\mathbf{p}]_{\mathcal{E}}$ when $[\mathbf{p}]_{\mathcal{B}} = (2, 2, -1)$.
- (b) Find $[\mathbf{q}]_{\mathcal{E}}$ when $[\mathbf{q}]_{\mathcal{B}} = (-5, 0, -2)$.
- (c) Find $[\mathbf{r}]_{\mathcal{E}}$ when $[\mathbf{r}]_{\mathcal{B}} = (0, 3, 3)$.
- (d) Find $[\mathbf{s}]_{\mathcal{E}}$ when $[\mathbf{s}]_{\mathcal{B}} = (-1, 5, 0)$.
- (e) Find $[\mathbf{t}]_{\mathcal{B}}$ when $[\mathbf{t}]_{\mathcal{E}} = (-31, 26, -5, 19, -14)$.
- (f) Find $[\mathbf{u}]_{\mathcal{B}}$ when $[\mathbf{u}]_{\mathcal{E}} = (-1, 6, 4, 14, -5)$.
- (g) Find $[\mathbf{v}]_{\mathcal{B}}$ when $[\mathbf{v}]_{\mathcal{E}} = (-21, 18, -7, 9, -8)$.
- (h) Find $[\mathbf{w}]_{\mathcal{B}}$ when $[\mathbf{w}]_{\mathcal{E}} = (-0.2, -0.6, 1.8, 1.3, -0.7)$.
- (i) Find $[\mathbf{x}]_{\mathcal{B}}$ when $[\mathbf{x}]_{\mathcal{E}} = (0.7, -0.4, -0.3, -0.3, 1.2)$.
- (j) Find $[\mathbf{y}]_{\mathcal{B}}$ when $[\mathbf{y}]_{\mathcal{E}} = (4.8, -3.8, 0.8, -2.6, 2.1)$.

Exercise 7.2.15. Repeat Exercise 7.2.14 but with basis vectors $\mathbf{b}_1 = (-3, 8, -9, -1, 1)$, $\mathbf{b}_2 = (10, -20, 14, -7, 2)$ and $\mathbf{b}_3 = (-1, -2, 5, 3, -2)$ (specified in the standard basis \mathcal{E}).



7.3 Diagonalisation identifies the transformation

Section Contents

7.3.1	Solve systems of differential equations	648
7.3.2	Exercises	659

Population modelling Recall that this Chapter 7 started by introducing the dynamics of two interacting species of animals. Recall we let $y(t)$ and $z(t)$ be the number of female animals in each of the species at time t (years). Modelling might deduce the populations interact according to the rule that the population one year later is $y(t+1) = 2y(t) - 4z(t)$ and $z(t+1) = -y(t) + 2z(t)$. Then seeking solutions proportional to λ^t led to the eigen-problem

$$\begin{bmatrix} 2 & -4 \\ -1 & 2 \end{bmatrix} \mathbf{x} = \lambda \mathbf{x}.$$

This section introduces an alternate equivalent approach.

The alternate approach invokes non-orthogonal coordinates. Start by writing the population model as a system in terms of vector $\mathbf{y}(t) = (y(t), z(t))$, namely

$$\mathbf{y}(t+1) = \begin{bmatrix} y(t+1) \\ z(t+1) \end{bmatrix} = \begin{bmatrix} 2y - 4z \\ -y + 2z \end{bmatrix} = \begin{bmatrix} 2 & -4 \\ -1 & 2 \end{bmatrix} \mathbf{y}.$$

Now let's ask if there is a basis $\mathcal{P} = \{\mathbf{p}_1, \mathbf{p}_2\}$ for the yz -plane that simplifies the system? In such a basis every vector may be written as $\mathbf{y} = Y_1 \mathbf{p}_1 + Y_2 \mathbf{p}_2$ for some components Y_1 and Y_2 (where $(Y_1, Y_2) = \mathbf{Y} = [\mathbf{y}]_{\mathcal{P}}$, but to simplify writing we use the symbol \mathbf{Y} not $[\mathbf{y}]_{\mathcal{P}}$). Write this relation as the matrix-vector product $\mathbf{y} = P\mathbf{Y}$ where matrix $P = [\mathbf{p}_1 \ \mathbf{p}_2]$ and vector $\mathbf{Y} = (Y_1, Y_2)$. The populations \mathbf{y} depends upon time t , and hence so does \mathbf{Y} ; that is, $\mathbf{y}(t) = P\mathbf{Y}(t)$. Substitute this into the system of equations:

$$\mathbf{y}(t+1) = P\mathbf{Y}(t+1) = \begin{bmatrix} 2 & -4 \\ -1 & 2 \end{bmatrix} P\mathbf{Y}.$$

Multiply both sides by P^{-1} (which exists by linear independence of the columns, Theorem 7.2.31k) to give

$$\mathbf{Y}(t+1) = \underbrace{P^{-1} \begin{bmatrix} 1 & -4 \\ -1 & 1 \end{bmatrix}}_{P^{-1}AP} P\mathbf{Y}.$$

The question then becomes, for a given square matrix A , such as this, can we find a matrix P such that $P^{-1}AP$ is somehow simple? The answer is yes: using eigenvalues and eigenvectors, most of the time the product $P^{-1}AP$ can be made into a simple diagonal matrix.

Recall that (section 4.2.2) for a symmetric matrix A we could always factor $A = VDV^T = VDV^{-1}$ for orthogonal matrix V and diagonal matrix D : thus a symmetric matrix is always orthogonally diagonalisable. For non-symmetric matrices, a diagonalisation mostly (although not always) can be done with the difference being we need an invertible matrix, typically called P , instead of the orthogonal matrix V . Such a matrix is termed ‘diagonalisable’ instead of ‘orthogonally diagonalisable’.

Definition 7.3.1. An $n \times n$ square matrix A is **diagonalisable** if there exists a diagonal matrix D and an invertible matrix P such that $A = PDP^{-1}$, equivalently $AP = PD$ or $P^{-1}AP = D$.

Example 7.3.2. (a) Show that $A = \begin{bmatrix} 0 & 1 \\ 2 & -1 \end{bmatrix}$ is diagonalisable by matrix $P = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}$.

Solution: First find the 2×2 inverse (Theorem 3.2.6)

$$P^{-1} = \frac{1}{3} \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix}.$$

Second, compute the product

$$P^{-1}AP = P^{-1} \begin{bmatrix} 1 & 2 \\ 1 & -4 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 3 & 0 \\ 0 & -6 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & -2 \end{bmatrix}.$$

Since this product is diagonal, $\text{diag}(1, -2)$, the matrix A is diagonalisable.

(b) $B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ is not diagonalisable.

Solution: Assume B is diagonalisable by the invertible matrix $P = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$. Being invertible, P has inverse $P^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$ (Theorem 3.2.6). Then the product

$$P^{-1}BP = P^{-1} \begin{bmatrix} c & d \\ 0 & 0 \end{bmatrix} = \frac{1}{ad-bc} \begin{bmatrix} cd & d^2 \\ -c^2 & -cd \end{bmatrix}.$$

For this matrix ($P^{-1}BP$) to be diagonal requires the off-diagonal elements to be zero: $d^2 = -c^2 = 0$. This requires both $c = d = 0$, but then the determinant $ad - bc = 0 - 0 = 0$ and so matrix P is not invertible (Theorem 3.2.6). This contradiction means that matrix B is not diagonalisable.

(c) Is matrix $C = \begin{bmatrix} 1.2 & 3.2 & 2.3 \\ 2.2 & -0.5 & -2.2 \end{bmatrix}$ diagonalisable?

Solution: No, as it is not a square matrix. (Perhaps an SVD could answer the needs of whatever problem led to this question.)

Example 7.3.3. Example 7.3.2a showed that matrix $P = \begin{bmatrix} 1 & -1 \\ 1 & 2 \end{bmatrix}$ diagonalises matrix $A = \begin{bmatrix} 0 & 1 \\ 2 & -1 \end{bmatrix}$ to matrix $D = \text{diag}(1, -2)$. As a prelude to the next Theorem 7.3.4, show that the columns of P are eigenvectors of A .

Solution: Invoke the original Definition 4.1.1 of an eigenvector for a matrix.

- The first column of P is $\mathbf{p}_1 = (1, 1)$. Multiplying $A\mathbf{p}_1 = (0+1, 2-1) = (1, 1) = 1\mathbf{p}_1$ so vector \mathbf{p}_1 is an eigenvector of A corresponding to the eigenvalue 1. Correspondingly, this eigenvalue is the first entry in the diagonal D .
- The second column of P is $\mathbf{p}_2 = (-1, 2)$. Multiplying $A\mathbf{p}_2 = (0+2, -2-2) = (2, -4) = -2\mathbf{p}_2$ so vector \mathbf{p}_2 is an eigenvector of A corresponding to the eigenvalue -2 . Correspondingly, this eigenvalue is the second entry in the diagonal D .

Theorem 7.3.4. An $n \times n$ square matrix A is diagonalisable if and only if A has n linearly independent eigenvectors. If A is diagonalisable, with diagonal matrix $D = P^{-1}AP$, then the diagonal entries of D are eigenvalues, and the columns of P are corresponding eigenvectors.

Proof. First, let matrix A be diagonalisable by invertible P and diagonal D . Write $P = [\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n]$ in terms of its columns, and let $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ in terms of its diagonal entries. Then $AP = PD$ becomes

$$A[\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n] = [\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n] \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix}.$$

Multiplying the matrix-column products on both sides gives

$$[A\mathbf{p}_1 \ A\mathbf{p}_2 \ \cdots \ A\mathbf{p}_n] = [\lambda_1\mathbf{p}_1 \ \lambda_2\mathbf{p}_2 \ \cdots \ \lambda_n\mathbf{p}_n].$$

Equating columns implies $A\mathbf{p}_1 = \lambda_1\mathbf{p}_1$, $A\mathbf{p}_2 = \lambda_2\mathbf{p}_2$, ..., $A\mathbf{p}_n = \lambda_n\mathbf{p}_n$. As the matrix P is invertible, all its columns must be non-zero (Theorem 6.2.5a). Hence $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$ are eigenvectors of matrix A corresponding to eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Since matrix P is invertible, Theorem 7.2.31k implies the columns vectors $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$ are linearly independent.

Second, suppose matrix A has n linearly independent eigenvectors $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$ with corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Then follow the above argument backwards to deduce $AP = PD$ for invertible matrix $P = [\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n]$, and hence A is diagonalisable. Lastly, in these arguments, P is the matrix of eigenvectors and the diagonal of D are the corresponding eigenvalues, as required. \square

Example 7.3.5. Recall Example 7.0.3 found the triangular matrix

$$A = \begin{bmatrix} -3 & 2 & 0 \\ 0 & -4 & 2 \\ 0 & 0 & 4 \end{bmatrix}$$

has eigenvalues $-3, -4$ and 4 (from its diagonal) and corresponding eigenvectors are proportional to $(1, 0, 0)$, $(-2, 1, 0)$ and $(\frac{1}{14}, \frac{1}{4}, 1)$. Is matrix A diagonalisable?

Solution: These three eigenvectors are linearly independent as they correspond to distinct eigenvalues (Theorem 7.2.10). Hence the matrix is diagonalisable.

The previous paragraph answer the question. But further, forming these eigenvectors into the columns of matrix

$$P = \begin{bmatrix} 1 & -2 & \frac{1}{14} \\ 0 & 1 & \frac{1}{4} \\ 0 & 0 & 1 \end{bmatrix},$$

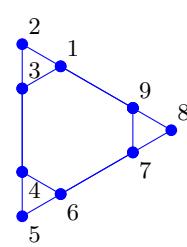
we know (Theorem 7.3.4) $P^{-1}AP = \text{diag}(-3, -4, 4)$ where the eigenvalues appear in the same order as that of the eigenvectors in P .

One may check this by hand or with Matlab/Octave. Enter the matrices with

```
A=[-3 2 0;0 -4 2;0 0 4]
P=[1 -2 1/14;0 1 1/4;0 0 1]
```

then compute $D = P^{-1}AP$ with $D=P\backslash A*P$ to find as required the following diagonal result

```
D =
-3.0000    0.0000    0.0000
 0.0000   -4.0000    0.0000
 0.0000    0.0000    4.0000
```



Example 7.3.6. Recall the Sierpinski network of Example 4.1.15 (shown in the margin). Is the 9×9 matrix encoding the network diagonalisable?

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

Solution: In that example we used Matlab/Octave command `[V,D]=eig(A)` to compute a matrix of eigenvectors V and the corresponding diagonal matrix of eigenvalues D where (2 d.p.)

```
V =
-0.41 0.51 -0.16 -0.21 -0.45 0.18 -0.40 0.06 0.33
 0.00 -0.13 0.28 0.63 0.13 -0.18 -0.58 -0.08 0.33
 0.41 -0.20 -0.49 -0.42 0.32 0.01 -0.36 -0.17 0.33
-0.41 -0.11 0.52 -0.42 0.32 0.01 0.14 -0.37 0.33
-0.00 -0.18 -0.26 0.37 -0.22 0.51 0.36 -0.46 0.33
 0.41 0.53 0.07 0.05 -0.10 -0.51 0.33 -0.23 0.33
-0.41 -0.39 -0.36 0.05 -0.10 -0.51 0.25 0.31 0.33
 0.00 0.31 -0.03 0.16 0.55 0.34 0.22 0.55 0.33
 0.41 -0.33 0.42 -0.21 -0.45 0.18 0.03 0.40 0.33

D =
-5.00 0 0 0 0 0 0 0 0
 0 -4.30 0 0 0 0 0 0 0
 0 0 -4.30 0 0 0 0 0 0
 0 0 0 -3.00 0 0 0 0 0
 0 0 0 0 -3.00 0 0 0 0
 0 0 0 0 0 -3.00 0 0 0
 0 0 0 0 0 0 -0.70 0 0
 0 0 0 0 0 0 0 -0.70 0
 0 0 0 0 0 0 0 0 -0.00
```

Since matrix A is symmetric, Matlab/Octave computes for us an orthogonal matrix V with columns eigenvectors. Since V is orthogonal its columns are orthonormal and hence its columns are linearly independent (Theorem 7.2.6). Since there exist nine linearly independent eigenvectors, the matrix A is diagonalisable. Further, the product $V^{-1}AV = D$ for the above diagonal matrix D of eigenvalues in the order of the eigenvectors in V . (Also, since V is orthogonal, $V^TAV = D$.)

■

Example 7.3.7. Recall Example 7.1.9 found eigenvalues and corresponding eigenspaces for various matrices. Revisit these cases and show none of the matrices are diagonalisable.

- (a) Matrix $A = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix}$ had one eigenvalue $\lambda = 3$ with multiplicity two and corresponding eigenspace $\mathbb{E}_3 = \text{span}\{(1, 0)\}$. This

matrix is not diagonalisable as it has only one linearly independent eigenvector, such as $(1, 0)$ or any non-zero multiple, and it needs two.

(b) Matrix $B = \begin{bmatrix} -1 & 1 & -2 \\ -1 & 0 & -1 \\ 0 & -3 & 1 \end{bmatrix}$ has eigenvalues $\lambda = -2$ (multiplicity one) and $\lambda = 1$ (multiplicity two). The corresponding eigenspaces are $\mathbb{E}_{-2} = \text{span}\{(1, 1, 1)\}$ and $\mathbb{E}_1 = \text{span}\{(-1, 0, 1)\}$. Thus the matrix has only two linearly independent eigenvectors, one from each eigenspace, and it needs three to be diagonalisable.

(c) Matrix $C = \begin{bmatrix} -1 & 0 & -2 \\ 0 & -3 & 2 \\ 0 & -2 & 1 \end{bmatrix}$ has only the eigenvalue $\lambda = -1$ with multiplicity three. The corresponding eigenspace $\mathbb{E}_{-1} = \text{span}\{(1, 0, 0)\}$. With only one linearly independent eigenvector, the matrix is not diagonalisable. ■

Example 7.3.8. Use the results of Example 7.1.10 to show the following matrix is diagonalisable:

$$A = \begin{bmatrix} 0 & 3 & 0 & 0 & 0 \\ 1 & 0 & 3 & 0 & 0 \\ 0 & 1 & 0 & 3 & 0 \\ 0 & 0 & 1 & 0 & 3 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

Solution: Example 7.1.10 derived the five eigenvalues are $\lambda = 0, \pm\sqrt{3}, \pm 3$, all of multiplicity one. Further, the corresponding eigenspaces are

$$\begin{aligned} \mathbb{E}_0 &= \text{span}\{(9, 0, -3, 0, 1)\}, \\ \mathbb{E}_{\pm\sqrt{3}} &= \text{span}\{(-9, \mp 3\sqrt{3}, 0, \pm\sqrt{3}, 1)\}, \\ \mathbb{E}_{\pm 3} &= \text{span}\{(9, \pm 9, 6, \pm 3, 1)\}. \end{aligned}$$

Here there are five linearly independent eigenvectors, one from each distinct eigenspace (Theorem 7.2.10). Since A is a 5×5 matrix it is thus diagonalisable. Further, Theorem 7.3.4 establishes that the matrix formed from the columns of the five eigenvectors will be a possible matrix ⁷

$$P = \begin{bmatrix} 9 & -9 & -9 & 9 & 9 \\ 0 & -3\sqrt{3} & 3\sqrt{3} & 9 & -9 \\ -3 & 0 & 0 & 6 & 6 \\ 0 & \sqrt{3} & -\sqrt{3} & 3 & -3 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}.$$

⁷ One could also scale each column of P by a different arbitrary non-zero constant, and the diagonalisation still holds.
© AJ Roberts, ORCID:0000-0001-8930-1552, July 26, 2016

Lastly, Theorem 7.3.4 establishes the diagonal matrix $P^{-1}AP = D = \text{diag}(0, \sqrt{3}, -\sqrt{3}, 3, -3)$ is that of the eigenvalues in the order corresponding to the eigenvectors in P . ■

These examples illustrate a widely useful property. The 5×5 matrix in Example 7.3.8 has five distinct eigenvalues whose corresponding eigenvectors are necessarily linearly independent (Theorem 7.2.10) and so diagonalise the matrix (Theorem 7.3.4). The 3×3 matrix in Example 7.3.5 has three distinct eigenvalues whose corresponding eigenvectors are necessarily linearly independent (Theorem 7.2.10) and so diagonalise the matrix (Theorem 7.3.4). However, the matrices of Examples 7.3.6 and 7.3.7 have repeated eigenvalues—eigenvalues of multiplicity two or more—and these matrices may (Example 7.3.6) or may not (Example 7.3.7) be diagonalisable. The following theorem confirms that matrices with as many *distinct* eigenvalues as the size of the matrix are always diagonalisable.

Theorem 7.3.9. *If an $n \times n$ square matrix A has n distinct eigenvalues, then A is diagonalisable. Consequently, and allowing complex eigenvalues, a non-diagonalisable matrix must be non-symmetric and must have at least one repeated eigenvalue (an eigenvalue with multiplicity two or more).*

Proof. First, let $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ be eigenvectors corresponding to the n distinct eigenvalues of matrix A . (Recall that Theorem 7.1.1 establishes that there cannot be more than n eigenvalues.) As the corresponding eigenvalues are distinct, Theorem 7.2.10 establishes that these eigenvectors are linearly independent. Theorem 7.3.4 then establishes the matrix A is diagonalisable.

Second, the converse of the first statement in the theorem then also holds. Since an $n \times n$ matrix has n eigenvalues, when counted accordingly to multiplicity and allowing for complex eigenvalues (Procedure 7.1.8), a non-diagonalisable matrix must have at least one repeated eigenvalue. But further, by Theorem 4.2.18 a symmetric matrix is always diagonalisable: hence a non-diagonalisable matrix must also be non-symmetric. □

Example 7.3.10. From the given information, are the matrices diagonalisable?

- (a) The only eigenvalues of a 4×4 matrix are 1.8, -3, 0.4 and 3.2.

Solution: Theorem 7.3.9 implies the matrix must be diagonalisable.

- (b) The only eigenvalues of a 5×5 matrix are 1.8, -3, 0.4 and 3.2.

Solution: Here there are only four distinct eigenvalues of the 5×5 matrix. Theorem 7.3.9 does not apply as the

precondition that there be five distinct eigenvalues is not met: the matrix may or may not be diagonalisable—it is unknowable on this information.

- (c) The only eigenvalues of a 3×3 matrix are 1.8, -3, 0.4 and 3.2.

Solution: An error has been made in determining the eigenvalues as a 3×3 matrix has at most three distinct eigenvalues (Theorem 7.1.1). Because of the error, we cannot answer.

■

Example 7.3.11. Matlab/Octave computes the eigenvalues of matrix

$$A = \begin{bmatrix} -1 & 2 & -2 & 1 & -2 \\ -3 & -1 & -2 & 5 & 6 \\ 3 & 1 & 6 & -2 & -1 \\ 1 & 1 & 2 & 1 & -1 \\ 7 & 5 & -3 & 0 & 0 \end{bmatrix}$$

via `eig(A)` and reports them to be (2 d.p.)

```
ans =
-3.45 + 3.50i
-3.45 - 3.50i
5.00
5.00
1.91
```

Is the matrix diagonalisable?

Solution: The matrix appears to have only four distinct eigenvalues (two of them complex valued), and so on the given information Theorem 7.3.9 cannot determine whether the matrix is diagonalisable or not.

However, upon reporting the eigenvalues to four decimal places we find the two eigenvalues of 5.00 (2 d.p.) more precisely are two separate eigenvalues of 5.0000 and 4.9961. Hence this matrix has five distinct eigenvalues and so Theorem 7.3.9 implies the matrix is diagonalisable.⁸

■

Theorem 7.3.12. *The dimension of the eigenspace \mathbb{E}_{λ_j} corresponding to a given eigenvalue λ_j of a matrix A is less than or equal to the*

⁸ Nonetheless, in an application where errors are significant then the matrix may be effectively non-diagonalisable. Such effective non-diagonalisability is indicated by poor conditioning of the matrix of eigenvectors which here has the poor `rcond` of 0.0004 (Procedure 2.4).

multiplicity of λ_j ; that is, $1 \leq \dim \mathbb{E}_{\lambda_j} \leq \text{multiplicity of } \lambda_j$. (Recall that, from Definition 4.1.11, for a symmetric matrix, $\dim \mathbb{E}_{\lambda_j} = \text{multiplicity of } \lambda_j$.)

Proof. Suppose λ_j is an eigenvalue of matrix A and $\dim \mathbb{E}_{\lambda_j} = p < n$ (the case $p = n$ is proved by Exercise 7.3.7). Then choose p orthonormal vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p$ to span \mathbb{E}_{λ_j} : vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p$ are eigenvectors as they are in the eigenspace. Let

$$P = [\mathbf{v}_1 \ \cdots \ \mathbf{v}_p \ \mathbf{w}_{p+1} \ \cdots \ \mathbf{w}_n]$$

be any orthogonal matrix with $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p$ as its first p columns. Equivalently, write as the partitioned matrix $P = [V \ W]$ for corresponding $n \times p$ and $n \times (n-p)$ matrices. Since P is orthogonal,

its inverse $P^{-1} = P^T = \begin{bmatrix} V^T \\ W^T \end{bmatrix}$. Since the columns of V are eigenvectors corresponding to eigenvalue λ_j , $AV = A[\mathbf{v}_1 \ \cdots \ \mathbf{v}_p] =$

$$[A\mathbf{v}_1 \ \cdots \ A\mathbf{v}_p] = [\lambda_j \mathbf{v}_1 \ \cdots \ \lambda_j \mathbf{v}_p] = \lambda_j [\mathbf{v}_1 \ \cdots \ \mathbf{v}_p] = \lambda_j V.$$

Now consider

$$\begin{aligned} P^{-1}AP &= \begin{bmatrix} V^T \\ W^T \end{bmatrix} A [V \ W] = \begin{bmatrix} V^T AV & V^T AW \\ W^T AV & W^T AW \end{bmatrix} \\ &= \begin{bmatrix} \lambda_j V^T V & V^T AW \\ \lambda_j W^T V & W^T AW \end{bmatrix} = \begin{bmatrix} \lambda_j I_p & V^T AW \\ O & W^T AW \end{bmatrix} \end{aligned}$$

where the last equality follows from the orthonormality of columns of $P = [V \ W]$. Then the characteristic polynomial of matrix A becomes

$$\begin{aligned} &\det(A - \lambda I_n) \\ &= \det(PP^{-1}APP^{-1} - \lambda PP^{-1}) \quad (\text{as } PP^{-1} = I_n) \\ &= \det[P(P^{-1}AP - \lambda I_n)P^{-1}] \\ &= \det P \det(P^{-1}AP - \lambda I_n) \det(P^{-1}) \quad (\text{product Thm. 6.1.12}) \\ &= \det P \det(P^{-1}AP - \lambda I_n) \frac{1}{\det P} \quad (\text{inverse Thm. 6.1.23}) \\ &= \det(P^{-1}AP - \lambda I_n) \\ &= \det \begin{bmatrix} (\lambda_j - \lambda)I_p & V^T AW \\ O & W^T AW - \lambda I_{n-p} \end{bmatrix} \quad (\text{by above } P^{-1}AP) \\ &= (\lambda_j - \lambda)^p \det(W^T AW - \lambda I_{n-p}) \end{aligned}$$

by p successive first column expansions of the determinant (Theorem 6.2.20). Because of the factor $(\lambda_j - \lambda)^p$ in the characteristic polynomial of A , the eigenvalue λ_j must have multiplicity of at least $p = \dim \mathbb{E}_{\lambda_j}$ (there may be more factors of $(\lambda_j - \lambda)$ hidden within $\det(W^T AW - \lambda I_{n-p})$). \square

Example 7.3.13. Show the following matrix has one eigenvalue of multiplicity three, and the corresponding eigenspace has dimension two:

$$A = \begin{bmatrix} 0 & 5 & 6 \\ -8 & 22 & 24 \\ 6 & -15 & -16 \end{bmatrix}$$

Solution: Find eigenvalues via the characteristic polynomial

$$\begin{aligned}\det(A - \lambda I) &= \begin{vmatrix} -\lambda & 5 & 6 \\ -8 & 22 - \lambda & 24 \\ 6 & -15 & -16 - \lambda \end{vmatrix} \\ &= -\lambda(22 - \lambda)(-16 - \lambda) + 5 \cdot 24 \cdot 6 + 6(-8)(-15) \\ &\quad - 6(22 - \lambda)6 + \lambda 24(-15) - 5(-8)(-16 - \lambda) \\ &= \dots \\ &= -\lambda^3 + 6\lambda^2 - 12\lambda + 8 \\ &= -(\lambda - 2)^3.\end{aligned}$$

This characteristic polynomial is zero only for eigenvalue $\lambda = 2$ which is of multiplicity three.

The corresponding eigenspace comes from solving $(A - \lambda I)\mathbf{x} = \mathbf{0}$ which here is

$$\begin{bmatrix} -2 & 5 & 6 \\ -8 & 20 & 24 \\ 6 & -15 & -18 \end{bmatrix} \mathbf{x} = \mathbf{0}.$$

Observe that the second row is just four times the first row, and the third row is $(-3) \times$ the first row, hence all three equations in this system are equivalent to just the one from the first row, namely $-2x_1 + 5x_2 + 6x_3 = 0$. A general solution of this equation is $x_1 = \frac{5}{2}x_2 - 3x_3$. That is all solutions are $\mathbf{x} = (\frac{5}{2}x_2 - 3x_3, x_2, x_3) = x_2(\frac{5}{2}, 1, 0) + x_3(-3, 0, 1)$. Hence all solutions form the two dimensional eigenspace $E_2 = \text{span}\{(\frac{5}{2}, 1, 0), (-3, 0, 1)\}$.

■

Example 7.3.14. Use Matlab/Octave to find the eigenvalues and the dimension of the eigenspaces of the matrix

$$B = \begin{bmatrix} 344 & -1165 & -149 & -1031 & 1065 & -2816 \\ 90 & -306 & -38 & -272 & 280 & -742 \\ -45 & 140 & 12 & 117 & -115 & 302 \\ 135 & -470 & -70 & -421 & 445 & -1175 \\ -165 & 555 & 67 & 493 & -506 & 1338 \\ -105 & 360 & 48 & 322 & -335 & 886 \end{bmatrix}.$$

Solution: In Matlab/Octave enter the matrix with

```
B=[344 -1165 -149 -1031 1065 -2816
    90 -306 -38 -272 280 -742
   -45  140  12  117 -115  302
   135 -470 -70 -421 445 -1175
  -165  555  67  493 -506 1338
  -105  360  48  322 -335  886]
```



Then $[V, D] = \text{eig}(B)$ computes something like the following (2 d.p.)

```

V =
-0.19  0.19 -0.45  0.75 -0.20  0.15
-0.38  0.38  0.12  0.26 -0.03  0.00
-0.58  0.58  0.08 -0.00 -0.54  0.56
-0.00 -0.00 -0.83  0.28 -0.50  0.52
  0.58 -0.58 -0.29 -0.45 -0.64  0.63
  0.38 -0.38  0.08 -0.29 -0.04  0.03

D =
  4.00    0    0    0    0    0
    0  4.00    0    0    0    0
    0    0  4.00    0    0    0
    0    0    0 -1.00    0    0
    0    0    0    0 -1.00    0
    0    0    0    0    0 -1.00

```

Evidently, the matrix B has two eigenvalues, $\lambda = 4$ and $\lambda = -1$, both of multiplicity three (although due to round-off error Matlab/Octave will report these with errors of about 10^{-5} (section 7.1.2)—possibly complex errors in which case ignore small complex parts). For each eigenvalue Matlab/Octave reports three corresponds columns of V containing corresponding eigenvectors. These eigenvectors do span the eigenspace, but are not necessarily linearly independent.

- For eigenvalue $\lambda = 4$ the first two columns of V are clearly the negative of each other, and so are essentially the same eigenvector. The third column of V is clearly not proportional to the first two columns and so is linearly independent (Theorem 7.2.8). Thus Matlab/Octave has computed only two linearly independent eigenvectors, either the first and third column, or the second and third column. Consequently, the dimension $\dim \mathbb{E}_4 = 2$.

One can confirm this dimensionality by computing the singular values of the first three columns of V with `svd(V(:,1:3))` to find they are 1.4278, 0.9806 and 0.0000. The two non-zero singular values indicate the dimension of the span is two (Theorem 3.4.18).

- For eigenvalue $\lambda = -1$ the last three columns of V look linearly independent and so we suspect the eigenspace dimension $\dim \mathbb{E}_{-1} = 3$.

To confirm the eigenspace dimensionality (Theorem 3.4.18) compute `svd(V(:,4:6))` to find the three singular values are 1.4136, 1.0001 and 0.0414. Since all three singular values are non-zero, $\dim \mathbb{E}_{-1} = 3$.

Matlab/Octave may produce for you a quite different matrix V of eigenvectors (possibly with complex parts). As discussed by section 7.1.2, repeated eigenvalues are very sensitive and this sensitivity

means small variations in the hidden Matlab/Octave algorithm may produce quite large changes in the matrix V for repeated eigenvalues. However, each eigenspace spanned by the appropriate columns of V is robust. ■

7.3.1 Solve systems of differential equations

Population modelling The population modelling seen so far (section 7.1.3) expressed the changes of the population over discrete intervals in time. One such example is to describe the population numbers year by year. The alternative is to model the changes in the population *continuously* in time. This alternative invokes and analyses differential equations.

Let's start with a continuous time version of the population modelling discussed at the start of this Chapter 7. Let two species interact continuously in time with populations $y(t)$ and $z(t)$ at time t (years). Suppose they interact according to differential equations $dy/dt = y - 4z$ and $dz/dt = -y + z$ (instead of discrete time equations $y(t+1) = \dots$ and $z(t+1) = \dots$). Analogous to the start of this Section 7.3, we now ask the following question: is there a matrix transformation to new variables, the vector $\mathbf{Y}(t)$, such that $(y, z) = \mathbf{y} = P\mathbf{Y}$ where the differential equations for \mathbf{Y} are simple?

- First, form the differential equations into a matrix-vector system:

$$\begin{bmatrix} dy/dt \\ dz/dt \end{bmatrix} = \begin{bmatrix} y - 4z \\ -y + z \end{bmatrix} = \begin{bmatrix} 1 & -4 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix}.$$

So using vector $\mathbf{y} = (y, z)$, this system is

$$\frac{d\mathbf{y}}{dt} = A\mathbf{y} \quad \text{for matrix } A = \begin{bmatrix} 1 & -4 \\ -1 & 1 \end{bmatrix}.$$

- Second, see what happens when we transform to some, as yet unknown, new variables $\mathbf{Y}(t)$ such that $\mathbf{y} = P\mathbf{Y}$ for some invertible matrix P . Under such a transform: $\frac{d\mathbf{y}}{dt} = \frac{d}{dt}P\mathbf{Y} = P\frac{d\mathbf{Y}}{dt}$; and $A\mathbf{y} = AP\mathbf{Y}$. Hence substituting such an assumed transformation into the differential equations leads to

$$P\frac{d\mathbf{Y}}{dt} = AP\mathbf{Y}, \quad \text{that is } \frac{d\mathbf{Y}}{dt} = (P^{-1}AP)\mathbf{Y}.$$

To simplify this system for \mathbf{Y} , we diagonalise the matrix on the right-hand side. The procedure is to choose the columns of P to be eigenvectors of the matrix (Theorem 7.3.4).

- Third, here the matrix $A = \begin{bmatrix} 1 & -4 \\ -1 & 1 \end{bmatrix}$ has characteristic polynomial $\det(A - \lambda I) = (1 - \lambda)^2 - 4$. This is zero for $(1 - \lambda)^2 = 4$, that is, $(1 - \lambda) = \pm 2$. Hence the eigenvalues $\lambda = 1 \pm 2 = 3, -1$.

- For eigenvalue $\lambda_1 = 3$ the corresponding eigenvectors satisfy

$$(A - \lambda_1 I) \mathbf{p}_1 = \begin{bmatrix} -2 & -4 \\ -1 & -2 \end{bmatrix} \mathbf{p}_1 = \mathbf{0},$$

with general solution $\mathbf{p}_1 \propto (2, -1)$.

- For eigenvalue $\lambda_2 = -1$ the corresponding eigenvectors satisfy

$$(A - \lambda_2 I) \mathbf{p}_2 = \begin{bmatrix} 2 & -4 \\ -1 & 2 \end{bmatrix} \mathbf{p}_2 = \mathbf{0},$$

with general solution $\mathbf{p}_2 \propto (2, 1)$.

Thus setting transformation matrix (any scalar multiple of the two columns would also work)

$$P = \begin{bmatrix} 2 & 2 \\ -1 & 1 \end{bmatrix} \implies \frac{d\mathbf{Y}}{dt} = \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix} \mathbf{Y}.$$

- Fourth, having diagonalised the matrix, expand this diagonalised set of differential equations to write this system in terms of components:

$$\frac{dY_1}{dt} = 3Y_1 \quad \text{and} \quad \frac{dY_2}{dt} = -Y_2.$$

Each of these differential equations have well-known exponential solutions, respectively $Y_1 = c_1 e^{3t}$ and $Y_2 = c_2 e^{-t}$, for any constants c_1 and c_2 .

- Lastly, what does this mean for the original problem? From the relation

$$\begin{bmatrix} y \\ z \end{bmatrix} = \mathbf{y} = P\mathbf{Y} = \begin{bmatrix} 2 & 2 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} c_1 e^{3t} \\ c_2 e^{-t} \end{bmatrix} = \begin{bmatrix} 2c_1 e^{3t} + 2c_2 e^{-t} \\ -c_1 e^{3t} + c_2 e^{-t} \end{bmatrix}.$$

That is, a general solution of the original system of differential equations is $y(t) = 2c_1 e^{3t} + 2c_2 e^{-t}$ and $z(t) = -c_1 e^{3t} + c_2 e^{-t}$.

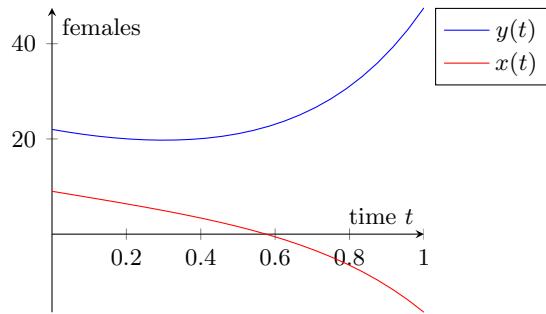
The diagonalisation of the matrix empowers us to solve complicated systems of differential equations as a set of simple systems.

Such a general solution makes predictions. For example, suppose at time zero ($t = 0$) the initial population of female y -animals is 22 and the population of female z -animals is 9. From the above general solution we then know that at time $t = 0$

$$\begin{bmatrix} 22 \\ 9 \end{bmatrix} = \begin{bmatrix} y(0) \\ z(0) \end{bmatrix} = \begin{bmatrix} 2c_1 e^{3 \cdot 0} + 2c_2 e^{-0} \\ -c_1 e^{3 \cdot 0} + c_2 e^{-0} \end{bmatrix} = \begin{bmatrix} 2c_1 + 2c_2 \\ -c_1 + c_2 \end{bmatrix}$$

This determines the coefficients: $2c_1 + 2c_2 = 22$ and $-c_1 + c_2 = 9$. Adding the first to twice the second gives $4c_2 = 40$, that is, $c_2 = 10$. Then either equation determines $c_1 = 1$. Consequently, the particular solution from this initial population is

$$\begin{bmatrix} y \\ z \end{bmatrix} = \begin{bmatrix} 2 \cdot 1e^{3t} + 2 \cdot 10e^{-t} \\ -1e^{3t} + 10e^{-t} \end{bmatrix} = \begin{bmatrix} 2e^{3t} + 20e^{-t} \\ -e^{3t} + 10e^{-t} \end{bmatrix}.$$



The above graph of this solution shows that the population of y -animals grows in time, whereas the population of z -animals crashes and becomes extinct at about time 0.6 years.⁹

The next theorem confirms that the same approach solves general systems of differential equations: it corresponds to Theorem 7.1.16 for discrete dynamics.

Theorem 7.3.15. *Let $n \times n$ square matrix A be diagonalisable by matrix $P = [\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n]$ whose columns are eigenvectors corresponding to eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. Then a general solution $\mathbf{x}(t)$ to the differential equation system $d\mathbf{x}/dt = A\mathbf{x}$ is the linear combination*

$$\mathbf{x}(t) = c_1 \mathbf{p}_1 e^{\lambda_1 t} + c_2 \mathbf{p}_2 e^{\lambda_2 t} + \cdots + c_n \mathbf{p}_n e^{\lambda_n t} \quad (7.4)$$

for arbitrary constants c_1, c_2, \dots, c_n .

Proof. First, instead of finding solutions for $\mathbf{x}(t)$ directly, let's write the differential equations in terms of the alternate basis for \mathbb{R}^n , basis $\mathcal{P} = \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n\}$ (as $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$ are linearly independent). That is, solve for the coordinates $\mathbf{X}(t) = [\mathbf{x}(t)]_{\mathcal{P}}$ with respect to basis \mathcal{P} . Now recall (Theorem 7.2.27) that $\mathbf{X} = [\mathbf{x}]_{\mathcal{P}}$ means that $\mathbf{x} = X_1 \mathbf{p}_1 + X_2 \mathbf{p}_2 + \cdots + X_n \mathbf{p}_n = P\mathbf{X}$. Substitute this into the differential equation $d\mathbf{x}/dt = A\mathbf{x}$ requires $\frac{d}{dt}(P\mathbf{X}) = A(P\mathbf{X})$ which is the same as $P \frac{d\mathbf{X}}{dt} = AP\mathbf{X}$. Since matrix P is invertible, this equation is the same as $\frac{d\mathbf{X}}{dt} = P^{-1}AP\mathbf{X}$. Because the columns of matrix P are eigenvectors, the product $P^{-1}AP$ is the diagonal matrix $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$, hence the system becomes $\frac{d\mathbf{X}}{dt} = D\mathbf{X}$. Because matrix D is diagonal, this is a much simpler system of differential equations. All the rows of the system are

$$\frac{dX_1}{dt} = \lambda_1 X_1, \quad \frac{dX_2}{dt} = \lambda_2 X_2, \quad \dots, \quad \frac{dX_n}{dt} = \lambda_n X_n.$$

Each of these have general solution

$$X_1 = c_1 e^{\lambda_1 t}, \quad X_2 = c_2 e^{\lambda_2 t}, \quad \dots, \quad X_n = c_n e^{\lambda_n t},$$

⁹ After time 0.6 years the differential equation model and its predictions becomes meaningless as there is no biological meaning to a negative number of animals z .

where c_1, c_2, \dots, c_n are arbitrary constants. To rewrite this solution for the original coordinates \mathbf{x} use

$$\begin{aligned}\mathbf{x} &= P\mathbf{X} \\ &= [\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n] \begin{bmatrix} c_1 e^{\lambda_1 t} \\ c_2 e^{\lambda_2 t} \\ \vdots \\ c_n e^{\lambda_n t} \end{bmatrix} \\ &= c_1 \mathbf{p}_1 e^{\lambda_1 t} + c_2 \mathbf{p}_2 e^{\lambda_2 t} + \cdots + c_n \mathbf{p}_n e^{\lambda_n t}\end{aligned}$$

to derive the solution (7.4).

Second, being able to use the constants $\mathbf{c} = (c_1, c_2, \dots, c_n)$ to match any given initial condition shows formula (7.4) is a general solution. Suppose the value of $\mathbf{x}(0)$ is given. Recalling $e^0 = 1$, formula (7.4) evaluated at $t = 0$ requires

$$\begin{aligned}\mathbf{x}(0) &= c_1 \mathbf{p}_1 e^{\lambda_1 0} + c_2 \mathbf{p}_2 e^{\lambda_2 0} + \cdots + c_n \mathbf{p}_n e^{\lambda_n 0} \\ &= c_1 \mathbf{p}_1 + c_2 \mathbf{p}_2 + \cdots + c_n \mathbf{p}_n \\ &= [\mathbf{p}_1 \ \mathbf{p}_2 \ \cdots \ \mathbf{p}_n] \mathbf{c} \\ &= P\mathbf{c}.\end{aligned}$$

If matrix P is invertible, then choose constants $\mathbf{c} = P^{-1}\mathbf{x}(0)$ for any given $\mathbf{x}(0)$. Theorems 7.3.4 and 7.2.31 establish that such an invertible P exists for diagonalisable matrices A . \square

Example 7.3.16. Find (by hand) a general solution to the system of differential equations $\frac{du}{dt} = -2u + 2v$, $\frac{dv}{dt} = u - 2v + w$, and $\frac{dw}{dt} = 2v - 2w$.

Solution: Let vector $\mathbf{u} = (u, v, w)$, and then form the differential equations into the matrix-vector system

$$\begin{aligned}\frac{d\mathbf{u}}{dt} &= \begin{bmatrix} \frac{du}{dt} \\ \frac{dv}{dt} \\ \frac{dw}{dt} \end{bmatrix} = \begin{bmatrix} -2u + 2v \\ u - 2v + w \\ 2v - 2w \end{bmatrix} \\ &= \begin{bmatrix} -2 & 2 & 0 \\ 1 & -2 & 1 \\ 0 & 2 & -2 \end{bmatrix} \begin{bmatrix} u \\ v \\ w \end{bmatrix} = \underbrace{\begin{bmatrix} -2 & 2 & 0 \\ 1 & -2 & 1 \\ 0 & 2 & -2 \end{bmatrix}}_A \mathbf{u}.\end{aligned}$$

To use Theorem 7.3.15 we need eigenvalues and eigenvectors of the matrix A . Here the characteristic polynomial of A is, using (6.1),

$$\begin{aligned}\det(A - \lambda I) &= \det \begin{bmatrix} -2 - \lambda & 2 & 0 \\ 1 & -2 - \lambda & 1 \\ 0 & 2 & -2 - \lambda \end{bmatrix} \\ &= -(2 + \lambda)^3 + 0 + 0 - 0 + 2(2 + \lambda) + 2(2 + \lambda) \\ &= (2 + \lambda)[-(2 + \lambda)^2 + 4]\end{aligned}$$

$$\begin{aligned} &= (2 + \lambda)(-\lambda^2 - 4\lambda) \\ &= -\lambda(\lambda + 2)(\lambda + 4). \end{aligned}$$

This determinant is only zero for eigenvalues $\lambda = 0, -2, -4$.

- For eigenvalue $\lambda = 0$, corresponding eigenvectors \mathbf{p} satisfy

$$(A - 0I)\mathbf{p} = \begin{bmatrix} -2 & 2 & 0 \\ 1 & -2 & 1 \\ 0 & 2 & -2 \end{bmatrix} \mathbf{p} = \mathbf{0}.$$

The last row of this equation requires $p_3 = p_2$, and the first row requires $p_1 = p_2$. Hence all solutions may be written as $\mathbf{p} = (p_2, p_2, p_2)$. Choose any one, say $\mathbf{p} = (1, 1, 1)$.

- For eigenvalue $\lambda = -2$, corresponding eigenvectors \mathbf{p} satisfy

$$(A + 2I)\mathbf{p} = \begin{bmatrix} 0 & 2 & 0 \\ 1 & 0 & 1 \\ 0 & 2 & 0 \end{bmatrix} \mathbf{p} = \mathbf{0}.$$

The first and last rows of this equation require $p_2 = 0$, and the second row requires $p_3 = -p_1$. Hence all solutions may be written as $\mathbf{p} = (p_1, 0, -p_1)$. Choose any one, say $\mathbf{p} = (1, 0, -1)$.

- For eigenvalue $\lambda = -4$, corresponding eigenvectors \mathbf{p} satisfy

$$(A + 4I)\mathbf{p} = \begin{bmatrix} 2 & 2 & 0 \\ 1 & 2 & 1 \\ 0 & 2 & 2 \end{bmatrix} \mathbf{p} = \mathbf{0}.$$

The last row of this equation requires $p_3 = -p_2$, and the first row requires $p_1 = -p_2$. Hence all solutions may be written as $\mathbf{p} = (-p_2, p_2, -p_2)$. Choose any one, say $\mathbf{p} = (-1, 1, -1)$.

With these three distinct eigenvalues, corresponding eigenvectors are linearly independent, and so Theorem 7.3.15 gives a general solution of the differential equations as

$$\begin{bmatrix} u \\ v \\ w \end{bmatrix} = c_1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} e^{0t} + c_2 \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} e^{-2t} + c_3 \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} e^{-4t}.$$

That is, $u(t) = c_1 + c_2 e^{-2t} - c_3 e^{-4t}$, $v(t) = c_1 + c_3 e^{-4t}$, and $w(t) = c_1 - c_2 e^{2t} - c_3 e^{-4t}$ for any constants c_1 , c_2 and c_3 .

■

Example 7.3.17. Use the general solution derived in Example 7.3.16 to predict the solution of the differential equations $\frac{du}{dt} = -2u + 2v$, $\frac{dv}{dt} = u - 2v + w$, and $\frac{dw}{dt} = 2v - 2w$ given the initial conditions that $u(0) = v(0) = 0$ and $w(0) = 4$.

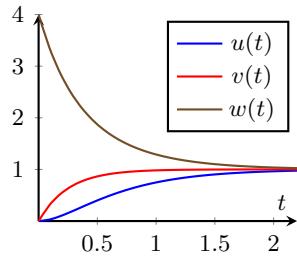
Solution: Evaluating the general solution

$$\begin{bmatrix} u \\ v \\ w \end{bmatrix} = c_1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + c_2 \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} e^{-2t} + c_3 \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} e^{-4t}$$

at time $t = 0$ gives, using the initial conditions and $e^0 = 1$,

$$\begin{bmatrix} 0 \\ 0 \\ 4 \end{bmatrix} = c_1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} + c_2 \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} + c_3 \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} c_1 + c_2 - c_3 \\ c_1 + c_3 \\ c_1 - c_2 - c_3 \end{bmatrix}.$$

Solving by hand, the second row requires $c_3 = -c_1$, so the first row then requires $c_1 + c_2 + c_1 = 0$, that is, $c_2 = -2c_1$. Putting both of these into the third row requires $c_1 + 2c_1 + c_1 = 4$, that is, $c_1 = 1$. Then $c_2 = -2$ and $c_3 = -1$. Consequently, as drawn in the margin, the particular solution is



$$\begin{aligned} \begin{bmatrix} u \\ v \\ w \end{bmatrix} &= \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} - 2 \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix} e^{-2t} - \begin{bmatrix} -1 \\ 1 \\ -1 \end{bmatrix} e^{-4t} \\ &= \begin{bmatrix} 1 - 2e^{-2t} + e^{-4t} \\ 1 - e^{-4t} \\ 1 + 2e^{-2t} + e^{-4t} \end{bmatrix}. \end{aligned}$$

■

Example 7.3.18. Use Matlab/Octave to find a general solution to the system of differential equations

$$\begin{aligned} dx_1/dt &= -\frac{1}{2}x_1 - \frac{1}{2}x_2 + x_3 + 2x_4, \\ dx_2/dt &= -\frac{1}{2}x_1 - \frac{1}{2}x_2 + 2x_3 + x_4, \\ dx_3/dt &= x_1 + 2x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4, \\ dx_4/dt &= 2x_1 + x_2 - \frac{1}{2}x_3 - \frac{1}{2}x_4. \end{aligned}$$

What is the particular solution that satisfies the initial conditions $x_1(0) = -5$, $x_2(0) = -1$ and $x_3(0) = x_4(0) = 0$? Record your commands and give reasons.

Solution: Write the system in matrix-vector form $\frac{d\mathbf{x}}{dt} = A\mathbf{x}$ for vector

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \quad \text{and matrix } A = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2} & 1 & 2 \\ -\frac{1}{2} & -\frac{1}{2} & 2 & 1 \\ 1 & 2 & -\frac{1}{2} & -\frac{1}{2} \\ 2 & 1 & -\frac{1}{2} & -\frac{1}{2} \end{bmatrix}.$$

Enter the matrix into Matlab/Octave and then find its eigenvalues and eigenvectors as follows



```
A=[-1/2 -1/2 1 2
   -1/2 -1/2 2 1
   1 2 -1/2 -1/2
   2 1 -1/2 -1/2]
[V,D]=eig(A)
```

Matlab/Octave tells us the eigenvectors and eigenvalues:

```
V =
-0.5000 0.5000 -0.5000 -0.5000
-0.5000 -0.5000 0.5000 -0.5000
0.5000 0.5000 0.5000 -0.5000
0.5000 -0.5000 -0.5000 -0.5000

D =
-4.0000 0 0 0
0 -1.0000 0 0
0 0 1.0000 0
0 0 0 2.0000
```

Then Theorem 7.3.15 gives that a general solution of the differential equations is

$$\mathbf{x} = c_1 \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{bmatrix} e^{-4t} + c_2 \begin{bmatrix} \frac{1}{2} \\ -\frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} e^{-t} + c_3 \begin{bmatrix} -\frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} e^t + c_4 \begin{bmatrix} -\frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} e^{2t}.$$

Given the specified initial conditions at $t = 0$, when all the above exponentials reduce to $e^0 = 1$, we just need to find the linear combination of the eigenvectors that equals the initial vector $\mathbf{x}(0) = (-5, -1, 0, 0)$; that is, we solve $V\mathbf{c} = \mathbf{x}(0)$. In Matlab/Octave compute $\mathbf{c}=V\backslash[-5;-1;0;0]$ to find the vector of coefficients is $\mathbf{c} = (3, -2, 2, 3)$. Hence the particular solution is

$$\mathbf{x} = \begin{bmatrix} -\frac{3}{2} \\ -\frac{3}{2} \\ \frac{3}{2} \\ \frac{3}{2} \end{bmatrix} e^{-4t} + \begin{bmatrix} -1 \\ 1 \\ -1 \\ 1 \end{bmatrix} e^{-t} + \begin{bmatrix} -1 \\ 1 \\ 1 \\ -1 \end{bmatrix} e^t + \begin{bmatrix} -\frac{3}{2} \\ -\frac{3}{2} \\ -\frac{3}{2} \\ -\frac{3}{2} \end{bmatrix} e^{2t}.$$

■

Example 7.3.19. Find (by hand) a general solution to the system of differential equations $\frac{dy}{dt} = z$ and $\frac{dz}{dt} = -4y$.

Solution: Let vector $\mathbf{y} = (y, z)$, and then form the differential equations into the matrix-vector system

$$\frac{d\mathbf{y}}{dt} = \begin{bmatrix} \frac{dy}{dt} \\ \frac{dz}{dt} \end{bmatrix} = \begin{bmatrix} z \\ -4y \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -4 & 0 \end{bmatrix} \begin{bmatrix} y \\ z \end{bmatrix} = \underbrace{\begin{bmatrix} 0 & 1 \\ -4 & 0 \end{bmatrix}}_A \mathbf{y}.$$

To use Theorem 7.3.15 we need eigenvalues and eigenvectors of the matrix A . Here the characteristic polynomial of A is

$$\det(A - \lambda I) = \det \begin{bmatrix} -\lambda & 1 \\ -4 & -\lambda \end{bmatrix} = \lambda^2 + 4.$$

This determinant is only zero for $\lambda^2 = -4$, that is, $\lambda = \pm 2i$ — a pair of complex conjugate eigenvalues.

- For eigenvalue $\lambda = +2i$ the corresponding eigenvectors \mathbf{p} satisfy

$$(A - \lambda I)\mathbf{p} = \begin{bmatrix} -2i & 1 \\ -4 & -2i \end{bmatrix} \mathbf{p} = \mathbf{0}.$$

The second row of this matrix is $-2i$ times the first row so we just need to satisfy the first row equation $[-2i \ 1] \mathbf{p} = 0$. This equation is $-2ip_1 + p_2 = 0$, that is, $p_2 = 2ip_1$. Hence all eigenvectors are of the form $\mathbf{p} = (1, 2i)p_1$. Choose any one, say $\mathbf{p} = (1, 2i)$.

- Similarly, for eigenvalue $\lambda = -2i$ the corresponding eigenvectors \mathbf{p} satisfy

$$(A - \lambda I)\mathbf{p} = \begin{bmatrix} 2i & 1 \\ -4 & 2i \end{bmatrix} \mathbf{p} = \mathbf{0}.$$

The second row of this matrix is $2i$ times the first row so we just need to satisfy the first row equation $[2i \ 1] \mathbf{p} = 0$. This equation is $2ip_1 + p_2 = 0$, that is, $p_2 = -2ip_1$. Hence all eigenvectors are of the form $\mathbf{p} = (1, -2i)p_1$. Choose any one, say $\mathbf{p} = (1, -2i)$.

With these two distinct eigenvalues, corresponding eigenvectors are linearly independent, and so Theorem 7.3.15 gives a general solution of the differential equations as

$$\begin{bmatrix} y \\ z \end{bmatrix} = c_1 \begin{bmatrix} 1 \\ 2i \end{bmatrix} e^{i2t} + c_2 \begin{bmatrix} 1 \\ -2i \end{bmatrix} e^{-i2t}.$$

That is, $y(t) = c_1 e^{i2t} + c_2 e^{-i2t}$ and $z(t) = 2ic_1 e^{i2t} - 2ic_2 e^{-i2t}$ for any constants c_1 and c_2 . These formulas answer the exercise. The next examples show that because of the complex exponentials, this solution describes oscillations in time t .

■

Example 7.3.20. Further consider Example 7.3.19. Suppose we additionally know that $y(0) = 3$ and $z(0) = 0$. Find the particular solution that satisfies these two initial conditions.

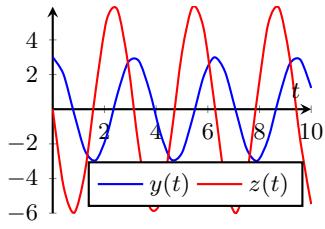
Solution: Use the derived general solution that $y(t) = c_1 e^{i2t} + c_2 e^{-i2t}$ and $z(t) = 2ic_1 e^{i2t} - 2ic_2 e^{-i2t}$. We find the constants c_1

and c_2 so that these satisfy the conditions $y(0) = 3$ and $z(0) = 0$. Substitute $t = 0$ into the general solution to require, using $e^0 = 1$,

$$\begin{aligned}y(0) &= 3 = c_1 e^{i2 \cdot 0} + c_2 e^{-i2 \cdot 0} = c_1 + c_2, \\z(0) &= 0 = 2ic_1 e^{i2 \cdot 0} - 2ic_2 e^{-i2 \cdot 0} = 2ic_1 - 2ic_2,\end{aligned}$$

The second of these equations requires $2ic_1 = 2ic_2$, that is, $c_1 = c_2$. The first, $c_1 + c_2 = 3$, then requires that $2c_1 = 3$, that is, $c_1 = 3/2$ and so $c_2 = 3/2$. Hence the particular solution is

$$y(t) = \frac{3}{2}e^{i2t} + \frac{3}{2}e^{-i2t} \quad \text{and} \quad z(t) = 3ie^{i2t} - 3ie^{-i2t}.$$



But recall Euler's formula that $e^{i\theta} = \cos \theta + i \sin \theta$ for any θ . Invoking Euler's formula the above particular solution simplifies:

$$\begin{aligned}y(t) &= \frac{3}{2}[\cos 2t + i \sin 2t] + \frac{3}{2}[\cos(-2t) + i \sin(-2t)] \\&= \frac{3}{2} \cos 2t + \frac{3}{2}i \sin 2t + \frac{3}{2} \cos 2t - \frac{3}{2}i \sin 2t \\&= 3 \cos 2t, \\z(t) &= 3i[\cos 2t + i \sin 2t] - 3i[\cos(-2t) + i \sin(-2t)] \\&= 3i \cos 2t - 3 \sin 2t - 3i \cos 2t - 3 \sin 2t \\&= -6 \sin 2t.\end{aligned}$$

Because $y(t)$ and $z(t)$ are just trigonometric functions of t , they oscillate in time t , as illustrated in the margin.

■

Example 7.3.21. In a real application the complex numbers of the general solution to Example 7.3.19 are usually inconvenient. Instead we often express the solution solely in terms of real quantities as just done in the previous Example 7.3.20. Use Euler's formula, that $e^{i\theta} = \cos \theta + i \sin \theta$ for any θ , to rewrite the general solution of Example 7.3.19 in terms of real functions.

Solution: Here use Euler's formula in the expression for

$$\begin{aligned}y(t) &= c_1 e^{i2t} + c_2 e^{-i2t} \\&= c_1 [\cos 2t + i \sin 2t] + c_2 [\cos(-2t) + i \sin(-2t)] \\&= c_1 \cos 2t + ic_1 \sin 2t + c_2 \cos 2t - ic_2 \sin 2t \\&= (c_1 + c_2) \cos 2t + (ic_1 - ic_2) \sin 2t \\&= C_1 \cos 2t + C_2 \sin 2t\end{aligned}$$

for constants $C_1 = c_1 + c_2$ and $C_2 = i(c_1 - c_2)$. Let's view the arbitrariness in c_1 and c_2 as being 'transferred' to C_1 and C_2 , then $y(t) = C_1 \cos 2t + C_2 \sin 2t$ is a general solution to the differential equations, and is expressed purely in real factors. In such a real form

we explicitly see the oscillations in time t through the trigonometric functions $\cos 2t$ and $\sin 2t$.

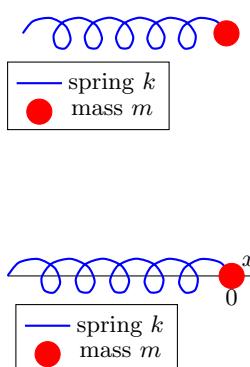
The function $z(t) = 2ic_1e^{i2t} - 2ic_2e^{-i2t}$ has a corresponding form found by replacing c_1 and c_2 in terms of C_1 and C_2 . From above, the constants $C_1 = c_1 + c_2$ and $C_2 = i(c_1 - c_2)$. Adding the first to $\pm i$ times the second determines $C_1 - iC_2 = 2c_1$ and $C_1 + iC_2 = 2c_2$, so that $c_1 = (C_1 - iC_2)/2$ and $c_2 = (C_1 + iC_2)/2$. Then the expression for

$$\begin{aligned} z(t) &= 2ic_1e^{i2t} - 2ic_2e^{-i2t} \\ &= 2i \frac{C_1 - iC_2}{2} [\cos 2t + i \sin 2t] \\ &\quad - 2i \frac{C_1 + iC_2}{2} [\cos 2t - i \sin 2t] \\ &= (iC_1 + C_2)[\cos 2t + i \sin 2t] \\ &\quad + (-iC_1 + C_2)[\cos 2t - i \sin 2t] \\ &= iC_1 \cos 2t - C_1 \sin 2t + C_2 \cos 2t + iC_2 \sin 2t \\ &\quad - iC_1 \cos 2t - C_1 \sin 2t + C_2 \cos 2t - iC_2 \sin 2t \\ &= -2C_1 \sin 2t + 2C_2 \cos 2t. \end{aligned}$$

That is, the corresponding general solution $z(t) = -2C_1 \sin 2t + 2C_2 \cos 2t$ is now also expressed in real factors.

■

Example 7.3.22 (oscillating applications). A huge variety of vibrating systems are analogous to the basic oscillations of a mass on a spring, illustrated schematically in the margin. The mass generally will oscillate to and fro. Describe such a system mathematically with two differential equations, and solve the differential equations to confirm it oscillates.



Solution: At any time t , let the position of the mass relative to its rest position be denoted by $x(t)$: that is, we put an x -axis on the picture with $x = 0$ where the mass would stay at rest, as illustrated. At any time t let the mass be moving with velocity $v(t)$ (positive to the right, negative to the left). Then we know one differential equation, that $\frac{dx}{dt} = v$.

Newton's law, that mass \times acceleration = force, provides another differential equation. Here the mass is denoted by m , and the acceleration is $\frac{dv}{dt}$. The force on the mass come from the spring: typically the force by the spring is proportional to the stretching of the spring, namely to x . Different springs give different strength forces so let's denote the constant of proportionality by k —a constant that varies from spring to spring. Then the force will be $-kx$ as springs try to pull/push the mass back towards $x = 0$. Consequently Newton's law gives us the differential equation, $m \frac{dv}{dt} = -kx$. Divide by mass m and the differential equation is $\frac{dv}{dt} = -\frac{k}{m}x$.

Write these two differential equations together as a matrix-vector system:

$$\begin{bmatrix} \frac{dx}{dt} \\ \frac{dv}{dt} \end{bmatrix} = \begin{bmatrix} v \\ -\frac{k}{m}x \end{bmatrix}, \quad \text{that is, } \frac{d}{dt} \begin{bmatrix} x \\ v \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{k}{m} & 0 \end{bmatrix} \begin{bmatrix} x \\ v \end{bmatrix}.$$

Theorem 7.3.15 asserts a general solution comes from the eigenvalues and eigenvectors of the matrix.

- Here the characteristic polynomial is

$$\det \begin{bmatrix} -\lambda & 1 \\ -\frac{k}{m} & -\lambda \end{bmatrix} = \lambda^2 + \frac{k}{m}.$$

This polynomial is zero only when the eigenvalues $\lambda = \pm\sqrt{-k/m} = \pm i\sqrt{k/m}$: these are pure imaginary eigenvalues.

- The corresponding eigenvectors \mathbf{p} satisfy

$$\begin{bmatrix} \mp i\sqrt{k/m} & 1 \\ -k/m & \mp i\sqrt{k/m} \end{bmatrix} \mathbf{p} = \mathbf{0}.$$

The two rows of this equation are satisfied by the corresponding eigenvectors $\mathbf{p} \propto (1, \pm i\sqrt{k/m})$.

Then a general solution to the system of differential equations is

$$\begin{bmatrix} x \\ v \end{bmatrix} = c_1 \begin{bmatrix} 1 \\ i\sqrt{k/m} \end{bmatrix} e^{i\sqrt{k/m}t} + c_2 \begin{bmatrix} 1 \\ -i\sqrt{k/m} \end{bmatrix} e^{-i\sqrt{k/m}t}.$$

This formula shows that the mass on the spring generally oscillates as the complex exponentials are oscillatory.

However, in real applications we usually prefer a real algebraic expression. Just as in Example 7.3.21, we make the above formula real by changing from (complex) arbitrary constants c_1 and c_2 to new (real) arbitrary constants C_1 and C_2 where $c_1 = (C_1 - iC_2)/2$ and $c_2 = (C_1 + iC_2)/2$. Substitute these relations into the above general solution, and using Euler's formula, gives the position

$$\begin{aligned} x(t) &= c_1 e^{i\sqrt{k/m}t} + c_2 e^{-i\sqrt{k/m}t} \\ &= \frac{C_1 - iC_2}{2} [\cos(\sqrt{k/m}t) + i \sin(\sqrt{k/m}t)] \\ &\quad + \frac{C_1 + iC_2}{2} [\cos(\sqrt{k/m}t) - i \sin(\sqrt{k/m}t)] \\ &= \frac{C_1}{2} \cos(\sqrt{k/m}t) + i \frac{C_1}{2} \sin(\sqrt{k/m}t) \\ &\quad - i \frac{C_2}{2} \cos(\sqrt{k/m}t) + \frac{C_2}{2} \sin(\sqrt{k/m}t) \\ &\quad + \frac{C_1}{2} \cos(\sqrt{k/m}t) - i \frac{C_1}{2} \sin(\sqrt{k/m}t) \\ &\quad + i \frac{C_2}{2} \cos(\sqrt{k/m}t) + \frac{C_2}{2} \sin(\sqrt{k/m}t) \end{aligned}$$

$$= C_1 \cos(\sqrt{k/m}t) + C_2 \sin(\sqrt{k/m}t).$$

For all values of the arbitrary constants C_1 and C_2 , this formula describes the position $x(t)$ of the mass as purely real oscillations in time. Similarly for the velocity $v(t)$ (Exercise 7.3.12). ■

7.3.2 Exercises

Exercise 7.3.1. Which of the following matrices diagonalise the matrix $Z = \begin{bmatrix} 7 & 12 \\ -2 & -3 \end{bmatrix}$? Show your working.

$$(a) P_a = \begin{bmatrix} -2 & 3 \\ 1 & -1 \end{bmatrix}$$

$$(b) P_b = \begin{bmatrix} 3 & -2 \\ -1 & 1 \end{bmatrix}$$

$$(c) P_c = \begin{bmatrix} 1 & -1 \\ -2 & 3 \end{bmatrix}$$

$$(d) P_d = \begin{bmatrix} 1 & 3 \\ 1 & 2 \end{bmatrix}$$

$$(e) P_e = \begin{bmatrix} 4 & 3 \\ -2 & -1 \end{bmatrix}$$

$$(f) P_f = \begin{bmatrix} -2 & 3 \\ 2 & -2 \end{bmatrix}$$

$$(g) P_g = \begin{bmatrix} -1 & 1 \\ 3 & -2 \end{bmatrix}$$

$$(h) P_h = \begin{bmatrix} 3 & 1 \\ 2 & 1 \end{bmatrix}$$

Exercise 7.3.2. Redo Exercise 7.3.1 by finding which matrices P_a, \dots, P_h diagonalise each of the following matrices.

$$(a) A = \begin{bmatrix} 5 & 12 \\ -2 & -5 \end{bmatrix}$$

$$(b) B = \begin{bmatrix} -3 & -12 \\ 2 & 7 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} -1 & -1 \\ 6 & 4 \end{bmatrix}$$

$$(d) D = \begin{bmatrix} 4 & 6 \\ -1 & -1 \end{bmatrix}$$

Exercise 7.3.3. Following Example 7.3.2b, prove that the matrix $\begin{bmatrix} k & 1 \\ 0 & k \end{bmatrix}$ is not diagonalisable for any scalar k .

Exercise 7.3.4. In each of the following cases, you are given three linearly independent eigenvectors and corresponding eigenvalues for some 3×3 matrix A . Write down three different matrices P that will diagonalise the matrix A , and for each write down the corresponding diagonal matrix $D = P^{-1}AP$.

$$(a) \lambda_1 = -1, \mathbf{p}_1 = (3, 2, -1); \lambda_2 = 1, \mathbf{p}_2 = (-4, -2, 2); \lambda_3 = 3, \mathbf{p}_3 = (-1, 0, 2).$$

$$(b) \lambda_1 = -1, \mathbf{p}_1 = (2, 1, 2); \lambda_2 = -1, \mathbf{p}_2 = (0, 3, 1); \lambda_3 = 2, \mathbf{p}_3 = (4, -2, -2).$$

- (c) $\lambda_1 = 1$, $\mathbf{p}_1 = (3, -7, 2)$; $\lambda_2 = -2$, $\mathbf{p}_2 = (-4, -5, 1)$; $\lambda_3 = 4$, $\mathbf{p}_3 = (1, 2, 3)$.
- (d) $\lambda_1 = 3$, $\mathbf{p}_1 = (2, -3, 0)$; $\lambda_2 = 1$, $\mathbf{p}_2 = (-1, 2, -6)$; $\lambda_3 = -4$, $\mathbf{p}_3 = (-2, -1, -3)$.

Exercise 7.3.5. From the given information, are each of the matrices diagonalisable? Give reasons.

- (a) The only eigenvalues of a 2×2 matrix are 2.2 and 0.1.
- (b) The only eigenvalues of a 4×4 matrix are 2.2, 1.9, -1.8 and -1.
- (c) The only eigenvalue of a 2×2 matrix is 0.7.
- (d) The only eigenvalues of a 6×6 matrix are -1.6, 0.3, 0.1 and -2.3.
- (e) The only eigenvalues of a 5×5 matrix are -1.7, 1.4, 1.3, 2.4, 0.5 and -2.3.
- (f) The only eigenvalues of a 6×6 matrix are 1.2, -0.9, -0.8, 2.2, 0.2 and -0.2.
- (g) The Matlab/Octave function `eig(A)` returns the result

```
ans =
2.6816
-0.1445
0.0798
0.3844
```

- (h) The Matlab/Octave function `eig(A)` returns the result

```
ans =
3.0821 + 0.0000i
-2.7996 + 0.0000i
-0.7429 + 1.6123i
-0.7429 - 1.6123i
```

- (i) The Matlab/Octave function `eig(A)` returns the result

```
ans =
-1.0000
1.0000
2.0000
-1.0000
```

Exercise 7.3.6. For each of the following 3×3 matrices, show with hand algebra that each matrix has one eigenvalue of multiplicity three, and then determine the dimension of the corresponding eigenspace. Also compute the eigenvalue and eigenvectors with Matlab/Octave: comment on any limitations in the computed ‘eigenvalues’ and ‘eigenvectors’.

$$(a) A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix} \quad (b) B = \begin{bmatrix} 2 & 1 & -2 \\ -8 & -4 & 8 \\ -2 & -1 & 2 \end{bmatrix}$$

$$(c) C = \begin{bmatrix} -2 & 1 & 0 \\ 0 & -2 & 0 \\ 0 & 1 & -2 \end{bmatrix} \quad (d) D = \begin{bmatrix} 7 & 1 & 22 \\ -6 & -1 & -20 \\ -1 & 0 & -3 \end{bmatrix}$$

$$(e) E = \begin{bmatrix} 3 & 5 & -1 \\ -4 & -6 & 1 \\ -5 & -6 & 0 \end{bmatrix} \quad (f) F = \begin{bmatrix} -6 & -5 & 3 \\ -4 & -7 & 3 \\ -12 & -15 & 7 \end{bmatrix}$$

$$(g) G = \begin{bmatrix} -3 & 2 & 7 \\ 2 & -3 & -9 \\ -1 & 1 & 3 \end{bmatrix} \quad (h) H = \begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix}$$

Exercise 7.3.7. For a given $n \times n$ square matrix A , suppose λ_1 is an eigenvalue with n corresponding linearly independent eigenvectors $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$. Adapt parts of the proof of Theorem 7.3.12 to prove that the characteristic polynomial of matrix A is $\det(A - \lambda I) = (\lambda_1 - \lambda)^n$. Then deduce that λ_1 is the only eigenvalue of A , and is of multiplicity n .

Exercise 7.3.8. For each of the following matrices, use Matlab/Octave to find the eigenvalues, their multiplicity, and the dimension of the eigenspaces. Give reasons (remember computational error).



$$(a) A = \begin{bmatrix} -30.5 & 25.5 & 22 & -5 & -48.5 \\ -111 & 88 & 76 & -20 & -173 \\ 87.5 & -70.5 & -61 & 16 & 136.5 \\ 51 & -43 & -36 & 9 & 81 \\ -4.5 & 2.5 & 2 & -1 & -6.5 \end{bmatrix}$$



$$(b) B = \begin{bmatrix} 2929 & 1140 & -1359 & 2352 & 406 \\ -2441 & -950 & 1132 & -1960 & -338 \\ -1070 & -416 & 495 & -858 & -148 \\ -2929 & -1140 & 1359 & -2352 & -406 \\ -895 & -347 & 412 & -715 & -124 \end{bmatrix}$$



$$(c) C = \begin{bmatrix} -168 & 115 & 305 & -120 & 70 \\ -4710 & 3202 & 8510 & -3360 & 1990 \\ 450 & -305 & -813 & 320 & -190 \\ -4545 & 3090 & 8205 & -3243 & 1920 \\ -2415 & 1640 & 4355 & -1720 & 1017 \end{bmatrix}$$



$$(d) D = \begin{bmatrix} -8 & 744 & -564 & 270 & 321 \\ 1 & -43 & 33 & -15 & -19 \\ 1 & 29 & -21 & 12 & 12 \\ 5 & -431 & 327 & -156 & -186 \\ -5 & 535 & -405 & 195 & 231 \end{bmatrix}$$



$$(e) E = \begin{bmatrix} 15.6 & 31.4 & -11.6 & -14.6 & -10.6 \\ -16.6 & -32.4 & 11.6 & 14.6 & 10.6 \\ 33.6 & 56.4 & -19.6 & -27.6 & -15.6 \\ 38.6 & 75.4 & -28.6 & -35.6 & -26.6 \\ -122 & -226 & 82 & 106 & 73 \end{bmatrix}$$

$$(f) F = \begin{bmatrix} -208 & 420 & -518 & 82 & 264 \\ 655 & -1336 & 1642 & -260 & -838 \\ 229 & -467 & 574 & -91 & -294 \\ -171 & 348 & -428 & 66 & 218 \\ -703 & 1431 & -1760 & 279 & 896 \end{bmatrix}$$

Exercise 7.3.9. For each of the following systems of differential equations, find eigenvalues and eigenvectors to derive a general solution of the system of differential equations. Show your working.

$$(a) \frac{dx}{dt} = x - 1.5y, \quad \frac{dy}{dt} = 4x - 4y$$

$$(b) \frac{dx}{dt} = x, \frac{dy}{dt} = -12x + 5y$$

$$(c) \frac{dx}{dt} = 7x - 3y$$

$$(d) \frac{du}{dt} = 2.8u - 3.6v, \quad \frac{dv}{dt} = -0.6u + 2.2v$$

$$(e) \frac{dp}{dt} = 14p + 16q, \quad \frac{dq}{dt} = -8p - 10q$$

$$(f) \frac{dx}{dt} = 6.5x - 0.6y - 5.7z, \quad \frac{dy}{dt} = -3x + 4.4y + 7.8z$$

$$(g) \frac{dx}{dt} = -31x + 26y - 24z, \quad \frac{dy}{dt} = -48x + 39y - 36z, \quad \frac{dz}{dt} = -14x + 10y - 9z$$

$$(h) \frac{dx}{dt} = 0.2x + 1.2z, \frac{dy}{dt} = -x, \quad \frac{dz}{dt} = 1.8x + 0.8z$$

$$(i) \frac{du}{dt} = 4.5u + 7.5v + 7.5w, \quad \frac{dv}{dt} = 3u + 4v + 5w, \quad \frac{dw}{dt} = -7.5u - 11.5v - 12.5w$$

$$(j) \frac{dp}{dt} = -13p + 30q + 6r, \quad \frac{dq}{dt} = -32p + 69q + 14r, \quad \frac{dr}{dt} = 125p - 265q - 54r$$

Exercise 7.3.10. In each of the following, a general solution to a differential equation is given. Find the particular solution that satisfies the specified initial conditions. Show your working.

$$(a) (x, y) = c_1(0, 1)e^{-t} + c_2(1, 3)e^{2t} \text{ where } x(0) = 2 \text{ and } y(0) = 1$$

$$(b) (x, y) = c_1(0, 1)e^{-2t} + c_2(1, 3)e^t \text{ where } x(0) = 0 \text{ and } y(0) = 2$$

$$(c) x = 3c_1e^t + c_2e^{-t}, y = 5c_1e^t + 2c_2e^{-t} \text{ where } x(0) = 0 \text{ and } y(0) = 2$$

$$(d) x = 3c_1e^t + c_2e^{2t}, y = 5c_1e^t + 2c_2e^{2t} \text{ where } x(0) = 1 \text{ and } y(0) = 3$$

(e) $(x, y, z) = c_1(0, 0, 1) + c_2(1, -1, 1)e^{2t} + c_3(-2, 3, 1)e^{-2t}$ where
 $x(0) = 3, y(0) = -4$ and $z(0) = 0$

(f) $(x, y, z) = c_1(0, 0, 1)e^{-3t} + c_2(1, -1, 1) + c_3(-2, 3, 1)e^{-t}$ where
 $x(0) = 3, y(0) = -4$ and $z(0) = 2$

Exercise 7.3.11. For each of the following systems of differential equations, use Matlab/Octave to find eigenvalues and eigenvectors and hence derive a general solution of the system of differential equations. Record your working.

(a) $\frac{dx_1}{dt} = 15.8x_1 + 17.1x_2 - 119.7x_3 + 153.9x_4,$
 $\frac{dx_2}{dt} = 1.4x_1 + 0.1x_2 + 12.9x_3 - 17.1x_4,$
 $\frac{dx_3}{dt} = 6.2x_1 + 6.2x_2 - 43x_3 + 57.6x_4,$
 $\frac{dx_4}{dt} = 3.4x_1 + 3.4x_2 - 25x_3 + 34x_4.$

(b) $\frac{dx_1}{dt} = -12.2x_1 + 53.7x_2 + 50.1x_3 - 22.8x_4,$
 $\frac{dx_2}{dt} = 0.6x_1 - 2.3x_2 - 2.7x_3 + 1.2x_4,$
 $\frac{dx_3}{dt} = -20.4x_1 + 93.8x_2 + 90.2x_3 - 40.8x_4,$
 $\frac{dx_4}{dt} = -38.2x_1 + 177.9x_2 + 170.7x_3 - 77.2x_4.$

(c) $\frac{dx_1}{dt} = x_1 + 29.4x_2 - 3.2x_3 - 12.9x_4,$
 $\frac{dx_2}{dt} = 1.4x_1 - 38.4x_2 + 5.6x_3 + 18.2x_4,$
 $\frac{dx_3}{dt} = 2.3x_1 - 80.3x_2 + 12.3x_3 + 36.7x_4,$
 $\frac{dx_4}{dt} = 2.4x_1 - 65x_2 + 9.2x_3 + 31.1x_4.$

(d) $\frac{dx_1}{dt} = -50.2x_1 - 39.5x_2 - 20.2x_3 + 68.9x_4,$
 $\frac{dx_2}{dt} = 62.8x_1 + 50.2x_2 + 28.4x_3 - 85.2x_4,$
 $\frac{dx_3}{dt} = -17.3x_1 - 13.7x_2 - 8.9x_3 + 22.7x_4,$
 $\frac{dx_4}{dt} = -6.4x_1 - 4.6x_2 - 1.4x_3 + 9x_4.$

Exercise 7.3.12. Recall the general complex solution that Example 7.3.22 derives for the oscillations of a mass on a spring. Show that substituting $c_1 = (C_1 - iC_2)/2$ and $c_2 = (C_1 + iC_2)/2$ for real C_1 and C_2 results in the velocity $v(t)$ being expressed algebraically in purely real terms.



Answers to selected exercises

7.0.1b : real

7.0.1d : complex

7.0.1f : complex

7.0.1h : complex

7.0.2b : $\mathbb{E}_2 = \text{span}\{(0, 1)\}$.

7.0.2d : $\mathbb{E}_{-2} = \text{span}\{(0, 0, 1)\}$, $\mathbb{E}_{-4} = \text{span}\{(0, 2, -5)\}$, $\mathbb{E}_0 = \text{span}\{(8, -6, -11)\}$.

7.0.2f : This matrix is not triangular so we cannot answer (as yet).

7.0.2h : $\mathbb{E}_{-2} = \text{span}\{(1, 0, 0, 0)\}$, $\mathbb{E}_{-3} = \text{span}\{(6, -1, 1, 0)\}$, $\mathbb{E}_2 = \text{span}\{(6, 9, 1, 5)\}$.

7.0.2j : $\mathbb{E}_0 = \text{span}\{(1, 0, 0, 0)\}$, $\mathbb{E}_7 = \text{span}\{(-2, 7, 0, 0)\}$, $\mathbb{E}_3 = \text{span}\{(-22, 3, 12, 0)\}$.

7.1.1b : yes

7.1.1d : yes

7.1.2a : $\lambda^2 + 5\lambda - 12$

7.1.2c : $-\lambda^3 + 0\lambda^2 + \dots + 4$

7.1.2e : $-\lambda^3 - \lambda^2 + \dots - 228$

7.1.2g : $\lambda^4 - \lambda^3 + \dots + 24$

7.1.3a : $2 \times 2, -5, -6$

7.1.3c : $2 \times 2, 1, 0$

7.1.3e : $3 \times 3, 0, -199$

7.1.3g : $4 \times 4, -3, -56$

7.1.4a : 1 once, -3 once

7.1.4c : 3 twice

7.1.4e : $1/2$ twice

7.1.4g : 0 once, -2 once, 1 once

7.1.4i : 2 once, $-1 \pm i$ once

7.1.4k : -1 once, 0 twice

7.1.5a : -1 once, -2.3 twice

7.1.5c : -0.9 once, -0.1 once, 0.1 twice

7.1.5e : 0.6 four times

7.1.5g : $1.7 \pm 1.3i$ once, -3.9 thrice

7.1.6a : $\mathbb{E}_3 = \text{span}\{(2, 3)\}$

7.1.6c : $\mathbb{E}_{-1} = \text{span}\{(1, 1)\}$

7.1.6e : $\mathbb{E}_7 = \text{span}\{(7, -8, 0), (-2, 2, 1)\}$

7.1.6g : $\mathbb{E}_{-2} = \text{span}\{(-2, 2, 1)\}$

7.1.7a : $\lambda = -3$, once, $(0, -0.89, -0.45)$; $\lambda = 2$, twice, $(0.27, 0.53, 0.80)$.

7.1.7c : $\lambda = 5$, once, $(-0.82, 0, -0.41, 0.41)$; $\lambda = 3$, once, $(0.55, 0.28, 0.55, 0.55)$; $\lambda = -4$, twice, $(0.67, 0, 0.33, -0.67)$.

7.1.7e : $\lambda = -3.94$, once, $(-0.33, -0.41, 0.72, 0.44, -0.15)$; $\lambda = 1$, once, $(0, 0.8, -0.53, 0.27, 0)$; $\lambda = 7.49 \pm 0.39i$, once, $(-0.24 \pm 0.55i, 0.23 \pm 0.02i, 0.043 \pm 0.01i, 0.72, 0.24 \mp 0.11i)$; $\lambda = 4.95$, once, $(-0.14, 0.45, 0.62, 0.46, -0.43)$.

7.1.8a : Sensitive.

7.1.8c : Not sensitive.

7.1.9b : $\lambda = 3$ sensitive.

7.1.9d : neither sensitive, symmetric.

7.1.9f : -0.9 and -0.1 not sensitive, 0.1 sensitive

7.1.9h : 0.6 mixed sensitivity!

7.1.10b : $(0, -12), (12, -24)$ and $(24, -96)$

7.1.10d : $(3, 2), (4, 2)$ and $(2, 0)$

7.1.10f : $(-7, 16, -6), (79, -64, 126)$ and $(-367, 256, -606)$

7.1.10h : $(-12, 6, -39), (-63, -30, 0)$ and $(-168, -96, 39)$

7.1.17b : $\lambda = 0, \pm 13, (0, \frac{12}{13}, \frac{5}{13})$ and $(\pm 1, -\frac{5}{13}, \frac{12}{13})$

7.1.17d : $\lambda = \pm 1, \pm 4, (1, 0, 0, \pm 1)$ and $(0, 1, \mp 1, 0)$

7.2.1b : lin. dep.

7.2.1d : lin. indep.

7.2.1f : lin. indep.

7.2.1h : lin. dep.

7.2.1j : lin. indep.

7.2.2b : lin. dep.

7.2.2d : lin. dep.

7.2.2f : lin. dep.

7.2.2h : lin. dep.

7.2.6b : Two possibilities are $\{(-3/4, 1, 1/4)\}$ and $\{(3, -4, -1)\}$.

7.2.6d : Two possibilities are $\{(0, 1, 1)\}$ and $\{(0, -2, -2)\}$.

- 7.2.6f : Two possibilities are $\{(0,1,0),(2,0,1)\}$ and $\{(2,-4,1),(4,1,2)\}$.
- 7.2.6h : Does not have a solution subspace.
- 7.2.7b : $\{(0.67, -0.57, -0.47)\}$ (2 d.p.)
- 7.2.7d : $\{(-0.50, 0.50, 0.50, -0.50)\}$ (2 d.p.)
- 7.2.7f : $\{(0.61, 0.61, -0.51, -0.10)\}$ (2 d.p.)
- 7.2.7h : $\{(0.62, 0.45, -0.62, -0.17), (-0.13, 0.63, 0.13, 0.76)\}$ (2 d.p.)
- 7.2.9b : $(-1.5, 0), (0.5, -0.5), (-1, 2.5), (2, 0.5)$
- 7.2.9d : $(-1, -0.5), (1, 1), (-1, 0.5), (0.5, 0)$
- 7.2.9f : $(-0.9, 0.3), (0.6, 0.2), (-0.2, -1.1), (0.7, -0.6)$
- 7.2.10a : $[p]_{\mathcal{E}} = (-2, 11, 9)$
- 7.2.10c : $[r]_{\mathcal{E}} = (-4, -1, -11)$
- 7.2.10e : $[t]_{\mathcal{E}} = (-1, 11/2, 9/2)$
- 7.2.10g : $[v]_{\mathcal{E}} = (1/2, -3, -5/2)$
- 7.2.10i : $[x]_{\mathcal{E}} = (-0.3, 4.0, 4.3)$
- 7.2.12a : $[p]_{\mathcal{B}} = (-1, 1)$
- 7.2.12c : $[r]_{\mathcal{B}} = (-3, -1)$
- 7.2.12e : not in \mathcal{B}
- 7.2.12g : not in \mathcal{B}
- 7.2.12i : $[x]_{\mathcal{B}} = (-0.1, -0.2)$
- 7.2.14a : $[p]_{\mathcal{E}} = (-2, -14, 18, 2, -5)$
- 7.2.14c : $[r]_{\mathcal{E}} = (-18, 12, -6, 3, -6)$
- 7.2.14e : $[t]_{\mathcal{B}} = (3, 2, 6)$
- 7.2.14g : $[v]_{\mathcal{B}} = (1, 2, 4)$
- 7.2.14i : not in \mathbb{B}
- 7.3.1a : yes
- 7.3.1c : no
- 7.3.1e : yes
- 7.3.1g : no
- 7.3.5a : yes
- 7.3.5c : unknown
- 7.3.5e : error
- 7.3.5g : yes
- 7.3.5i : unknown

- 7.3.6b : $\lambda = 0$, two; errors 10^{-8} and two eigenvectors effectively the same.
- 7.3.6d : $\lambda = 1$, one; errors 10^{-5} , all three eigenvectors effectively the same
- 7.3.6f : $\lambda = 0$, two; errors 10^{-8} and two eigenvectors effectively the same.
- 7.3.6h : $\lambda = -1$, three; all good
- 7.3.8b : $\lambda = -1$ twice, $\dim \mathbb{E}_{-1} = 1$; $\lambda = 0$ thrice, $\dim \mathbb{E}_0 = 2$
- 7.3.8d : $\lambda = 0$ twice, $\dim \mathbb{E}_0 = 2$; $\lambda = 1$ thrice, $\dim \mathbb{E}_1 = 2$
- 7.3.8f : $\lambda = -1 \pm i$ once each, $\dim \mathbb{E}_{-1\pm i} = 1$; $\lambda = -2$ thrice, $\dim \mathbb{E}_{-2} = 2$
- 7.3.9b : $(x, y) = c_1(0, 1)e^{5t} + c_2(1, 3)e^t$
- 7.3.9d : $(u, v) = c_1(3, -1)e^{4t} + c_2(2, 1)e^t$
- 7.3.9f : Not possible: only two equations for three unknowns.
- 7.3.9h : $(x, y, z) = c_1(0, 1, 0) + c_2(1, 1, -1)e^{-t} + c_3(-2, 1, 3)e^{2t}$
- 7.3.9j : $(p, q, r) = c_1(2, 2, -5)e^{2t} + c_2(3, 1, 2)e^t + c_3(0, -1, 5)e^{-t}$
- 7.3.10b : $(x, y) = (0, 2)e^{-2t}$
- 7.3.10d : $x = -3e^t + 4e^{2t}, y = -5e^t + 8e^{2t}$
- 7.3.10f : $(x, y, z) = (0, 0, 2)e^{-3t} + (1, -1, 1) + (2, -3, -1)e^{-t}$
- 7.3.11b : (2 d.p.) $\mathbf{x} = c_1(0.63, 0.63, -0.42, 0.21)e^{0.5t} + c_2(-0.23, -0.23, -0.23, -0.92)e^{0.4t} + c_3(0.64, 0, 0.43, 0.64)e^{-1.6t} + c_4(-0.89, 0, 0, 0.45)e^{0.8t}$
- 7.3.11d : (2 d.p.) $\mathbf{x} = c_1(-0.78, 0.38 + 0.3i, 0.095 + 0.076i, -0.32 + 0.21i)e^{(0.1-1.1i)t} + c_2(-0.78, 0.38 - 0.3i, 0.095 - 0.076i, -0.32 - 0.21i)e^{(0.1+1.1i)t} + c_3(0, -0.87, 0.22, -0.44)e^{0.5t} + c_4(-0.5, -0.5, 0.5, -0.5)e^{-0.6t}$

Bibliography

- Alpers, B., Demlova, M., Fant, C.-H., Gustafsson, T., Lawson, D., Mustoe, L., Olsson-Lehtonen, B., Robinson, C. & Velichova, D. (2013), A framework for mathematics curricula in engineering education, Technical report, European Society for Engineering Education (SEFI).
<http://sefi.htw-aalen.de/curriculum.htm>
- Anton, H. & Rorres, C. (1991), *Elementary linear algebra. Applications version*, 6th edn, Wiley.
- Arnold, V. I. (2014), *Mathematical understanding of nature*, Amer. Math. Soc.
- Berry, M. W., Dumais, S. T. & O'Brien, G. W. (1995), 'Using linear algebra for intelligent information retrieval', *SIAM Review* **37**(4), 573–595.
<http://pubs.siam.org/doi/abs/10.1137/1037127>
- Bliss, K., Fowler, K., Galluzzo, B., Garfunkel, S., Giordano, F., Godbold, L., Gould, H., Levy, R., Libertini, J., Long, M., Malkevitch, J., Montgomery, M., Pollak, H., Teague, D., van der Kooij, H. & Zbiek, R. (2016), GAIMME—Guidelines for Assessment and Instruction in Mathematics Modeling Education, Technical report, SIAM and COMAP.
http://www.siam.org/reports/gaimme.php?_ga=1.192190826.137003948.1386128871
- Bressoud, D. M., Friedlander, E. M. & Levermore, C. D. (2014), 'Meeting the challenges of improved post-secondary education in the mathematical sciences', *Notices of the AMS* **61**(5), 502–3.
- Chartier, T. (2015), *When life is linear: from computer graphics to bracketology*, Math Assoc Amer.
<http://www.maa.org/press/books/when-life-is-linear-from-computer-graphics-to-bracketology>
- Cowen, C. C. (1997), On the centrality of linear algebra in the curriculum, Technical report, Mathematical Association of America.
<http://www.maa.org/centrality-of-linear-algebra>
- Cuyt, A. (2015), Approximation theory, in N. J. Higham, M. R. Dennis, P. Glendinning, P. A. Martin, F. Santosa & J. Tanner, eds, 'Princeton Companion to Applied Mathematics', Princeton, chapter IV.9, pp. 248–262.
- Davis, B. & Uhl, J. (1999), *Matrices, Geometry and Mathematica*, Wolfram Research.

- Donoho, D. L. & Stodden, V. (2015), Reproducible research in the mathematical sciences, *in* N. J. Higham, M. R. Dennis, P. Glendinning, P. A. Martin, F. Santosa & J. Tanner, eds, ‘Princeton Companion to Applied Mathematics’, Princeton, chapter VIII.5, pp. 916–925.
- Driscoll, T. A. & Maki, K. L. (2007), ‘Searching for rare growth factors using multicanonical monte carlo methods’, *SIAM Review* **49**(4), 673–692.
<http://link.aip.org/link/?SIR/49/673/1>
- Gorodetski, V. I., Popack, L. J. & Samoilov, V. (2001), SVD-based approach to transparent embedding data into digital images, *in* V. I. Gorodetski, V. A. Skormin & L. J. Popack, eds, ‘Information assurance in computer networks: methods, models, and architectures for network security’, Vol. 2052 of *Lecture Notes in Computer Science*, Springer, pp. 263–274.
- Halpern, D. F. & Hakel, M. D. (2003), ‘Applying the science of learning to the university and beyond: Teaching for long-term retention and transfer’, *Change: The Magazine of Higher Learning* **35**(4), 36–41.
- Hannah, J. (1996), ‘A geometric approach to determinants’, *The American Mathematical Monthly* **103**(5), 401–409.
<http://www.jstor.org/stable/2974931>
- Higham, N. J. (1996), *Accuracy and stability of numerical algorithms*, SIAM.
- Higham, N. J. (2015), Numerical linear algebra and matrix analysis, *in* N. J. Higham, M. R. Dennis, P. Glendinning, P. A. Martin, F. Santosa & J. Tanner, eds, ‘Princeton Companion to Applied Mathematics’, Princeton, chapter IV.10, pp. 263–281.
- Holt, J. (2013), *Linear algebra with applications*, Freeman.
- Hughes-Hallett, D., Gleason, A. M. & McCallum, et al., W. G. (2013), *Calculus: single and multivariable*, 6th edn, Wiley.
- Kleiber, M. (1947), ‘Body size and metabolic rate’, *Physiological Reviews* **27**, 511–541.
- Kress, R. (2015), Integral equations, *in* N. J. Higham, M. R. Dennis, P. Glendinning, P. A. Martin, F. Santosa & J. Tanner, eds, ‘Princeton Companion to Applied Mathematics’, Princeton, chapter IV.4, pp. 200–208.
- Larson, R. (2013), *Elementary linear algebra*, 7th edn, Brooks/Cole Cengage learning.
- Lay, D. C. (2012), *Linear Algebra and its Applications*, 4th edn, Addison–Wesley.

- Lichman, M. (2013), ‘UCI machine learning repository’, [online].
<http://archive.ics.uci.edu/ml>
- Mandelbrot, B. B. (1982), *The fractal geometry of nature*, W. H. Freeman.
- Moody, D. L. (2009), ‘The “physics” of notations: Towards a scientific basis for constructing visual notations in software engineering’, *IEEE Trans. Soft. Engrg.* **35**(5), 756–778.
- Nakos, G. & Joyner, D. (1998), *Linear algebra with applications*, Brooks/Cole.
- Poole, D. (2015), *Linear algebra: A modern introduction*, 4th edn, Cengage Learning.
- Roulstone, I. & Norbury, J. (2013), *Invisible in the storm: the role of mathematics in understanding weather*, Princeton.
- Schonefeld, S. (1995), ‘Eigenpictures: Picturing the eigenvector problem’, *The College Mathematics Journal* **26**(4), 316–319.
<http://www.jstor.org/stable/2687037>
- Schumacher, C. S., Siegel, M. J. & Zorn, P. (2015), 2015 CUPM curriculum guide to majors in the mathematical sciences, Technical report, The Mathematical Association of America.
<http://www.maa.org/programs/faculty-and-departments/curriculum-department-guidelines-recommendations/cupm>
- Trefethen, L. N. & Bau, III, D. (1997), *Numerical linear algebra*, SIAM.
- Turner, P. R., Crowley, J. M., Humpherys, J., Levy, R., Socha, K. & Wasserstein, R. (2015), Modeling across the curriculum II. report on the second SIAM-NSF workshop, Alexandria, VA, Technical report, [http://www.siam.org/reports/modeling_14.pdf].
- Uhlig, F. (2002), A new unified, balanced, and conceptual approach to teaching linear algebra, Technical report, Department of Mathematics, Auburn University, <http://www.auburn.edu/~uhligfd/TLA/download/tlateach.pdf>.
- Will, T. (2004), Introduction to the singular value decomposition, Technical report, [<http://www.uwlax.edu/faculty/will/svd>].

Index

- $|\cdot|$, **16**
 $()$, **73**
 $+,-,*$, **73**, **75**, **147**
 $.*$, **176**
 $./$, **176**
 $.^$, **98**
 $/$, **73**
 $=$, **73**, **98**
 $[...]$, **73**, **98**
2-norm, **445**
2 d.p., **99**
 \setminus , **98**
`acos()`, **36**
`acos()`, **73**
addition, **23**, **137**, **149**
adjacency matrix, **401**
adolescent, **578–580**
adult, **578–580**, **590**
age structure, **577–592**
age structured population, **140**
angle, **35**, **36**, **38**, **51**, **60**, **182**, **200**
`ans`, **73**
approximate solution, **279**
`arc-cos`, **73**
area, **61**
Arrow, Kenneth, **283**
artificial intelligence, **325**
associative law, **28**, **29**
augmented matrix, **107**, **107**, **112**
average, **444**, **463**
`axis`, **444**

Babbage, Charles, **72**, **79**, **478**
basin area, **325**
basis, **567**, **568**, **616**, **619**, **620**, **623**, **625**, **628**,
 635, **636**
best straight line, **283**
body-centered cubic, **37**
bond angles, **36**
bounded, **512**
brackets, **97**, **134**
bulls eye, **438**, **451**

 \mathbb{C} , **11**

canonical form, **423**, **431**
Cardano, Gerolamo, **11**
Cauchy–Schwarz inequality, **43**, **44**, **45**, **55**,
 448
characteristic equation, **392**, **399**, **552**, **559**,
 567
characteristic polynomial, **559**, **560–564**, **568**,
 595, **596**, **661**
chemistry, **36**
closed, **242**
coastlines, **326**
coefficient, **87**, **122**, **561**
`colormap()`, **293**
column space, **245**, **247**, **260**, **270**, **272**, **274**,
 278, **297**, **304**, **313**
column vector, **97**, **101**, **135**, **148**, **186**, **245**,
 267, **631**
commutative law, **21**, **28**, **29**, **42**, **63**
complex conjugate, **409**
complex eigenvalue, **408**, **553**, **557**, **568**, **585**,
 587
complex eigenvector, **553**
complex numbers, **11**, **395**, **409**
components, **14**, **18**, **135**
composition, **353**, **363**
computed tomography, **292**, **328**, **504**
computer algebra, **73**
`cond`, **217**, **224**, **237**, **238**, **289**, **352**
condition number, **98**, **217**, **217–225**, **236**,
 352, **493**, **502**
conic section, **418**, **430**, **548**
consistent, **91**, **112**, **124**
constant term, **87**, **114**, **561**
contradiction, **321**, **409**, **415**
coordinate axes, **430**
coordinate system, **12**
coordinate vector, **625**
coordinates, **625**
cosine, **36**
cosine rule, **34**, **39**, **53**
cross product, **59**, **57–71**
cross product direction, **60**
`csvread()`, **444**, **463**

- CT scan, 292, 494
`cumprod()`, 479
- data mining, 325
 data reduction, 456
 De Moivre's theorem, 587
 decimal places, 99
 Descartes, 28
`det()`, 523
 determinant, 66, 167, 391, 506, 509, 509, 515, 518, 523, 533, 539, 595
`diag`, 175–179, 181, 193, 197, 210, 218, 329, 362, 367, 384, 390, 416, 435, 452, 639
`diag()`, 176
 diagonal entries, 175
 diagonal matrix, 175, 174–182, 193, 380, 416, 509, 533, 637–639
 diagonalisable, 638, 639, 642, 643, 650
 diagonalisation, 637–663
 difference, 23, 137, 149
 differential equation, 648, 650
 differential equations, 648–659
`dim`, 260, 316, 317, 619, 645, 667
 dimension, 260, 260, 262, 273, 383, 456, 619, 644, 660
 direction vector, 26
 discrepancy principle, 501
 discriminant, 43, 45
 displacement vector, 12, 12, 18
 distance, 25, 437
 distributive law, 28, 29, 42, 63
 dolphin, 603
 dot product, 35, 36, 73, 182, 182, 469, 474
`dot()`, 73
 double subscript, 135
 Duhem, Pierre, 276
- e_j , 24
`eig()`, 384, 566
 eigen-problem, 552, 583
 eigenspace, 381, 383, 387, 392, 400, 552, 558, 567, 568, 597, 644, 660, 661
 eigenvalue, 377, 380, 381, 383, 384, 387, 392, 398, 400, 403, 405, 409, 410, 414–416, 423, 424, 552, 554, 558, 559, 561, 564, 566–568, 576, 584, 593, 595, 596, 603, 605, 611, 639, 643, 644, 650
- eigenvector, 377, 381, 384, 392, 398, 405, 410, 415, 423, 552, 554, 567, 568, 584, 593, 603, 605, 611, 639, 650
 El Nino, 254, 260
 elementary row operation, 107, 112
 elements, 135
 elephant, 600
 ellipsis, 11
 ensemble of simulations, 259
 entries, 135
 equal, 14, 136
 equation of the plane, 49, 49
 Error using, 149, 152
`Error using`, 77
`error:`, 77
 error: operator, 149, 152
 Euler, 401, 441
 Euler's formula, 396, 656
`exp()`, 293
 exponential, 293
`eye()`, 147
- factorial, 542
 factorisation, 205
 female, 578, 580
 Feynman, Richard, 483
 floating point, 73
`format long`, 574
 fractal, 477
 free variable, 110, 110, 112, 216, 291, 622
- Galileo Galilei, 240
 Gauss–Jordan elimination, 112, 113
 general solution, 584, 585, 600, 605, 649, 650, 662, 663
 giant mouse lemur, 602
 global positioning system, 88, 93, 104, 120
 GPS, 88, 93, 104, 120
- Hack's exponent, 325, 326
 Hankel matrix, 256
`hankel()`, 255, 274
 Hawking, Stephen, 179
`hilb()`, 487
 Hilbert matrix, 487, 503, 504
 homogeneous, 114, 116, 220, 248, 267, 381, 623, 631
 hyper-cube, 509, 511
 hyper-volume, 509

- i, 24**
- idempotent, 522
- identical columns, 525
- identical rows, 525
- identity matrix, 135, 140, 147
- identity transformation, 356
- image compression, 438–456
- `imagesc()`, 293
- Impossibility Theorem, 283
- `imread()`, 444
- in, \in , 241
- inconsistent, 91, 216, 278, 279, 299, 484
- inconsistent equations, 276–335
- inconsistent system, 297–299, 305
- induction, 227, 228, 415, 416
- `Inf`, 98
- infant, 578
- inference, 325
- inferring, 83
- infinitely many solutions, 91, 114, 116
- infinity, 98
- initial condition, 652, 653, 655, 662
- inner product, 35
- integer, 10
- integral equations, 139
- interchanging two columns, 525
- interchanging two rows, 108, 525
- intersection, 313
- `inv()`, 192
- inverse, 165, 166, 169, 364, 523
- inverse cosine, 73
- inverse matrix, 192
- inverse transformation, 356
- invertible, 165, 166, 167, 169–171, 177, 186, 202, 220, 237, 267, 356, 358, 391, 403, 405, 506, 517, 584, 631, 638
- Iris, 458
- j, 24**
- jewel in the crown, 4, 205
- jpeg, 438
- juvenile, 578, 580, 590
- k, 24**
- Kepler's law, 286
- kitten, 590
- Kleiber's power law, 325
- knowledge discovery, 325
- Lagrange multiplier, 329
- latent semantic indexing, 466–475
- leading one, 110, 110
- least square, 279, 304, 322, 346, 500, 502
- left triangular, 533
- Leibnitz, Gottfried, 95
- length, 16, 17, 20, 25, 43, 44, 61, 73, 182, 200, 437, 444
- Leslie matrix, 141, 144, 163, 172, 407
- life expectancy, 283
- linear combination, 122, 122–125, 127, 128, 584, 606, 607, 610, 625, 650
- linear dependence, 605–636
- linear equation, 86, 91, 96, 106, 107, 110, 112, 114, 116, 124, 169, 201, 231, 291, 483, 484
- linear independence, 568, 605–636
- linear transform, 336–360
- linear transformation, 337, 342, 346, 353, 358–361, 518, 519
- linearly dependent, 608, 605–636
- linearly independent, 553, 608, 605–636, 639, 643
- linguistic vector, 15
- `log()`, 176, 293
- `log10()`, 176, 288
- logarithm, 293
- Lovelace, Ada, 478
- lower triangular, 533
- lower-left triangular, 533
- magnitude, 16, 20, 73
- Markov chain, 141
- Matlab, 2, 8, 10, 73, 72–80, 95, 97, 98, 98, 100, 108, 112, 117, 121, 134, 138, 147, 147, 148, 150–152, 160, 162, 163, 169, 174, 176, 176, 178, 191–193, 211, 214, 217, 218, 221, 222, 233, 235–237, 255, 256, 270, 272, 274, 283, 286, 289, 290, 292, 293, 293, 295, 306, 327, 331, 333, 362, 384, 384, 387, 390, 399, 401, 413, 414, 417, 428, 429, 431, 444, 444, 446, 447, 462, 466, 468, 472, 473, 477–479, 503, 523, 559, 566, 568, 571, 573, 574, 596–598, 613, 633, 636, 644, 646–648, 653, 660, 661, 663
- matrix, 96, 134

- matrix norm, 444, 445, 445–448, 450–452, 455, 456, 462
 matrix power, 144
 matrix product, 143
 matrix-vector form, 96, 123
 matrix-vector product, 139
 maximum, 228
 mean, 463
`mean()`, 444, 463
 minor, 530, 539
 Molar, Cleve, 290
 Moore–Penrose inverse, 346
 multiplicity, 383, 384, 387, 400, 564, 564, 566, 568, 596, 597, 643, 645, 661
NaN, 98
 natural numbers, 10
 n D-cube, 509, 509, 511
 n D-parallelepiped, 511
 n D-volume, 509, 509, 511, 512
 negative, 23
 nilpotent, 522
 no solution, 91, 114, 216
 non-diagonalisable matrix, 643
 non-homogeneous, 114
 non-orthogonal coordinates, 637
 nonconformant arguments, 77
 nonlinear equation, 86, 87, 88
`norm()`, 72, 73, 444, 446, 447
 normal equation, 289, 498
 normal vector, 48, 49, 55, 58, 59, 62, 65, 68
 not a number, 98
 null, 248, 313
 nullity, 262, 262, 264, 267, 273, 631
 nullspace, 248, 262, 272–274, 313, 623
 Occam’s razor, 494
 Octave, 2, 8, 10, 73, 72–80, 95, 97, 98, 98, 100, 108, 112, 117, 121, 134, 138, 147, 147, 148, 150–152, 160, 162, 163, 169, 174, 176, 176, 178, 191–193, 211, 211, 214, 217, 218, 221, 222, 233, 235–237, 255, 256, 270, 272, 274, 283, 286, 289, 290, 292, 293, 293–296, 306, 327, 331, 333, 362, 384, 384, 387, 390, 399, 401, 413, 414, 417, 428, 429, 431, 444, 444, 446, 447, 462, 466, 468, 472, 473, 477–479, 503, 523, 559, 566, 568, 571, 573, 574, 596–598, 613, 633, 636, 644, 646–648, 653, 660, 661, 663
 ones matrix, 147
`ones()`, 98, 147
 orangutan, 579
 orbital period, 286
 orthogonal, 46, 47, 48, 60, 182, 183, 410, 423, 554
 orthogonal basis, 400
 orthogonal complement, 311, 311, 313, 332, 333, 335
 orthogonal decomposition, 321, 335
 orthogonal matrix, 184, 184–190, 200, 202, 203, 209, 223, 384, 416, 448, 509, 638
 orthogonal projection, 298, 300, 296–322, 330, 331
 orthogonal set, 183, 183, 197–199
 orthogonally diagonalisable, 416, 417, 429, 638
 orthonormal, 209, 423
 orthonormal basis, 249, 251, 257, 260, 272, 274, 300, 306, 331, 468, 472, 616, 618, 620, 623, 633
 orthonormal set, 183, 183–184, 186, 197–199, 202, 249, 610
 parallelepiped, 340
 parallelepiped volume, 65–67
 parallelogram area, 57, 60, 61, 67
 parameter, 26, 51
 parametric equation, 26, 26, 27, 43, 51, 50–53, 123, 128
 parentheses, 135
 partial fraction, 101
 particular solution, 649, 653, 655, 662
 PCA, 456
 perp, 317, 318, 320, 335
 perpendicular, 46
 perpendicular component, 318, 320, 334
 pigeonhole principle, 415
 pixel pattern, 453
 planets, 286
 player rating, 278, 279, 305
 Poincaré, Henri, 225
 polar form, 587
`poly()`, 559
 population, 577–592

- population model, 590
population modelling, 551, 637
position vector, 12, 12, 14, 18, 26, 51
precision, 73, 162
principal axes, 424, 431
principal component analysis, 456, 467, 479
principal components, 461, 456–468
principal vector, 461, 456–466, 468
proj, 298–301, 304–308, 317, 318, 320, 322, 329, 330, 341, 371
pseudo-inverse, 346, 345–352, 362, 363
Punch, John, 494
Pythagoras, 10
- QR factorisation, 211, 568
qr-code, 74, 98
quadratic equation, 418
quadratic form, 423, 423, 424
quit, 73
quote, 147
- \mathbb{R} , 11
`randn()`, 147, 575
random matrix, 147
random perturbation, 598
range, 245
rank, 218, 218–220, 224, 236, 238, 251, 260, 264, 267, 275, 279, 351, 363, 452, 620, 623, 631
rank theorem, 264, 273
rational function, 121
rational numbers, 10
`rcond`, 217
`rcond()`, 97, 98
real numbers, 11
Recorde, Robert, 14
reduced row echelon form, 4, 109, 110, 110–112, 131
reflection, 182, 184
regularisation parameter, 498, 501, 502
relative error, 502
repeated eigenvalue, 564, 574–577, 598, 643
`reshape()`, 293
rhombus, 340
Richardson, L. F., 326
right triangular, 533
right-angles, 14, 24
river length, 325
 \mathbb{R}^n , 14
- rotation, 182, 184
rotation and/or reflection, 184, 186
row space, 246, 260, 270, 272, 274, 313, 314
row vector, 135, 186, 246, 267, 461, 631
- scalar, 11, 23, 26, 29, 35, 41, 51
scalar multiplication, 23
scalar product, 138
scalar triple product, 65, 66, 70
`scatter()`, 444, 459
`scatter3()`, 444
searching, 470
`semilogy()`, 444
sensitivity, 574–577, 598
serval, 590
Seven Bridges of Königsberg, 401
shear transformation, 357
Sierpinski carpet, 477
Sierpinski network, 387
Sierpinski triangle, 477
significant digits, 73
Singular Spectrum Analysis, 254, 274
singular value, 209, 216–218, 220, 251, 267, 415, 440, 444, 445, 452, 482, 484, 500, 501, 516, 593, 620, 623, 631
singular value decomposition, 4, 133, 205, 208, 209, 209–212, 214, 216, 218, 219, 225, 227, 228, 230–234, 236, 238, 239, 250, 251, 254, 256, 260–262, 264, 272, 274–277, 279, 283, 286, 290, 292, 295, 297, 302, 305, 323, 327–329, 331, 348, 362, 375, 376, 415–417, 428, 432, 437, 438, 440, 447, 450, 452, 455, 456, 461, 466, 468, 470, 472, 475–479, 481, 484, 485, 494, 502, 504, 506, 515, 516, 559, 568, 592, 593, 603, 620, 623, 632
singular vector, 209, 251, 593
size, 14, 14, 134, 148
`size()`, 73, 147
smallest change, 277, 278, 281, 282
smallest solution, 290, 291, 500
Southern Oscillation Index, 254
span, 125, 125–127, 129, 131, 244–246, 248, 251, 252, 260, 262, 270, 297, 300, 305, 308, 311, 313, 314, 317, 319, 321, 330, 331, 333–335, 372, 383, 392, 410, 433, 567, 616, 617, 619,

- 620, 641, 642, 664, 665
 species, 458
 spectral norm, 445
 square matrix, 96, 97, 98, 134, 147, 165, 220,
 267, 377, 381, 384, 392, 403, 405,
 416, 417, 509, 513, 514, 516, 517,
 533, 559, 567, 584, 631, 637–639,
 643, 650
 square-root, 429
 standard basis, 625, 627, 628, 635, 636
 standard coordinate system, 13, 14
 standard deviation, 444, 463, 464, 466
 standard matrix, 342, 346, 353, 358, 360,
 361
 standard unit vector, 23, 24, 46, 47, 59, 64,
 67, 69, 183, 259, 342, 625, 627, 635
`std()`, 444, 464
 steganography, 478
 stereo pair, 22, 179, 184, 189, 198, 203, 266
 subspace, 240, 242, 244–246, 248, 249, 257,
 260, 268, 270, 300, 306, 311, 313,
 317, 318, 320, 321, 381, 616, 618–
 620, 623, 625
 subtraction, 137, 149
 sum, 23, 137, 149
 SVD, 4, 133, 205, 208–212, 214, 216, 218,
 219, 225, 227, 228, 230–234, 236,
 238, 239, 250, 251, 254, 256, 260–
 262, 264, 272, 274–277, 279, 283,
 286, 290, 292, 295, 297, 302, 305,
 323, 327–329, 331, 348, 362, 375,
 376, 415–417, 428, 432, 437, 438,
 440, 447, 450, 452, 455, 456, 461,
 466, 468, 470, 472, 475–479, 481,
 484, 485, 494, 502, 504, 506, 515,
 516, 559, 568, 592, 593, 603, 620,
 623, 632
`svd()`, 211
`svds()`, 444, 464, 466, 468, 472
 symmetric matrix, 146, 163, 383, 409, 410,
 414, 415, 576, 638, 645
 system, 87, 96, 107, 112
 table tennis, 278
 Tasmanian Devil, 600
 Tasmanian Devils, 162
 Tikhonov regularisation, 498, 498–501, 504
 tomography, 292
 trace, 561, 595
 transformation, 336, 337
 transpose, 145, 146, 147, 152, 163
 triangle inequality, 43, 44, 45, 55, 447
 triangular matrix, 532, 533, 533, 554, 558,
 640
 triple product, scalar, 65
 unique solution, 91, 97, 114, 169, 217, 220,
 267, 631
 unit cube, 340
 unit square, 338
 unit vector, 16, 48, 182, 183, 376, 423, 424
 upper triangular, 533
 upper-right triangular, 533
 van de Snepscheut, Jan L. A., 97
 variance, 466
 vector, 14, 16
 vector product, 59
 vectors, 12
 velocity vector, 13
 volume, parallelepiped, 65, 70
 walking gait, 274
 weather forecasts, 259
 Wikipedia, 11, 120, 283, 286, 292, 498, 579,
 590
 Wiles, Andrew, 225
 wine recognition, 462
 word vector, 20, 38, 54, 467, 471
 work, 40
 zero, 114, 291, 403, 415, 484
 zero column, 525
 zero matrix, 135, 147
 zero row, 525
 zero vector, 14, 16, 17, 313, 381
`zeros()`, 147