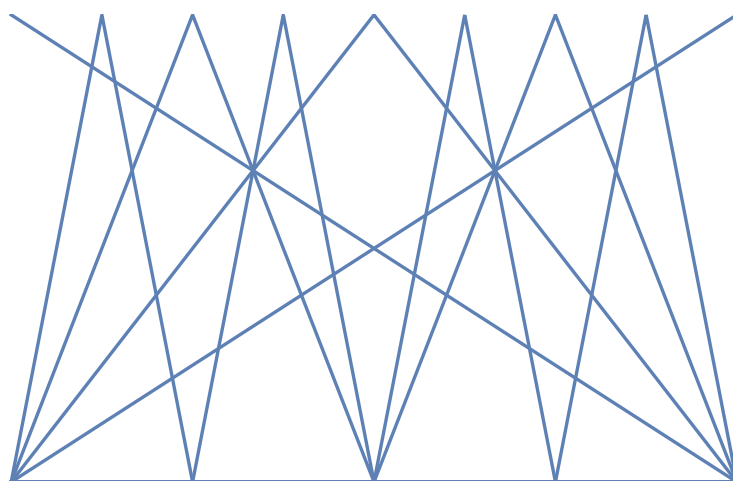


# Topics in Linear and Nonlinear Functional Analysis

Gerald Teschl



Graduate Studies  
in Mathematics  
Volume (to appear)



American Mathematical Society  
Providence, Rhode Island

Gerald Teschl  
Fakultät für Mathematik  
Oskar-Mogenstern-Platz 1  
Universität Wien  
1090 Wien, Austria

*E-mail:* Gerald.Teschl@univie.ac.at

*URL:* <http://www.mat.univie.ac.at/~gerald/>

---

*2010 Mathematics subject classification.* 46-01, 46E30, 47H10, 47H11, 58Exx, 76D05

---

**Abstract.** This manuscript provides a brief introduction to linear and non-linear Functional Analysis. It covers basic Hilbert and Banach space theory as well as some advanced topics like operator semigroups, mapping degrees and fixed point theorems.

*Keywords and phrases.* Functional Analysis, Banach space, Hilbert space, operator semigroup, mapping degree, fixed point theorem, differential equation, Navier–Stokes equation.

Typeset by  $\mathcal{A}\mathcal{M}\mathcal{S}$ - $\text{\LaTeX}$  and Makeindex.

Version: December 23, 2021

Copyright © 1998–2020 by Gerald Teschl



---

# Contents

Preface	vii
<b>Part 1. Functional Analysis</b>	
Chapter 1. A first look at Banach and Hilbert spaces	3
§1.1. Introduction: Linear partial differential equations	3
§1.2. The Banach space of continuous functions	7
§1.3. The geometry of Hilbert spaces	17
§1.4. Completeness	24
§1.5. Compactness	25
§1.6. Bounded operators	28
§1.7. Sums and quotients of Banach spaces	33
§1.8. Spaces of continuous and differentiable functions	37
Chapter 2. Hilbert spaces	41
§2.1. Orthonormal bases	41
§2.2. The projection theorem and the Riesz representation theorem	48
§2.3. Operators defined via forms	51
§2.4. Orthogonal sums and tensor products	56
§2.5. Applications to Fourier series	58
Chapter 3. Compact operators	65
§3.1. Compact operators	65
§3.2. The spectral theorem for compact symmetric operators	68
§3.3. Applications to Sturm–Liouville operators	74

---

§3.4. Estimating eigenvalues	82
§3.5. Singular value decomposition of compact operators	85
§3.6. Hilbert–Schmidt and trace class operators	89
Chapter 4. The main theorems about Banach spaces	97
§4.1. The Baire theorem and its consequences	97
§4.2. The Hahn–Banach theorem and its consequences	104
§4.3. The adjoint operator	114
§4.4. Weak convergence	119
Chapter 5. Bounded linear operators	129
§5.1. Banach algebras	129
§5.2. The $C^*$ algebra of operators and the spectral theorem	140
§5.3. Spectral measures	143
 <b>Part 2. Advanced Functional Analysis</b>	
Chapter 6. More on convexity	151
§6.1. The geometric Hahn–Banach theorem	151
§6.2. Convex sets and the Krein–Milman theorem	155
§6.3. Weak topologies	160
§6.4. Beyond Banach spaces: Locally convex spaces	165
§6.5. Uniformly convex spaces	172
Chapter 7. Advanced Spectral theory	177
§7.1. Spectral theory for compact operators	177
§7.2. Fredholm operators	184
§7.3. The Gelfand representation theorem	191
Chapter 8. Unbounded operators	199
§8.1. Closed operators	199
§8.2. Spectral theory for unbounded operators	209
§8.3. Reducing subspaces and spectral projections	213
§8.4. Relatively bounded and relatively compact operators	217
§8.5. Unbounded Fredholm operators	222
 <b>Part 3. Nonlinear Functional Analysis</b>	
Chapter 9. Analysis in Banach spaces	229
§9.1. Single variable calculus in Banach spaces	229

---

§9.2. Multivariable calculus in Banach spaces	232
§9.3. Minimizing nonlinear functionals via calculus	244
§9.4. Minimizing nonlinear functionals via compactness	249
§9.5. Contraction principles	257
§9.6. Ordinary differential equations	261
§9.7. Bifurcation theory	267
Chapter 10. Operator semigroups	273
§10.1. Uniformly continuous operator groups	273
§10.2. Strongly continuous semigroups	276
§10.3. Generator theorems	284
Chapter 11. The nonlinear Schrödinger equation	301
§11.1. Local well-posedness in $H^r$ for $r > \frac{n}{2}$	301
§11.2. Strichartz estimates	304
§11.3. Well-posedness in $L^2$ and $H^1$	308
§11.4. Blowup in $H^1$	314
§11.5. Standing waves	316
Chapter 12. The Brouwer mapping degree	319
§12.1. Introduction	319
§12.2. Definition of the mapping degree and the determinant formula	321
§12.3. Extension of the determinant formula	325
§12.4. The Brouwer fixed point theorem	332
§12.5. Kakutani's fixed point theorem and applications to game theory	334
§12.6. Further properties of the degree	337
§12.7. The Jordan curve theorem	340
Chapter 13. The Leray–Schauder mapping degree	341
§13.1. The mapping degree on finite dimensional Banach spaces	341
§13.2. Compact maps	342
§13.3. The Leray–Schauder mapping degree	344
§13.4. The Leray–Schauder principle and the Schauder fixed point theorem	346
§13.5. Applications to integral and differential equations	348
Chapter 14. Monotone maps	355
§14.1. Monotone maps	355

§14.2. The nonlinear Lax–Milgram theorem	357
§14.3. The main theorem of monotone maps	358
Appendix A. Some set theory	361
Appendix B. Metric and topological spaces	369
§B.1. Basics	369
§B.2. Convergence and completeness	375
§B.3. Functions	379
§B.4. Product topologies	381
§B.5. Compactness	384
§B.6. Separation	390
§B.7. Connectedness	394
§B.8. Continuous functions on metric spaces	398
Bibliography	405
Glossary of notation	409
Index	413

---

# Preface

The present manuscript was written for my course *Functional Analysis* given at the University of Vienna in winter 2004 and 2009. The second part are the notes for my course *Nonlinear Functional Analysis* held at the University of Vienna in Summer 1998, 2001, and 2018. The two parts are essentially independent. In particular, the first part does not assume any knowledge from measure theory (at the expense of hardly mentioning  $L^p$  spaces). However, there is an accompanying part on Real Analysis [48], where these topics are covered.

It is updated whenever I find some errors and extended from time to time. Hence you might want to make sure that you have the most recent version, which is available from

<http://www.mat.univie.ac.at/~gerald/ftp/book-fa/>

**Please do not redistribute this file or put a copy on your personal webpage but link to the page above.**

## *Goals*

The main goal of the present book is to give students a concise introduction which gets to some interesting results without much ado while using a sufficiently general approach suitable for further studies. Still I have tried to always start with some interesting special cases and then work my way up to the general theory. While this unavoidably leads to some duplications, it usually provides much better motivation and implies that the core material always comes first (while the more general results are then optional). Moreover, this book is *not* written under the assumption that it will be



read linearly starting with the first chapter and ending with the last. Consequently, I have tried to separate core and optional materials as much as possible while keeping the optional parts as independent as possible.

Furthermore, my aim is not to present an encyclopedic treatment but to provide the reader with a versatile toolbox for further study. Moreover, in contradistinction to many other books, I do not have a particular direction in mind and hence I am trying to give a broad introduction which should prepare you for diverse fields such as spectral theory, partial differential equations, or probability theory. This is related to the fact that I am working in mathematical physics, an area where you never know what mathematical theory you will need next.

I have tried to keep a balance between verbosity and clarity in the sense that I have tried to provide sufficient detail for being able to follow the arguments but without drowning the key ideas in boring details. In particular, you will find a *show this* from time to time encouraging the reader to check the claims made (these tasks typically involve only simple routine calculations). Moreover, to make the presentation student friendly, I have tried to include many worked out examples within the main text. Some of them are standard counterexamples pointing out the limitations of theorems (and explaining why the assumptions are important). Others show how to use the theory in the investigation of practical examples.

### *Preliminaries*

The present manuscript is intended to be gentle when it comes to required background. Of course I assume basic familiarity with analysis (real and complex numbers, limits, differentiation, basic (Riemann) integration, open sets) and linear algebra (finite dimensional vector spaces, matrices).

Apart from this natural assumptions I also expect some familiarity with metric spaces and point set topology. However, only a few basic things are required to begin with. This and much more is collected in the Appendix and I will refer you there from time to time such that you can refresh your memory should need arise. Moreover, you can always go there if you are unsure about a certain term (using the extensive index) or if there should be a need to clarify notation or conventions. I prefer this over referring you to several other books which might not always be readily available. For convenience, the Appendix contains full proofs in case one needs to fill some gaps. As some things are only outlined (or outsourced to exercises), it will require extra effort in case you see all this for the first time.

On the other hand I do not assume familiarity with Lebesgue integration and consequently  $L^p$  spaces will only be briefly mentioned as the completion

of continuous functions with respect to the corresponding integral norms in the first part. At a few places I also assume some basic results from complex analysis but it will be sufficient to just believe them.

The second part of course requires basic familiarity with functional analysis and measure theory (Lebesgue and Sobolev spaces). But apart from this it is again independent from the first two parts.

## *Content*

Below follows a short description of each chapter together with some hints which parts can be skipped.

**Chapter 1.** The first part starts with Fourier's treatment of the heat equation which led to the theory of Fourier analysis as well as the development of spectral theory which drove much of the development of functional analysis around the turn of the last century. In particular, the first chapter tries to introduce and motivate some of the key concepts and should be covered in detail except for Section 1.8 which introduces some interesting examples for later use.

**Chapter 2** discusses basic Hilbert space theory and should be considered core material except for the last section discussing applications to Fourier series. They will only be used in some examples and could be skipped in case they are covered in a different course.

**Chapter 3** develops basic spectral theory for compact self-adjoint operators. The first core result is the spectral theorem for compact symmetric (self-adjoint) operators which is then applied to Sturm–Liouville problems. Of course this application could be skipped, but this would reduce the didactical concept to absurdity. Nevertheless it is clearly possible to shorten the material as none of it (including the follow-up section which touches upon some more tools from spectral theory) will be required in later chapters. The last two sections on singular value decompositions as well as Hilbert–Schmidt and trace class operators cover important topics for applications, but will again not be required later on.

**Chapter 4** discusses what is typically considered as the core results from Banach space theory. In order to keep the topological requirements to a minimum some advanced topics are shifted to the following chapters.

**Chapter 5** develops spectral theory for bounded self-adjoint operators via the framework of  $C^*$  algebras. The last section contains some optional results establishing the connection with the measure theoretic formulation of the spectral theorem.

The next chapters contain selected advanced topics.

**Chapter 6** centers around convexity. Except for the geometric Hahn–Banach theorem, which is a prerequisite for the other sections, the remaining sections are independent of each other to simplify the selection of topics.

**Chapter 7** presents some advanced topics from spectral theory: The Gelfand representation theorem, spectral theory for compact operators in Banach spaces and Fredholm theory. Again these sections are independent of each other except for the fact that Section 7.1, which contains the spectral theorem for compact operators, and hence the Fredholm alternative for compact perturbations of the identity, is of course used to identify compact perturbations of the identity as premier examples of Fredholm operators.

**Chapter 8** touches upon unbounded operators starting with the basic results about closed operators. Since unbounded operators play an increasing role in applications I felt it is appropriate to discuss at least some basics.

Finally, there is a part on nonlinear functional analysis.

**Chapter 9** discusses analysis in Banach spaces (with a view towards applications in the calculus of variations and infinite dimensional dynamical systems).

**Chapter 10** finally gives a brief introduction to operator semigroups.

**Chapter 11** applies the results obtained so far to an ubiquitous example, the nonlinear Schrödinger equation.

**Chapter 12 and 13** cover degree theory and fixed point theorems in finite and infinite dimensional spaces. Several applications to integral equations, ordinary differential equations and to the stationary Navier–Stokes equation are given.

**Chapter 14** provides some basics about monotone maps.

Sometimes also the historic development of the subject is of interest. This is however not covered in the present book and we refer to [28, 42, 43] as good starting points.

### *To the teacher*

There are a couple of courses to be taught from this book. First of all there is of course a basic functional analysis course: Chapters 1 to 4 (skipping some optional material as discussed above) and perhaps adding some material from Chapter 5 or 6. If one wants to cover Lebesgue spaces, this can be easily done by including Chapters 1, 2, and 3 from [48]. In this case one could cover Section 1.2 (Section 1.1 contains just motivation), give an outline of Section 1.3 (by covering Dynkin’s  $\pi$ - $\lambda$  theorem, the uniqueness theorem for measures, and then quoting the existence theorem for Lebesgue measure), cover Section 1.5. The core material from Chapter 2 are the

first two sections and from Chapter 3 the first three sections. I think that this gives a well-balanced introduction to functional analysis which contains several optional topics to choose from depending on personal preferences and time constraints.

The remaining material from the first part could then be used for a course on advanced functional analysis. Typically one could also add some further topics from the second part or some material from unbounded operators in Hilbert spaces following [47] (where one can start with Chapter 2) or from unbounded operators in Banach spaces following the book by Kato [25] (e.g. Sections 3.4, 3.5 and 4.1).

The third part gives a short basis for a course on nonlinear functional analysis.

Problems relevant for the main text are marked with a "\*". A Solutions Manual will be available electronically for instructors only.

### *Acknowledgments*

I wish to thank my readers, Kerstin Ammann, Phillip Bachler, Batuhan Bayır, Alexander Beigl, Mikhail Botchkarev, Ho Boon Suan, Peng Du, Christian Ekstrand, Damir Ferizović, Michael Fischer, Raffaello Giulietti, Melanie Graf, Josef Greilhuber, Julian Grüber, Matthias Hammerl, Jona Marie Hasenbach, Nobuya Kakehashi, Jerzy Knopik, Nikolas Knotz, Florian Kogelbauer, Helge Krüger, Reinhold Küstner, Oliver Leingang, Juho Leppäkanigas, Joris Mestdagh, Alice Mikikits-Leitner, Claudiu Mîndrilă, Jakob Möller, Caroline Moosmüller, Matthias Ostermann, Piotr Owczarek, Martina Pflegpeter, Mateusz Piorkowski, Tobias Preinerstorfer, Maximilian H. Rüp, Tidhar Sarel, Chiara Schindler, Christian Schmid, Laura Shou, Bertram Tschiderer, Liam Urban, Vincent Valmorin, David Wallauch, Richard Welke, David Wimmesberger, Gunter Wirthumer, Song Xiaojun, Markus Youssef, Rudolf Zeidler, and colleagues Pierre-Antoine Absil, Nils C. Framstad, Fritz Gesztesy, Heinz Hanßmann, Günther Hörmann, Aleksey Kostenko, Wallace Lam, Daniel Lenz, Johanna Michor, Viktor Qvarfordt, Alex Strohmaier, David C. Ullrich, Hendrik Vogt, Marko Stautz, Maxim Zinchenko who have pointed out several typos and made useful suggestions for improvements. Moreover, I am most grateful to Iryna Karpenko who read several parts of the manuscript, provided long lists of typos, and also contributed some of the problems. I am also grateful to Volker Enß for making his lecture notes on nonlinear Functional Analysis available to me.

**Finally, no book is free of errors. So if you find one, or if you have comments or suggestions (no matter how small), please let me know.**

Gerald Teschl

Vienna, Austria  
January, 2019

---

*Part 1*

# Functional Analysis



# A first look at Banach and Hilbert spaces

Functional analysis is an important tool in the investigation of all kind of problems in pure mathematics, physics, biology, economics, etc.. In fact, it is hard to find a branch in science where functional analysis is not used.

The main objects are (infinite dimensional) vector spaces with different concepts of convergence. The classical theory focuses on linear operators (i.e., functions) between these spaces but nonlinear operators are of course equally important. However, since one of the most important tools in investigating nonlinear mappings is linearization (differentiation), linear functional analysis will be our first topic in any case.

## 1.1. Introduction: Linear partial differential equations

Rather than listing an overwhelming number of classical examples I want to focus on one: linear partial differential equations. We will use this example as a guide throughout our first three chapters and will develop all necessary tools for a successful treatment of our particular problem.

In his investigation of heat conduction Fourier was led to the (one dimensional) **heat** or diffusion equation

$$\frac{\partial}{\partial t}u(t, x) = \frac{\partial^2}{\partial x^2}u(t, x). \quad (1.1)$$

Here  $u : \mathbb{R} \times [0, 1] \rightarrow \mathbb{R}$  is the temperature distribution in a thin rod at time  $t \in \mathbb{R}$  at the point  $x \in [0, 1]$ . It is usually assumed, that the temperature at  $x = 0$  and  $x = 1$  is fixed, say  $u(t, 0) = a$  and  $u(t, 1) = b$ . By considering  $u(t, x) \rightarrow u(t, x) - a - (b - a)x$  it is clearly no restriction to assume  $a = b = 0$ .



Moreover, the initial temperature distribution  $u(0, x) = u_0(x)$  is assumed to be known as well.

Since finding the solution seems at first sight unfeasable, we could try to find at least some solutions of (1.1). For example, we could make an ansatz for  $u(t, x)$  as a product of two functions, each of which depends on only one variable, that is,

$$u(t, x) := w(t)y(x). \quad (1.2)$$

Plugging this ansatz into the heat equation we arrive at

$$\dot{w}(t)y(x) = y''(x)w(t), \quad (1.3)$$

where the dot refers to differentiation with respect to  $t$  and the prime to differentiation with respect to  $x$ . Bringing all  $t$ ,  $x$  dependent terms to the left, right side, respectively, we obtain

$$\frac{\dot{w}(t)}{w(t)} = \frac{y''(x)}{y(x)}. \quad (1.4)$$

Accordingly, this ansatz is called **separation of variables**.

Now if this equation should hold for all  $t$  and  $x$ , the quotients must be equal to a constant  $-\lambda$  (we choose  $-\lambda$  instead of  $\lambda$  for convenience later on). That is, we are led to the equations

$$-\dot{w}(t) = \lambda w(t) \quad (1.5)$$

and

$$-y''(x) = \lambda y(x), \quad y(0) = y(1) = 0, \quad (1.6)$$

which can easily be solved. The first one gives

$$w(t) = c_1 e^{-\lambda t} \quad (1.7)$$

and the second one

$$y(x) = c_2 \cos(\sqrt{\lambda}x) + c_3 \sin(\sqrt{\lambda}x). \quad (1.8)$$

However,  $y(x)$  must also satisfy the boundary conditions  $y(0) = y(1) = 0$ . The first one  $y(0) = 0$  is satisfied if  $c_2 = 0$  and the second one yields ( $c_3$  can be absorbed by  $w(t)$ )

$$\sin(\sqrt{\lambda}) = 0, \quad (1.9)$$

which holds if  $\lambda = (\pi n)^2$ ,  $n \in \mathbb{N}$  (in the case  $\lambda < 0$  we get  $\sinh(\sqrt{-\lambda}) = 0$ , which cannot be satisfied and explains our choice of sign above). In summary, we obtain the solutions

$$u_n(t, x) := c_n e^{-(\pi n)^2 t} \sin(n\pi x), \quad n \in \mathbb{N}. \quad (1.10)$$

So we have found a large number of solutions, but we still have not dealt with our initial condition  $u(0, x) = u_0(x)$ . This can be done using the superposition principle which holds since our equation is linear: Any finite linear combination of the above solutions will again be a solution. Moreover,

under suitable conditions on the coefficients we can even consider infinite linear combinations. In fact, choosing

$$u(t, x) := \sum_{n=1}^{\infty} c_n e^{-(\pi n)^2 t} \sin(n\pi x), \quad (1.11)$$

where the coefficients  $c_n$  decay sufficiently fast (e.g. absolutely summable), we obtain further solutions of our equation. Moreover, these solutions satisfy

$$u(0, x) = \sum_{n=1}^{\infty} c_n \sin(n\pi x) \quad (1.12)$$

and expanding the initial conditions into a Fourier sine series

$$u_0(x) = \sum_{n=1}^{\infty} \hat{u}_{0,n} \sin(n\pi x), \quad (1.13)$$

we see that the solution of our original problem is given by (1.11) if we choose  $c_n = \hat{u}_{0,n}$  (cf. Problem 1.2).

Of course for this last statement to hold we need to ensure that the series in (1.11) converges and that we can interchange summation and differentiation. You are asked to do so in Problem 1.1.

In fact, many equations in physics can be solved in a similar way:

• **Reaction-Diffusion equation:**

$$\begin{aligned} \frac{\partial}{\partial t} u(t, x) - \frac{\partial^2}{\partial x^2} u(t, x) + q(x) u(t, x) &= 0, \\ u(0, x) &= u_0(x), \\ u(t, 0) = u(t, 1) &= 0. \end{aligned} \quad (1.14)$$

Here  $u(t, x)$  could be the density of some gas in a pipe and  $q(x) > 0$  describes that a certain amount per time is removed (e.g., by a chemical reaction).

• **Wave equation:**

$$\begin{aligned} \frac{\partial^2}{\partial t^2} u(t, x) - \frac{\partial^2}{\partial x^2} u(t, x) &= 0, \\ u(0, x) &= u_0(x), \quad \frac{\partial u}{\partial t}(0, x) = v_0(x) \\ u(t, 0) = u(t, 1) &= 0. \end{aligned} \quad (1.15)$$

Here  $u(t, x)$  is the displacement of a vibrating string which is fixed at  $x = 0$  and  $x = 1$ . Since the equation is of second order in time, both the initial displacement  $u_0(x)$  and the initial velocity  $v_0(x)$  of the string need to be known.

• **Schrödinger equation:**

$$\begin{aligned} i \frac{\partial}{\partial t} u(t, x) &= -\frac{\partial^2}{\partial x^2} u(t, x) + q(x)u(t, x), \\ u(0, x) &= u_0(x), \\ u(t, 0) &= u(t, 1) = 0. \end{aligned} \quad (1.16)$$

Here  $|u(t, x)|^2$  is the probability distribution of a particle trapped in a box  $x \in [0, 1]$  and  $q(x)$  is a given external potential which describes the forces acting on the particle.

All these problems (and many others) lead to the investigation of the following problem

$$Ly(x) = \lambda y(x), \quad L := -\frac{d^2}{dx^2} + q(x), \quad (1.17)$$

subject to the **boundary conditions**

$$y(a) = y(b) = 0. \quad (1.18)$$

Such a problem is called a **Sturm–Liouville boundary value problem**. Our example shows that we should prove the following facts about Sturm–Liouville problems:

- (i) The Sturm–Liouville problem has a countable number of eigenvalues  $E_n$  with corresponding eigenfunctions  $u_n$ , that is,  $u_n$  satisfies the boundary conditions and  $Lu_n = E_n u_n$ .
- (ii) The eigenfunctions  $u_n$  are complete, that is, any *nice* function  $u$  can be expanded into a generalized Fourier series

$$u(x) = \sum_{n=1}^{\infty} c_n u_n(x).$$

This problem is very similar to the eigenvalue problem of a matrix and we are looking for a generalization of the well-known fact that every symmetric matrix has an orthonormal basis of eigenvectors. However, our linear operator  $L$  is now acting on some space of functions which is not finite dimensional and it is not at all clear what (e.g.) orthogonal should mean in this context. Moreover, since we need to handle infinite series, we need convergence and hence we need to define the distance of two functions as well.

Hence our program looks as follows:

- What is the distance of two functions? This automatically leads us to the problem of convergence and completeness.
- If we additionally require the concept of orthogonality, we are led to Hilbert spaces which are the proper setting for our eigenvalue problem.

- Finally, the spectral theorem for compact symmetric operators will be the solution of our above problem.

**Problem 1.1.** Suppose  $\sum_{n=1}^{\infty} |c_n| < \infty$ . Show that (1.11) is continuous for  $(t, x) \in [0, \infty) \times [0, 1]$  and solves the heat equation for  $(t, x) \in (0, \infty) \times [0, 1]$ . (Hint: Weierstraß M-test. When can you interchange the order of summation and differentiation?)

**Problem 1.2.** Show that for  $n, m \in \mathbb{N}$  we have

$$2 \int_0^1 \sin(n\pi x) \sin(m\pi x) dx = \begin{cases} 1 & , n = m, \\ 0, & n \neq m. \end{cases}$$

Conclude that the Fourier sine coefficients are given by

$$\hat{u}_{0,n} = 2 \int_0^1 \sin(n\pi x) u_0(x) dx$$

provided the sum in (1.13) converges uniformly. Conclude that in this case the solution can be expressed as

$$u(t, x) = \int_0^1 K(t, x, y) u_0(y) dy, \quad t > 0,$$

where

$$\begin{aligned} K(t, x, y) &:= 2 \sum_{n=1}^{\infty} e^{-(\pi n)^2 t} \sin(n\pi x) \sin(n\pi y) \\ &= \frac{1}{2} \left( \vartheta\left(\frac{x-y}{2}, i\pi t\right) - \vartheta\left(\frac{x+y}{2}, i\pi t\right) \right). \end{aligned}$$

Here

$$\vartheta(z, \tau) := \sum_{n \in \mathbb{Z}} e^{i\pi n^2 \tau + 2\pi i n z} = 1 + 2 \sum_{n \in \mathbb{N}} e^{i\pi n^2 \tau} \cos(2\pi n z), \quad \text{Im}(\tau) > 0,$$

is the **Jacobi theta function**.

## 1.2. The Banach space of continuous functions

Our point of departure will be the set of continuous functions  $C(I)$  on a compact interval  $I := [a, b] \subset \mathbb{R}$ . Since we want to handle both real and complex models, we will formulate most results for the more general complex case only. In fact, most of the time there will be no difference but we will add a remark in the rare case where the real and complex case do indeed differ.

One way of declaring a distance, well-known from calculus, is the **maximum norm** of a function  $f \in C(I)$ :

$$\|f\|_{\infty} := \max_{x \in I} |f(x)|. \quad (1.19)$$

It is not hard to see that with this definition  $C(I)$  becomes a normed vector space:

A **normed vector space**  $X$  is a vector space  $X$  over  $\mathbb{C}$  (or  $\mathbb{R}$ ) with a nonnegative function (the **norm**)  $\|\cdot\| : X \rightarrow [0, \infty)$  such that

- $\|f\| > 0$  for  $f \neq 0$  (**positive definiteness**),
- $\|\alpha f\| = |\alpha| \|f\|$  for all  $\alpha \in \mathbb{C}$ ,  $f \in X$  (**positive homogeneity**), and
- $\|f + g\| \leq \|f\| + \|g\|$  for all  $f, g \in X$  (**triangle inequality**).

If positive definiteness is dropped from the requirements, one calls  $\|\cdot\|$  a **seminorm**.

From the triangle inequality we also get the **inverse triangle inequality** (Problem 1.3)

$$|||f\| - \|g\|| \leq \|f - g\|, \quad (1.20)$$

which shows that the norm is continuous.

Also note that norms are closely related to convexity. To this end recall that a subset  $C \subseteq X$  is called **convex** if for every  $x, y \in C$  we also have  $\lambda f + (1 - \lambda)g \in C$  whenever  $\lambda \in (0, 1)$ . Moreover, a mapping  $F : C \rightarrow \mathbb{R}$  is called **convex** if  $F(\lambda f + (1 - \lambda)g) \leq \lambda F(f) + (1 - \lambda)F(g)$  whenever  $\lambda \in (0, 1)$  and  $f, g \in C$ . In our case the triangle inequality plus homogeneity imply that every norm is convex:

$$\|\lambda f + (1 - \lambda)g\| \leq \lambda \|f\| + (1 - \lambda)\|g\|, \quad \lambda \in [0, 1]. \quad (1.21)$$

Moreover, choosing  $\lambda = \frac{1}{2}$  we get back the triangle inequality upon using homogeneity. In particular, the triangle inequality could be replaced by convexity in the definition.

Once we have a norm, we have a **distance**  $d(f, g) := \|f - g\|$  and hence we know when a sequence of vectors  $f_n$  **converges** to a vector  $f$  (namely if  $d(f_n, f) \rightarrow 0$ ). We will write  $f_n \rightarrow f$  or  $\lim_{n \rightarrow \infty} f_n = f$ , as usual, in this case. Moreover, a mapping  $F : X \rightarrow Y$  between two normed spaces is called **continuous** if for every convergent sequence  $f_n \rightarrow f$  from  $X$  we have  $F(f_n) \rightarrow F(f)$  (with respect to the norm of  $X$  and  $Y$ , respectively). In fact, the norm, vector addition, and multiplication by scalars are continuous (Problem 1.4).

Two normed spaces  $X$  and  $Y$  are called **isomorphic** if there exists a linear bijection  $T : X \rightarrow Y$  such that  $T$  and its inverse  $T^{-1}$  are continuous. We will write  $X \cong Y$  in this case. They are called **isometrically isomorphic** if in addition,  $T$  is an **isometry**,  $\|T(f)\| = \|f\|$  for every  $f \in X$ .

In addition to the concept of convergence, we also have the concept of a **Cauchy sequence** and hence the concept of completeness: A normed space

is called **complete** if every Cauchy sequence has a limit. A complete normed space is called a **Banach space**.

**Example 1.1.** By completeness of the real numbers,  $\mathbb{R}$  as well as  $\mathbb{C}$  with the absolute value as norm are Banach spaces.  $\diamond$

**Example 1.2.** The space  $\ell^1(\mathbb{N})$  of all complex-valued sequences  $a = (a_j)_{j=1}^\infty$  for which the norm

$$\|a\|_1 := \sum_{j=1}^{\infty} |a_j| \quad (1.22)$$

is finite is a Banach space.

To show this, we need to verify three things: (i)  $\ell^1(\mathbb{N})$  is a vector space, that is, closed under addition and scalar multiplication, (ii)  $\|\cdot\|_1$  satisfies the three requirements for a norm, and (iii)  $\ell^1(\mathbb{N})$  is complete.

First of all, observe

$$\sum_{j=1}^k |a_j + b_j| \leq \sum_{j=1}^k |a_j| + \sum_{j=1}^k |b_j| \leq \|a\|_1 + \|b\|_1$$

for every finite  $k$ . Letting  $k \rightarrow \infty$ , we conclude that  $\ell^1(\mathbb{N})$  is closed under addition and that the triangle inequality holds. That  $\ell^1(\mathbb{N})$  is closed under scalar multiplication together with homogeneity as well as definiteness are straightforward. It remains to show that  $\ell^1(\mathbb{N})$  is complete. Let  $a^n = (a_j^n)_{j=1}^\infty$  be a Cauchy sequence; that is, for given  $\varepsilon > 0$  we can find an  $N_\varepsilon$  such that  $\|a^m - a^n\|_1 \leq \varepsilon$  for  $m, n \geq N_\varepsilon$ . This implies, in particular,  $|a_j^m - a_j^n| \leq \varepsilon$  for every fixed  $j$ . Thus  $a_j^n$  is a Cauchy sequence for fixed  $j$  and, by completeness of  $\mathbb{C}$ , it has a limit:  $a_j := \lim_{n \rightarrow \infty} a_j^n$ . Now consider  $\sum_{j=1}^k |a_j^m - a_j^n| \leq \varepsilon$  and take  $m \rightarrow \infty$ :

$$\sum_{j=1}^k |a_j - a_j^n| \leq \varepsilon.$$

Since this holds for all finite  $k$ , we even have  $\|a - a^n\|_1 \leq \varepsilon$ . Hence  $(a - a^n) \in \ell^1(\mathbb{N})$  and since  $a^n \in \ell^1(\mathbb{N})$ , we finally conclude  $a = a^n + (a - a^n) \in \ell^1(\mathbb{N})$ . By our estimate  $\|a - a^n\|_1 \leq \varepsilon$ , our candidate  $a$  is indeed the limit of  $a^n$ .  $\diamond$

**Example 1.3.** The previous example can be generalized by considering the space  $\ell^p(\mathbb{N})$  of all complex-valued sequences  $a = (a_j)_{j=1}^\infty$  for which the norm

$$\|a\|_p := \left( \sum_{j=1}^{\infty} |a_j|^p \right)^{1/p}, \quad p \in [1, \infty), \quad (1.23)$$

is finite. By  $|a_j + b_j|^p \leq 2^p \max(|a_j|, |b_j|)^p = 2^p \max(|a_j|^p, |b_j|^p) \leq 2^p(|a_j|^p + |b_j|^p)$  it is a vector space, but the triangle inequality is only easy to see in the

case  $p = 1$ . (It is also not hard to see that it fails for  $p < 1$ , which explains our requirement  $p \geq 1$ . See also Problem 1.15.)

To prove the triangle inequality we need Young's inequality (Problem 1.8)

$$\alpha^{1/p} \beta^{1/q} \leq \frac{1}{p} \alpha + \frac{1}{q} \beta, \quad \frac{1}{p} + \frac{1}{q} = 1, \quad \alpha, \beta \geq 0, \quad (1.24)$$

which implies **Hölder's inequality**

$$\|ab\|_1 \leq \|a\|_p \|b\|_q \quad (1.25)$$

for  $a \in \ell^p(\mathbb{N})$ ,  $b \in \ell^q(\mathbb{N})$ . In fact, by homogeneity of the norm it suffices to prove the case  $\|a\|_p = \|b\|_q = 1$ . But this case follows by choosing  $\alpha = |a_j|^p$  and  $\beta = |b_j|^q$  in (1.24) and summing over all  $j$ .

Now using  $|a_j + b_j|^p \leq |a_j| |a_j + b_j|^{p-1} + |b_j| |a_j + b_j|^{p-1}$ , we obtain from Hölder's inequality (note  $(p-1)q = p$ )

$$\begin{aligned} \|a + b\|_p^p &\leq \|a\|_p \|(a + b)^{p-1}\|_q + \|b\|_p \|(a + b)^{p-1}\|_q \\ &= (\|a\|_p + \|b\|_p) \|a + b\|_p^{p-1}. \end{aligned}$$

Hence  $\ell^p(\mathbb{N})$  is a normed space. That it is complete can be shown as in the case  $p = 1$  (Problem 1.9).

The unit ball with respect to these norms in  $\mathbb{R}^2$  is depicted in Figure 1.1. One sees that for  $p < 1$  the unit ball is not convex (explaining once more our restriction  $p \geq 1$ ). Moreover, for  $1 < p < \infty$  it is even strictly convex (that is, the line segment joining two distinct points is always in the interior). This is related to the question of equality in the triangle inequality and will be discussed in Problems 1.12 and 1.13.  $\diamond$

**Example 1.4.** The space  $\ell^\infty(\mathbb{N})$  of all complex-valued bounded sequences  $a = (a_j)_{j=1}^\infty$  together with the norm

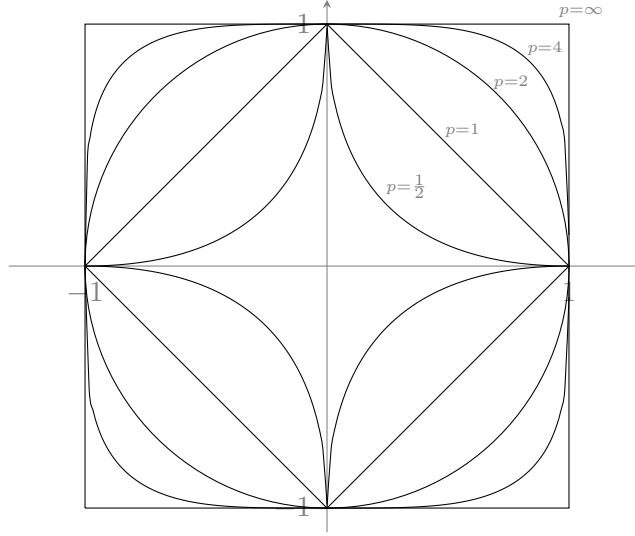
$$\|a\|_\infty := \sup_{j \in \mathbb{N}} |a_j| \quad (1.26)$$

is a Banach space (Problem 1.10). Note that with this definition, Hölder's inequality (1.25) remains true for the cases  $p = 1$ ,  $q = \infty$  and  $p = \infty$ ,  $q = 1$ . The reason for the notation is explained in Problem 1.14.  $\diamond$

**Example 1.5.** Every closed subspace of a Banach space is again a Banach space. For example, the space  $c_0(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$  of all sequences converging to zero is a closed subspace. In fact, if  $a \in \ell^\infty(\mathbb{N}) \setminus c_0(\mathbb{N})$ , then  $\limsup_{j \rightarrow \infty} |a_j| = \varepsilon > 0$  and thus  $\|a - b\|_\infty \geq \varepsilon$  for every  $b \in c_0(\mathbb{N})$ .  $\diamond$

Now what about completeness of  $C(I)$ ? A sequence of functions  $f_n$  converges to  $f$  if and only if

$$\lim_{n \rightarrow \infty} \|f - f_n\|_\infty = \lim_{n \rightarrow \infty} \max_{x \in I} |f(x) - f_n(x)| = 0. \quad (1.27)$$



**Figure 1.1.** Unit balls for  $\|\cdot\|_p$  in  $\mathbb{R}^2$

That is, in the language of real analysis,  $f_n$  converges uniformly to  $f$ . Now let us look at the case where  $f_n$  is only a Cauchy sequence. Then  $f_n(x)$  is clearly a Cauchy sequence of complex numbers for every fixed  $x \in I$ . In particular, by completeness of  $\mathbb{C}$ , there is a limit  $f(x)$  for each  $x$ . Thus we get a limiting function  $f(x)$ . Moreover, letting  $m \rightarrow \infty$  in

$$|f_m(x) - f_n(x)| \leq \varepsilon \quad \forall m, n > N_\varepsilon, x \in I, \quad (1.28)$$

we see

$$|f(x) - f_n(x)| \leq \varepsilon \quad \forall n > N_\varepsilon, x \in I; \quad (1.29)$$

that is,  $f_n(x)$  converges uniformly to  $f(x)$ . However, up to this point we do not know whether  $f$  is in our vector space  $C(I)$ , that is, whether it is continuous. Fortunately, there is a well-known result from real analysis which tells us that the uniform limit of continuous functions is again continuous: Fix  $x \in I$  and  $\varepsilon > 0$ . To show that  $f$  is continuous we need to find a  $\delta$  such that  $|x - y| < \delta$  implies  $|f(x) - f(y)| < \varepsilon$ . Pick  $n$  so that  $\|f_n - f\|_\infty < \varepsilon/3$  and  $\delta$  so that  $|x - y| < \delta$  implies  $|f_n(x) - f_n(y)| < \varepsilon/3$ . Then  $|x - y| < \delta$  implies

$$|f(x) - f(y)| \leq |f(x) - f_n(x)| + |f_n(x) - f_n(y)| + |f_n(y) - f(y)| < \frac{\varepsilon}{3} + \frac{\varepsilon}{3} + \frac{\varepsilon}{3} = \varepsilon$$

as required. Hence  $f \in C(I)$  and thus every Cauchy sequence in  $C(I)$  converges. Or, in other words,

**Theorem 1.1.** *Let  $I \subset \mathbb{R}$  be a compact interval, then the continuous functions  $C(I)$  with the maximum norm is a Banach space.*



For finite dimensional vector spaces the concept of a basis plays a crucial role. In the case of infinite dimensional vector spaces one could define a basis as a maximal set of linearly independent vectors (known as a **Hamel basis**; Problem 1.7). Such a basis has the advantage that it only requires finite linear combinations. However, the price one has to pay is that such a basis will be way too large (typically uncountable, cf. Problems 1.6 and 4.2). Since we have the notion of convergence, we can handle countable linear combinations and try to look for *countable bases*. We start with a few definitions.

The set of all finite linear combinations of a set of vectors  $\{u_n\}_{n \in \mathcal{N}} \subset X$  is called the **span** of  $\{u_n\}_{n \in \mathcal{N}}$  and denoted by

$$\text{span}\{u_n\}_{n \in \mathcal{N}} := \left\{ \sum_{j=1}^m \alpha_j u_{n_j} \mid n_j \in \mathcal{N}, \alpha_j \in \mathbb{C}, m \in \mathbb{N} \right\}. \quad (1.30)$$

A set of vectors  $\{u_n\}_{n \in \mathcal{N}} \subset X$  is called **linearly independent** if every finite subset is. If  $\{u_n\}_{n=1}^N \subset X$ ,  $N \in \mathbb{N} \cup \{\infty\}$ , is countable, we can throw away all elements which can be expressed as linear combinations of the previous ones to obtain a subset of linearly independent vectors which have the same span.

We will call a countable sequence of vectors  $(u_n)_{n=1}^N \subset X$  a **Schauder basis** if every element  $f \in X$  can be uniquely written as a countable linear combination of the basis elements:

$$f = \sum_{n=1}^N \alpha_n u_n, \quad \alpha_n = \alpha_n(f) \in \mathbb{C}, \quad (1.31)$$

where the sum has to be understood as a limit if  $N = \infty$  (the sum is not required to converge unconditionally and hence the order of the basis elements is important). Since we have assumed the coefficients  $\alpha_n(f)$  to be uniquely determined, the vectors are necessarily linearly independent. Moreover, one can show that the coordinate functionals  $f \mapsto \alpha_n(f)$  are continuous (cf. Problem 4.7). A Schauder basis and its corresponding coordinate functionals  $u_n^* : X \rightarrow \mathbb{C}$ ,  $f \mapsto \alpha_n(f)$  form a so-called **biorthogonal system**:  $u_m^*(u_n) = \delta_{m,n}$ , where

$$\delta_{n,m} := \begin{cases} 1, & n = m, \\ 0, & n \neq m, \end{cases} \quad (1.32)$$

is the **Kronecker delta**.

**Example 1.6.** The sequence of vectors  $\delta^n = (\delta_m^n := \delta_{n,m})_{m \in \mathbb{N}}$  is a Schauder basis for the Banach space  $\ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ .

Let  $a = (a_j)_{j=1}^\infty \in \ell^p(\mathbb{N})$  be given and set  $a^m := \sum_{n=1}^m a_n \delta^n$ . Then

$$\|a - a^m\|_p = \left( \sum_{j=m+1}^\infty |a_j|^p \right)^{1/p} \rightarrow 0$$

since  $a_j^m = a_j$  for  $1 \leq j \leq m$  and  $a_j^m = 0$  for  $j > m$ . Hence

$$a = \sum_{n=1}^\infty a_n \delta^n$$

and  $(\delta^n)_{n=1}^\infty$  is a Schauder basis (uniqueness of the coefficients is left as an exercise).

Note that  $(\delta^n)_{n=1}^\infty$  is also Schauder basis for  $c_0(\mathbb{N})$  but not for  $\ell^\infty(\mathbb{N})$  (try to approximate a constant sequence).  $\diamond$

A set whose span is dense is called **total**, and if we have a countable total set, we also have a countable dense set (consider only linear combinations with rational coefficients — show this). A normed vector space containing a countable dense set is called **separable**.

Warning: Some authors use the term total in a slightly different way — see the warning on page 112.

**Example 1.7.** Every Schauder basis is total and thus every Banach space with a Schauder basis is separable (the converse puzzled mathematicians for quite some time and was eventually shown to be false by Per Enflo). In particular, the Banach space  $\ell^p(\mathbb{N})$  is separable for  $1 \leq p < \infty$ .

However,  $\ell^\infty(\mathbb{N})$  is not separable (Problem 1.11)!  $\diamond$

While we will not give a Schauder basis for  $C(I)$  (Problem 1.17), we will at least show that  $C(I)$  is separable. We will do this by showing that every continuous function can be approximated by polynomials, a result which is of independent interest. But first we need a lemma.

**Lemma 1.2** (Smoothing). *Let  $u_n$  be a sequence of nonnegative continuous functions on  $[-1, 1]$  such that*

$$\int_{|x| \leq 1} u_n(x) dx = 1 \quad \text{and} \quad \int_{\delta \leq |x| \leq 1} u_n(x) dx \rightarrow 0, \quad \delta > 0. \quad (1.33)$$

(In other words,  $u_n$  has mass one and concentrates near  $x = 0$  as  $n \rightarrow \infty$ .)

Then for every  $f \in C[-\frac{1}{2}, \frac{1}{2}]$  which vanishes at the endpoints,  $f(-\frac{1}{2}) = f(\frac{1}{2}) = 0$ , we have that

$$f_n(x) := \int_{-1/2}^{1/2} u_n(x-y) f(y) dy \quad (1.34)$$

converges uniformly to  $f(x)$ .

**Proof.** Since  $f$  is uniformly continuous, for given  $\varepsilon$  we can find a  $\delta < 1/2$  (independent of  $x$ ) such that  $|f(x) - f(y)| \leq \varepsilon$  whenever  $|x - y| \leq \delta$ . Moreover, we can choose  $n$  such that  $\int_{\delta \leq |y| \leq 1} u_n(y) dy \leq \varepsilon$ . Now abbreviate  $M := \max_{x \in [-1/2, 1/2]} \{1, |f(x)|\}$  and note

$$|f(x) - \int_{-1/2}^{1/2} u_n(x-y)f(y)dy| = |f(x)| \left| 1 - \int_{-1/2}^{1/2} u_n(x-y)dy \right| \leq M\varepsilon.$$

In fact, either the distance of  $x$  to one of the boundary points  $\pm \frac{1}{2}$  is smaller than  $\delta$  and hence  $|f(x)| \leq \varepsilon$  or otherwise  $[-\delta, \delta] \subset [x - 1/2, x + 1/2]$  and the difference between one and the integral is smaller than  $\varepsilon$ .

Using this, we have

$$\begin{aligned} |f_n(x) - f(x)| &\leq \int_{-1/2}^{1/2} u_n(x-y)|f(y) - f(x)|dy + M\varepsilon \\ &= \int_{|y| \leq 1/2, |x-y| \leq \delta} u_n(x-y)|f(y) - f(x)|dy \\ &\quad + \int_{|y| \leq 1/2, |x-y| \geq \delta} u_n(x-y)|f(y) - f(x)|dy + M\varepsilon \\ &\leq \varepsilon + 2M\varepsilon + M\varepsilon = (1 + 3M)\varepsilon, \end{aligned}$$

which proves the claim.  $\square$

Note that  $f_n$  will be as smooth as  $u_n$ , hence the title smoothing lemma. Moreover,  $f_n$  will be a polynomial if  $u_n$  is. The same idea is used to approximate noncontinuous functions by smooth ones (of course the convergence will no longer be uniform in this case).

Now we are ready to show:

**Theorem 1.3** (Weierstraß). *Let  $I \subset \mathbb{R}$  be a compact interval. Then the set of polynomials is dense in  $C(I)$ .*

**Proof.** Let  $f \in C(I)$  be given. By considering  $f(x) - f(a) - \frac{f(b)-f(a)}{b-a}(x-a)$  it is no loss to assume that  $f$  vanishes at the boundary points. Moreover, without restriction, we only consider  $I = [-\frac{1}{2}, \frac{1}{2}]$  (why?).

Now the claim follows from Lemma 1.2 using the **Landau kernel**

$$u_n(x) := \frac{1}{I_n}(1-x^2)^n,$$

where (using integration by parts)

$$\begin{aligned} I_n &:= \int_{-1}^1 (1-x^2)^n dx = \frac{n}{n+1} \int_{-1}^1 (1-x)^{n-1} (1+x)^{n+1} dx \\ &= \dots = \frac{n!}{(n+1) \cdots (2n+1)} 2^{2n+1} = \frac{(n!)^2 2^{2n+1}}{(2n+1)!} = \frac{n!}{\frac{1}{2}(\frac{1}{2}+1) \cdots (\frac{1}{2}+n)}. \end{aligned}$$

Indeed, the first part of (1.33) holds by construction, and the second part follows from the elementary estimate

$$\frac{1}{\frac{1}{2}+n} < I_n < 2,$$

which shows  $\int_{\delta \leq |x| \leq 1} u_n(x) dx \leq 2u_n(\delta) < (2n+1)(1-\delta^2)^n \rightarrow 0$ .  $\square$

**Corollary 1.4.** *The monomials are total and hence  $C(I)$  is separable.*

Note that while the proof of Theorem 1.3 provides an explicit way of constructing a sequence of polynomials  $f_n(x)$  which will converge uniformly to  $f(x)$ , this method still has a few drawbacks from a practical point of view: Suppose we have approximated  $f$  by a polynomial of degree  $n$  but our approximation turns out to be insufficient for the intended purpose. First of all, since our polynomial will not be optimal in general, we could try to find another polynomial of the same degree giving a better approximation. However, as this is by no means straightforward, it seems more feasible to simply increase the degree. However, if we do this, all coefficients will change and we need to start from scratch. This is in contradistinction to a Schauder basis where we could just add one new element from the basis (and where it suffices to compute one new coefficient).

In particular, note that this shows that the monomials are no Schauder basis for  $C(I)$  since adding monomials incrementally to the expansion gives a uniformly convergent power series whose limit must be analytic. This observation emphasizes that a Schauder basis is more than a set of linearly independent vectors whose span is dense.

We will see in the next section that the concept of orthogonality resolves these problems.

**Problem\* 1.3.** *Show that  $|||f|| - ||g||| \leq ||f - g||$ .*

**Problem\* 1.4.** *Let  $X$  be a Banach space. Show that the norm, vector addition, and multiplication by scalars are continuous. That is, if  $f_n \rightarrow f$ ,  $g_n \rightarrow g$ , and  $\alpha_n \rightarrow \alpha$ , then  $||f_n|| \rightarrow ||f||$ ,  $f_n + g_n \rightarrow f + g$ , and  $\alpha_n g_n \rightarrow \alpha g$ .*

**Problem 1.5.** *Let  $X$  be a Banach space. Show that  $\sum_{j=1}^{\infty} ||f_j|| < \infty$  implies that*

$$\sum_{j=1}^{\infty} f_j = \lim_{n \rightarrow \infty} \sum_{j=1}^n f_j$$

exists. The series is called **absolutely convergent** in this case. Conversely, show that a normed space is complete if every absolutely convergent series converges.

**Problem 1.6.** While  $\ell^1(\mathbb{N})$  is separable, it still has room for an uncountable set of linearly independent vectors. Show this by considering vectors of the form

$$a^\alpha = (1, \alpha, \alpha^2, \dots), \quad \alpha \in (0, 1).$$

(Hint: Recall the Vandermonde determinant. See Problem 4.2 for a generalization.)

**Problem 1.7.** A **Hamel basis** is a maximal set of linearly independent vectors. Show that every vector space  $X$  has a Hamel basis  $\{u_\alpha\}_{\alpha \in A}$ . Show that given a Hamel basis, every  $x \in X$  can be written as a finite linear combination  $x = \sum_{j=1}^n c_j u_{\alpha_j}$ , where the vectors  $u_{\alpha_j}$  and the constants  $c_j$  are uniquely determined. (Hint: Use Zorn's lemma, see Appendix A, to show existence.)

**Problem\* 1.8.** Prove Young's inequality (1.24). Show that equality occurs precisely if  $\alpha = \beta$ . (Hint: Take logarithms on both sides.)

**Problem\* 1.9.** Show that  $\ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ , is complete.

**Problem\* 1.10.** Show that  $\ell^\infty(\mathbb{N})$  is a Banach space.

**Problem\* 1.11.** Show that  $\ell^\infty(\mathbb{N})$  is not separable. (Hint: Consider sequences which take only the value one and zero. How many are there? What is the distance between two such sequences?)

**Problem\* 1.12.** Show that there is equality in the Hölder inequality (1.25) for  $1 < p < \infty$  if and only if either  $a = 0$  or  $|b_j|^p = \alpha |a_j|^q$  for all  $j \in \mathbb{N}$ . Show that we have equality in the triangle inequality for  $\ell^1(\mathbb{N})$  if and only if  $a_j b_j^* \geq 0$  for all  $j \in \mathbb{N}$  (here the  $*$  denotes complex conjugation). Show that we have equality in the triangle inequality for  $\ell^p(\mathbb{N})$  with  $1 < p < \infty$  if and only if  $a = 0$  or  $b = \alpha a$  with  $\alpha \geq 0$ .

**Problem\* 1.13.** Let  $X$  be a normed space. Show that the following conditions are equivalent.

- (i) If  $\|x + y\| = \|x\| + \|y\|$  then  $y = \alpha x$  for some  $\alpha \geq 0$  or  $x = 0$ .
- (ii) If  $\|x\| = \|y\| = 1$  and  $x \neq y$  then  $\|\lambda x + (1 - \lambda)y\| < 1$  for all  $0 < \lambda < 1$ .
- (iii) If  $\|x\| = \|y\| = 1$  and  $x \neq y$  then  $\frac{1}{2}\|x + y\| < 1$ .
- (iv) The function  $x \mapsto \|x\|^2$  is strictly convex.

A norm satisfying one of them is called **strictly convex**.

Show that  $\ell^p(\mathbb{N})$  is strictly convex for  $1 < p < \infty$  but not for  $p = 1, \infty$ .

**Problem 1.14.** Show that  $p_0 \leq p$  implies  $\ell^{p_0}(\mathbb{N}) \subset \ell^p(\mathbb{N})$  and  $\|a\|_p \leq \|a\|_{p_0}$ . Moreover, show

$$\lim_{p \rightarrow \infty} \|a\|_p = \|a\|_\infty.$$

**Problem 1.15.** Formally extend the definition of  $\ell^p(\mathbb{N})$  to  $p \in (0, 1)$ . Show that  $\|\cdot\|_p$  does not satisfy the triangle inequality. However, show that it is a **quasinormed space**, that is, it satisfies all requirements for a normed space except for the triangle inequality which is replaced by

$$\|a + b\| \leq K(\|a\| + \|b\|)$$

with some constant  $K \geq 1$ . Show, in fact,

$$\|a + b\|_p \leq 2^{1/p-1}(\|a\|_p + \|b\|_p), \quad p \in (0, 1).$$

Moreover, show that  $\|\cdot\|_p^p$  satisfies the triangle inequality in this case, but of course it is no longer homogeneous (but at least you can get an honest metric  $d(a, b) = \|a - b\|_p^p$  which gives rise to the same topology). (Hint: Show  $\alpha + \beta \leq (\alpha^p + \beta^p)^{1/p} \leq 2^{1/p-1}(\alpha + \beta)$  for  $0 < p < 1$  and  $\alpha, \beta \geq 0$ .)

**Problem 1.16.** Let  $I$  be a compact interval. Show that the set  $Y := \{f \in C(I) \mid f(x) > 0\}$  is open in  $X := C(I)$ . Compute its closure.

**Problem 1.17.** Show that the following set of functions is a Schauder basis for  $C[0, 1]$ : We start with  $u_1(t) = t$ ,  $u_2(t) = 1 - t$  and then split  $[0, 1]$  into  $2^n$  intervals of equal length and let  $u_{2^n+k+1}(t)$ , for  $1 \leq k \leq 2^n$ , be a piecewise linear peak of height 1 supported in the  $k$ 'th subinterval:  $u_{2^n+k+1}(t) = \max(0, 1 - |2^{n+1}t - 2k + 1|)$  for  $n \in \mathbb{N}_0$  and  $1 \leq k \leq 2^n$ .

### 1.3. The geometry of Hilbert spaces

So far it looks like  $C(I)$  has all the properties we want. However, there is still one thing missing: How should we define orthogonality in  $C(I)$ ? In Euclidean space, two vectors are called **orthogonal** if their scalar product vanishes, so we would need a scalar product:

Suppose  $\mathfrak{H}$  is a vector space. A map  $\langle \cdot, \cdot \rangle : \mathfrak{H} \times \mathfrak{H} \rightarrow \mathbb{C}$  is called a **sesquilinear form** if it is conjugate linear in the first argument and linear in the second; that is,

$$\begin{aligned} \langle \alpha_1 f_1 + \alpha_2 f_2, g \rangle &= \alpha_1^* \langle f_1, g \rangle + \alpha_2^* \langle f_2, g \rangle, \\ \langle f, \alpha_1 g_1 + \alpha_2 g_2 \rangle &= \alpha_1 \langle f, g_1 \rangle + \alpha_2 \langle f, g_2 \rangle, \end{aligned} \quad \alpha_1, \alpha_2 \in \mathbb{C}, \quad (1.35)$$

where  $*$  denotes complex conjugation. A symmetric

$$\langle f, g \rangle = \langle g, f \rangle^* \quad (\text{symmetry})$$

sesquilinear form is also called a **Hermitian form** and a positive definite

$$\langle f, f \rangle > 0 \text{ for } f \neq 0 \quad (\text{positive definite}),$$

Hermitian form is called an **inner product** or **scalar product**. Note that positivity already implies symmetry in the complex case (Problem 1.21). Associated with every scalar product is a norm

$$\|f\| := \sqrt{\langle f, f \rangle}. \quad (1.36)$$

Only the triangle inequality is nontrivial. It will follow from the Cauchy–Schwarz inequality below. Until then, just regard (1.36) as a convenient short hand notation.

Warning: There is no common agreement whether a sesquilinear form (scalar product) should be linear in the first or in the second argument and different authors use different conventions.

The pair  $(\mathfrak{H}, \langle \cdot, \cdot \rangle)$  is called an **inner product space**. If  $\mathfrak{H}$  is complete (with respect to the norm (1.36)), it is called a **Hilbert space**.

**Example 1.8.** Clearly,  $\mathbb{C}^n$  with the usual scalar product

$$\langle a, b \rangle := \sum_{j=1}^n a_j^* b_j \quad (1.37)$$

is a (finite dimensional) Hilbert space.  $\diamond$

**Example 1.9.** A somewhat more interesting example is the Hilbert space  $\ell^2(\mathbb{N})$ , that is, the set of all complex-valued sequences

$$\left\{ (a_j)_{j=1}^\infty \mid \sum_{j=1}^\infty |a_j|^2 < \infty \right\} \quad (1.38)$$

with scalar product

$$\langle a, b \rangle := \sum_{j=1}^\infty a_j^* b_j. \quad (1.39)$$

That this sum is (absolutely) convergent (and thus well-defined) for  $a, b \in \ell^2(\mathbb{N})$  follows from Hölder's inequality (1.25) in the case  $p = q = 2$ .

Observe that the norm  $\|a\| = \sqrt{\langle a, a \rangle}$  is identical to the norm  $\|a\|_2$  defined in the previous section. In particular,  $\ell^2(\mathbb{N})$  is complete and thus indeed a Hilbert space.  $\diamond$

A vector  $f \in \mathfrak{H}$  is called **normalized** or a **unit vector** if  $\|f\| = 1$ . Two vectors  $f, g \in \mathfrak{H}$  are called **orthogonal** or **perpendicular** ( $f \perp g$ ) if  $\langle f, g \rangle = 0$  and **parallel** if one is a multiple of the other.

If  $f$  and  $g$  are orthogonal, we have the **Pythagorean theorem**:

$$\|f + g\|^2 = \|f\|^2 + \|g\|^2, \quad f \perp g, \quad (1.40)$$

which is one line of computation (do it!).

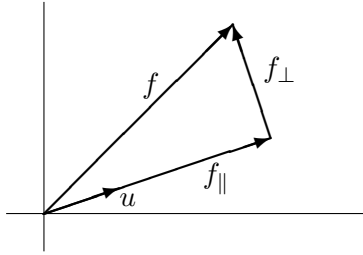
Suppose  $u$  is a unit vector. Then the projection of  $f$  in the direction of  $u$  is given by

$$f_{\parallel} := \langle u, f \rangle u, \quad (1.41)$$

and  $f_{\perp}$ , defined via

$$f_{\perp} := f - \langle u, f \rangle u, \quad (1.42)$$

is perpendicular to  $u$  since  $\langle u, f_{\perp} \rangle = \langle u, f - \langle u, f \rangle u \rangle = \langle u, f \rangle - \langle u, f \rangle \langle u, u \rangle = 0$ .



Taking any other vector parallel to  $u$ , we obtain from (1.40)

$$\|f - \alpha u\|^2 = \|f_{\perp} + (f_{\parallel} - \alpha u)\|^2 = \|f_{\perp}\|^2 + |\langle u, f \rangle - \alpha|^2 \quad (1.43)$$

and hence  $f_{\parallel}$  is the unique vector parallel to  $u$  which is closest to  $f$ .

As a first consequence we obtain the **Cauchy–Bunyakovsky–Schwarz** inequality:

**Theorem 1.5** (Cauchy–Bunyakovsky–Schwarz). *Let  $\mathfrak{H}_0$  be an inner product space. Then for every  $f, g \in \mathfrak{H}_0$  we have*

$$|\langle f, g \rangle| \leq \|f\| \|g\| \quad (1.44)$$

with equality if and only if  $f$  and  $g$  are parallel.

**Proof.** It suffices to prove the case  $\|g\| = 1$ . But then the claim follows from  $\|f\|^2 = |\langle g, f \rangle|^2 + \|f_{\perp}\|^2$ .  $\square$

We will follow common practice and refer to (1.44) simply as Cauchy–Schwarz inequality. Note that the Cauchy–Schwarz inequality implies that the scalar product is continuous in both variables; that is, if  $f_n \rightarrow f$  and  $g_n \rightarrow g$ , we have  $\langle f_n, g_n \rangle \rightarrow \langle f, g \rangle$ .

As another consequence we infer that the map  $\|\cdot\|$  is indeed a norm. In fact,

$$\|f + g\|^2 = \|f\|^2 + \langle f, g \rangle + \langle g, f \rangle + \|g\|^2 \leq (\|f\| + \|g\|)^2. \quad (1.45)$$

But let us return to  $C(I)$ . Can we find a scalar product which has the maximum norm as associated norm? Unfortunately the answer is no! The reason is that the maximum norm does not satisfy the parallelogram law (Problem 1.20).



**Theorem 1.6** (Jordan–von Neumann). *A norm is associated with a scalar product if and only if the **parallelogram law***

$$\|f + g\|^2 + \|f - g\|^2 = 2\|f\|^2 + 2\|g\|^2 \quad (1.46)$$

*holds.*

*In this case the scalar product can be recovered from its norm by virtue of the **polarization identity***

$$\langle f, g \rangle = \frac{1}{4} (\|f + g\|^2 - \|f - g\|^2 + i\|f - ig\|^2 - i\|f + ig\|^2). \quad (1.47)$$

**Proof.** If an inner product space is given, verification of the parallelogram law and the polarization identity is straightforward (Problem 1.21).

To show the converse, we define

$$s(f, g) := \frac{1}{4} (\|f + g\|^2 - \|f - g\|^2 + i\|f - ig\|^2 - i\|f + ig\|^2).$$

Then  $s(f, f) = \|f\|^2$  and  $s(f, g) = s(g, f)^*$  are straightforward to check. Moreover, another straightforward computation using the parallelogram law shows

$$s(f, g) + s(f, h) = 2s(f, \frac{g+h}{2}).$$

Now choosing  $h = 0$  (and using  $s(f, 0) = 0$ ) shows  $s(f, g) = 2s(f, \frac{g}{2})$  and thus  $s(f, g) + s(f, h) = s(f, g+h)$ . Furthermore, by induction we infer  $\frac{m}{2^n} s(f, g) = s(f, \frac{m}{2^n} g)$ ; that is,  $\alpha s(f, g) = s(f, \alpha g)$  for a dense set of positive rational numbers  $\alpha$ . By continuity (which follows from continuity of the norm) this holds for all  $\alpha \geq 0$  and  $s(f, -g) = -s(f, g)$ , respectively,  $s(f, ig) = i s(f, g)$ , finishes the proof.  $\square$

In the case of a real Hilbert space, the polarization identity of course simplifies to  $\langle f, g \rangle = \frac{1}{4} (\|f + g\|^2 - \|f - g\|^2)$ .

Note that the parallelogram law and the polarization identity even hold for sesquilinear forms (Problem 1.21).

But how do we define a scalar product on  $C(I)$ ? One possibility is

$$\langle f, g \rangle := \int_a^b f^*(x)g(x)dx. \quad (1.48)$$

The corresponding inner product space is denoted by  $\mathcal{L}_{cont}^2(I)$ . Note that we have

$$\|f\| \leq \sqrt{|b-a|} \|f\|_\infty \quad (1.49)$$

and hence the maximum norm is stronger than the  $\mathcal{L}_{cont}^2$  norm.

Suppose we have two norms  $\|\cdot\|_1$  and  $\|\cdot\|_2$  on a vector space  $X$ . Then  $\|\cdot\|_2$  is said to be **stronger** than  $\|\cdot\|_1$  if there is a constant  $m > 0$  such that

$$\|f\|_1 \leq m\|f\|_2. \quad (1.50)$$

It is straightforward to check the following.

**Lemma 1.7.** *If  $\|\cdot\|_2$  is stronger than  $\|\cdot\|_1$ , then every  $\|\cdot\|_2$  Cauchy sequence is also a  $\|\cdot\|_1$  Cauchy sequence.*

Hence if a function  $F : X \rightarrow Y$  is continuous in  $(X, \|\cdot\|_1)$ , it is also continuous in  $(X, \|\cdot\|_2)$ , and if a set is dense in  $(X, \|\cdot\|_2)$ , it is also dense in  $(X, \|\cdot\|_1)$ .

In particular,  $\mathcal{L}_{cont}^2$  is separable since the polynomials are dense. But is it also complete? Unfortunately the answer is no:

**Example 1.10.** Take  $I = [0, 2]$  and define

$$f_n(x) := \begin{cases} 0, & 0 \leq x \leq 1 - \frac{1}{n}, \\ 1 + n(x - 1), & 1 - \frac{1}{n} \leq x \leq 1, \\ 1, & 1 \leq x \leq 2. \end{cases}$$

Then  $f_n(x)$  is a Cauchy sequence in  $\mathcal{L}_{cont}^2$ , but there is no limit in  $\mathcal{L}_{cont}^2$ ! Clearly, the limit should be the step function which is 0 for  $0 \leq x < 1$  and 1 for  $1 \leq x \leq 2$ , but this step function is discontinuous (Problem 1.24)!  $\diamond$

**Example 1.11.** The previous example indicates that we should consider (1.48) on a larger class of functions, for example on the class of Riemann integrable functions

$$\mathcal{R}(I) := \{f : I \rightarrow \mathbb{C} \mid f \text{ is Riemann integrable}\}$$

such that the integral makes sense. While this seems natural it implies another problem: Any function which vanishes outside a set which is negligible for the integral (e.g. finitely many points) has norm zero! That is,  $\|f\|_2 := (\int_I |f(x)|^2 dx)^{1/2}$  is only a seminorm on  $\mathcal{R}(I)$  (Problem 1.23). To get a norm we consider  $\mathcal{N}(I) := \{f \in \mathcal{R}(I) \mid \|f\|_2 = 0\}$ . By homogeneity and the triangle inequality  $\mathcal{N}(I)$  is a subspace and we can consider equivalence classes of functions which differ by a negligible function from  $\mathcal{N}(I)$ :

$$\mathcal{L}_{Ri}^2 := \mathcal{R}(I)/\mathcal{N}(I).$$

Since  $\|f\|_2 = \|g\|_2$  for  $f - g \in \mathcal{N}(I)$  we have a norm on  $\mathcal{L}_{Ri}^2$ . Moreover, since this norm inherits the parallelogram law we even have an inner product space. However, this space will not be complete unless we replace the Riemann by the Lebesgue integral. Hence we will not pursue this further until we have the Lebesgue integral at our disposal.  $\diamond$

This shows that in infinite dimensional vector spaces, different norms will give rise to different convergent sequences. In fact, the key to solving problems in infinite dimensional spaces is often finding the right norm! This is something which cannot happen in the finite dimensional case.

**Theorem 1.8.** *If  $X$  is a finite dimensional vector space, then all norms are equivalent. That is, for any two given norms  $\|\cdot\|_1$  and  $\|\cdot\|_2$ , there are positive constants  $m_1$  and  $m_2$  such that*

$$\frac{1}{m_2}\|f\|_1 \leq \|f\|_2 \leq m_1\|f\|_1. \quad (1.51)$$

**Proof.** Choose a basis  $\{u_j\}_{1 \leq j \leq n}$  such that every  $f \in X$  can be written as  $f = \sum_j \alpha_j u_j$ . Since equivalence of norms is an equivalence relation (check this!), we can assume that  $\|\cdot\|_2$  is the usual Euclidean norm:  $\|f\|_2 := \|\sum_j \alpha_j u_j\|_2 = (\sum_j |\alpha_j|^2)^{1/2}$ . Then by the triangle and Cauchy–Schwarz inequalities,

$$\|f\|_1 \leq \sum_j |\alpha_j| \|u_j\|_1 \leq \sqrt{\sum_j \|u_j\|_1^2} \|f\|_2$$

and we can choose  $m_2 = \sqrt{\sum_j \|u_j\|_1^2}$ .

In particular, if  $f_n$  is convergent with respect to  $\|\cdot\|_2$ , it is also convergent with respect to  $\|\cdot\|_1$ . Thus  $\|\cdot\|_1$  is continuous with respect to  $\|\cdot\|_2$  and attains its minimum  $m > 0$  on the unit sphere  $S := \{u \mid \|u\|_2 = 1\}$  (which is compact by the Heine–Borel theorem, Theorem B.22). Now choose  $m_1 = 1/m$ .  $\square$

Finally, I remark that a real Hilbert space can always be embedded into a complex Hilbert space. In fact, if  $\mathfrak{H}$  is a real Hilbert space, then  $\mathfrak{H} \times \mathfrak{H}$  is a complex Hilbert space if we define

$$(f_1, f_2) + (g_1, g_2) = (f_1 + g_1, f_2 + g_2), \quad (\alpha + i\beta)(f_1, f_2) = (\alpha f_1 - \beta f_2, \alpha f_2 + \beta f_1) \quad (1.52)$$

and

$$\langle (f_1, f_2), (g_1, g_2) \rangle = \langle f_1, g_1 \rangle + \langle f_2, g_2 \rangle + i(\langle f_1, g_2 \rangle - \langle f_2, g_1 \rangle). \quad (1.53)$$

Here you should think of  $(f_1, f_2)$  as  $f_1 + if_2$ . Note that we have a conjugate linear map  $C : \mathfrak{H} \times \mathfrak{H} \rightarrow \mathfrak{H} \times \mathfrak{H}$ ,  $(f_1, f_2) \mapsto (f_1, -f_2)$  which satisfies  $C^2 = \mathbb{I}$  and  $\langle Cf, Cg \rangle = \langle g, f \rangle$ . In particular, we can get our original Hilbert space back if we consider  $\text{Re}(f) = \frac{1}{2}(f + Cf) = (f_1, 0)$ .

**Problem 1.18.** *Show that the norm in a Hilbert space satisfies  $\|f + g\| = \|f\| + \|g\|$  if and only if  $f = \alpha g$ ,  $\alpha \geq 0$ , or  $g = 0$ . Hence Hilbert spaces are strictly convex (cf. Problem 1.13).*

**Problem 1.19** (Generalized parallelogram law). *Show that, in a Hilbert space,*

$$\sum_{1 \leq j < k \leq n} \|x_j - x_k\|^2 + \sum_{1 \leq j \leq n} \|x_j\|^2 = n \sum_{1 \leq j \leq n} \|x_j\|^2$$

for every  $n \in \mathbb{N}$ . The case  $n = 2$  is (1.46).

**Problem 1.20.** Show that the maximum norm on  $C[0, 1]$  does not satisfy the parallelogram law.

**Problem\* 1.21.** Suppose  $\mathfrak{Q}$  is a complex vector space. Let  $s(f, g)$  be a sesquilinear form on  $\mathfrak{Q}$  and  $q(f) := s(f, f)$  the associated quadratic form. Prove the **parallelogram law**

$$q(f + g) + q(f - g) = 2q(f) + 2q(g) \quad (1.54)$$

and the **polarization identity**

$$s(f, g) = \frac{1}{4} (q(f + g) - q(f - g) + i q(f - ig) - i q(f + ig)). \quad (1.55)$$

Show that  $s(f, g)$  is symmetric if and only if  $q(f)$  is real-valued.

Note, that if  $\mathfrak{Q}$  is a real vector space, then the parallelogram law is unchanged but the polarization identity in the form  $s(f, g) = \frac{1}{4}(q(f + g) - q(f - g))$  will only hold if  $s(f, g)$  is symmetric.

**Problem 1.22.** A sesquilinear form on a complex inner product space is called **bounded** if

$$\|s\| := \sup_{\|f\|=\|g\|=1} |s(f, g)|$$

is finite. Similarly, the associated quadratic form  $q$  is **bounded** if

$$\|q\| := \sup_{\|f\|=1} |q(f)|$$

is finite. Show

$$\|q\| \leq \|s\| \leq 2\|q\|$$

with  $\|q\| = \|s\|$  if  $s$  is symmetric. (Hint: Use the polarization identity from the previous problem. For the symmetric case look at the real part.)

**Problem\* 1.23.** Suppose  $\mathfrak{Q}$  is a vector space. Let  $s(f, g)$  be a sesquilinear form on  $\mathfrak{Q}$  and  $q(f) := s(f, f)$  the associated quadratic form. Show that the Cauchy–Schwarz inequality

$$|s(f, g)| \leq q(f)^{1/2} q(g)^{1/2}$$

holds if  $q(f) \geq 0$ . In this case  $q(\cdot)^{1/2}$  satisfies the triangle inequality and hence is a seminorm.

(Hint: Consider  $0 \leq q(f + \alpha g) = q(f) + 2\operatorname{Re}(\alpha s(f, g)) + |\alpha|^2 q(g)$  and choose  $\alpha = t s(f, g)^* / |s(f, g)|$  with  $t \in \mathbb{R}$ .)

**Problem\* 1.24.** Prove the claims made about  $f_n$  in Example 1.10.

### 1.4. Completeness

Since  $\mathcal{L}_{cont}^2$  is not complete, how can we obtain a Hilbert space from it? Well, the answer is simple: take the **completion**.

If  $X$  is an (incomplete) normed space, consider the set of all Cauchy sequences  $\mathcal{X}$ . Call two Cauchy sequences equivalent if their difference converges to zero and denote by  $\bar{X}$  the set of all equivalence classes. It is easy to see that  $\bar{X}$  (and  $\mathcal{X}$ ) inherit the vector space structure from  $X$ . Moreover,

**Lemma 1.9.** *If  $x_n$  is a Cauchy sequence in  $X$ , then  $\|x_n\|$  is also a Cauchy sequence and thus converges.*

Consequently, the norm of an equivalence class  $[(x_n)_{n=1}^\infty]$  can be defined by  $\|[(x_n)_{n=1}^\infty]\| := \lim_{n \rightarrow \infty} \|x_n\|$  and is independent of the representative (show this!). Thus  $\bar{X}$  is a normed space.

**Theorem 1.10.**  *$\bar{X}$  is a Banach space containing  $X$  as a dense subspace if we identify  $x \in X$  with the equivalence class of all sequences converging to  $x$ .*

**Proof.** (Outline) It remains to show that  $\bar{X}$  is complete. Let  $\xi_n = [(x_{n,j})_{j=1}^\infty]$  be a Cauchy sequence in  $\bar{X}$ . Without loss of generality (by dropping terms) we can choose the representatives  $x_{n,j}$  such that  $|x_{n,j} - x_{n,k}| \leq \frac{1}{n}$  for  $j, k \geq n$ . Then it is not hard to see that  $\xi = [(x_{j,j})_{j=1}^\infty]$  is its limit.  $\square$

Notice that the completion  $\bar{X}$  is unique. More precisely, every other complete space which contains  $X$  as a dense subset is isomorphic to  $\bar{X}$ . This can for example be seen by showing that the identity map on  $X$  has a unique extension to  $\bar{X}$  (compare Theorem 1.16 below).

In particular, it is no restriction to assume that a normed vector space or an inner product space is complete (note that by continuity of the norm the parallelogram law holds for  $\bar{X}$  if it holds for  $X$ ).

**Example 1.12.** The completion of the space  $\mathcal{L}_{cont}^2(I)$  is denoted by  $L^2(I)$ . While this defines  $L^2(I)$  uniquely (up to isomorphisms) it is often inconvenient to work with equivalence classes of Cauchy sequences. A much more convenient characterization can be given with the help of the Lebesgue integral (see Chapter 3 from [48] if you are familiar with basic Lebesgue integration; Theorem 3.18 from [48] will establish equivalence of both approaches).

Similarly, we define  $L^p(I)$ ,  $1 \leq p < \infty$ , as the completion of  $C(I)$  with respect to the norm

$$\|f\|_p := \left( \int_a^b |f(x)|^p dx \right)^{1/p}.$$

The only requirement for a norm which is not immediate is the triangle inequality (except for  $p = 1, 2$ ) but this can be shown as for  $\ell^p$  (cf. Problem 1.27).  $\diamond$

**Problem 1.25.** Provide a detailed proof of Theorem 1.10.

**Problem 1.26.** For every  $f \in L^1(I)$  we can define its integral

$$\int_c^d f(x)dx$$

as the (unique) extension of the corresponding linear functional from  $C(I)$  to  $L^1(I)$  (by Theorem 1.16 below). Show that this integral is linear and satisfies

$$\int_c^e f(x)dx = \int_c^d f(x)dx + \int_d^e f(x)dx, \quad \left| \int_c^d f(x)dx \right| \leq \int_c^d |f(x)|dx.$$

**Problem\* 1.27.** Show the **Hölder inequality**

$$\|fg\|_1 \leq \|f\|_p \|g\|_q, \quad \frac{1}{p} + \frac{1}{q} = 1, \quad 1 \leq p, q \leq \infty,$$

for  $f \in L^p(I)$ ,  $g \in L^q(I)$  and conclude that  $\|\cdot\|_p$  is a norm on  $C(I)$ . Also conclude that  $L^p(I) \subseteq L^1(I)$ .

## 1.5. Compactness

In finite dimensions relatively compact sets are easily identified as they are precisely the bounded sets by the Heine–Borel theorem (Theorem B.22). In the infinite dimensional case the situation is more complicated. Before we look into this, please recall that for a subset  $U$  of a Banach space the following are equivalent (see Corollary B.20 and Lemma B.26):

- $U$  is relatively compact (i.e. its closure is compact)
- every sequence from  $U$  has a convergent subsequence
- $U$  is totally bounded (i.e. it has a finite  $\varepsilon$ -cover for every  $\varepsilon > 0$ )

**Example 1.13.** Consider the bounded sequence  $(\delta^n)_{n=1}^\infty$  in  $\ell^p(\mathbb{N})$ . Since  $\|\delta^n - \delta^m\|_p = 2^{1/p}$  for  $n \neq m$ , there is no way to extract a convergent subsequence.  $\diamond$

In particular, the Heine–Borel theorem fails for  $\ell^p(\mathbb{N})$ . In fact, it turns out that it fails in any infinite dimensional space as we will see in Theorem 4.27 below. Hence one needs criteria when a given subset is relatively compact. Our strategy will be based on total boundedness and can be outlined as follows: Project the original set to some finite dimensional space such that the information loss can be made arbitrarily small (by increasing the dimension of the finite dimensional space) and apply Heine–Borel to the finite dimensional space. This idea is formalized in the following lemma.

**Lemma 1.11.** *Let  $X$  be a metric space and  $K$  some subset. Assume that for every  $\varepsilon > 0$  there is a metric space  $Y_n$ , a surjective map  $P_n : X \rightarrow Y_n$ , and some  $\delta > 0$  such that  $P_n(K)$  is totally bounded and  $d(x, y) < \varepsilon$  whenever  $x, y \in K$  with  $d(P_n(x), P_n(y)) < \delta$ . Then  $K$  is totally bounded.*

*In particular, if  $X$  is a Banach space the claim holds if  $P_n$  can be chosen a linear map onto a finite dimensional subspace  $Y_n$  such that  $\|P_n\| \leq C$ ,  $P_n K$  is bounded, and  $\|(1 - P_n)x\| \leq \varepsilon$  for  $x \in K$ .*

**Proof.** Fix  $\varepsilon > 0$ . Then by total boundedness of  $P_n(K)$  we can find a  $\delta$ -cover  $\{B_\delta(y_j)\}_{j=1}^m$  for  $P_n(K)$ . Now if we choose  $x_j \in P_n^{-1}(\{y_j\}) \cap K$ , then  $\{B_\varepsilon(x_j)\}_{j=1}^n$  is an  $\varepsilon$ -cover for  $K$  since  $P_n^{-1}(B_\delta(y_j)) \cap K \subseteq B_\varepsilon(x_j)$ .

For the last claim take  $P_n$  corresponding to  $\varepsilon/3$  and note that  $\|x - y\| \leq \|(1 - P_n)x\| + \|P_n(x - y)\| + \|(1 - P_n)y\| < \varepsilon$  for  $\delta := \varepsilon/3$ .  $\square$

The first application will be to  $\ell^p(\mathbb{N})$ .

**Theorem 1.12** (Fréchet). *Consider  $\ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ , and let  $P_n a = (a_1, \dots, a_n, 0, \dots)$  be the projection onto the first  $n$  components. A subset  $\mathcal{K} \subseteq \ell^p(\mathbb{N})$  is relatively compact if and only if*

- (i) *it is pointwise bounded,  $\sup_{a \in \mathcal{K}} |a_j| \leq M_j$  for all  $j \in \mathbb{N}$ , and*
- (ii) *for every  $\varepsilon > 0$  there is some  $n$  such that  $\|(1 - P_n)a\|_p \leq \varepsilon$  for all  $a \in \mathcal{K}$ .*

*In the case  $p = \infty$  conditions (i) and (ii) still imply that  $\mathcal{K}$  is relatively compact, but the converse only holds for  $\mathcal{K} \subseteq c_0(\mathbb{N})$ .*

**Proof.** Clearly (i) and (ii) is what is needed for Lemma 1.11.

Conversely, if  $\mathcal{K}$  is relatively compact it is bounded. Moreover, given  $\delta$  we can choose a finite  $\delta$ -cover  $\{B_\delta(a^j)\}_{j=1}^m$  for  $\mathcal{K}$  and some  $n$  such that  $\|(1 - P_n)a^j\|_p \leq \delta$  for all  $1 \leq j \leq m$ . Now given  $a \in \mathcal{K}$  we have  $a \in B_\delta(a^j)$  for some  $j$  and hence  $\|(1 - P_n)a\|_p \leq \|(1 - P_n)(a - a^j)\|_p + \|(1 - P_n)a^j\|_p \leq 2\delta$  as required.  $\square$

**Example 1.14.** Fix  $a \in \ell^p(\mathbb{N})$  if  $1 \leq p < \infty$  or  $a \in c_0(\mathbb{N})$  else. Then  $\mathcal{K} := \{b \mid |b_j| \leq |a_j|\} \subset \ell^p(\mathbb{N})$  is compact.  $\diamond$

The second application will be to  $C(I)$ . A family of functions  $F \subset C(I)$  is called (pointwise) **equicontinuous** if for every  $\varepsilon > 0$  and every  $x \in I$  there is a  $\delta > 0$  such that

$$|f(y) - f(x)| \leq \varepsilon \quad \text{whenever} \quad |y - x| < \delta, \quad \forall f \in F. \quad (1.56)$$

That is, in this case  $\delta$  is required to be independent of the function  $f \in F$ .

**Theorem 1.13** (Arzelà–Ascoli). *Let  $F \subset C(I)$  be a family of continuous functions. Then every sequence from  $F$  has a uniformly convergent subsequence if and only if  $F$  is equicontinuous and the set  $\{f(x_0) | f \in F\}$  is bounded for one  $x_0 \in I$ . In this case  $F$  is even bounded.*

**Proof.** Suppose  $F$  is equicontinuous and pointwise bounded. Fix  $\varepsilon > 0$ . By compactness of  $I$  there are finitely many points  $x_1, \dots, x_n \in I$  such that the balls  $B_{\delta_j}(x_j)$  cover  $I$ , where  $\delta_j$  is the  $\delta$  corresponding to  $x_j$  as in the definition of equicontinuity. Now first of all note that, since  $I$  is connected and since  $x_0 \in B_{\delta_j}(x_j)$  for some  $j$ , we see that  $F$  is bounded:  $|f(x)| \leq \sup_{f \in F} |f(x_0)| + n\varepsilon$ .

Next consider  $P : C[0, 1] \rightarrow \mathbb{C}^n$ ,  $P(f) = (f(x_1), \dots, f(x_n))$ . Then  $P(F)$  is bounded and  $\|f - g\|_\infty \leq 3\varepsilon$  whenever  $\|P(f) - P(g)\|_\infty < \varepsilon$ . Indeed, just note that for every  $x$  there is some  $j$  such that  $x \in B_{\delta_j}(x_j)$  and thus  $|f(x) - g(x)| \leq |f(x) - f(x_j)| + |f(x_j) - g(x_j)| + |g(x_j) - g(x)| \leq 3\varepsilon$ . Hence  $F$  is relatively compact by Lemma 1.11.

Conversely, suppose  $F$  is relatively compact. Then  $F$  is totally bounded and hence bounded. To see equicontinuity fix  $x \in I$ ,  $\varepsilon > 0$  and choose a corresponding  $\varepsilon$ -cover  $\{B_\varepsilon(f_j)\}_{j=1}^n$  for  $F$ . Pick  $\delta > 0$  such that  $y \in B_\delta(x)$  implies  $|f_j(y) - f_j(x)| < \varepsilon$  for all  $1 \leq j \leq n$ . Then  $f \in B_\varepsilon(f_j)$  for some  $j$  and hence  $|f(y) - f(x)| \leq |f(y) - f_j(y)| + |f_j(y) - f_j(x)| + |f_j(x) - f(x)| \leq 3\varepsilon$ , proving equicontinuity.  $\square$

**Example 1.15.** Consider the solution  $y_n(x)$  of the initial value problem

$$y' = \sin(ny), \quad y(0) = 1.$$

(Assuming this solution exists — it can in principle be found using separation of variables.) Then  $|y'_n(x)| \leq 1$  and hence the mean value theorem shows that the family  $\{y_n\} \subseteq C([0, 1])$  is equicontinuous. Hence there is a uniformly convergent subsequence.  $\diamond$

**Problem 1.28.** *Show that a subset  $F \subset c_0(\mathbb{N})$  is relatively compact if and only if there is a nonnegative sequence  $a \in c_0(\mathbb{N})$  such that  $|b_n| \leq a_n$  for all  $n \in \mathbb{N}$  and all  $b \in F$ .*

**Problem 1.29.** *Find a sequence in  $C[0, 1]$  which is bounded but has no convergent subsequence.*

**Problem 1.30.** *Find a family in  $C[0, 1]$  that is equicontinuous but not bounded.*

**Problem 1.31.** *Which of the following families are relatively compact in  $C[0, 1]$ ?*

$$(i) \ F = \{f \in C^1[0, 1] \mid \|f\|_\infty \leq 1\}$$



- (ii)  $F = \{f \in C^1[0, 1] \mid \|f'\|_\infty \leq 1\}$
- (iii)  $F = \{f \in C^1[0, 1] \mid \|f\|_\infty \leq 1, \|f'\|_2 \leq 1\}$

## 1.6. Bounded operators

Given two normed spaces  $X$  and  $Y$ , a linear map

$$A : \mathfrak{D}(A) \subseteq X \rightarrow Y \quad (1.57)$$

will be called a **(linear) operator**. The linear subspace  $\mathfrak{D}(A)$  on which  $A$  is defined is called the **domain** of  $A$  and is frequently required to be dense. The **kernel** (also **null space**)

$$\text{Ker}(A) := \{f \in \mathfrak{D}(A) \mid Af = 0\} \subseteq X \quad (1.58)$$

and **range**

$$\text{Ran}(A) := \{Af \mid f \in \mathfrak{D}(A)\} = A\mathfrak{D}(A) \subseteq Y \quad (1.59)$$

are again linear subspaces. Note that a linear map  $A$  will be continuous if and only if it is continuous at 0, that is,  $x_n \in \mathfrak{D}(A) \rightarrow 0$  implies  $Ax_n \rightarrow 0$ .

The operator  $A$  is called **bounded** if the operator norm

$$\|A\| := \sup_{f \in \mathfrak{D}(A), \|f\|_X=1} \|Af\|_Y \quad (1.60)$$

is finite. This says that  $A$  is bounded if the image of the closed unit ball  $\bar{B}_1(0) \subset X$  is contained in some closed ball  $\bar{B}_r(0) \subset Y$  of finite radius  $r$  (with the smallest radius being the operator norm). Hence  $A$  is bounded if and only if it maps bounded sets to bounded sets.

Note that if you replace the norm on  $X$  or  $Y$  then the operator norm will of course also change in general. However, if the norms are equivalent so will be the operator norms.

By construction, a bounded operator satisfies

$$\|Af\|_Y \leq \|A\| \|f\|_X, \quad f \in \mathfrak{D}(A), \quad (1.61)$$

and hence is Lipschitz continuous, that is,  $\|Af - Ag\|_Y \leq \|A\| \|f - g\|_X$  for  $f, g \in \mathfrak{D}(A)$ . In particular, it is continuous. The converse is also true:

**Theorem 1.14.** *A linear operator  $A$  is bounded if and only if it is continuous.*

**Proof.** Suppose  $A$  is continuous but not bounded. Then there is a sequence of unit vectors  $u_n \in \mathfrak{D}(A)$  such that  $\|Au_n\|_Y \geq n$ . Then  $f_n := \frac{1}{n}u_n$  converges to 0 but  $\|Af_n\|_Y \geq 1$  does not converge to 0.  $\square$

Of course it suffices to check continuity at one point in  $X$ , say at 0, since continuity at all other points will then follow by a simple translation.

If  $X$  is finite dimensional, then every operator is bounded.

**Lemma 1.15.** *Let  $X, Y$  be normed spaces with  $X$  finite dimensional. Then every linear operator  $A : X \rightarrow Y$  is bounded.*

**Proof.** Choose a basis  $\{x_j\}_{j=1}^n$  for  $X$  such that every  $x \in X$  can be written as  $x = \sum_{j=1}^n \alpha_j x_j$ . By Theorem 1.8 we can assume  $\|x\|_X = (\sum_{j=1}^n |\alpha_j|^2)^{1/2}$  without loss of generality. Then

$$\|Ax\|_Y \leq \sum_{j=1}^n |\alpha_j| \|Ax_j\|_Y \leq \sqrt{\sum_{j=1}^n \|Ax_j\|_Y^2} \|x\|_X$$

and thus  $\|A\| \leq (\sum_{j=1}^n \|Ax_j\|_Y^2)^{1/2}$ .  $\square$

In the infinite dimensional case an operator can be unbounded. Moreover, one and the same operation might be bounded (i.e. continuous) or unbounded, depending on the norm chosen.

**Example 1.16.** Let  $X := \ell^p(\mathbb{N})$  and  $a \in \ell^\infty(\mathbb{N})$ . Consider the multiplication operator  $A : X \rightarrow X$  defined by

$$(Ab)_j := a_j b_j.$$

Then  $|(Ab)_j| \leq \|a\|_\infty |b_j|$  shows  $\|A\| \leq \|a\|_\infty$ . In fact, we even have  $\|A\| = \|a\|_\infty$  (show this). If  $a$  is unbounded we need a domain  $\mathfrak{D}(A) := \{b \in \ell^p(\mathbb{N}) \mid (a_j b_j)_{j \in \mathbb{N}} \in \ell^p(\mathbb{N})\}$  and  $A$  will be unbounded (show this).  $\diamond$

**Example 1.17.** Let  $X := C[0, 1]$ ,  $Y := \mathbb{C}$  and consider the operator  $A : X \rightarrow Y$  given by  $Ax := \int_0^1 x(t) dt$ . Then  $\|A\| = 1$  since  $|Ax| \leq \int_0^1 |x(t)| dt \leq \|x\|_\infty$  with equality for constant functions. If we replace  $X$  by  $X_0 := \{x \in C[0, 1] \mid x(0) = 0\}$  (check that this is a closed subspace), then  $A$  restricted to  $X_0$  still has norm one (show this), but the sup in (1.60) is no longer attained.  $\diamond$

**Example 1.18.** Consider the vector space of differentiable functions  $X := C^1[0, 1]$  and equip it with the norm (cf. Problem 1.35)

$$\|f\|_{\infty, 1} := \max_{x \in [0, 1]} |f(x)| + \max_{x \in [0, 1]} |f'(x)|.$$

Let  $Y := C[0, 1]$  and observe that the differential operator  $A = \frac{d}{dx} : X \rightarrow Y$  is bounded since

$$\|Af\|_\infty = \max_{x \in [0, 1]} |f'(x)| \leq \max_{x \in [0, 1]} |f(x)| + \max_{x \in [0, 1]} |f'(x)| = \|f\|_{\infty, 1}.$$

However, if we consider  $A = \frac{d}{dx} : \mathfrak{D}(A) \subseteq Y \rightarrow Y$  defined on  $\mathfrak{D}(A) = C^1[0, 1]$ , then we have an unbounded operator. Indeed, choose  $u_n(x) := \sin(n\pi x)$  which is normalized,  $\|u_n\|_\infty = 1$ , and observe that

$$Au_n(x) = u'_n(x) = n\pi \cos(n\pi x)$$

is unbounded,  $\|Au_n\|_\infty = n\pi$ . Note that  $\mathfrak{D}(A)$  contains the set of polynomials and thus is dense by the Weierstraß approximation theorem (Theorem 1.3).  $\diamond$

If  $A$  is bounded and densely defined, it is no restriction to assume that it is defined on all of  $X$ .

**Theorem 1.16** (extension principle). *Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  be a bounded linear operator between a normed space  $X$  and a Banach space  $Y$ . If  $\mathfrak{D}(A)$  is dense, there is a unique (continuous) extension of  $A$  to  $X$  which has the same operator norm.*

**Proof.** Since a bounded operator maps Cauchy sequences to Cauchy sequences, this extension can only be given by

$$\overline{A}f := \lim_{n \rightarrow \infty} Af_n, \quad f_n \in \mathfrak{D}(A), \quad f \in X.$$

To show that this definition is independent of the sequence  $f_n \rightarrow f$ , let  $g_n \rightarrow f$  be a second sequence and observe

$$\|Af_n - Ag_n\| = \|A(f_n - g_n)\| \leq \|A\|\|f_n - g_n\| \rightarrow 0.$$

Since for  $f \in \mathfrak{D}(A)$  we can choose  $f_n = f$ , we see that  $\overline{A}f = Af$  in this case, that is,  $\overline{A}$  is indeed an extension. From continuity of vector addition and scalar multiplication it follows that  $\overline{A}$  is linear. Finally, from continuity of the norm we conclude that the operator norm does not increase.  $\square$

The set of all bounded linear operators from  $X$  to  $Y$  is denoted by  $\mathcal{L}(X, Y)$ . If  $X = Y$ , we write  $\mathcal{L}(X) := \mathcal{L}(X, X)$ . An operator in  $\mathcal{L}(X, \mathbb{C})$  is called a **bounded linear functional**, and the space  $X^* := \mathcal{L}(X, \mathbb{C})$  is called the **dual space** of  $X$ . The dual space takes the role of coordinate functions in a Banach space.

**Example 1.19.** Let  $X$  be a finite dimensional space and  $\{x_j\}_{j=1}^n$  a basis. Then every  $x \in X$  can be uniquely written as  $x = \sum_{j=1}^n \alpha_j x_j$  and we can define linear functionals via  $\ell_j(x) := \alpha_j$  for  $1 \leq j \leq n$ . The functionals  $\{\ell_j\}_{j=1}^n$  are called a **dual basis** since  $\ell_k(x_j) = \delta_{j,k}$  and since any other linear functional  $\ell \in X^*$  can be written as  $\ell = \sum_{j=1}^n \ell(x_j)\ell_j$ . In particular,  $X$  and  $X^*$  have the same dimension.  $\diamond$

**Example 1.20.** Let  $X := \ell^p(\mathbb{N})$ . Then the coordinate functions

$$\ell_j(a) := a_j$$

are bounded linear functionals:  $|\ell_j(a)| = |a_j| \leq \|a\|_p$  and hence  $\|\ell_j\| = 1$ . More general, let  $b \in \ell^q(\mathbb{N})$  where  $\frac{1}{p} + \frac{1}{q} = 1$ . Then

$$\ell_b(a) := \sum_{j=1}^n b_j a_j$$

is a bounded linear functional satisfying  $\|\ell_b\| \leq \|b\|_q$  by Hölder's inequality. In fact, we even have  $\|\ell_b\| = \|b\|_q$  (Problem 4.14). Note that the first example is a special case of the second one upon choosing  $b = \delta^j$ .  $\diamond$

**Example 1.21.** Consider  $X := C(I)$ . Then for every  $x_0 \in I$  the point evaluation  $\ell_{x_0}(f) := f(x_0)$  is a bounded linear functional. In fact,  $\|\ell_{x_0}\| = 1$  (show this).

However, note that  $\ell_{x_0}$  is unbounded on  $\mathcal{L}_{cont}^2(I)$ ! To see this take  $f_n(x) := \sqrt{\frac{3n}{2}} \max(0, 1 - n|x - x_0|)$  which is a triangle shaped peak supported on  $[x_0 - n^{-1}, x_0 + n^{-1}]$  and normalized according to  $\|f_n\|_2 = 1$  for  $n$  sufficiently large such that the support is contained in  $I$ . Then  $\ell_{x_0}(f) = f_n(x_0) = \sqrt{\frac{3n}{2}} \rightarrow \infty$ . This implies that  $\ell_{x_0}$  cannot be extended to the completion of  $\mathcal{L}_{cont}^2(I)$  in a natural way and reflects the fact that the integral cannot see individual points (changing the value of a function at one point does not change its integral).  $\diamond$

**Example 1.22.** Consider  $X := C(I)$  and let  $g$  be some continuous function. Then

$$\ell_g(f) := \int_a^b g(x)f(x)dx$$

is a linear functional with norm  $\|\ell_g\| = \|g\|_1$ . Indeed, first of all note that

$$|\ell_g(f)| \leq \int_a^b |g(x)f(x)|dx \leq \|f\|_\infty \int_a^b |g(x)|dx$$

shows  $\|\ell_g\| \leq \|g\|_1$ . To see that we have equality consider  $f_\varepsilon = g^*/(|g| + \varepsilon)$  and note

$$|\ell_g(f_\varepsilon)| = \int_a^b \frac{|g(x)|^2}{|g(x)| + \varepsilon} dx \geq \int_a^b \frac{|g(x)|^2 - \varepsilon^2}{|g(x)| + \varepsilon} dx = \|g\|_1 - (b-a)\varepsilon.$$

Since  $\|f_\varepsilon\| \leq 1$  and  $\varepsilon > 0$  is arbitrary this establishes the claim.  $\diamond$

**Theorem 1.17.** *The space  $\mathcal{L}(X, Y)$  together with the operator norm (1.60) is a normed space. It is a Banach space if  $Y$  is.*

**Proof.** That (1.60) is indeed a norm is straightforward. If  $Y$  is complete and  $A_n$  is a Cauchy sequence of operators, then  $A_n f$  converges to an element  $g$  for every  $f$ . Define a new operator  $A$  via  $Af = g$ . By continuity of the vector operations,  $A$  is linear and by continuity of the norm  $\|Af\| = \lim_{n \rightarrow \infty} \|A_n f\| \leq (\lim_{n \rightarrow \infty} \|A_n\|)\|f\|$ , it is bounded. Furthermore, given  $\varepsilon > 0$ , there is some  $N$  such that  $\|A_n - A_m\| \leq \varepsilon$  for  $n, m \geq N$  and thus  $\|A_n f - A_m f\| \leq \varepsilon\|f\|$ . Taking the limit  $m \rightarrow \infty$ , we see  $\|A_n f - Af\| \leq \varepsilon\|f\|$ ; that is,  $A_n \rightarrow A$ .  $\square$

The Banach space of bounded linear operators  $\mathcal{L}(X)$  even has a multiplication given by composition. Clearly, this multiplication is distributive

$$(A+B)C = AC + BC, \quad A(B+C) = AB + BC, \quad A, B, C \in \mathcal{L}(X) \quad (1.62)$$

and associative

$$(AB)C = A(BC), \quad \alpha(AB) = (\alpha A)B = A(\alpha B), \quad \alpha \in \mathbb{C}. \quad (1.63)$$

Moreover, it is easy to see that we have

$$\|AB\| \leq \|A\|\|B\|. \quad (1.64)$$

In other words,  $\mathcal{L}(X)$  is a so-called **Banach algebra**. However, note that our multiplication is not commutative (unless  $X$  is one-dimensional). We even have an **identity**, the identity operator  $\mathbb{I}$ , satisfying  $\|\mathbb{I}\| = 1$ .

**Problem 1.32.** Show that two norms on  $X$  are equivalent if and only if they give rise to the same convergent sequences.

**Problem 1.33.** Consider  $X = \mathbb{C}^n$  and let  $A \in \mathcal{L}(X)$  be a matrix. Equip  $X$  with the norm (show that this is a norm)

$$\|x\|_\infty := \max_{1 \leq j \leq n} |x_j|$$

and compute the operator norm  $\|A\|$  with respect to this norm in terms of the matrix entries. Do the same with respect to the norm

$$\|x\|_1 := \sum_{1 \leq j \leq n} |x_j|.$$

**Problem 1.34.** Show that the integral operator

$$(Kf)(x) := \int_0^1 K(x, y)f(y)dy,$$

where  $K(x, y) \in C([0, 1] \times [0, 1])$ , defined on  $\mathfrak{D}(K) := C[0, 1]$ , is a bounded operator both in  $X := C[0, 1]$  (max norm) and  $X := \mathcal{L}_{\text{cont}}^2(0, 1)$ . Show that the norm in the  $X = C[0, 1]$  case is given by

$$\|K\| = \max_{x \in [0, 1]} \int_0^1 |K(x, y)|dy.$$

**Problem\* 1.35.** Let  $I$  be a compact interval. Show that the set of differentiable functions  $C^1(I)$  becomes a Banach space if we set  $\|f\|_{\infty, 1} := \max_{x \in I} |f(x)| + \max_{x \in I} |f'(x)|$ .

**Problem\* 1.36.** Show that  $\|AB\| \leq \|A\|\|B\|$  for every  $A, B \in \mathcal{L}(X)$ . Conclude that the multiplication is continuous:  $A_n \rightarrow A$  and  $B_n \rightarrow B$  imply  $A_n B_n \rightarrow AB$ .

**Problem 1.37.** Let  $A \in \mathcal{L}(X)$  be a bijection. Show

$$\|A^{-1}\|^{-1} = \inf_{f \in X, \|f\|=1} \|Af\|.$$

**Problem\* 1.38.** Suppose  $B \in \mathcal{L}(X)$  with  $\|B\| < 1$ . Then  $\mathbb{I} + B$  is invertible with

$$(\mathbb{I} + B)^{-1} = \sum_{n=0}^{\infty} (-1)^n B^n.$$

Consequently for  $A, B \in \mathcal{L}(X, Y)$ ,  $A + B$  is invertible if  $A$  is invertible and  $\|B\| < \|A^{-1}\|^{-1}$ .

**Problem\* 1.39.** Let

$$f(z) := \sum_{j=0}^{\infty} f_j z^j, \quad |z| < R,$$

be a convergent power series with radius of convergence  $R > 0$ . Suppose  $X$  is a Banach space and  $A \in \mathcal{L}(X)$  is a bounded operator with  $\limsup_n \|A^n\|^{1/n} < R$  (note that by  $\|A^n\| \leq \|A\|^n$  the limsup is finite). Show that

$$f(A) := \sum_{j=0}^{\infty} f_j A^j$$

exists and defines a bounded linear operator. Moreover, if  $f$  and  $g$  are two such functions and  $\alpha \in \mathbb{C}$ , then

$$(f + g)(A) = f(A) + g(A), \quad (\alpha f)(A) = \alpha f(A), \quad (fg)(A) = f(A)g(A).$$

(Hint: Problem 1.5.)

**Problem\* 1.40.** Show that a linear map  $\ell : X \rightarrow \mathbb{C}$  is continuous if and only if its kernel is closed. (Hint: If  $\ell$  is not continuous, we can find a sequence of normalized vectors  $x_n$  with  $|\ell(x_n)| \rightarrow \infty$  and a vector  $y$  with  $\ell(y) = 1$ .)

## 1.7. Sums and quotients of Banach spaces

Given two Banach spaces  $X_1$  and  $X_2$  we can define their **(direct) sum**  $X := X_1 \oplus X_2$  as the Cartesian product  $X_1 \times X_2$  together with the norm  $\|(x_1, x_2)\| := \|x_1\| + \|x_2\|$ . Clearly  $X$  is again a Banach space and a sequence in  $X$  converges if and only if the components converge in  $X_1$  and  $X_2$ , respectively. In fact, since all norms on  $\mathbb{R}^2$  are equivalent (Theorem 1.8), we could as well take  $\|(x_1, x_2)\|_p := (\|x_1\|^p + \|x_2\|^p)^{1/p}$  or  $\|(x_1, x_2)\|_{\infty} := \max(\|x_1\|, \|x_2\|)$ . We will write  $X_1 \oplus_p X_2$  if we want to emphasize the norm used. In particular, in the case of Hilbert spaces the choice  $p = 2$  will ensure that  $X$  is again a Hilbert space.

Note that  $X_1$  and  $X_2$  can be regarded as closed subspaces of  $X_1 \times X_2$  by virtue of the obvious embeddings  $x_1 \mapsto (x_1, 0)$  and  $x_2 \mapsto (0, x_2)$ . It is straightforward to generalize this concept to finitely many spaces (Problem 1.41).

If  $A_j : \mathfrak{D}(A_j) \subseteq X_j \rightarrow Y_j$ ,  $j = 1, 2$ , are linear operators, then  $A_1 \oplus A_2 : \mathfrak{D}(A_1) \times \mathfrak{D}(A_2) \subseteq X_1 \oplus X_2 \rightarrow Y_1 \oplus Y_2$  is defined as  $A_1 \oplus A_2(x_1, x_2) = (A_1x_1, A_2x_2)$ . Clearly  $A_1 \oplus A_2$  will be bounded if and only if both  $A_1$  and  $A_2$  are bounded and  $\|A_1 \oplus A_2\| = \max(\|A_1\|, \|A_2\|)$ .

Note that if  $A_j : X_j \rightarrow Y$ ,  $j = 1, 2$ , there is another natural way of defining an associated operator  $X_1 \oplus X_2 \rightarrow Y$  given by  $A_1 \hat{\oplus} A_2(x_1, x_2) := A_1x_1 + A_2x_2$ . In particular, in the case  $Y = \mathbb{C}$  we get that  $(X_1 \oplus_p X_2)^* \cong X_1^* \oplus_q X_2^*$  for  $\frac{1}{p} + \frac{1}{q} = 1$  via the identification  $(\ell_1, \ell_2) \in X_1^* \oplus_q X_2^* \mapsto \ell_1 \hat{\oplus} \ell_2 \in (X_1 \oplus_p X_2)^*$ . Clearly this identification is bijective and preserves the norm (to see this relate it to Hölder's inequality in  $\mathbb{C}^2$  and note that equality is attained).

Given two subspaces  $M, N \subseteq X$  of a vector space, we can define their sum as usual:  $M + N := \{x + y | x \in M, y \in N\}$ . In particular, the decomposition  $x + y$  with  $x \in M$ ,  $y \in N$  is unique iff  $M \cap N = \{0\}$  and we will write  $M \dot{+} N$  in this case. It is important to observe, that  $M \dot{+} N$  is in general different from  $M \oplus N$  since both have different norms. In fact,  $M \dot{+} N$  might not even be closed (no problems occur if one of the spaces is finite dimensional — see Corollary 1.19 below).

**Example 1.23.** Consider  $X := \ell^p(\mathbb{N})$ . Let  $M = \{a \in X | a_{2n} = 0\}$  and  $N = \{a \in X | a_{2n+1} = n^3 a_{2n}\}$ . Then both subspaces are closed and  $M \cap N = \{0\}$ . Moreover,  $M \dot{+} N$  is dense since it contains all sequences with finite support. However, it is not all of  $X$  since  $a_n = \frac{1}{n^2} \notin M \dot{+} N$ . Indeed, if we could write  $a = b + c \in M \dot{+} N$ , then  $c_{2n} = \frac{1}{4n^2}$  and hence  $c_{2n+1} = \frac{n}{4}$  contradicting  $c \in N \subseteq X$ .  $\diamond$

A closed subspace  $M$  is called **complemented** if we can find another closed subspace  $N$  with  $M \cap N = \{0\}$  and  $M \dot{+} N = X$ . In this case every  $x \in X$  can be uniquely written as  $x = x_1 + x_2$  with  $x_1 \in M$ ,  $x_2 \in N$  and we can define a projection  $P : X \rightarrow M$ ,  $x \mapsto x_1$ . By definition  $P^2 = P$  and we have a complementary projection  $Q := \mathbb{I} - P$  with  $Q : X \rightarrow N$ ,  $x \mapsto x_2$ . Moreover, it is straightforward to check  $M = \text{Ker}(Q) = \text{Ran}(P)$  and  $N = \text{Ker}(P) = \text{Ran}(Q)$ . Of course one would like  $P$  (and hence also  $Q$ ) to be continuous. If we consider the linear bijection  $\phi : M \oplus N \rightarrow X$ ,  $(x_1, x_2) \mapsto x_1 + x_2$  then this is equivalent to the question if  $\phi^{-1}$  is continuous. By the triangle inequality  $\phi$  is continuous with  $\|\phi\| \leq 1$  and the inverse mapping theorem (Theorem 4.6) will answer this question affirmative. In summary, we have  $M \oplus N \cong X$ .

It is important to emphasize, that it is precisely the requirement that  $N$  is closed which makes  $P$  continuous (conversely observe that  $N = \text{Ker}(P)$  is closed if  $P$  is continuous). Without this requirement we can always find  $N$  by a simple application of Zorn's lemma (order the subspaces which have trivial intersection with  $M$  by inclusion and note that a maximal element has the required properties). Moreover, the question which closed subspaces can be complemented is a highly nontrivial one. If  $M$  is finite (co)dimensional, then it can be complemented (see Problems 1.47 and 4.20).

Given a subspace  $M$  of a linear space  $X$  we can define the **quotient space**  $X/M$  as the set of all equivalence classes  $[x] = x + M$  with respect to the equivalence relation  $x \equiv y$  if  $x - y \in M$ . It is straightforward to see that  $X/M$  is a vector space when defining  $[x] + [y] = [x + y]$  and  $\alpha[x] = [\alpha x]$  (show that these definitions are independent of the representative of the equivalence class). The **quotient map**  $\pi : X \rightarrow X/M$ ,  $x \mapsto [x]$  is a linear surjective map with  $\text{Ker}(\pi) = M$ . In particular, for a linear operator  $A : X \rightarrow Y$  the linear space  $\text{Coker}(A) := Y/\text{Ran}(A)$  is known as the **cokernel** of  $A$ . The dimension of  $X/M$  is known as the **codimension** of  $M$ .

**Lemma 1.18.** *Let  $M$  be a closed subspace of a Banach space  $X$ . Then  $X/M$  together with the norm*

$$\|[x]\| := \text{dist}(x, M) = \inf_{y \in M} \|x + y\| \quad (1.65)$$

*is a Banach space.*

**Proof.** First of all we need to show that (1.65) is indeed a norm. If  $\|[x]\| = 0$  we must have a sequence  $y_j \in M$  with  $y_j \rightarrow -x$  and since  $M$  is closed we conclude  $x \in M$ , that is  $[x] = [0]$  as required. To see  $\|\alpha[x]\| = |\alpha|\|[x]\|$  we use again the definition

$$\begin{aligned} \|\alpha[x]\| &= \|[\alpha x]\| = \inf_{y \in M} \|\alpha x + y\| = \inf_{y \in M} \|\alpha x + \alpha y\| \\ &= |\alpha| \inf_{y \in M} \|x + y\| = |\alpha|\|[x]\|. \end{aligned}$$

The triangle inequality follows with a similar argument and is left as an exercise.

Thus (1.65) is a norm and it remains to show that  $X/M$  is complete. To this end let  $[x_n]$  be a Cauchy sequence. Since it suffices to show that some subsequence has a limit, we can assume  $\|[x_{n+1}] - [x_n]\| < 2^{-n}$  without loss of generality. Moreover, by definition of (1.65) we can choose the representatives  $x_n$  such that  $\|x_{n+1} - x_n\| < 2^{-n}$  (start with  $x_1$  and then choose the remaining ones inductively). By construction  $x_n$  is a Cauchy sequence which has a limit  $x \in X$  since  $X$  is complete. Moreover, by  $\|[x_n] - [x]\| = \|[x_n - x]\| \leq \|x_n - x\|$  we see that  $[x]$  is the limit of  $[x_n]$ .  $\square$



Observe that  $\|[x]\| = \text{dist}(x, M) = 0$  whenever  $x \in \overline{M}$  and hence we only get a semi-norm if  $M$  is not closed.

**Example 1.24.** If  $X := C[0, 1]$  and  $M := \{f \in X \mid f(0) = 0\}$  then  $X/M \cong \mathbb{C}$ .  $\diamond$

**Example 1.25.** If  $X := c(\mathbb{N})$ , the convergent sequences and  $M := c_0(\mathbb{N})$  the sequences converging to 0, then  $X/M \cong \mathbb{C}$ . In fact, note that every sequence  $x \in c(\mathbb{N})$  can be written as  $x = y + \alpha e$  with  $y \in c_0(\mathbb{N})$ ,  $e = (1, 1, 1, \dots)$ , and  $\alpha \in \mathbb{C}$  its limit.  $\diamond$

Note that by  $\|[x]\| \leq \|x\|$  the quotient map  $\pi : X \rightarrow X/M$ ,  $x \mapsto [x]$  is bounded with norm at most one. As a small application we note:

**Corollary 1.19.** Let  $X$  be a Banach space and let  $M, N \subseteq X$  be two closed subspaces with one of them, say  $N$ , finite dimensional. Then  $M + N$  is also closed.

**Proof.** If  $\pi : X \rightarrow X/M$  denotes the quotient map, then  $M + N = \pi^{-1}(\pi(N))$ . Moreover, since  $\pi(N)$  is finite dimensional it is closed and hence  $\pi^{-1}(\pi(N))$  is closed by continuity.  $\square$

**Problem\* 1.41.** Let  $X_j$ ,  $j = 1, \dots, n$ , be Banach spaces. Let  $X := \bigoplus_{j=1}^n X_j$  be the Cartesian product  $X_1 \times \dots \times X_n$  together with the norm

$$\|(x_1, \dots, x_n)\|_p := \begin{cases} \left( \sum_{j=1}^n \|x_j\|^p \right)^{1/p}, & 1 \leq p < \infty, \\ \max_{j=1, \dots, n} \|x_j\|, & p = \infty. \end{cases}$$

Show that  $X$  is a Banach space. Show that all norms are equivalent and that this sum is associative  $(X_1 \oplus_p X_2) \oplus_p X_3 = X_1 \oplus_p (X_2 \oplus_p X_3)$ .

**Problem 1.42.** Let  $X_j$ ,  $j \in \mathbb{N}$ , be Banach spaces. Let  $X := \bigoplus_{j \in \mathbb{N}} X_j$  be the set of all elements  $x = (x_j)_{j \in \mathbb{N}}$  of the Cartesian product for which the norm

$$\|x\|_p := \begin{cases} \left( \sum_{j \in \mathbb{N}} \|x_j\|^p \right)^{1/p}, & 1 \leq p < \infty, \\ \max_{j \in \mathbb{N}} \|x_j\|, & p = \infty, \end{cases}$$

is finite. Show that  $X$  is a Banach space. Show that for  $1 \leq p < \infty$  the elements with finitely many nonzero terms are dense and conclude that  $X$  is separable if all  $X_j$  are.

**Problem 1.43.** Let  $\ell$  be a nontrivial linear functional. Then its kernel has codimension one.

**Problem 1.44.** Compute  $\|[e]\|$  in  $\ell^\infty(\mathbb{N})/c_0(\mathbb{N})$ , where  $e = (1, 1, 1, \dots)$ .

**Problem 1.45** (Complexification). Given a real normed space  $X$  its **complexification** is given by  $X_{\mathbb{C}} := X \times X$  together with the (complex) scalar

multiplication  $\alpha(x, y) = (\operatorname{Re}(\alpha)x - \operatorname{Im}(\alpha)y, \operatorname{Re}(\alpha)y + \operatorname{Im}(\alpha)x)$ . By virtue of the embedding  $x \mapsto (x, 0)$  you should of course think of  $(x, y)$  as  $x + iy$ .

Show that

$$\|x + iy\|_{\mathbb{C}} := \max_{0 \leq t \leq \pi} \|\cos(t)x + \sin(t)y\|,$$

defines a norm on  $X_{\mathbb{C}}$  which satisfies  $\|x\|_{\mathbb{C}} = \|x\|$  and

$$\max(\|x\|, \|y\|) \leq \|x + iy\|_{\mathbb{C}} \leq (\|x\|^2 + \|y\|^2)^{1/2}$$

In particular, this norm is equivalent to the product norm on  $X \oplus X$ .

If  $X$  is a Hilbert space, then the above norm will in general not give rise to a scalar product. However, any bilinear form  $s : X \times X \rightarrow \mathbb{R}$  gives rise to a sesquilinear form  $s_{\mathbb{C}}(x_1 + iy_1, x_2 + iy_2) := s(x_1, x_2) + s(y_1, y_2) + i(s(x_1, y_2) - s(y_1, x_2))$ . If  $s$  is symmetric or positive definite, so will be  $s_{\mathbb{C}}$ . The corresponding norm satisfies  $\langle x + iy, x + iy \rangle_{\mathbb{C}} = \|x\|^2 + \|y\|^2$  and is equivalent to the above one since  $\frac{1}{2}(\|x\|^2 + \|y\|^2) \leq \|x + iy\|_{\mathbb{C}}^2 \leq \|x\|^2 + \|y\|^2$ .

Given two real normed spaces  $X, Y$ , every linear operator  $A : X \rightarrow Y$  gives rise to a linear operator  $A_{\mathbb{C}} : X_{\mathbb{C}} \rightarrow Y_{\mathbb{C}}$  via  $A_{\mathbb{C}}(x + iy) = Ax + iAy$ . Show  $\|A_{\mathbb{C}}\| = \|A\|$ .

**Problem\* 1.46.** Suppose  $A \in \mathcal{L}(X, Y)$ . Show that  $\operatorname{Ker}(A)$  is closed. Suppose  $M \subseteq \operatorname{Ker}(A)$  is a closed subspace. Show that the induced map  $\tilde{A} : X/M \rightarrow Y$ ,  $[x] \mapsto Ax$  is a well-defined operator satisfying  $\|\tilde{A}\| = \|A\|$  and  $\operatorname{Ker}(\tilde{A}) = \operatorname{Ker}(A)/M$ . In particular,  $\tilde{A}$  is injective for  $M = \operatorname{Ker}(A)$ .

**Problem\* 1.47.** Show that if a closed subspace  $M$  of a Banach space  $X$  has finite codimension, then it can be complemented. (Hint: Start with a basis  $\{[x_j]\}$  for  $X/M$  and choose a corresponding dual basis  $\{\ell_k\}$  with  $\ell_k([x_j]) = \delta_{j,k}$ .)

## 1.8. Spaces of continuous and differentiable functions

In this section we introduce a few further sets of continuous and differentiable functions which are of interest in applications. Let  $I$  be some compact interval, then we can make  $C^1(I)$  into a Banach space by (Problem 1.35) by introducing the norm  $\|f\|_{1,\infty} := \|f\|_{\infty} + \|f'\|_{\infty}$ . By a straightforward extension we can even get (cf. Problem 1.49)

**Theorem 1.20.** Let  $I \subseteq \mathbb{R}$  be some interval. The space  $C_b^k(I)$  of all functions whose partial derivatives up to order  $k$  are bounded and continuous form a Banach space with norm

$$\|f\|_{k,\infty} := \sum_{|\alpha| \leq k} \sup_{x \in I} |f^{(\alpha)}(x)|. \quad (1.66)$$

Note that the space  $C_b^k(I)$  could be further refined by requiring the highest derivatives to be Hölder continuous. Recall that a function  $f : I \rightarrow \mathbb{C}$  is called uniformly **Hölder continuous** with exponent  $\gamma \in (0, 1]$  if

$$[f]_\gamma := \sup_{x \neq y \in I} \frac{|f(x) - f(y)|}{|x - y|^\gamma} \quad (1.67)$$

is finite. Clearly, any Hölder continuous function is uniformly continuous and, in the special case  $\gamma = 1$ , we obtain the **Lipschitz continuous** functions. Note that for  $\gamma = 0$  the Hölder condition boils down to boundedness and also the case  $\gamma > 1$  is not very interesting (Problem 1.48).

**Example 1.26.** By the mean value theorem every function  $f \in C_b^1(I)$  is Lipschitz continuous with  $[f]_\gamma \leq \|f'\|_\infty$ .  $\diamond$

**Example 1.27.** The prototypical example of a Hölder continuous function is of course  $f(x) = x^\gamma$  on  $[0, \infty)$  with  $\gamma \in (0, 1]$ . In fact, without loss of generality we can assume  $0 \leq x < y$  and set  $t = \frac{x}{y} \in [0, 1)$ . Then we have

$$\frac{y^\gamma - x^\gamma}{(y - x)^\gamma} \leq \frac{1 - t^\gamma}{(1 - t)^\gamma} \leq \frac{1 - t}{1 - t} = 1.$$

From this one easily gets further examples since the composition of two Hölder continuous functions is again Hölder continuous (the exponent being the product).  $\diamond$

It is easy to verify that this is a seminorm and that the corresponding space is complete.

**Theorem 1.21.** *Let  $I \subseteq \mathbb{R}$  be an interval. The space  $C_b^{k,\gamma}(I)$  of all functions whose partial derivatives up to order  $k$  are bounded and Hölder continuous with exponent  $\gamma \in (0, 1]$  form a Banach space with norm*

$$\|f\|_{k,\gamma,\infty} := \|f\|_{k,\infty} + [f^{(k)}]_\gamma. \quad (1.68)$$

As already noted before, in the case  $\gamma = 0$  we get a norm which is equivalent to  $\|f\|_{\infty,k}$  and we will set  $C_b^{k,0}(I) := C_b^k(I)$  for notational convenience later on.

Note that by the mean value theorem all derivatives up to order lower than  $k$  are automatically Lipschitz continuous. Moreover, every Hölder continuous function is uniformly continuous and hence has a unique extension to the closure  $\bar{I}$  (cf. Theorem B.39). In this sense, the spaces  $C_b^{0,\gamma}(I)$  and  $C_b^{0,\gamma}(\bar{I})$  are naturally isomorphic. Finally, since Hölder continuous functions on a bounded domain are automatically bounded, we can drop the subscript  $b$  in this situation.

**Theorem 1.22.** *Suppose  $I \subset \mathbb{R}$  is a compact interval. Then  $C^{0,\gamma_2}(I) \subseteq C^{0,\gamma_1}(I) \subseteq C(I)$  for  $0 < \gamma_1 < \gamma_2 \leq 1$  with the embeddings being compact.*

**Proof.** That we have continuous embeddings follows since  $|x - y|^{-\gamma_1} = |x - y|^{-\gamma_2 + (\gamma_2 - \gamma_1)} \leq (2r)^{\gamma_2 - \gamma_1} |x - y|^{-\gamma_2}$  if  $r$  denotes the length of  $I$ . Moreover, that the embedding  $C^{0, \gamma_1}(I) \subseteq C(I)$  is compact follows from the Arzelà–Ascoli theorem (Theorem B.40). To see the remaining claim let  $f_m$  be a bounded sequence in  $C^{0, \gamma_1}(I)$ , explicitly  $\|f_m\|_\infty \leq C$  and  $[f]_{\gamma_1} \leq C$ . Hence by the Arzelà–Ascoli theorem we can assume that  $f_m$  converges uniformly to some  $f \in C(I)$ . Moreover, taking the limit in  $|f_m(x) - f_m(y)| \leq C|x - y|^{\gamma_1}$  we see that we even have  $f \in C^{0, \gamma_1}(I)$ . To see that  $f$  is the limit of  $f_m$  in  $C^{0, \gamma_2}(I)$  we need to show  $[g_m]_{\gamma_2} \rightarrow 0$ , where  $g_m := f_m - f$ . Now observe that

$$\begin{aligned} [g_m]_{\gamma_2} &= \sup_{x \neq y \in I: |x-y| \geq \varepsilon} \frac{|g_m(x) - g_m(y)|}{|x - y|^{\gamma_2}} + \sup_{x \neq y \in I: |x-y| < \varepsilon} \frac{|g_m(x) - g_m(y)|}{|x - y|^{\gamma_2}} \\ &\leq 2\|g_m\|_\infty \varepsilon^{-\gamma_2} + [g_m]_{\gamma_1} \varepsilon^{\gamma_1 - \gamma_2} \leq 2\|g_m\|_\infty \varepsilon^{-\gamma_2} + 2C\varepsilon^{\gamma_1 - \gamma_2}, \end{aligned}$$

implying  $\limsup_{m \rightarrow \infty} [g_m]_{\gamma_2} \leq 2C\varepsilon^{\gamma_1 - \gamma_2}$  and since  $\varepsilon > 0$  is arbitrary this establishes the claim.  $\square$

As pointed out in the example before, the embedding  $C_b^1(I) \subseteq C_b^{0,1}(I)$  is continuous and combining this with the previous result immediately gives

**Corollary 1.23.** *Suppose  $I \subset \mathbb{R}$  is a compact interval,  $k_1, k_2 \in \mathbb{N}_0$ , and  $0 \leq \gamma_1, \gamma_2 \leq 1$ . Then  $C^{k_2, \gamma_2}(I) \subseteq C^{k_1, \gamma_1}(I)$  for  $k_1 + \gamma_1 \leq k_2 + \gamma_2$  with the embeddings being compact if the inequality is strict.*

For now continuous functions on intervals will be sufficient for our purpose. However, once we delve deeper into the subject we will also need continuous functions on topological spaces  $X$ . Luckily most of the results extend to this case in a more or less straightforward way. If you are not familiar with these extensions you can find them in Section B.8.

**Problem 1.48.** *Let  $I$  be an interval. Suppose  $f : I \rightarrow \mathbb{C}$  is Hölder continuous with exponent  $\gamma > 1$ . Show that  $f$  is constant.*

**Problem\* 1.49.** *Suppose  $X$  is a vector space and  $\|\cdot\|_j$ ,  $1 \leq j \leq n$ , is a finite family of seminorms. Show that  $\|x\| := \sum_{j=1}^n \|x\|_j$  is a seminorm. It is a norm if and only if  $\|x\|_j = 0$  for all  $j$  implies  $x = 0$ .*

**Problem 1.50.** *Let  $I$ . Show that the product of two bounded Hölder continuous functions is again Hölder continuous with*

$$[fg]_\gamma \leq \|f\|_\infty [g]_\gamma + [f]_\gamma \|g\|_\infty.$$



# Hilbert spaces

The additional geometric structure of Hilbert spaces allows for an intuitive geometric solution of many problems. In fact, in many situations, e.g. in Quantum Mechanics, Hilbert spaces occur naturally. This makes them the weapon of choice whenever possible. Throughout this chapter  $\mathfrak{H}$  will be a (complex) Hilbert space.

## 2.1. Orthonormal bases

In this section we will investigate orthonormal series and you will notice hardly any difference between the finite and infinite dimensional cases. As our first task, let us generalize the projection into the direction of one vector.

A set of vectors  $\{u_j\}$  is called an **orthonormal set** if  $\langle u_j, u_k \rangle = 0$  for  $j \neq k$  and  $\langle u_j, u_j \rangle = 1$ . Note that every orthonormal set is linearly independent (show this).

**Lemma 2.1.** *Suppose  $\{u_j\}_{j=1}^n$  is a finite orthonormal set in a Hilbert space  $\mathfrak{H}$ . Then every  $f \in \mathfrak{H}$  can be written as*

$$f = f_{\parallel} + f_{\perp}, \quad f_{\parallel} := \sum_{j=1}^n \langle u_j, f \rangle u_j, \quad (2.1)$$

where  $f_{\parallel}$  and  $f_{\perp}$  are orthogonal. Moreover,  $\langle u_j, f_{\perp} \rangle = 0$  for all  $1 \leq j \leq n$ . In particular,

$$\|f\|^2 = \sum_{j=1}^n |\langle u_j, f \rangle|^2 + \|f_{\perp}\|^2. \quad (2.2)$$

Moreover, every  $\hat{f}$  in the span of  $\{u_j\}_{j=1}^n$  satisfies

$$\|f - \hat{f}\| \geq \|f_{\perp}\| \quad (2.3)$$

with equality holding if and only if  $\hat{f} = f_{\parallel}$ . In other words,  $f_{\parallel}$  is uniquely characterized as the vector in the span of  $\{u_j\}_{j=1}^n$  closest to  $f$ .

**Proof.** A straightforward calculation shows  $\langle u_j, f - f_{\parallel} \rangle = 0$  and hence  $f_{\parallel}$  and  $f_{\perp} := f - f_{\parallel}$  are orthogonal. The formula for the norm follows by applying (1.40) iteratively.

Now, fix a vector  $\hat{f} := \sum_{j=1}^n \alpha_j u_j$  in the span of  $\{u_j\}_{j=1}^n$ . Then one computes

$$\begin{aligned} \|f - \hat{f}\|^2 &= \|f_{\parallel} + f_{\perp} - \hat{f}\|^2 = \|f_{\perp}\|^2 + \|f_{\parallel} - \hat{f}\|^2 \\ &= \|f_{\perp}\|^2 + \sum_{j=1}^n |\alpha_j - \langle u_j, f \rangle|^2 \end{aligned}$$

from which the last claim follows.  $\square$

From (2.2) we obtain **Bessel's inequality**

$$\sum_{j=1}^n |\langle u_j, f \rangle|^2 \leq \|f\|^2 \quad (2.4)$$

with equality holding if and only if  $f$  lies in the span of  $\{u_j\}_{j=1}^n$ .

Of course, since we cannot assume  $\mathfrak{H}$  to be a finite dimensional vector space, we need to generalize Lemma 2.1 to arbitrary orthonormal sets  $\{u_j\}_{j \in J}$ . We start by assuming that  $J$  is countable. Then Bessel's inequality (2.4) shows that

$$\sum_{j \in J} |\langle u_j, f \rangle|^2 \quad (2.5)$$

converges absolutely. Moreover, for any finite subset  $K \subset J$  we have

$$\left\| \sum_{j \in K} \langle u_j, f \rangle u_j \right\|^2 = \sum_{j \in K} |\langle u_j, f \rangle|^2 \quad (2.6)$$

by the Pythagorean theorem and thus  $\sum_{j \in J} \langle u_j, f \rangle u_j$  is a Cauchy sequence if and only if  $\sum_{j \in J} |\langle u_j, f \rangle|^2$  is. Now let  $J$  be arbitrary. Again, Bessel's inequality shows that for any given  $\varepsilon > 0$  there are at most finitely many  $j$  for which  $|\langle u_j, f \rangle| \geq \varepsilon$  (namely at most  $\|f\|/\varepsilon$ ). Hence there are at most countably many  $j$  for which  $|\langle u_j, f \rangle| > 0$ . Thus it follows that

$$\sum_{j \in J} |\langle u_j, f \rangle|^2 \quad (2.7)$$

is well defined (as a countable sum over the nonzero terms) and (by completeness) so is

$$\sum_{j \in J} \langle u_j, f \rangle u_j. \quad (2.8)$$

Furthermore, it is also independent of the order of summation.

In particular, by continuity of the scalar product we see that Lemma 2.1 can be generalized to arbitrary orthonormal sets.

**Theorem 2.2.** *Suppose  $\{u_j\}_{j \in J}$  is an orthonormal set in a Hilbert space  $\mathfrak{H}$ . Then every  $f \in \mathfrak{H}$  can be written as*

$$f = f_{\parallel} + f_{\perp}, \quad f_{\parallel} := \sum_{j \in J} \langle u_j, f \rangle u_j, \quad (2.9)$$

where  $f_{\parallel}$  and  $f_{\perp}$  are orthogonal. Moreover,  $\langle u_j, f_{\perp} \rangle = 0$  for all  $j \in J$ . In particular,

$$\|f\|^2 = \sum_{j \in J} |\langle u_j, f \rangle|^2 + \|f_{\perp}\|^2. \quad (2.10)$$

Furthermore, every  $\hat{f} \in \overline{\text{span}\{u_j\}_{j \in J}}$  satisfies

$$\|f - \hat{f}\| \geq \|f_{\perp}\| \quad (2.11)$$

with equality holding if and only if  $\hat{f} = f_{\parallel}$ . In other words,  $f_{\parallel}$  is uniquely characterized as the vector in  $\overline{\text{span}\{u_j\}_{j \in J}}$  closest to  $f$ .

**Proof.** The first part follows as in Lemma 2.1 using continuity of the scalar product. The same is true for the last part except for the fact that every  $f \in \overline{\text{span}\{u_j\}_{j \in J}}$  can be written as  $f = \sum_{j \in J} \alpha_j u_j$  (i.e.,  $f = f_{\parallel}$ ). To see this, let  $f_n \in \text{span}\{u_j\}_{j \in J}$  converge to  $f$ . Then  $\|f - f_n\|^2 = \|f_{\parallel} - f_n\|^2 + \|f_{\perp}\|^2 \rightarrow 0$  implies  $f_n \rightarrow f_{\parallel}$  and  $f_{\perp} = 0$ .  $\square$

Note that from Bessel's inequality (which of course still holds), it follows that the map  $f \rightarrow f_{\parallel}$  is continuous.

Of course we are particularly interested in the case where every  $f \in \mathfrak{H}$  can be written as  $\sum_{j \in J} \langle u_j, f \rangle u_j$ . In this case we will call the orthonormal set  $\{u_j\}_{j \in J}$  an **orthonormal basis** (ONB).

If  $\mathfrak{H}$  is separable it is easy to construct an orthonormal basis. In fact, if  $\mathfrak{H}$  is separable, then there exists a countable total set  $\{f_j\}_{j=1}^N$ . Here  $N \in \mathbb{N}$  if  $\mathfrak{H}$  is finite dimensional and  $N = \infty$  otherwise. After throwing away some vectors, we can assume that  $f_{n+1}$  cannot be expressed as a linear combination of the vectors  $f_1, \dots, f_n$ . Now we can construct an orthonormal set as follows: We begin by normalizing  $f_1$ :

$$u_1 := \frac{f_1}{\|f_1\|}. \quad (2.12)$$

Next we take  $f_2$  and remove the component parallel to  $u_1$  and normalize again:

$$u_2 := \frac{f_2 - \langle u_1, f_2 \rangle u_1}{\|f_2 - \langle u_1, f_2 \rangle u_1\|}. \quad (2.13)$$



Proceeding like this, we define recursively

$$u_n := \frac{f_n - \sum_{j=1}^{n-1} \langle u_j, f_n \rangle u_j}{\|f_n - \sum_{j=1}^{n-1} \langle u_j, f_n \rangle u_j\|}. \quad (2.14)$$

This procedure is known as **Gram–Schmidt orthogonalization**. Hence we obtain an orthonormal set  $\{u_j\}_{j=1}^N$  such that  $\text{span}\{u_j\}_{j=1}^n = \text{span}\{f_j\}_{j=1}^n$  for any finite  $n$  and thus also for  $n = N$  (if  $N = \infty$ ). Since  $\{f_j\}_{j=1}^N$  is total, so is  $\{u_j\}_{j=1}^N$ . Now suppose there is some  $f = f_{\parallel} + f_{\perp} \in \mathfrak{H}$  for which  $f_{\perp} \neq 0$ . Since  $\{u_j\}_{j=1}^N$  is total, we can find a  $\hat{f}$  in its span such that  $\|f - \hat{f}\| < \|f_{\perp}\|$ , contradicting (2.11). Hence we infer that  $\{u_j\}_{j=1}^N$  is an orthonormal basis.

**Theorem 2.3.** *Every separable Hilbert space has a countable orthonormal basis.*

**Example 2.1.** The vectors  $\{\delta^n\}_{n \in \mathbb{N}}$  form an orthonormal basis for  $\ell^2(\mathbb{N})$ .  $\diamond$

**Example 2.2.** In  $\mathcal{L}_{\text{cont}}^2(-1, 1)$ , we can orthogonalize the monomials  $f_n(x) = x^n$  (which are total by the Weierstraß approximation theorem — Theorem 1.3). The resulting polynomials are up to a normalization known as **Legendre polynomials**

$$P_0(x) = 1, \quad P_1(x) = x, \quad P_2(x) = \frac{3x^2 - 1}{2}, \quad \dots \quad (2.15)$$

(which are normalized such that  $P_n(1) = 1$ ).  $\diamond$

**Example 2.3.** The set of functions

$$u_n(x) = \frac{1}{\sqrt{2\pi}} e^{inx}, \quad n \in \mathbb{Z}, \quad (2.16)$$

forms an orthonormal basis for  $\mathfrak{H} = \mathcal{L}_{\text{cont}}^2(0, 2\pi)$ . The corresponding orthogonal expansion is just the ordinary Fourier series. We will discuss this example in detail in Section 2.5.  $\diamond$

The following equivalent properties also characterize a basis.

**Theorem 2.4.** *For an orthonormal set  $\{u_j\}_{j \in J}$  in a Hilbert space  $\mathfrak{H}$ , the following conditions are equivalent:*

- (i)  $\{u_j\}_{j \in J}$  is a maximal orthogonal set.
- (ii) For every vector  $f \in \mathfrak{H}$  we have

$$f = \sum_{j \in J} \langle u_j, f \rangle u_j. \quad (2.17)$$

- (iii) For every vector  $f \in \mathfrak{H}$  we have **Parseval's relation**

$$\|f\|^2 = \sum_{j \in J} |\langle u_j, f \rangle|^2. \quad (2.18)$$

(iv)  $\langle u_j, f \rangle = 0$  for all  $j \in J$  implies  $f = 0$ .

**Proof.** We will use the notation from Theorem 2.2.

(i)  $\Rightarrow$  (ii): If  $f_\perp \neq 0$ , then we can normalize  $f_\perp$  to obtain a unit vector  $\tilde{f}_\perp$  which is orthogonal to all vectors  $u_j$ . But then  $\{u_j\}_{j \in J} \cup \{\tilde{f}_\perp\}$  would be a larger orthonormal set, contradicting the maximality of  $\{u_j\}_{j \in J}$ .

(ii)  $\Rightarrow$  (iii): This follows since (ii) implies  $f_\perp = 0$ .

(iii)  $\Rightarrow$  (iv): If  $\langle f, u_j \rangle = 0$  for all  $j \in J$ , we conclude  $\|f\|^2 = 0$  and hence  $f = 0$ .

(iv)  $\Rightarrow$  (i): If  $\{u_j\}_{j \in J}$  were not maximal, there would be a unit vector  $g$  such that  $\{u_j\}_{j \in J} \cup \{g\}$  is a larger orthonormal set. But  $\langle u_j, g \rangle = 0$  for all  $j \in J$  implies  $g = 0$  by (iv), a contradiction.  $\square$

By continuity of the norm it suffices to check (iii), and hence also (ii), for  $f$  in a dense set. In fact, by the inverse triangle inequality for  $\ell^2(\mathbb{N})$  and the Bessel inequality we have

$$\left| \sum_{j \in J} |\langle u_j, f \rangle|^2 - \sum_{j \in J} |\langle u_j, g \rangle|^2 \right| \leq \sqrt{\sum_{j \in J} |\langle u_j, f - g \rangle|^2} \sqrt{\sum_{j \in J} |\langle u_j, f + g \rangle|^2} \\ \leq \|f - g\| \|f + g\| \quad (2.19)$$

implying  $\sum_{j \in J} |\langle u_j, f_n \rangle|^2 \rightarrow \sum_{j \in J} |\langle u_j, f \rangle|^2$  if  $f_n \rightarrow f$ .

It is not surprising that if there is one countable basis, then it follows that every other basis is countable as well.

**Theorem 2.5.** *In a Hilbert space  $\mathfrak{H}$  every orthonormal basis has the same cardinality.*

**Proof.** Let  $\{u_j\}_{j \in J}$  and  $\{v_k\}_{k \in K}$  be two orthonormal bases. We first look at the case where one of them, say the first, is finite dimensional:  $J = \{1, \dots, n\}$ . Suppose the other basis has at least  $n$  elements  $\{1, \dots, n\} \subseteq K$ . Then  $v_k = \sum_{j=1}^n U_{k,j} u_j$ , where  $U_{k,j} = \langle u_j, v_k \rangle$ . By  $\delta_{j,k} = \langle v_j, v_k \rangle = \sum_{l=1}^n U_{j,l}^* U_{k,l}$  we see  $u_j = \sum_{k=1}^n U_{k,j}^* v_k$  showing that  $K$  cannot have more than  $n$  elements.

Now let us turn to the case where both  $J$  and  $K$  are infinite. Set  $K_j = \{k \in K | \langle v_k, u_j \rangle \neq 0\}$ . Since these are the expansion coefficients of  $u_j$  with respect to  $\{v_k\}_{k \in K}$ , this set is countable (and nonempty). Hence the set  $\tilde{K} = \bigcup_{j \in J} K_j$  satisfies  $|\tilde{K}| \leq |J \times \mathbb{N}| = |J|$  (Theorem A.9) But  $k \in K \setminus \tilde{K}$  implies  $v_k = 0$  and hence  $\tilde{K} = K$ . So  $|K| \leq |J|$  and reversing the roles of  $J$  and  $K$  shows  $|K| = |J|$ .  $\square$

The cardinality of an orthonormal basis is also called the Hilbert space **dimension** of  $\mathfrak{H}$ .

It even turns out that, up to unitary equivalence, there is only one separable infinite dimensional Hilbert space:

A bijective linear operator  $U \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  is called **unitary** if  $U$  preserves scalar products:

$$\langle Ug, Uf \rangle_2 = \langle g, f \rangle_1, \quad g, f \in \mathfrak{H}_1. \quad (2.20)$$

By the polarization identity, (1.47) this is the case if and only if  $U$  preserves norms:  $\|Uf\|_2 = \|f\|_1$  for all  $f \in \mathfrak{H}_1$  (note a norm preserving linear operator is automatically injective). The two Hilbert spaces  $\mathfrak{H}_1$  and  $\mathfrak{H}_2$  are called **unitarily equivalent** in this case.

Let  $\mathfrak{H}$  be a separable infinite dimensional Hilbert space and let  $\{u_j\}_{j \in \mathbb{N}}$  be any orthogonal basis. Then the map  $U : \mathfrak{H} \rightarrow \ell^2(\mathbb{N})$ ,  $f \mapsto (\langle u_j, f \rangle)_{j \in \mathbb{N}}$  is unitary. Indeed by Theorem 2.4 (iii) it is norm preserving and hence injective. To see that it is onto, let  $a \in \ell^2(\mathbb{N})$  and observe that by  $\|\sum_{j=m}^n a_j u_j\|^2 = \sum_{j=m}^n |a_j|^2$  the vector  $f := \sum_{j \in \mathbb{N}} a_j u_j$  is well defined and satisfies  $a_j = \langle u_j, f \rangle$ . In particular,

**Theorem 2.6.** *Any separable infinite dimensional Hilbert space is unitarily equivalent to  $\ell^2(\mathbb{N})$ .*

Of course the same argument shows that every finite dimensional Hilbert space of dimension  $n$  is unitarily equivalent to  $\mathbb{C}^n$  with the usual scalar product.

Finally we briefly turn to the case where  $\mathfrak{H}$  is not separable.

**Theorem 2.7.** *Every Hilbert space has an orthonormal basis.*

**Proof.** To prove this we need to resort to Zorn's lemma (see Appendix A): The collection of all orthonormal sets in  $\mathfrak{H}$  can be partially ordered by inclusion. Moreover, every linearly ordered chain has an upper bound (the union of all sets in the chain). Hence Zorn's lemma implies the existence of a maximal element, that is, an orthonormal set which is not a proper subset of every other orthonormal set.  $\square$

Hence, if  $\{u_j\}_{j \in J}$  is an orthogonal basis, we can show that  $\mathfrak{H}$  is unitarily equivalent to  $\ell^2(J)$  and, by prescribing  $J$ , we can find a Hilbert space of any given dimension. Here  $\ell^2(J)$  is the set of all complex-valued functions  $(a_j)_{j \in J}$  where at most countably many values are nonzero and  $\sum_{j \in J} |a_j|^2 < \infty$ .

**Example 2.4.** Define the set of **almost periodic functions**  $AP(\mathbb{R})$  as the closure of the set of trigonometric polynomials

$$f(t) = \sum_{k=1}^n \alpha_k e^{i\theta_k t}, \quad \alpha_k \in \mathbb{C}, \theta_k \in \mathbb{R},$$

with respect to the sup norm. In particular  $AP(\mathbb{R}) \subset C_b(\mathbb{R})$  is a Banach space when equipped with the sup norm. Since the trigonometric polynomials form an algebra, it is even a Banach algebra. Using the Stone–Weierstraß theorem one can verify that every periodic function is almost periodic (make the approximation on one period and note that you get the rest of  $\mathbb{R}$  for free from periodicity) but the converse is not true (e.g.  $e^{it} + e^{i\sqrt{2}t}$  is not periodic).

It is not difficult to show that

$$\lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T}^T e^{i\theta t} dt = \begin{cases} 1, & \theta = 0, \\ 0, & \theta \neq 0, \end{cases}$$

and hence one can conclude that every almost periodic function has a mean value

$$M(f) := \lim_{T \rightarrow \infty} \frac{1}{2T} \int_{-T}^T f(t) dt.$$

Next one can show that

$$\langle f, g \rangle := M(f^* g)$$

defines a scalar product on  $AP(\mathbb{R})$  (only positivity is nontrivial and it will not be shown here). Note that  $\|f\| \leq \|f\|_\infty$ . Abbreviating  $e_\theta(t) = e^{i\theta t}$  we see that  $\{e_\theta\}_{\theta \in \mathbb{R}}$  is an uncountable orthonormal set and

$$f(t) \mapsto \hat{f}(\theta) := \langle e_\theta, f \rangle = M(e_{-\theta} f)$$

maps  $AP(\mathbb{R})$  isometrically (with respect to  $\|\cdot\|$ ) into  $\ell^2(\mathbb{R})$ . This map is however not surjective (take e.g. a Fourier series which converges in mean square but not uniformly — see later).  $\diamond$

**Problem 2.1.** Given some vectors  $f_1, \dots, f_n$  we define their **Gram determinant** as

$$\Gamma(f_1, \dots, f_n) := \det(\langle f_j, f_k \rangle)_{1 \leq j, k \leq n}.$$

Show that the Gram determinant is nonzero if and only if the vectors are linearly independent. Moreover, show that in this case

$$\text{dist}(g, \text{span}\{f_1, \dots, f_n\})^2 = \frac{\Gamma(f_1, \dots, f_n, g)}{\Gamma(f_1, \dots, f_n)}$$

and

$$\Gamma(f_1, \dots, f_n) \leq \prod_{j=1}^n \|f_j\|^2.$$

with equality if the vectors are orthogonal. (Hint: How does  $\Gamma$  change when you apply the Gram–Schmidt procedure?)

**Problem 2.2.** Let  $\{u_j\}$  be some orthonormal basis. Show that a bounded linear operator  $A$  is uniquely determined by its matrix elements  $A_{jk} := \langle u_j, Au_k \rangle$  with respect to this basis.

**Problem 2.3.** Give an example of a nonempty closed bounded subset of a Hilbert space which does not contain an element with minimal norm. Can this happen in finite dimensions? (Hint: Look for a discrete set.)

## 2.2. The projection theorem and the Riesz representation theorem

Let  $M \subseteq \mathfrak{H}$  be a subset. Then  $M^\perp := \{f \mid \langle g, f \rangle = 0, \forall g \in M\}$  is called the **orthogonal complement** of  $M$ . By continuity of the scalar product it follows that  $M^\perp$  is a closed linear subspace and by linearity that  $(\overline{\text{span}(M)})^\perp = M^\perp$ . For example, we have  $\mathfrak{H}^\perp = \{0\}$  since any vector in  $\mathfrak{H}^\perp$  must be in particular orthogonal to all vectors in some orthonormal basis.

**Theorem 2.8** (Projection theorem). *Let  $M$  be a closed linear subspace of a Hilbert space  $\mathfrak{H}$ . Then every  $f \in \mathfrak{H}$  can be uniquely written as  $f = f_\parallel + f_\perp$  with  $f_\parallel \in M$  and  $f_\perp \in M^\perp$ , where  $f_\parallel$  is uniquely characterized as the vector in  $M$  closest to  $f$ . One writes*

$$M \oplus M^\perp = \mathfrak{H} \quad (2.21)$$

in this situation.

**Proof.** Since  $M$  is closed, it is a Hilbert space and has an orthonormal basis  $\{u_j\}_{j \in J}$ . Hence the existence part follows from Theorem 2.2. To see uniqueness, suppose there is another decomposition  $f = \tilde{f}_\parallel + \tilde{f}_\perp$ . Then  $f_\parallel - \tilde{f}_\parallel = \tilde{f}_\perp - f_\perp \in M \cap M^\perp = \{0\}$  (since  $g \in M \cap M^\perp$  implies  $\|g\|^2 = \langle g, g \rangle = 0$ ).  $\square$

**Corollary 2.9.** *Every orthogonal set  $\{u_j\}_{j \in J}$  can be extended to an orthogonal basis.*

**Proof.** Just add an orthogonal basis for  $(\{u_j\}_{j \in J})^\perp$ .  $\square$

The operator  $P_M f := f_\parallel$  is called the **orthogonal projection** corresponding to  $M$ . Note that we have

$$P_M^2 = P_M \quad \text{and} \quad \langle P_M g, f \rangle = \langle g, P_M f \rangle \quad (2.22)$$

since  $\langle P_M g, f \rangle = \langle g_\parallel, f_\parallel \rangle = \langle g, P_M f \rangle$ . Clearly we have  $P_{M^\perp} f = f - P_M f = f_\perp$ . Furthermore, (2.22) uniquely characterizes orthogonal projections (Problem 2.6).

Moreover, if  $M$  is a closed subspace, we have  $P_{M^{\perp\perp}} = \mathbb{I} - P_{M^\perp} = \mathbb{I} - (\mathbb{I} - P_M) = P_M$ ; that is,  $M^{\perp\perp} = M$ . If  $M$  is an arbitrary subset, we have at least

$$M^{\perp\perp} = \overline{\text{span}(M)}. \quad (2.23)$$

Note that by  $\mathfrak{H}^\perp = \{0\}$  we see that  $M^\perp = \{0\}$  if and only if  $M$  is total.

Next we turn to **linear functionals**, that is, to operators  $\ell : \mathfrak{H} \rightarrow \mathbb{C}$ . By the Cauchy–Schwarz inequality we know that  $\ell_g : f \mapsto \langle g, f \rangle$  is a bounded linear functional (with norm  $\|g\|$ ). It turns out that, in a Hilbert space, every bounded linear functional can be written in this way.

**Theorem 2.10** (Riesz representation theorem). *Suppose  $\ell$  is a bounded linear functional on a Hilbert space  $\mathfrak{H}$ . Then there is a unique vector  $g \in \mathfrak{H}$  such that  $\ell(f) = \langle g, f \rangle$  for all  $f \in \mathfrak{H}$ .*

*In other words, a Hilbert space is equivalent to its own dual space  $\mathfrak{H}^* \cong \mathfrak{H}$  via the map  $f \mapsto \langle f, \cdot \rangle$  which is a conjugate linear isometric bijection between  $\mathfrak{H}$  and  $\mathfrak{H}^*$ .*

**Proof.** If  $\ell \equiv 0$ , we can choose  $g = 0$ . Otherwise  $\text{Ker}(\ell) = \{f \mid \ell(f) = 0\}$  is a proper subspace and we can find a unit vector  $\tilde{g} \in \text{Ker}(\ell)^\perp$ . For every  $f \in \mathfrak{H}$  we have  $\ell(f)\tilde{g} - \ell(\tilde{g})f \in \text{Ker}(\ell)$  and hence

$$0 = \langle \tilde{g}, \ell(f)\tilde{g} - \ell(\tilde{g})f \rangle = \ell(f) - \ell(\tilde{g})\langle \tilde{g}, f \rangle.$$

In other words, we can choose  $g = \ell(\tilde{g})^* \tilde{g}$ . To see uniqueness, let  $g_1, g_2$  be two such vectors. Then  $\langle g_1 - g_2, f \rangle = \langle g_1, f \rangle - \langle g_2, f \rangle = \ell(f) - \ell(f) = 0$  for every  $f \in \mathfrak{H}$ , which shows  $g_1 - g_2 \in \mathfrak{H}^\perp = \{0\}$ .  $\square$

In particular, this shows that  $\mathfrak{H}^*$  is again a Hilbert space whose scalar product (in terms of the above identification) is given by  $\langle \langle f, \cdot \rangle, \langle g, \cdot \rangle \rangle_{\mathfrak{H}^*} = \langle f, g \rangle^*$ .

We can even get a unitary map between  $\mathfrak{H}$  and  $\mathfrak{H}^*$  but such a map is not unique. To this end note that every Hilbert space has a conjugation  $C$  which generalizes taking the complex conjugate of every coordinate. In fact, choosing an orthonormal basis (and different choices will produce different maps in general) we can set

$$Cf := \sum_{j \in J} \langle u_j, f \rangle^* u_j = \sum_{j \in J} \langle f, u_j \rangle u_j.$$

Then  $C$  is conjugate linear, isometric  $\|Cf\| = \|f\|$ , and idempotent  $C^2 = \mathbb{I}$ . Note also  $\langle Cf, Cg \rangle = \langle f, g \rangle^*$ . As promised, the map  $f \mapsto \langle Cf, \cdot \rangle$  is a unitary map from  $\mathfrak{H}$  to  $\mathfrak{H}^*$ .

Finally, we remark that projections can not only be defined for subspaces but also for closed convex sets (of course they will no longer be linear in this case).

**Theorem 2.11** (Hilbert projection theorem). *Let  $\mathfrak{H}$  be a Hilbert space and  $K$  a nonempty closed convex subset. Then for every  $f \in \mathfrak{H} \setminus K$  there is a unique  $P_K(f) \in K$  such that  $\|P_K(f) - f\| = \inf_{g \in K} \|f - g\|$ . If we extend  $P_K : \mathfrak{H} \rightarrow K$  by setting  $P_K(g) = g$  for  $g \in K$  then  $P_K$  will be Lipschitz continuous:  $\|P_K(f) - P_K(g)\| \leq \|f - g\|$ ,  $f, g \in \mathfrak{H}$ .*

**Proof.** Fix  $f \in \mathfrak{H} \setminus K$  and choose a sequence  $f_n \in K$  with  $\|f_n - f\| \rightarrow d := \inf_{g \in K} \|f - g\|$ . Then applying the parallelogram law to the vectors  $f_n - f$  and  $f_m - f$  we obtain

$$\begin{aligned} \|f_n - f_m\|^2 &= 2(\|f - f_n\|^2 + \|f - f_m\|^2) - 4\|f - \tfrac{1}{2}(f_n + f_m)\|^2 \\ &\leq 2(\|f - f_n\|^2 + \|f - f_m\|^2) - 4d^2, \end{aligned}$$

which shows that  $f_n$  is Cauchy and hence converges to some point in  $K$  which we call  $P(f)$ . By construction  $\|P(f) - f\| = d$ . If there would be another point  $\tilde{P}(f)$  with the same property, we could apply the parallelogram law to  $P(f) - f$  and  $\tilde{P}(f) - f$  giving  $\|P(f) - \tilde{P}(f)\|^2 \leq 0$  and hence  $P(f)$  is uniquely defined.

Next, let  $f \in \mathfrak{H}$ ,  $g \in K$  and consider  $\tilde{g} = (1-t)P(f) + tg \in K$ ,  $t \in [0, 1]$ . Then

$$0 \geq \|f - P(f)\|^2 - \|f - \tilde{g}\|^2 = 2t\operatorname{Re}(\langle f - P(f), g - P(f) \rangle) - t^2\|g - P(f)\|^2$$

for arbitrary  $t \in [0, 1]$  shows  $\operatorname{Re}(\langle f - P(f), P(f) - g \rangle) \geq 0$ . Consequently we have  $\operatorname{Re}(\langle f - P(f), P(f) - P(g) \rangle) \geq 0$  for all  $f, g \in \mathfrak{H}$ . Now reverse to roles of  $f, g$  and add the two inequalities to obtain  $\|P(f) - P(g)\|^2 \leq \operatorname{Re}(\langle f - g, P(f) - P(g) \rangle) \leq \|f - g\|\|P(f) - P(g)\|$ . Hence Lipschitz continuity follows.  $\square$

If  $K$  is a closed subspace then this projection will of course coincide with the orthogonal projection defined before. By inspection of the proof, note that  $P_K(f)$  is alternatively characterized by  $\operatorname{Re}(\langle f - P_K(f), g - P_K(f) \rangle) \leq 0$  for all  $g \in K$ .

**Problem 2.4.** Suppose  $U : \mathfrak{H} \rightarrow \mathfrak{H}$  is unitary and  $M \subseteq \mathfrak{H}$ . Show that  $UM^\perp = (UM)^\perp$ .

**Problem 2.5.** Show that an orthogonal projection  $P_M \neq 0$  has norm one.

**Problem\* 2.6.** Suppose  $P \in \mathcal{L}(\mathfrak{H})$  satisfies

$$P^2 = P \quad \text{and} \quad \langle Pf, g \rangle = \langle f, Pg \rangle$$

and set  $M = \operatorname{Ran}(P)$ . Show

- $Pf = f$  for  $f \in M$  and  $M$  is closed,
- $g \in M^\perp$  implies  $Pg \in M^\perp$  and thus  $Pg = 0$ ,

and conclude  $P = P_M$ . In particular

$$\mathfrak{H} = \operatorname{Ker}(P) \oplus \operatorname{Ran}(P), \quad \operatorname{Ker}(P) = (\mathbb{I} - P)\mathfrak{H}, \quad \operatorname{Ran}(P) = P\mathfrak{H}.$$

### 2.3. Operators defined via forms

One of the key results about linear maps is that they are uniquely determined once we know the images of some basis vectors. In fact, the matrix elements with respect to some basis uniquely determine a linear map. Clearly this raises the question how this results extends to the infinite dimensional setting. As a first result we show that the Riesz lemma, Theorem 2.10, implies that a bounded operator  $A$  is uniquely determined by its associated sesquilinear form  $\langle g, Af \rangle$ . In fact, there is a one-to-one correspondence between bounded operators and bounded sesquilinear forms:

**Lemma 2.12.** *Suppose  $s : \mathfrak{H}_2 \times \mathfrak{H}_1 \rightarrow \mathbb{C}$  is a bounded sesquilinear form; that is,*

$$|s(g, f)| \leq C \|g\|_{\mathfrak{H}_2} \|f\|_{\mathfrak{H}_1}. \quad (2.24)$$

*Then there is a unique bounded operator  $A \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  such that*

$$s(g, f) = \langle g, Af \rangle_{\mathfrak{H}_2}. \quad (2.25)$$

*Moreover, the norm of  $A$  is given by*

$$\|A\| = \sup_{\|g\|_{\mathfrak{H}_2}=\|f\|_{\mathfrak{H}_1}=1} |\langle g, Af \rangle_{\mathfrak{H}_2}| \leq C. \quad (2.26)$$

**Proof.** For every  $f \in \mathfrak{H}_1$  we have an associated bounded linear functional  $\ell_f(g) := s(g, f)^*$  on  $\mathfrak{H}_2$ . By Theorem 2.10 there is a corresponding  $h \in \mathfrak{H}_2$  (depending on  $f$ ) such that  $\ell_f(g) = \langle h, g \rangle_{\mathfrak{H}_2}$ , that is  $s(g, f) = \langle g, h \rangle_{\mathfrak{H}_2}$  and we can define  $A$  via  $Af := h$ . It is not hard to check that  $A$  is linear and from

$$\|Af\|_{\mathfrak{H}_2}^2 = \langle Af, Af \rangle_{\mathfrak{H}_2} = s(Af, f) \leq C \|Af\|_{\mathfrak{H}_2} \|f\|_{\mathfrak{H}_1}$$

we infer  $\|Af\|_{\mathfrak{H}_2} \leq C \|f\|_{\mathfrak{H}_1}$ , which shows that  $A$  is bounded with  $\|A\| \leq C$ . Equation (2.26) is left as an exercise (Problem 2.9).  $\square$

Note that if  $\{u_k\}_{k \in K} \subseteq \mathfrak{H}_1$  and  $\{v_j\}_{j \in J} \subseteq \mathfrak{H}_2$  are some orthogonal bases, then the matrix elements  $A_{j,k} := \langle v_j, Au_k \rangle_{\mathfrak{H}_2}$  for all  $(j, k) \in J \times K$  uniquely determine  $\langle g, Af \rangle_{\mathfrak{H}_2}$  for arbitrary  $f \in \mathfrak{H}_1$ ,  $g \in \mathfrak{H}_2$  (just expand  $f, g$  with respect to these bases) and thus  $A$  by our theorem.

**Example 2.5.** Consider  $\ell^2(\mathbb{N})$  and let  $A \in \mathcal{L}(\ell^2(\mathbb{N}))$  be some bounded operator. Let  $A_{jk} = \langle \delta^j, A\delta^k \rangle$  be its matrix elements such that

$$(Aa)_j = \sum_{k=1}^{\infty} A_{jk} a_k.$$

Here the sum converges in  $\ell^2(\mathbb{N})$  and hence, in particular, for every fixed  $j$ . Moreover, choosing  $a_k^n = \alpha_n A_{jk}$  for  $k \leq n$  and  $a_k^n = 0$  for  $k > n$  with  $\alpha_n = (\sum_{j=1}^n |A_{jk}|^2)^{1/2}$  we see  $\alpha_n = |(Aa^n)_j| \leq \|A\| \|a^n\| = \|A\|$ . Thus  $\sum_{j=1}^{\infty} |A_{jk}|^2 \leq \|A\|^2$  and the sum is even absolutely convergent.  $\diamond$



Moreover, for  $A \in \mathcal{L}(\mathfrak{H})$  the polarization identity (Problem 1.21) implies that  $A$  is already uniquely determined by its quadratic form  $q_A(f) := \langle f, Af \rangle$ .

As a first application we introduce the **adjoint operator** via Lemma 2.12 as the operator associated with the sesquilinear form  $s(f, g) := \langle Af, g \rangle_{\mathfrak{H}_2}$ .

**Theorem 2.13.** *For every bounded operator  $A \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  there is a unique bounded operator  $A^* \in \mathcal{L}(\mathfrak{H}_2, \mathfrak{H}_1)$  defined via*

$$\langle f, A^*g \rangle_{\mathfrak{H}_1} = \langle Af, g \rangle_{\mathfrak{H}_2}. \quad (2.27)$$

A bounded operator  $A \in \mathcal{L}(\mathfrak{H})$  satisfying  $A^* = A$  is called **self-adjoint**. Note that  $q_{A^*}(f) = \langle Af, f \rangle = q_A(f)^*$  and hence a bounded operator is self-adjoint if and only if its quadratic form is real-valued.

**Example 2.6.** If  $\mathfrak{H} := \mathbb{C}^n$  and  $A := (a_{jk})_{1 \leq j, k \leq n}$ , then  $A^* = (a_{kj}^*)_{1 \leq j, k \leq n}$ .  $\diamond$

**Example 2.7.** If  $\mathbb{I} \in \mathcal{L}(\mathfrak{H})$  is the identity, then  $\mathbb{I}^* = \mathbb{I}$ .  $\diamond$

**Example 2.8.** Consider the linear functional  $\ell : \mathfrak{H} \rightarrow \mathbb{C}$ ,  $f \mapsto \langle g, f \rangle$ . Then by the definition  $\langle f, \ell^* \alpha \rangle = \ell(f)^* \alpha = \langle f, \alpha g \rangle$  we obtain  $\ell^* : \mathbb{C} \rightarrow \mathfrak{H}$ ,  $\alpha \mapsto \alpha g$ .  $\diamond$

**Example 2.9.** Let  $\mathfrak{H} := \ell^2(\mathbb{N})$ ,  $a \in \ell^\infty(\mathbb{N})$  and consider the multiplication operator

$$(Ab)_j := a_j b_j.$$

Then

$$\langle Ab, c \rangle = \sum_{j=1}^{\infty} (a_j b_j)^* c_j = \sum_{j=1}^{\infty} b_j^* (a_j^* c_j) = \langle b, A^* c \rangle$$

with  $(A^* c)_j = a_j^* c_j$ , that is,  $A^*$  is the multiplication operator with  $a^*$ .  $\diamond$

**Example 2.10.** Let  $\mathfrak{H} := \ell^2(\mathbb{N})$  and consider the shift operators defined via

$$(S^\pm a)_j := a_{j \pm 1}$$

with the convention that  $a_0 = 0$ . That is,  $S^-$  shifts a sequence to the right and fills up the left most place by zero and  $S^+$  shifts a sequence to the left dropping the left most place:

$$S^-(a_1, a_2, a_3, \dots) = (0, a_1, a_2, \dots), \quad S^+(a_1, a_2, a_3, \dots) = (a_2, a_3, a_4, \dots).$$

Then

$$\langle S^- a, b \rangle = \sum_{j=2}^{\infty} a_{j-1}^* b_j = \sum_{j=1}^{\infty} a_j^* b_{j+1} = \langle a, S^+ b \rangle,$$

which shows that  $(S^-)^* = S^+$ . Using symmetry of the scalar product we also get  $\langle b, S^- a \rangle = \langle S^+ b, a \rangle$ , that is,  $(S^+)^* = S^-$ .

Note that  $S^+$  is a left inverse of  $S^-$ ,  $S^+ S^- = \mathbb{I}$ , but not a right inverse as  $S^- S^+ \neq \mathbb{I}$ . This is different from the finite dimensional case, where a left inverse is also a right inverse and vice versa.  $\diamond$

**Example 2.11.** Suppose  $U \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  is unitary. Then  $U^* = U^{-1}$ . This follows from Lemma 2.12 since  $\langle f, g \rangle_{\mathfrak{H}_1} = \langle Uf, Ug \rangle_{\mathfrak{H}_2} = \langle f, U^*Ug \rangle_{\mathfrak{H}_1}$  implies  $U^*U = \mathbb{I}_{\mathfrak{H}_1}$ . Since  $U$  is bijective we can multiply this last equation from the right with  $U^{-1}$  to obtain the claim. Of course this calculation shows that the converse is also true, that is  $U \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  is unitary if and only if  $U^* = U^{-1}$ .  $\diamond$

A few simple properties of taking adjoints are listed below.

**Lemma 2.14.** *Let  $A, B \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$ ,  $C \in \mathcal{L}(\mathfrak{H}_2, \mathfrak{H}_3)$ , and  $\alpha \in \mathbb{C}$ . Then*

- (i)  $(A + B)^* = A^* + B^*$ ,  $(\alpha A)^* = \alpha^* A^*$ ,
- (ii)  $A^{**} = A$ ,
- (iii)  $(CA)^* = A^* C^*$ ,
- (iv)  $\|A^*\| = \|A\|$  and  $\|A\|^2 = \|A^* A\| = \|AA^*\|$ .

**Proof.** (i) is obvious. (ii) follows from  $\langle g, A^{**}f \rangle_{\mathfrak{H}_2} = \langle A^*g, f \rangle_{\mathfrak{H}_1} = \langle g, Af \rangle_{\mathfrak{H}_2}$ . (iii) follows from  $\langle g, (CA)f \rangle_{\mathfrak{H}_3} = \langle C^*g, Af \rangle_{\mathfrak{H}_2} = \langle A^*C^*g, f \rangle_{\mathfrak{H}_1}$ . (iv) follows using (2.26) from

$$\begin{aligned} \|A^*\| &= \sup_{\|f\|_{\mathfrak{H}_1}=\|g\|_{\mathfrak{H}_2}=1} |\langle f, A^*g \rangle_{\mathfrak{H}_1}| = \sup_{\|f\|_{\mathfrak{H}_1}=\|g\|_{\mathfrak{H}_2}=1} |\langle Af, g \rangle_{\mathfrak{H}_2}| \\ &= \sup_{\|f\|_{\mathfrak{H}_1}=\|g\|_{\mathfrak{H}_2}=1} |\langle g, Af \rangle_{\mathfrak{H}_2}| = \|A\| \end{aligned}$$

and

$$\begin{aligned} \|A^* A\| &= \sup_{\|f\|_{\mathfrak{H}_1}=\|g\|_{\mathfrak{H}_2}=1} |\langle f, A^* Ag \rangle_{\mathfrak{H}_1}| = \sup_{\|f\|_{\mathfrak{H}_1}=\|g\|_{\mathfrak{H}_2}=1} |\langle Af, Ag \rangle_{\mathfrak{H}_2}| \\ &= \sup_{\|f\|_{\mathfrak{H}_1}=1} \|Af\|^2 = \|A\|^2, \end{aligned}$$

where we have used that  $|\langle Af, Ag \rangle_{\mathfrak{H}_2}|$  attains its maximum when  $Af$  and  $Ag$  are parallel (compare Theorem 1.5).  $\square$

Note that  $\|A\| = \|A^*\|$  implies that taking adjoints is a continuous operation. For later use also note that (Problem 2.11)

$$\text{Ker}(A^*) = \text{Ran}(A)^\perp. \quad (2.28)$$

For the remainder of this section we restrict to the case of one Hilbert space. A sesquilinear form  $s : \mathfrak{H} \times \mathfrak{H} \rightarrow \mathbb{C}$  is called nonnegative if  $s(f, f) \geq 0$  and we will call  $A \in \mathcal{L}(\mathfrak{H})$  **nonnegative**,  $A \geq 0$ , if its associated sesquilinear form is. We will write  $A \geq B$  if  $A - B \geq 0$ . Observe that nonnegative operators are self-adjoint (as their quadratic forms are real-valued — here it is important that the underlying space is complex; in case of a real space a nonnegative form is required to be symmetric).

**Example 2.12.** For any operator  $A$  the operators  $A^*A$  and  $AA^*$  are both nonnegative. In fact  $\langle f, A^*Af \rangle = \langle Af, Af \rangle = \|Af\|^2 \geq 0$  and similarly  $\langle f, AA^*f \rangle = \|A^*f\|^2 \geq 0$ .  $\diamond$

**Lemma 2.15.** Suppose  $A \in \mathcal{L}(\mathfrak{H})$  satisfies  $\|Af\| \geq \varepsilon\|f\|^2$  for some  $\varepsilon > 0$ . Then  $\text{Ran}(A)$  is closed and  $A : \mathfrak{H} \rightarrow \text{Ran}(A)$  is a bijection with bounded inverse,  $\|A^{-1}\| \leq \frac{1}{\varepsilon}$ . If we have the stronger condition  $|\langle f, Af \rangle| \geq \varepsilon\|f\|^2$ , then  $\text{Ran}(A) = \mathfrak{H}$ .

**Proof.** Since  $Af = 0$  implies  $f = 0$  our operator is injective and thus for every  $g \in \text{Ran}(A)$  there is a unique  $f = A^{-1}g$ . Moreover, by  $\|A^{-1}g\| = \|f\| \leq \varepsilon^{-1}\|Af\| = \varepsilon^{-1}\|g\|$  the operator  $A^{-1}$  is bounded. So if  $g_n \in \text{Ran}(A)$  converges to some  $g \in \mathfrak{H}$ , then  $f_n = A^{-1}g_n$  converges to some  $f$ . Taking limits in  $g_n = Af_n$  shows that  $g = Af$  is in the range of  $A$ , that is, the range of  $A$  is closed.

By  $\varepsilon\|f\|^2 \leq |\langle f, Af \rangle| \leq \|f\|\|Af\|$  the second condition implies the first. To show that  $\text{Ran}(A) = \mathfrak{H}$  we pick  $h \in \text{Ran}(A)^\perp$ . Then  $0 = \langle h, Ah \rangle \geq \varepsilon\|h\|^2$  shows  $h = 0$  and thus  $\text{Ran}(A)^\perp = \{0\}$ .  $\square$

Combining the last two results we obtain the famous Lax–Milgram theorem which plays an important role in theory of elliptic partial differential equations.

**Theorem 2.16** (Lax–Milgram). Let  $s : \mathfrak{H} \times \mathfrak{H} \rightarrow \mathbb{C}$  be a sesquilinear form which is

- bounded,  $|s(f, g)| \leq C\|f\|\|g\|$ , and
- satisfies  $|s(f, f)| \geq \varepsilon\|f\|^2$  for some  $\varepsilon > 0$ .

Then for every  $g \in \mathfrak{H}$  there is a unique  $f \in \mathfrak{H}$  such that

$$s(h, f) = \langle h, g \rangle, \quad \forall h \in \mathfrak{H}. \quad (2.29)$$

Moreover,  $\|f\| \leq \frac{1}{\varepsilon}\|g\|$ .

**Proof.** Let  $A$  be the operator associated with  $s$ . Then  $A$  is a bijection and  $f = A^{-1}g$ .  $\square$

Instead of the second condition one frequently requires that  $s$  is **coercive**, that is,

$$\text{Re}(s(f, f)) \geq \varepsilon\|f\|^2. \quad (2.30)$$

This condition is clearly weaker. Note that (2.29) can also be phrased as a minimizing problem if  $s$  is nonnegative — Problem 2.13.

**Example 2.13.** Consider  $\mathfrak{H} = \ell^2(\mathbb{N})$  and introduce the operator

$$(Aa)_j := -a_{j+1} + 2a_j - a_{j-1}$$

which is a discrete version of a second derivative (discrete one-dimensional Laplace operator). Here we use the convention  $a_0 = 0$ , that is,  $(Aa)_1 = -a_2 + 2a_1$ . In terms of the shift operators  $S^\pm$  we can write

$$A = -S^+ + 2 - S^- = (S^+ - 1)(S^- - 1)$$

and using  $(S^\pm)^* = S^\mp$  we obtain

$$s_A(a, b) = \langle (S^- - 1)a, (S^- - 1)b \rangle = \sum_{j=1}^{\infty} (a_{j-1} - a_j)^*(b_{j-1} - b_j).$$

In particular, this shows  $A \geq 0$ . Moreover, we have  $|s_A(a, b)| \leq 4\|a\|_2\|b\|_2$  or equivalently  $\|A\| \leq 4$ .

Next, let

$$(Qa)_j = q_j a_j$$

for some sequence  $q \in \ell^\infty(\mathbb{N})$ . Then

$$s_Q(a, b) = \sum_{j=1}^{\infty} q_j a_j^* b_j$$

and  $|s_Q(a, b)| \leq \|q\|_\infty \|a\|_2 \|b\|_2$  or equivalently  $\|Q\| \leq \|q\|_\infty$ . If in addition  $q_j \geq \varepsilon > 0$ , then  $s_{A+Q}(a, b) = s_A(a, b) + s_Q(a, b)$  satisfies the assumptions of the Lax–Milgram theorem and

$$(A + Q)a = b$$

has a unique solution  $a = (A + Q)^{-1}b$  for every given  $b \in \ell^2(\mathbb{Z})$ . Moreover, since  $(A + Q)^{-1}$  is bounded, this solution depends continuously on  $b$ .  $\diamond$

**Problem\* 2.7.** Let  $\mathfrak{H}_1, \mathfrak{H}_2$  be Hilbert spaces and let  $u \in \mathfrak{H}_1, v \in \mathfrak{H}_2$ . Show that the operator

$$Af := \langle u, f \rangle v$$

is bounded and compute its norm. Compute the adjoint of  $A$ .

**Problem 2.8.** Show that under the assumptions of Problem 1.39 one has  $f(A)^* = f^\#(A^*)$  where  $f^\#(z) = f(z^*)^*$ .

**Problem\* 2.9.** Prove (2.26). (Hint: Use  $\|f\| = \sup_{\|g\|=1} |\langle g, f \rangle|$  — compare Theorem 1.5.)

**Problem 2.10.** Suppose  $A \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  has a bounded inverse  $A^{-1} \in \mathcal{L}(\mathfrak{H}_2, \mathfrak{H}_1)$ . Show  $(A^{-1})^* = (A^*)^{-1}$ .

**Problem\* 2.11.** Show (2.28).

**Problem\* 2.12.** Show that every operator  $A \in \mathcal{L}(\mathfrak{H})$  can be written as the linear combination of two self-adjoint operators  $\operatorname{Re}(A) := \frac{1}{2}(A + A^*)$  and  $\operatorname{Im}(A) := \frac{1}{2i}(A - A^*)$ . Moreover, every self-adjoint operator can be written as a linear combination of two unitary operators. (Hint: For the last part consider  $f_\pm(z) = z \pm i\sqrt{1 - z^2}$  and Problems 1.39, 2.8.)

**Problem 2.13** (Abstract Dirichlet problem). *Show that the solution of (2.29) is also the unique minimizer of*

$$h \mapsto \operatorname{Re} \left( \frac{1}{2} s(h, h) - \langle h, g \rangle \right)$$

if  $s$  is nonnegative with  $s(w, w) \geq \varepsilon \|w\|^2$  for all  $w \in \mathfrak{H}$ .

## 2.4. Orthogonal sums and tensor products

Given two Hilbert spaces  $\mathfrak{H}_1$  and  $\mathfrak{H}_2$ , we define their **orthogonal sum**  $\mathfrak{H}_1 \oplus \mathfrak{H}_2$  to be the set of all pairs  $(f_1, f_2) \in \mathfrak{H}_1 \times \mathfrak{H}_2$  together with the scalar product

$$\langle (g_1, g_2), (f_1, f_2) \rangle := \langle g_1, f_1 \rangle_{\mathfrak{H}_1} + \langle g_2, f_2 \rangle_{\mathfrak{H}_2}. \quad (2.31)$$

It is left as an exercise to verify that  $\mathfrak{H}_1 \oplus \mathfrak{H}_2$  is again a Hilbert space. Moreover,  $\mathfrak{H}_1$  can be identified with  $\{(f_1, 0) | f_1 \in \mathfrak{H}_1\}$ , and we can regard  $\mathfrak{H}_1$  as a subspace of  $\mathfrak{H}_1 \oplus \mathfrak{H}_2$ , and similarly for  $\mathfrak{H}_2$ . With this convention we have  $\mathfrak{H}_1^\perp = \mathfrak{H}_2$ . It is also customary to write  $f_1 \oplus f_2$  instead of  $(f_1, f_2)$ . In the same way we can define the orthogonal sum  $\bigoplus_{j=1}^n \mathfrak{H}_j$  of any finite number of Hilbert spaces.

**Example 2.14.** For example we have  $\bigoplus_{j=1}^n \mathbb{C} = \mathbb{C}^n$  and hence we will write  $\bigoplus_{j=1}^n \mathfrak{H} =: \mathfrak{H}^n$ .  $\diamond$

More generally, let  $\mathfrak{H}_j$ ,  $j \in \mathbb{N}$ , be a countable collection of Hilbert spaces and define

$$\bigoplus_{j=1}^{\infty} \mathfrak{H}_j := \left\{ \bigoplus_{j=1}^{\infty} f_j \mid f_j \in \mathfrak{H}_j, \sum_{j=1}^{\infty} \|f_j\|_{\mathfrak{H}_j}^2 < \infty \right\}, \quad (2.32)$$

which becomes a Hilbert space with the scalar product

$$\left\langle \bigoplus_{j=1}^{\infty} g_j, \bigoplus_{j=1}^{\infty} f_j \right\rangle := \sum_{j=1}^{\infty} \langle g_j, f_j \rangle_{\mathfrak{H}_j}. \quad (2.33)$$

**Example 2.15.**  $\bigoplus_{j=1}^{\infty} \mathbb{C} = \ell^2(\mathbb{N})$ .  $\diamond$

Similarly, if  $\mathfrak{H}$  and  $\tilde{\mathfrak{H}}$  are two Hilbert spaces, we define their tensor product as follows: The elements should be products  $f \otimes \tilde{f}$  of elements  $f \in \mathfrak{H}$  and  $\tilde{f} \in \tilde{\mathfrak{H}}$ . Hence we start with the set of all finite linear combinations of elements of  $\mathfrak{H} \times \tilde{\mathfrak{H}}$

$$\mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}}) := \left\{ \sum_{j=1}^n \alpha_j (f_j, \tilde{f}_j) \mid (f_j, \tilde{f}_j) \in \mathfrak{H} \times \tilde{\mathfrak{H}}, \alpha_j \in \mathbb{C} \right\}. \quad (2.34)$$

Since we want  $(f_1 + f_2) \otimes \tilde{f} = f_1 \otimes \tilde{f} + f_2 \otimes \tilde{f}$ ,  $f \otimes (\tilde{f}_1 + \tilde{f}_2) = f \otimes \tilde{f}_1 + f \otimes \tilde{f}_2$ , and  $(\alpha f) \otimes \tilde{f} = f \otimes (\alpha \tilde{f}) = \alpha(f \otimes \tilde{f})$  we consider  $\mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}})/\mathcal{N}(\mathfrak{H}, \tilde{\mathfrak{H}})$ , where

$$\mathcal{N}(\mathfrak{H}, \tilde{\mathfrak{H}}) := \text{span}\left\{ \sum_{j,k=1}^n \alpha_j \beta_k (f_j, \tilde{f}_k) - \left( \sum_{j=1}^n \alpha_j f_j, \sum_{k=1}^n \beta_k \tilde{f}_k \right) \right\} \quad (2.35)$$

and write  $f \otimes \tilde{f}$  for the equivalence class of  $(f, \tilde{f})$ . By construction, every element in this quotient space is a linear combination of elements of the type  $f \otimes \tilde{f}$ .

Next, we want to define a scalar product such that

$$\langle f \otimes \tilde{f}, g \otimes \tilde{g} \rangle = \langle f, g \rangle_{\mathfrak{H}} \langle \tilde{f}, \tilde{g} \rangle_{\tilde{\mathfrak{H}}} \quad (2.36)$$

holds. To this end we set

$$s\left(\sum_{j=1}^n \alpha_j (f_j, \tilde{f}_j), \sum_{k=1}^n \beta_k (g_k, \tilde{g}_k)\right) = \sum_{j,k=1}^n \alpha_j^* \beta_k \langle f_j, g_k \rangle_{\mathfrak{H}} \langle \tilde{f}_j, \tilde{g}_k \rangle_{\tilde{\mathfrak{H}}}, \quad (2.37)$$

which is a symmetric sesquilinear form on  $\mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}})$ . Moreover, one verifies that  $s(f, g) = 0$  for arbitrary  $f \in \mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}})$  and  $g \in \mathcal{N}(\mathfrak{H}, \tilde{\mathfrak{H}})$  and thus

$$\left\langle \sum_{j=1}^n \alpha_j f_j \otimes \tilde{f}_j, \sum_{k=1}^n \beta_k g_k \otimes \tilde{g}_k \right\rangle = \sum_{j,k=1}^n \alpha_j^* \beta_k \langle f_j, g_k \rangle_{\mathfrak{H}} \langle \tilde{f}_j, \tilde{g}_k \rangle_{\tilde{\mathfrak{H}}} \quad (2.38)$$

is a symmetric sesquilinear form on  $\mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}})/\mathcal{N}(\mathfrak{H}, \tilde{\mathfrak{H}})$ . To show that this is in fact a scalar product, we need to ensure positivity. Let  $f = \sum_i \alpha_i f_i \otimes \tilde{f}_i \neq 0$  and pick orthonormal bases  $u_j, \tilde{u}_k$  for  $\text{span}\{f_i\}, \text{span}\{\tilde{f}_i\}$ , respectively. Then

$$f = \sum_{j,k} \alpha_{jk} u_j \otimes \tilde{u}_k, \quad \alpha_{jk} = \sum_i \alpha_i \langle u_j, f_i \rangle_{\mathfrak{H}} \langle \tilde{u}_k, \tilde{f}_i \rangle_{\tilde{\mathfrak{H}}} \quad (2.39)$$

and we compute

$$\langle f, f \rangle = \sum_{j,k} |\alpha_{jk}|^2 > 0. \quad (2.40)$$

The completion of  $\mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}})/\mathcal{N}(\mathfrak{H}, \tilde{\mathfrak{H}})$  with respect to the induced norm is called the **tensor product**  $\mathfrak{H} \otimes \tilde{\mathfrak{H}}$  of  $\mathfrak{H}$  and  $\tilde{\mathfrak{H}}$ .

**Lemma 2.17.** *If  $u_j, \tilde{u}_k$  are orthonormal bases for  $\mathfrak{H}, \tilde{\mathfrak{H}}$ , respectively, then  $u_j \otimes \tilde{u}_k$  is an orthonormal basis for  $\mathfrak{H} \otimes \tilde{\mathfrak{H}}$ .*

**Proof.** That  $u_j \otimes \tilde{u}_k$  is an orthonormal set is immediate from (2.36). Moreover, since  $\text{span}\{u_j\}, \text{span}\{\tilde{u}_k\}$  are dense in  $\mathfrak{H}, \tilde{\mathfrak{H}}$ , respectively, it is easy to see that  $u_j \otimes \tilde{u}_k$  is dense in  $\mathcal{F}(\mathfrak{H}, \tilde{\mathfrak{H}})/\mathcal{N}(\mathfrak{H}, \tilde{\mathfrak{H}})$ . But the latter is dense in  $\mathfrak{H} \otimes \tilde{\mathfrak{H}}$ .  $\square$

Note that this in particular implies  $\dim(\mathfrak{H} \otimes \tilde{\mathfrak{H}}) = \dim(\mathfrak{H}) \dim(\tilde{\mathfrak{H}})$ .

**Example 2.16.** We have  $\mathfrak{H} \otimes \mathbb{C}^n = \mathfrak{H}^n$ .  $\diamond$

**Example 2.17.** We have  $\ell^2(\mathbb{N}) \otimes \ell^2(\mathbb{N}) = \ell^2(\mathbb{N} \times \mathbb{N})$  by virtue of the identification  $(a_{jk}) \mapsto \sum_{jk} a_{jk} \delta^j \otimes \delta^k$  where  $\delta^j$  is the standard basis for  $\ell^2(\mathbb{N})$ . In fact, this follows from the previous lemma as in the proof of Theorem 2.6.  $\diamond$

It is straightforward to extend the tensor product to any finite number of Hilbert spaces. We even note

$$\left(\bigoplus_{j=1}^{\infty} \mathfrak{H}_j\right) \otimes \mathfrak{H} = \bigoplus_{j=1}^{\infty} (\mathfrak{H}_j \otimes \mathfrak{H}), \quad (2.41)$$

where equality has to be understood in the sense that both spaces are unitarily equivalent by virtue of the identification

$$\left(\sum_{j=1}^{\infty} f_j\right) \otimes f = \sum_{j=1}^{\infty} f_j \otimes f. \quad (2.42)$$

**Problem 2.14.** Show that  $f \otimes \tilde{f} = 0$  if and only if  $f = 0$  or  $\tilde{f} = 0$ .

**Problem 2.15.** We have  $f \otimes \tilde{f} = g \otimes \tilde{g} \neq 0$  if and only if there is some  $\alpha \in \mathbb{C} \setminus \{0\}$  such that  $f = \alpha g$  and  $\tilde{f} = \alpha^{-1} \tilde{g}$ .

**Problem\* 2.16.** Show (2.41).

## 2.5. Applications to Fourier series

We have already encountered the Fourier sine series during our treatment of the heat equation in Section 1.1. Given an integrable function  $f$  we can define its **Fourier series**

$$S(f)(x) := \frac{a_0}{2} + \sum_{k \in \mathbb{N}} (a_k \cos(kx) + b_k \sin(kx)), \quad (2.43)$$

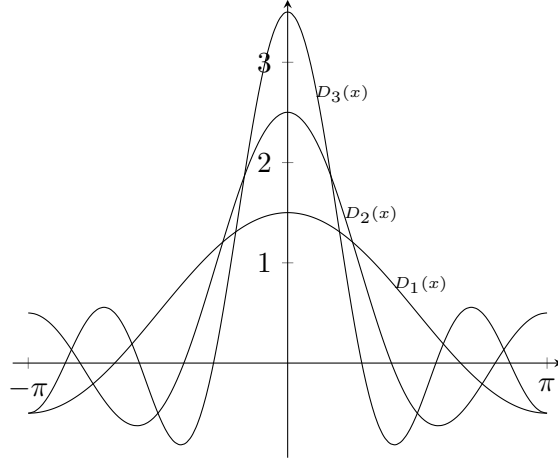
where the corresponding Fourier coefficients are given by

$$a_k := \frac{1}{\pi} \int_{-\pi}^{\pi} \cos(kx) f(x) dx, \quad b_k := \frac{1}{\pi} \int_{-\pi}^{\pi} \sin(kx) f(x) dx. \quad (2.44)$$

At this point (2.43) is just a formal expression and it was (and to some extent still is) a fundamental question in mathematics to understand in what sense the above series converges. For example, does it converge at a given point (e.g. at every point of continuity of  $f$ ) or when does it converge uniformly? We will give some first answers in the present section and then come back later to this when we have further tools at our disposal.

For our purpose the complex form

$$S(f)(x) = \sum_{k \in \mathbb{Z}} \hat{f}_k e^{ikx}, \quad \hat{f}_k := \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-iky} f(y) dy \quad (2.45)$$



**Figure 2.1.** The Dirichlet kernels  $D_1$ ,  $D_2$ , and  $D_3$

will be more convenient. The connection is given via  $\hat{f}_{\pm k} = \frac{a_k \mp ib_k}{2}$ ,  $k \in \mathbb{N}_0$  (with the convention  $b_0 = 0$ ). In this case the  $n$ 'th partial sum can be written as

$$S_n(f)(x) := \sum_{k=-n}^n \hat{f}_k e^{ikx} = \frac{1}{2\pi} \int_{-\pi}^{\pi} D_n(x-y) f(y) dy, \quad (2.46)$$

where

$$D_n(x) = \sum_{k=-n}^n e^{ikx} = \frac{\sin((n+1/2)x)}{\sin(x/2)} \quad (2.47)$$

is known as the **Dirichlet kernel** (to obtain the second form observe that the left-hand side is a geometric series). Note that  $D_n(-x) = D_n(x)$  and that  $|D_n(x)|$  has a global maximum  $D_n(0) = 2n+1$  at  $x=0$ . Moreover, by  $S_n(1) = 1$  we see that  $\int_{-\pi}^{\pi} D_n(x) dx = 1$ .

Since

$$\int_{-\pi}^{\pi} e^{-ikx} e^{ilx} dx = 2\pi \delta_{k,l} \quad (2.48)$$

the functions  $e_k(x) := (2\pi)^{-1/2} e^{ikx}$  are orthonormal in  $L^2(-\pi, \pi)$  and hence the Fourier series is just the expansion with respect to this orthogonal set. Hence we obtain

**Theorem 2.18.** *For every square integrable function  $f \in L^2(-\pi, \pi)$ , the Fourier coefficients  $\hat{f}_k$  are square summable*

$$\sum_{k \in \mathbb{Z}} |\hat{f}_k|^2 = \frac{1}{2\pi} \int_{-\pi}^{\pi} |f(x)|^2 dx \quad (2.49)$$

and the Fourier series converges to  $f$  in the sense of  $L^2$ . Moreover, this is a continuous bijection between  $L^2(-\pi, \pi)$  and  $\ell^2(\mathbb{Z})$ .



**Proof.** To show this theorem it suffices to show that the functions  $e_k$  form a basis. This will follow from Theorem 2.21 below (see the discussion after this theorem). It will also follow as a special case of Theorem 3.11 below (see the examples after this theorem) as well as from the Stone–Weierstraß theorem — Problem 2.21.  $\square$

This gives a satisfactory answer in the Hilbert space  $L^2(-\pi, \pi)$  but does not answer the question about pointwise or uniform convergence. The latter will be the case if the Fourier coefficients are summable. First of all we note that for integrable functions the Fourier coefficients will at least tend to zero.

**Lemma 2.19** (Riemann–Lebesgue lemma). *Suppose  $f \in L^1(-\pi, \pi)$ , then the Fourier coefficients  $\hat{f}_k$  converge to zero as  $|k| \rightarrow \infty$ .*

**Proof.** By our previous theorem this holds for continuous functions. But the map  $f \rightarrow \hat{f}$  is bounded from  $C[-\pi, \pi] \subset L^1(-\pi, \pi)$  to  $c_0(\mathbb{Z})$  (the sequences vanishing as  $|k| \rightarrow \infty$ ) since  $|\hat{f}_k| \leq (2\pi)^{-1} \|f\|_1$  and there is a unique extension to all of  $L^1(-\pi, \pi)$ .  $\square$

It turns out that this result is best possible in general and we cannot say more without additional assumptions on  $f$ . For example, if  $f$  is periodic of period  $2\pi$  and continuously differentiable, then integration by parts shows

$$\hat{f}_k = \frac{1}{2\pi i k} \int_{-\pi}^{\pi} e^{-ikx} f'(x) dx. \quad (2.50)$$

Then, since both  $k^{-1}$  and the Fourier coefficients of  $f'$  are square summable, we conclude that  $\hat{f}$  is absolutely summable and hence the Fourier series converges uniformly. So we have a simple sufficient criterion for summability of the Fourier coefficients, but can we do better? Of course continuity of  $f$  is a necessary condition for absolute summability but this alone will not even be enough for pointwise convergence as we will see in Example 4.3. Moreover, continuity will not tell us more about the decay of the Fourier coefficients than what we already know in the integrable case from the Riemann–Lebesgue lemma (see Example 4.4).

A few improvements are easy: (2.50) holds for any class of functions for which integration by parts holds, e.g., piecewise continuously differentiable functions or, slightly more general, absolutely continuous functions (cf. Lemma 4.30 from [48]) provided one assumes that the derivative is square integrable. However, for an arbitrary absolutely continuous function the Fourier coefficients might not be absolutely summable: For an absolutely continuous function  $f$  we have a derivative which is integrable (Theorem 4.29 from [48]) and hence the above formula combined with the Riemann–Lebesgue lemma implies  $\hat{f}_k = o(\frac{1}{k})$ . But on the other hand we

can choose an absolutely summable sequence  $c_k$  which does not obey this asymptotic requirement, say  $c_k = \frac{1}{k}$  for  $k = l^2$  and  $c_k = 0$  else. Then

$$f(x) := \sum_{k \in \mathbb{Z}} c_k e^{ikx} = \sum_{l \in \mathbb{N}} \frac{1}{l^2} e^{il^2 x} \quad (2.51)$$

is a function with absolutely summable Fourier coefficients  $\hat{f}_k = c_k$  (by uniform convergence we can interchange summation and integration) but which is not absolutely continuous. There are further criteria for absolute summability of the Fourier coefficients, but no simple necessary and sufficient one. A particularly simple sufficient one is:

**Theorem 2.20** (Bernstein). *Suppose that  $f \in C_{per}^{0,\gamma}[-\pi, \pi]$  is Hölder continuous (cf. (1.67)) of exponent  $\gamma > \frac{1}{2}$ , then*

$$\sum_{k \in \mathbb{Z}} |\hat{f}_k| \leq C_\gamma \|f\|_{0,\gamma}.$$

**Proof.** The proof starts with the observation that the Fourier coefficients of  $f_\delta(x) := f(x - \delta)$  are  $\hat{f}_k = e^{-ik\delta} \hat{f}_k$ . Now for  $\delta := \frac{2\pi}{3} 2^{-m}$  and  $2^m \leq |k| < 2^{m+1}$  we have  $|e^{ik\delta} - 1|^2 \geq 3$  implying

$$\begin{aligned} \sum_{2^m \leq |k| < 2^{m+1}} |\hat{f}_k|^2 &\leq \frac{1}{3} \sum_k |e^{ik\delta} - 1|^2 |\hat{f}_k|^2 = \frac{1}{6\pi} \int_{-\pi}^{\pi} |f_\delta(x) - f(x)|^2 dx \\ &\leq \frac{1}{3} [f]_\gamma^2 \delta^{2\gamma} \end{aligned}$$

Now the sum on the left has  $2 \cdot 2^m$  terms and hence Cauchy–Schwarz implies

$$\sum_{2^m \leq |k| < 2^{m+1}} |\hat{f}_k| \leq \frac{2^{(m+1)/2}}{\sqrt{3}} [f]_\gamma \delta^\gamma = \sqrt{\frac{2}{3}} \left(\frac{2\pi}{3}\right)^\gamma 2^{(1/2-\gamma)m} [f]_\gamma.$$

Summing over  $m$  shows

$$\sum_{k \neq 0} |\hat{f}_k| \leq C_\gamma [f]_\gamma$$

provided  $\gamma > \frac{1}{2}$  and establishes the claim since  $|\hat{f}_0| \leq \|f\|_\infty$ .  $\square$

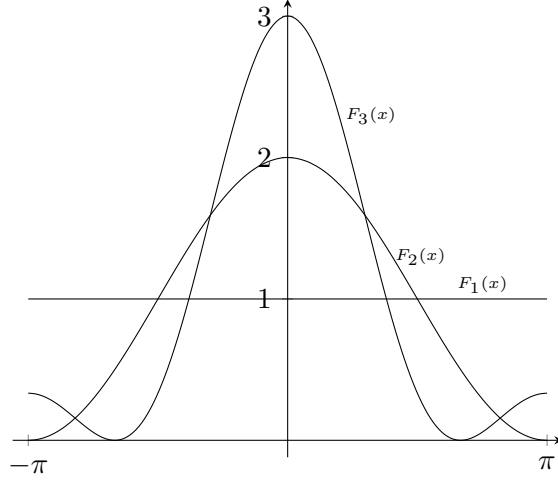
Note however, that the situation looks much brighter if one looks at mean values

$$\bar{S}_n(f)(x) := \frac{1}{n} \sum_{k=0}^{n-1} S_k(f)(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} F_n(x-y) f(y) dy, \quad (2.52)$$

where

$$F_n(x) = \frac{1}{n} \sum_{k=0}^{n-1} D_k(x) = \frac{1}{n} \left( \frac{\sin(nx/2)}{\sin(x/2)} \right)^2 \quad (2.53)$$

is the **Fejér kernel**. To see the second form we use the closed form for the



**Figure 2.2.** The Fejér kernels  $F_1$ ,  $F_2$ , and  $F_3$

Dirichlet kernel to obtain

$$\begin{aligned} nF_n(x) &= \sum_{k=0}^{n-1} \frac{\sin((k+1/2)x)}{\sin(x/2)} = \frac{1}{\sin(x/2)} \operatorname{Im} \sum_{k=0}^{n-1} e^{i(k+1/2)x} \\ &= \frac{1}{\sin(x/2)} \operatorname{Im} \left( e^{ix/2} \frac{e^{inx} - 1}{e^{ix} - 1} \right) = \frac{1 - \cos(nx)}{2 \sin(x/2)^2} = \left( \frac{\sin(nx/2)}{\sin(x/2)} \right)^2. \end{aligned}$$

The main difference to the Dirichlet kernel is positivity:  $F_n(x) \geq 0$ . Of course the property  $\int_{-\pi}^{\pi} F_n(x) dx = 1$  is inherited from the Dirichlet kernel.

**Theorem 2.21** (Fejér). *Suppose  $f$  is continuous and periodic with period  $2\pi$ . Then  $\bar{S}_n(f) \rightarrow f$  uniformly.*

**Proof.** Let us set  $F_n = 0$  outside  $[-\pi, \pi]$ . Then  $F_n(x) \leq \frac{1}{n \sin(\delta/2)^2}$  for  $\delta \leq |x| \leq \pi$  implies that a straightforward adaption of Lemma 1.2 to the periodic case is applicable.  $\square$

In particular, this shows that the functions  $\{e_k\}_{k \in \mathbb{Z}}$  are total in  $C_{per}[-\pi, \pi]$  (continuous periodic functions) and hence also in  $L^p(-\pi, \pi)$  for  $1 \leq p < \infty$  (Problem 2.20).

Note that for a given continuous function  $f$  this result shows that if  $S_n(f)(x)$  converges, then it must converge to  $\bar{S}_n(f)(x) = f(x)$ . We also remark that one can extend this result (see Lemma 3.21 from [48]) to show that for  $f \in L^p(-\pi, \pi)$ ,  $1 \leq p < \infty$ , one has  $\bar{S}_n(f) \rightarrow f$  in the sense of  $L^p$ . As a consequence note that the Fourier coefficients uniquely determine  $f$  for integrable  $f$  (for square integrable  $f$  this follows from Theorem 2.18).

Finally, we look at pointwise convergence.

**Theorem 2.22.** *Suppose*

$$\frac{f(x) - f(x_0)}{x - x_0} \quad (2.54)$$

*is integrable (e.g.  $f$  is Hölder continuous), then*

$$\lim_{m, n \rightarrow \infty} \sum_{k=-m}^n \hat{f}_k e^{ikx_0} = f(x_0). \quad (2.55)$$

**Proof.** Without loss of generality we can assume  $x_0 = 0$  (by shifting  $x \rightarrow x - x_0$  modulo  $2\pi$  implying  $\hat{f}_k \rightarrow e^{-ikx_0} \hat{f}_k$ ) and  $f(x_0) = 0$  (by linearity since the claim is trivial for constant functions). Then by assumption

$$g(x) := \frac{f(x)}{e^{ix} - 1}$$

is integrable and  $f(x) = (e^{ix} - 1)g(x)$  implies  $\hat{f}_k = \hat{g}_{k-1} - \hat{g}_k$  and hence

$$\sum_{k=-m}^n \hat{f}_k = \hat{g}_{-m-1} - \hat{g}_n.$$

Now the claim follows from the Riemann–Lebesgue lemma.  $\square$

If we look at symmetric partial sums  $S_n(f)$  we can do even better.

**Corollary 2.23** (Dirichlet–Dini criterion). *Suppose there is some  $\alpha$  such that*

$$\frac{f(x_0 + x) + f(x_0 - x) - 2\alpha}{x}$$

*is integrable. Then  $S_n(f)(x_0) \rightarrow \alpha$ .*

**Proof.** Without loss of generality we can assume  $x_0 = 0$ . Now observe (since  $D_n(-x) = D_n(x)$ )  $S_n(f)(0) = \alpha + S_n(g)(0)$ , where  $g(x) := \frac{1}{2}(f(x) + f(-x)) - \alpha$  and apply the previous result.  $\square$

**Problem 2.17.** *Compute the Fourier series of  $D_n$  and  $F_n$ .*

**Problem 2.18.** *Show  $|D_n(x)| \leq \min(2n + 1, \frac{\pi}{|x|})$  and  $F_n(x) \leq \min(n, \frac{\pi^2}{nx^2})$ .*

**Problem 2.19.** *Show that if  $f \in C_{per}^{0,\gamma}[-\pi, \pi]$  is Hölder continuous (cf. (1.67)), then*

$$|\hat{f}_k| \leq \frac{[f]_\gamma}{2} \left( \frac{\pi}{|k|} \right)^\gamma, \quad k \neq 0.$$

(Hint: What changes if you replace  $e^{-iky}$  by  $e^{-ik(y+\pi/k)}$  in (2.45)? Now make a change of variables  $y \rightarrow y - \pi/k$  in the integral.)

**Problem 2.20.** *Show that  $C_{per}[-\pi, \pi]$  is dense in  $L^p(-\pi, \pi)$  for  $1 \leq p < \infty$ .*

**Problem 2.21.** Show that the functions  $e_k(x) := \frac{1}{\sqrt{2\pi}}e^{ikx}$ ,  $k \in \mathbb{Z}$ , form an orthonormal basis for  $\mathfrak{H} = L^2(-\pi, \pi)$ . (Hint: Start with  $K = [-\pi, \pi]$  where  $-\pi$  and  $\pi$  are identified and use the Stone–Weierstraß theorem.)

# Compact operators

Typically, linear operators are much more difficult to analyze than matrices and many new phenomena appear which are not present in the finite dimensional case. So we have to be modest and slowly work our way up. A class of operators which still preserves some of the nice properties of matrices is the class of compact operators to be discussed in this chapter.

## 3.1. Compact operators

A linear operator  $A : X \rightarrow Y$  defined between normed spaces  $X, Y$  is called **compact** if every sequence  $Af_n$  has a convergent subsequence whenever  $f_n$  is bounded. Equivalently (cf. Corollary B.20),  $A$  is compact if it maps bounded sets to relatively compact ones. The set of all compact operators is denoted by  $\mathcal{K}(X, Y)$ . If  $X = Y$  we will just write  $\mathcal{K}(X) := \mathcal{K}(X, X)$  as usual.

**Example 3.1.** Every linear map between finite dimensional spaces is compact by the Bolzano–Weierstraß theorem. Slightly more general, a bounded operator is compact if its range is finite dimensional.  $\diamond$

The following elementary properties of compact operators are left as an exercise (Problem 3.1):

**Theorem 3.1.** *Let  $X, Y$ , and  $Z$  be normed spaces. Every compact linear operator is bounded,  $\mathcal{K}(X, Y) \subseteq \mathcal{L}(X, Y)$ . Linear combinations of compact operators are compact, that is,  $\mathcal{K}(X, Y)$  is a subspace of  $\mathcal{L}(X, Y)$ . Moreover, the product of a bounded and a compact operator is again compact, that is,  $A \in \mathcal{L}(X, Y)$ ,  $B \in \mathcal{K}(Y, Z)$  or  $A \in \mathcal{K}(X, Y)$ ,  $B \in \mathcal{L}(Y, Z)$  implies  $BA \in \mathcal{K}(X, Z)$ .*

In particular, the set of compact operators  $\mathcal{K}(X)$  is an ideal of the set of bounded operators. Moreover, if  $X$  is a Banach space this ideal is even closed:

**Theorem 3.2.** *Suppose  $X$  is a normed and  $Y$  a Banach space. Let  $A_n \in \mathcal{K}(X, Y)$  be a convergent sequence of compact operators. Then the limit  $A$  is again compact.*

**Proof.** Let  $f_j^0$  be a bounded sequence. Choose a subsequence  $f_j^1$  such that  $A_1 f_j^1$  converges. From  $f_j^1$  choose another subsequence  $f_j^2$  such that  $A_2 f_j^2$  converges and so on. Since there might be nothing left from  $f_j^n$  as  $n \rightarrow \infty$ , we consider the diagonal sequence  $f_j := f_j^j$ . By construction,  $f_j$  is a subsequence of  $f_j^n$  for  $j \geq n$  and hence  $A_n f_j$  is Cauchy for every fixed  $n$ . Now

$$\begin{aligned} \|A f_j - A f_k\| &= \|(A - A_n)(f_j - f_k) + A_n(f_j - f_k)\| \\ &\leq \|A - A_n\| \|f_j - f_k\| + \|A_n f_j - A_n f_k\| \end{aligned}$$

shows that  $A f_j$  is Cauchy since the first term can be made arbitrary small by choosing  $n$  large and the second by the Cauchy property of  $A_n f_j$ .  $\square$

**Example 3.2.** Let  $X := \ell^p(\mathbb{N})$  and consider the operator

$$(Qa)_j := q_j a_j$$

for some sequence  $q = (q_j)_{j=1}^\infty \in c_0(\mathbb{N})$  converging to zero. Let  $Q_n$  be associated with  $q_j^n = q_j$  for  $j \leq n$  and  $q_j^n = 0$  for  $j > n$ . Then the range of  $Q_n$  is finite dimensional and hence  $Q_n$  is compact. Moreover, by  $\|Q_n - Q\| = \sup_{j>n} |q_j|$  we see  $Q_n \rightarrow Q$  and thus  $Q$  is also compact by the previous theorem.  $\diamond$

**Example 3.3.** Let  $X := C^1[0, 1]$ ,  $Y := C[0, 1]$  (cf. Problem 1.35) then the embedding  $X \hookrightarrow Y$  is compact. Indeed, a bounded sequence in  $X$  has both the functions and the derivatives uniformly bounded. Hence by the mean value theorem the functions are equicontinuous and hence there is a uniformly convergent subsequence by the Arzelà–Ascoli theorem (Theorem 1.13). Of course the same conclusion holds if we take  $X := C^{0,\gamma}[0, 1]$  to be Hölder continuous functions (cf. Theorem 1.21).  $\diamond$

If  $A : X \rightarrow Y$  is a bounded operator there is a unique extension  $\bar{A} : \bar{X} \rightarrow \bar{Y}$  to the completion by Theorem 1.16. Moreover, if  $A \in \mathcal{K}(X, Y)$ , then  $A \in \mathcal{K}(X, \bar{Y})$  is immediate. That we also have  $\bar{A} \in \mathcal{K}(\bar{X}, \bar{Y})$  will follow from the next lemma. In particular, it suffices to verify compactness on a dense set.

**Lemma 3.3.** *Let  $X, Y$  be normed spaces and  $A \in \mathcal{K}(X, Y)$ . Let  $\bar{X}, \bar{Y}$  be the completion of  $X, Y$ , respectively. Then  $\bar{A} \in \mathcal{K}(\bar{X}, \bar{Y})$ , where  $\bar{A}$  is the unique extension of  $A$  (cf. Theorem 1.16).*

**Proof.** Let  $f_n \in \overline{X}$  be a given bounded sequence. We need to show that  $\overline{A}f_n$  has a convergent subsequence. Pick  $f_n^j \in X$  such that  $\|f_n^j - f_n\| \leq \frac{1}{j}$  and by compactness of  $A$  we can assume that  $Af_n^j \rightarrow g$ . But then  $\|\overline{A}f_n - g\| \leq \|A\|\|f_n - f_n^j\| + \|Af_n^j - g\|$  shows that  $\overline{A}f_n \rightarrow g$ .  $\square$

One of the most important examples of compact operators are integral operators. The proof will be based on the Arzelà–Ascoli theorem (Theorem 1.13).

**Lemma 3.4.** *Let  $X := C([a, b])$  or  $X := \mathcal{L}_{cont}^2(a, b)$ . The integral operator  $K : X \rightarrow X$  defined by*

$$(Kf)(x) := \int_a^b K(x, y)f(y)dy, \quad (3.1)$$

where  $K(x, y) \in C([a, b] \times [a, b])$ , is compact.

**Proof.** First of all note that  $K(., .)$  is continuous on  $[a, b] \times [a, b]$  and hence uniformly continuous. In particular, for every  $\varepsilon > 0$  we can find a  $\delta > 0$  such that  $|K(y, t) - K(x, t)| \leq \varepsilon$  for any  $t \in [a, b]$  whenever  $|y - x| \leq \delta$ . Moreover,  $\|K\|_\infty = \sup_{x, y \in [a, b]} |K(x, y)| < \infty$ .

We begin with the case  $X := \mathcal{L}_{cont}^2(a, b)$ . Let  $g := Kf$ . Then

$$|g(x)| \leq \int_a^b |K(x, t)| |f(t)| dt \leq \|K\|_\infty \int_a^b |f(t)| dt \leq \|K\|_\infty \|1\| \|f\|,$$

where we have used Cauchy–Schwarz in the last step (note that  $\|1\| = \sqrt{b - a}$ ). Similarly,

$$\begin{aligned} |g(x) - g(y)| &\leq \int_a^b |K(y, t) - K(x, t)| |f(t)| dt \\ &\leq \varepsilon \int_a^b |f(t)| dt \leq \varepsilon \|1\| \|f\|, \end{aligned}$$

whenever  $|y - x| \leq \delta$ . Hence, if  $f_n(x)$  is a bounded sequence in  $\mathcal{L}_{cont}^2(a, b)$ , then  $g_n := Kf_n$  is bounded and equicontinuous and hence has a uniformly convergent subsequence by the Arzelà–Ascoli theorem (Theorem 1.13). But a uniformly convergent sequence is also convergent in the norm induced by the scalar product. Therefore  $K$  is compact.

The case  $X := C([a, b])$  follows by the same argument upon observing  $\int_a^b |f(t)| dt \leq (b - a) \|f\|_\infty$ .  $\square$

Compact operators share many similarities with (finite) matrices as we will see in the next section.

**Problem\* 3.1.** *Show Theorem 3.1.*



**Problem 3.2.** Is the left shift  $(a_1, a_2, a_3, \dots) \mapsto (a_2, a_3, \dots)$  compact in  $\ell^2(\mathbb{N})$ ?

**Problem 3.3** (Ehrling's lemma). Let  $X$ ,  $Y$ , and  $Z$  be Banach spaces. Assume  $X$  is compactly embedded into  $Y$  and  $Y$  is continuously embedded into  $Z$ . Show that for every  $\varepsilon > 0$  there exists some  $C(\varepsilon)$  such that

$$\|x\|_Y \leq \varepsilon \|x\|_X + C(\varepsilon) \|x\|_Z.$$

**Problem 3.4.** Is the operator  $\frac{d}{dx} : C^k[0, 1] \rightarrow C[0, 1]$  compact for  $k = 1, 2$ ? (Hint: Problem 1.31 and Example 3.3.)

**Problem 3.5.** Is the multiplication operator  $M_t : C^k[0, 1] \rightarrow C[0, 1]$  with  $M_t f(t) = t f(t)$  compact for  $k = 0, 1$ ? (Hint: Problem 1.31 and Example 3.3.)

**Problem 3.6.** Let  $X := C([a, b])$  or  $X := \mathcal{L}_{\text{cont}}^2(a, b)$ . Show that the integral operator  $K : X \rightarrow X$  defined by

$$(Kf)(x) := \int_a^x K(x, y) f(y) dy,$$

where  $K(x, y) \in C([a, b] \times [a, b])$ , is compact.

**Problem\* 3.7.** Show that the adjoint of the integral operator  $K$  on  $\mathcal{L}_{\text{cont}}^2(a, b)$  from Lemma 3.4 is the integral operator with kernel  $K(y, x)^*$ :

$$(K^* f)(x) = \int_a^b K(y, x)^* f(y) dy.$$

(Hint: Fubini.)

### 3.2. The spectral theorem for compact symmetric operators

Let  $\mathfrak{H}$  be an inner product space. A linear operator  $A : \mathfrak{D}(A) \subseteq \mathfrak{H} \rightarrow \mathfrak{H}$  is called **symmetric** if its domain is dense and if

$$\langle g, Af \rangle = \langle Ag, f \rangle \quad f, g \in \mathfrak{D}(A). \quad (3.2)$$

If  $A$  is bounded (with  $\mathfrak{D}(A) = \mathfrak{H}$ ), then  $A$  is symmetric precisely if  $A = A^*$ , that is, if  $A$  is **self-adjoint**. However, for unbounded operators there is a subtle but important difference between symmetry and self-adjointness.

A number  $z \in \mathbb{C}$  is called **eigenvalue** of  $A$  if there is a nonzero vector  $u \in \mathfrak{D}(A)$  such that

$$Au = zu. \quad (3.3)$$

The vector  $u$  is called a corresponding **eigenvector** in this case. The set of all eigenvectors corresponding to  $z$  is called the **eigenspace**

$$\text{Ker}(A - z) \quad (3.4)$$

corresponding to  $z$ . Here we have used the shorthand notation  $A - z$  for  $A - z\mathbb{I}$ . An eigenvalue is called (geometrically) **simple** if there is only one linearly independent eigenvector.

**Example 3.4.** Let  $\mathfrak{H} := \ell^2(\mathbb{N})$  and consider the shift operators  $(S^\pm a)_j := a_{j\pm 1}$  (with  $a_0 := 0$ ). Suppose  $z \in \mathbb{C}$  is an eigenvalue, then the corresponding eigenvector  $u$  must satisfy  $u_{j\pm 1} = zu_j$ . For  $S^-$  the special case  $j = 1$  gives  $0 = u_0 = zu_1$ . So either  $z = 0$  and  $u = 0$  or  $z \neq 0$  and again  $u = 0$ . Hence there are no eigenvalues. For  $S^+$  we get  $u_j = z^j u_1$  and this will give an element in  $\ell^2(\mathbb{N})$  if and only if  $|z| < 1$ . Hence  $z$  with  $|z| < 1$  is an eigenvalue. All these eigenvalues are simple.  $\diamond$

**Example 3.5.** Let  $\mathfrak{H} := \ell^2(\mathbb{N})$  and consider the multiplication operator  $(Qa)_j := q_j a_j$  with a bounded sequence  $q \in \ell^\infty(\mathbb{N})$ . Suppose  $z \in \mathbb{C}$  is an eigenvalue, then the corresponding eigenvector  $u$  must satisfy  $(q_j - z)u_j = 0$ . Hence every value  $q_j$  is an eigenvalue with corresponding eigenvector  $u = \delta^j$ . If there is only one  $j$  with  $z = q_j$  the eigenvalue is simple (otherwise the numbers of independent eigenvectors equals the number of times  $z$  appears in the sequence  $q$ ). If  $z$  is different from all entries of the sequence then  $u = 0$  and  $z$  is no eigenvalue.  $\diamond$

Note that in the last example  $Q$  will be self-adjoint if and only if  $q$  is real-valued and hence if and only if all eigenvalues are real-valued. Moreover, the corresponding eigenfunctions are orthogonal. This has nothing to do with the simple structure of our operator and is in fact always true.

**Theorem 3.5.** *Let  $A$  be symmetric. Then all eigenvalues are real and eigenvectors corresponding to different eigenvalues are orthogonal.*

**Proof.** Suppose  $\lambda$  is an eigenvalue with corresponding normalized eigenvector  $u$ . Then  $\lambda = \langle u, Au \rangle = \langle Au, u \rangle = \lambda^*$ , which shows that  $\lambda$  is real. Furthermore, if  $Au_j = \lambda_j u_j$ ,  $j = 1, 2$ , we have

$$(\lambda_1 - \lambda_2)\langle u_1, u_2 \rangle = \langle Au_1, u_2 \rangle - \langle u_1, Au_2 \rangle = 0$$

finishing the proof.  $\square$

Note that while eigenvectors corresponding to the same eigenvalue  $\lambda$  will in general not automatically be orthogonal, we can of course replace each set of eigenvectors corresponding to  $\lambda$  by an set of orthonormal eigenvectors having the same linear span (e.g. using Gram–Schmidt orthogonalization).

**Example 3.6.** Let  $\mathfrak{H} := \ell^2(\mathbb{N})$  and consider the Jacobi operator  $J := \frac{1}{2}(S^+ + S^-)$ :

$$(Jc)_j := \frac{1}{2}(c_{j+1} + c_{j-1})$$

with the convention  $c_0 = 0$ . Recall that  $J^* = J$ . If we look for an eigenvalue  $Ju = zu$ , we need to solve the corresponding recursion  $u_{j+1} = 2zu_j - u_{j-1}$

starting from  $u_0 = 0$  (our convention) and  $u_1 = 1$  (normalization). Like an ordinary differential equation, a linear recursion relation with constant coefficients can be solved by an exponential ansatz  $u_j = k^j$  which leads to the characteristic polynomial  $k^2 = 2zk - 1$ . This gives two linearly independent solutions and our requirements lead us to

$$u_j(z) = \frac{k^j - k^{-j}}{k - k^{-1}}, \quad k = z - \sqrt{z^2 - 1}.$$

Note that  $k^{-1} = z + \sqrt{z^2 - 1}$  and in the case  $k = z = \pm 1$  the above expression has to be understood as its limit  $u_j(\pm 1) = (\pm 1)^{j+1}j$ . In fact,  $U_j(z) := u_{j+1}(z)$  are polynomials of degree  $j$  known as **Chebyshev polynomials** of the second kind.

Now for  $z \in \mathbb{R} \setminus [-1, 1]$  we have  $|k| < 1$  and  $u_j$  explodes exponentially. For  $z \in [-1, 1]$  we have  $|k| = 1$  and hence we can write  $k = e^{i\kappa}$  with  $\kappa \in \mathbb{R}$ . Thus  $u_j = \frac{\sin(\kappa j)}{\sin(\kappa)}$  is oscillating. So for no value of  $z \in \mathbb{R}$  our potential eigenvector  $u$  is square summable and thus  $J$  has no eigenvalues.  $\diamond$

The previous example shows that in the infinite dimensional case symmetry is not enough to guarantee existence of even a single eigenvalue. In order to always get this, we will need an extra condition. In fact, we will see that compactness provides a suitable extra condition to obtain an orthonormal basis of eigenfunctions. The crucial step is to prove existence of one eigenvalue, the rest then follows as in the finite dimensional case.

**Theorem 3.6.** *Let  $\mathfrak{H}$  be an inner product space. A symmetric compact operator  $A$  has an eigenvalue  $\alpha_1$  which satisfies  $|\alpha_1| = \|A\|$ .*

**Proof.** We set  $\alpha := \|A\|$  and assume  $\alpha \neq 0$  (i.e.,  $A \neq 0$ ) without loss of generality. Since

$$\|A\|^2 = \sup_{f:\|f\|=1} \|Af\|^2 = \sup_{f:\|f\|=1} \langle Af, Af \rangle = \sup_{f:\|f\|=1} \langle f, A^2 f \rangle$$

there exists a normalized sequence  $u_n$  such that

$$\lim_{n \rightarrow \infty} \langle u_n, A^2 u_n \rangle = \alpha^2.$$

Since  $A$  is compact, it is no restriction to assume that  $A^2 u_n$  converges, say  $\lim_{n \rightarrow \infty} A^2 u_n = \alpha^2 u$ . Now

$$\begin{aligned} \|(A^2 - \alpha^2)u_n\|^2 &= \|A^2 u_n\|^2 - 2\alpha^2 \langle u_n, A^2 u_n \rangle + \alpha^4 \\ &\leq 2\alpha^2(\alpha^2 - \langle u_n, A^2 u_n \rangle) \end{aligned}$$

(where we have used  $\|A^2 u_n\| \leq \|A\| \|Au_n\| \leq \|A\|^2 \|u_n\| = \alpha^2$ ) implies  $\lim_{n \rightarrow \infty} (A^2 u_n - \alpha^2 u_n) = 0$  and hence  $\lim_{n \rightarrow \infty} u_n = u$ . In addition,  $u$  is a normalized eigenvector of  $A^2$  since  $(A^2 - \alpha^2)u = 0$ . Factorizing this last equation according to  $(A - \alpha)u = v$  and  $(A + \alpha)v = 0$  shows that either

$v \neq 0$  is an eigenvector corresponding to  $-\alpha$  or  $v = 0$  and hence  $u \neq 0$  is an eigenvector corresponding to  $\alpha$ .  $\square$

Note that for a bounded operator  $A$ , there cannot be an eigenvalue with absolute value larger than  $\|A\|$ , that is, the set of eigenvalues is bounded by  $\|A\|$  (Problem 3.8).

Now consider a symmetric compact operator  $A$  with eigenvalue  $\alpha_1$  (as above) and corresponding normalized eigenvector  $u_1$ . Setting

$$\mathfrak{H}_1 := \{u_1\}^\perp = \{f \in \mathfrak{H} | \langle u_1, f \rangle = 0\} \quad (3.5)$$

we can restrict  $A$  to  $\mathfrak{H}_1$  since  $f \in \mathfrak{H}_1$  implies

$$\langle u_1, Af \rangle = \langle Au_1, f \rangle = \alpha_1 \langle u_1, f \rangle = 0 \quad (3.6)$$

and hence  $Af \in \mathfrak{H}_1$ . Denoting this restriction by  $A_1$ , it is not hard to see that  $A_1$  is again a symmetric compact operator. Hence we can apply Theorem 3.6 iteratively to obtain a sequence of eigenvalues  $\alpha_j$  with corresponding normalized eigenvectors  $u_j$ . Moreover, by construction,  $u_j$  is orthogonal to all  $u_k$  with  $k < j$  and hence the eigenvectors  $\{u_j\}$  form an orthonormal set. By construction we also have  $|\alpha_j| = \|A_j\| \leq \|A_{j-1}\| = |\alpha_{j-1}|$ . This procedure will not stop unless  $\mathfrak{H}$  is finite dimensional. However, note that  $\alpha_j = 0$  for  $j \geq n$  might happen if  $A_n = 0$ .

**Theorem 3.7** (Hilbert–Schmidt; Spectral theorem for compact symmetric operators). *Suppose  $\mathfrak{H}$  is an infinite dimensional Hilbert space and  $A : \mathfrak{H} \rightarrow \mathfrak{H}$  is a compact symmetric operator. Then there exists a sequence of real eigenvalues  $\alpha_j$  converging to 0. The corresponding normalized eigenvectors  $u_j$  form an orthonormal set and every  $f \in \mathfrak{H}$  can be written as*

$$f = \sum_{j=1}^{\infty} \langle u_j, f \rangle u_j + h, \quad (3.7)$$

where  $h$  is in the kernel of  $A$ , that is,  $Ah = 0$ .

*In particular, if 0 is not an eigenvalue, then the eigenvectors form an orthonormal basis (in addition,  $\mathfrak{H}$  need not be complete in this case).*

**Proof.** Existence of the eigenvalues  $\alpha_j$  and the corresponding eigenvectors  $u_j$  has already been established. Since the sequence  $|\alpha_j|$  is decreasing it has a limit  $\varepsilon \geq 0$  and we have  $|\alpha_j| \geq \varepsilon$ . If this limit is nonzero, then  $v_j = \alpha_j^{-1} u_j$  is a bounded sequence ( $\|v_j\| \leq \frac{1}{\varepsilon}$ ) for which  $Av_j$  has no convergent subsequence since  $\|Av_j - Av_k\|^2 = \|u_j - u_k\|^2 = 2$ , a contradiction.

Next, setting

$$f_n := \sum_{j=1}^n \langle u_j, f \rangle u_j,$$

we have

$$\|A(f - f_n)\| \leq |\alpha_{n+1}| \|f - f_n\| \leq |\alpha_{n+1}| \|f\|$$

since  $f - f_n \in \mathfrak{H}_n$  and  $\|A_n\| = |\alpha_{n+1}|$ . Letting  $n \rightarrow \infty$  shows  $A(f_\infty - f) = 0$  proving (3.7). Finally, note that without completeness  $f_\infty$  might not be well-defined unless  $h = 0$ .  $\square$

By applying  $A$  to (3.7) we obtain the following canonical form of compact symmetric operators.

**Corollary 3.8.** *Every compact symmetric operator  $A$  can be written as*

$$Af = \sum_{j=1}^N \alpha_j \langle u_j, f \rangle u_j, \quad (3.8)$$

where  $\alpha_j$  are the nonzero eigenvalues with corresponding eigenvectors  $u_j$  from the previous theorem.

**Remark:** There are two cases where our procedure might fail to construct an orthonormal basis of eigenvectors. One case is where there is an infinite number of nonzero eigenvalues. In this case  $\alpha_n$  never reaches 0 and all eigenvectors corresponding to 0 are missed. In the other case, 0 is reached, but there might not be a countable basis and hence again some of the eigenvectors corresponding to 0 are missed. In any case, by adding vectors from the kernel (which are automatically eigenvectors), one can always extend the eigenvectors  $u_j$  to an orthonormal basis of eigenvectors.

**Corollary 3.9.** *Every compact symmetric operator  $A$  has an associated orthonormal basis of eigenvectors  $\{u_j\}_{j \in J}$ . The corresponding unitary map  $U : \mathfrak{H} \rightarrow \ell^2(J)$ ,  $f \mapsto \{\langle u_j, f \rangle\}_{j \in J}$  diagonalizes  $A$  in the sense that  $UAU^{-1}$  is the operator which multiplies each basis vector  $\delta^j = Uu_j$  by the corresponding eigenvalue  $\alpha_j$ .*

**Example 3.7.** Let  $a, b \in c_0(\mathbb{N})$  be real-valued sequences and consider the operator

$$(Jc)_j := a_j c_{j+1} + b_j c_j + a_{j-1} c_{j-1}.$$

If  $A, B$  denote the multiplication operators by the sequences  $a, b$ , respectively, then we already know that  $A$  and  $B$  are compact. Moreover, using the shift operators  $S^\pm$  we can write

$$J = AS^+ + B + S^-A,$$

which shows that  $J$  is self-adjoint since  $A^* = A$ ,  $B^* = B$ , and  $(S^\pm)^* = S^\mp$ . Hence we can conclude that  $J$  has a countable number of eigenvalues converging to zero and a corresponding orthonormal basis of eigenvectors.  $\diamond$

In particular, in the new picture it is easy to define functions of our operator (thus extending the functional calculus from Problem 1.39). To this end set  $\Sigma := \{\alpha_j\}_{j \in J}$  and denote by  $B(K)$  the Banach algebra of bounded functions  $F : K \rightarrow \mathbb{C}$  together with the sup norm.

**Corollary 3.10** (Functional calculus). *Let  $A$  be a compact symmetric operator with associated orthonormal basis of eigenvectors  $\{u_j\}_{j \in J}$  and corresponding eigenvalues  $\{\alpha_j\}_{j \in J}$ . Suppose  $F \in B(\Sigma)$ , then*

$$F(A)f = \sum_{j \in J} F(\alpha_j) \langle u_j, f \rangle u_j \quad (3.9)$$

*defines a continuous algebra homomorphism from the Banach algebra  $B(\Sigma)$  to the algebra  $\mathcal{L}(\mathfrak{H})$  with  $1(A) = \mathbb{I}$  and  $\mathbb{I}(A) = A$ . Moreover  $F(A)^* = F^*(A)$ , where  $F^*$  is the function which takes complex conjugate values.*

**Proof.** This is straightforward to check for multiplication operators in  $\ell^2(J)$  and hence the result follows by the previous corollary.  $\square$

In many applications  $F$  will be given by a function on  $\mathbb{R}$  (or at least on  $[-\|A\|, \|A\|]$ ) and, since only the values  $F(\alpha_j)$  are used, two functions which agree on all eigenvalues will give the same result.

As a brief application we will say a few words about general spectral theory for bounded operators  $A \in \mathcal{L}(X)$  in a Banach space  $X$ . In the finite dimensional case, the spectrum is precisely the set of eigenvalues. In the infinite dimensional case one defines the **spectrum** as

$$\sigma(A) := \mathbb{C} \setminus \{z \in \mathbb{C} \mid \exists (A - z)^{-1} \in \mathcal{L}(X)\}. \quad (3.10)$$

It is important to emphasize that the inverse is required to exist as a bounded operator. Hence there are several ways in which this can fail: First of all,  $A - z$  could not be injective. In this case  $z$  is an eigenvalue and thus all eigenvalues belong to the spectrum. Secondly, it could not be surjective. And finally, even if it is bijective it could be unbounded. However, it will follow from the open mapping theorem that this last case cannot happen for a bounded operator. The inverse of  $A - z$  for  $z \in \mathbb{C} \setminus \sigma(A)$  is known as the **resolvent** of  $A$  and plays a crucial role in spectral theory. Using Problem 1.38 one can show that the complement of the spectrum is open, and hence the spectrum is closed. Since we will discuss this in detail in Chapter 5, we will not pursue this here but only look at our special case of symmetric compact operators.

To compute the inverse of  $A - z$  we will use the functional calculus and consider  $F(\alpha) = \frac{1}{\alpha - z}$ . Of course this function is unbounded on  $\mathbb{R}$  but if  $z$  is neither an eigenvalue nor zero it is bounded on  $\Sigma$  and hence satisfies our

requirements. Then

$$R_A(z)f := \sum_{j \in J} \frac{1}{\alpha_j - z} \langle u_j, f \rangle u_j \quad (3.11)$$

satisfies  $(A - z)R_A(z) = R_A(z)(A - z) = \mathbb{I}$ , that is,  $R_A(z) = (A - z)^{-1} \in \mathcal{L}(\mathfrak{H})$ . Of course, if  $z$  is an eigenvalue, then the above formula breaks down. However, in the infinite dimensional case it also breaks down if  $z = 0$  even if 0 is not an eigenvalue! In this case the above definition will still give an operator which is the inverse of  $A - z$ , however, since the sequence  $\alpha_j^{-1}$  is unbounded, so will be the corresponding multiplication operator in  $\ell^2(J)$  and the sum in (3.11) will only converge if  $\{\alpha_j^{-1} \langle u_j, f \rangle\}_{j \in J} \in \ell^2(J)$ . So in the infinite dimensional case 0 is in the spectrum even if it is not an eigenvalue. In particular,

$$\sigma(A) = \overline{\{\alpha_j\}_{j \in J}}. \quad (3.12)$$

Moreover, if we use  $\frac{1}{\alpha_j - z} = \frac{\alpha_j}{z(\alpha_j - z)} - \frac{1}{z}$  we can rewrite this as

$$R_A(z)f = \frac{1}{z} \left( \sum_{j=1}^N \frac{\alpha_j}{\alpha_j - z} \langle u_j, f \rangle u_j - f \right)$$

where it suffices to take the sum over all nonzero eigenvalues.

This is all we need and it remains to apply these results to Sturm–Liouville operators.

**Problem 3.8.** *Show that if  $A \in \mathcal{L}(\mathfrak{H})$ , then every eigenvalue  $\alpha$  satisfies  $|\alpha| \leq \|A\|$ .*

**Problem 3.9.** *Find the eigenvalues and eigenfunctions of the integral operator  $K \in \mathcal{L}(\mathcal{L}_{cont}^2(0, 1))$  given by*

$$(Kf)(x) := \int_0^1 u(x)v(y)f(y)dy,$$

where  $u, v \in C([0, 1])$  are some given continuous functions.

**Problem 3.10.** *Find the eigenvalues and eigenfunctions of the integral operator  $K \in \mathcal{L}(\mathcal{L}_{cont}^2(0, 1))$  given by*

$$(Kf)(x) := 2 \int_0^1 (2xy - x - y + 1)f(y)dy.$$

### 3.3. Applications to Sturm–Liouville operators

Now, after all this hard work, we can show that our Sturm–Liouville operator

$$L := -\frac{d^2}{dx^2} + q(x), \quad (3.13)$$

where  $q$  is continuous and real, defined on

$$\mathfrak{D}(L) := \{f \in C^2[0, 1] | f(0) = f(1) = 0\} \subset \mathcal{L}_{cont}^2(0, 1), \quad (3.14)$$

has an orthonormal basis of eigenfunctions.

The corresponding eigenvalue equation  $Lu = zu$  explicitly reads

$$-u''(x) + q(x)u(x) = zu(x). \quad (3.15)$$

It is a second order homogeneous linear ordinary differential equation and hence has two linearly independent solutions. In particular, specifying two initial conditions, e.g.  $u(0) = 0, u'(0) = 1$  determines the solution uniquely. Hence, if we require  $u(0) = 0$ , the solution is determined up to a multiple and consequently the additional requirement  $u(1) = 0$  cannot be satisfied by a nontrivial solution in general. However, there might be some  $z \in \mathbb{C}$  for which the solution corresponding to the initial conditions  $u(0) = 0, u'(0) = 1$  happens to satisfy  $u(1) = 0$  and these are precisely the eigenvalues we are looking for.

Note that the fact that  $\mathcal{L}_{cont}^2(0, 1)$  is not complete causes no problems since we can always replace it by its completion  $\mathfrak{H} = L^2(0, 1)$ . A thorough investigation of this completion will be given later, at this point this is not essential.

We first verify that  $L$  is symmetric:

$$\begin{aligned} \langle f, Lg \rangle &= \int_0^1 f(x)^* (-g''(x) + q(x)g(x)) dx \\ &= \int_0^1 f'(x)^* g'(x) dx + \int_0^1 f(x)^* q(x)g(x) dx \\ &= \int_0^1 -f''(x)^* g(x) dx + \int_0^1 f(x)^* q(x)g(x) dx \\ &= \langle Lf, g \rangle. \end{aligned} \quad (3.16)$$

Here we have used integration by parts twice (the boundary terms vanish due to our boundary conditions  $f(0) = f(1) = 0$  and  $g(0) = g(1) = 0$ ).

Of course we want to apply Theorem 3.7 and for this we would need to show that  $L$  is compact. But this task is bound to fail, since  $L$  is not even bounded (see Example 1.18)!

So here comes the trick: If  $L$  is unbounded its inverse  $L^{-1}$  might still be bounded. Moreover,  $L^{-1}$  might even be compact and this is the case here! Since  $L$  might not be injective (0 might be an eigenvalue), we consider  $R_L(z) := (L - z)^{-1}$ ,  $z \in \mathbb{C}$ , which is also known as the **resolvent** of  $L$ .

In order to compute the resolvent, we need to solve the inhomogeneous equation  $(L - z)f = g$ . This can be done using the variation of constants formula from ordinary differential equations which determines the solution



up to an arbitrary solution of the homogeneous equation. This homogeneous equation has to be chosen such that  $f \in \mathfrak{D}(L)$ , that is, such that  $f(0) = f(1) = 0$ .

Define

$$f(x) := \frac{u_+(z, x)}{W(z)} \left( \int_0^x u_-(z, t) g(t) dt \right) + \frac{u_-(z, x)}{W(z)} \left( \int_x^1 u_+(z, t) g(t) dt \right), \quad (3.17)$$

where  $u_{\pm}(z, x)$  are the solutions of the homogeneous differential equation  $-u_{\pm}''(z, x) + (q(x) - z)u_{\pm}(z, x) = 0$  satisfying the initial conditions  $u_-(z, 0) = 0$ ,  $u'_-(z, 0) = 1$  respectively  $u_+(z, 1) = 0$ ,  $u'_+(z, 1) = 1$  and

$$W(z) := W(u_+(z), u_-(z)) = u'_-(z, x)u_+(z, x) - u_-(z, x)u'_+(z, x) \quad (3.18)$$

is the Wronski determinant, which is independent of  $x$  (check this!).

Then clearly  $f(0) = 0$  since  $u_-(z, 0) = 0$  and similarly  $f(1) = 0$  since  $u_+(z, 1) = 0$ . Furthermore,  $f$  is differentiable and a straightforward computation verifies

$$f'(x) = \frac{u'_+(z, x)}{W(z)} \left( \int_0^x u_-(z, t) g(t) dt \right) + \frac{u'_-(z, x)}{W(z)} \left( \int_x^1 u_+(z, t) g(t) dt \right). \quad (3.19)$$

Thus we can differentiate once more giving

$$\begin{aligned} f''(x) &= \frac{u''_+(z, x)}{W(z)} \left( \int_0^x u_-(z, t) g(t) dt \right) \\ &\quad + \frac{u''_-(z, x)}{W(z)} \left( \int_x^1 u_+(z, t) g(t) dt \right) - g(x) \\ &= (q(x) - z)f(x) - g(x). \end{aligned} \quad (3.20)$$

In summary,  $f$  is in the domain of  $L$  and satisfies  $(L - z)f = g$ .

Note that  $z$  is an eigenvalue if and only if  $W(z) = 0$ . In fact, in this case  $u_+(z, x)$  and  $u_-(z, x)$  are linearly dependent and hence  $u_+(z, x) = c u_-(z, x)$  with  $c = u'_+(z, 0)$ . Evaluating this identity at  $x = 0$  shows  $u_+(z, 0) = c u_-(z, 0) = 0$  that  $u_+(z, x)$  satisfies both boundary conditions and is thus an eigenfunction.

Introducing the **Green function**

$$G(z, x, t) := \frac{1}{W(u_+(z), u_-(z))} \begin{cases} u_+(z, x)u_-(z, t), & x \geq t, \\ u_+(z, t)u_-(z, x), & x \leq t, \end{cases} \quad (3.21)$$

we see that  $(L - z)^{-1}$  is given by

$$(L - z)^{-1}g(x) = \int_0^1 G(z, x, t)g(t)dt. \quad (3.22)$$

Moreover, from  $G(z, x, t) = G(z, t, x)$  it follows that  $(L - z)^{-1}$  is symmetric for  $z \in \mathbb{R}$  (Problem 3.11) and from Lemma 3.4 it follows that it is compact. Hence Theorem 3.7 applies to  $(L - z)^{-1}$  once we show that we can find a real  $z$  which is not an eigenvalue.

**Theorem 3.11.** *The Sturm–Liouville operator  $L$  has a countable number of discrete and simple eigenvalues  $E_n$  which accumulate only at  $\infty$ . They are bounded from below and can hence be ordered as follows:*

$$\min_{x \in [0,1]} q(x) < E_0 < E_1 < \cdots. \quad (3.23)$$

*The corresponding normalized eigenfunctions  $u_n$  form an orthonormal basis for  $\mathcal{L}_{cont}^2(0, 1)$ , that is, every  $f \in \mathcal{L}_{cont}^2(0, 1)$  can be written as*

$$f(x) = \sum_{n=0}^{\infty} \langle u_n, f \rangle u_n(x). \quad (3.24)$$

*Moreover, for  $f \in \mathfrak{D}(L)$  this series is absolutely uniformly convergent.*

**Proof.** If  $E_j$  is an eigenvalue with corresponding normalized eigenfunction  $u_j$  we have

$$E_j = \langle u_j, Lu_j \rangle = \int_0^1 (|u_j'(x)|^2 + q(x)|u_j(x)|^2)dx < \min_{x \in [0,1]} q(x) \quad (3.25)$$

where we have used integration by parts as in (3.16). Note that equality could only occur if  $u_j$  is constant, which is incompatible with our boundary conditions. Hence the eigenvalues are bounded from below.

Now pick a value  $\lambda \in \mathbb{R}$  such that  $R_L(\lambda)$  exists ( $\lambda < \min_{x \in [0,1]} q(x)$  say). By Lemma 3.4  $R_L(\lambda)$  is compact and by Lemma 3.3 this remains true if we replace  $\mathcal{L}_{cont}^2(0, 1)$  by its completion. By Theorem 3.7 there are eigenvalues  $\alpha_n$  of  $R_L(\lambda)$  with corresponding eigenfunctions  $u_n$ . Moreover,  $R_L(\lambda)u_n = \alpha_n u_n$  is equivalent to  $Lu_n = (\lambda + \frac{1}{\alpha_n})u_n$ , which shows that  $E_n = \lambda + \frac{1}{\alpha_n}$  are eigenvalues of  $L$  with corresponding eigenfunctions  $u_n$ . Now everything follows from Theorem 3.7 except that the eigenvalues are simple. To show this, observe that if  $u_n$  and  $v_n$  are two different eigenfunctions corresponding to  $E_n$ , then  $u_n(0) = v_n(0) = 0$  implies  $W(u_n, v_n) = 0$  and hence  $u_n$  and  $v_n$  are linearly dependent.

To show that (3.24) converges uniformly if  $f \in \mathfrak{D}(L)$  we begin by writing  $f = R_L(\lambda)g$ ,  $g \in \mathcal{L}_{cont}^2(0, 1)$ , implying

$$\sum_{n=0}^{\infty} \langle u_n, f \rangle u_n(x) = \sum_{n=0}^{\infty} \langle R_L(\lambda)u_n, g \rangle u_n(x) = \sum_{n=0}^{\infty} \alpha_n \langle u_n, g \rangle u_n(x).$$

Moreover, the Cauchy-Schwarz inequality shows

$$\left| \sum_{j=m}^n |\alpha_j \langle u_j, g \rangle u_j(x)| \right|^2 \leq \sum_{j=m}^n |\langle u_j, g \rangle|^2 \sum_{j=m}^n |\alpha_j u_j(x)|^2.$$

Now, by (2.18),  $\sum_{j=0}^{\infty} |\langle u_j, g \rangle|^2 = \|g\|^2$  and hence the first term is part of a convergent series. Similarly, the second term can be estimated independent of  $x$  since

$$\alpha_n u_n(x) = R_L(\lambda)u_n(x) = \int_0^1 G(\lambda, x, t)u_n(t)dt = \langle u_n, G(\lambda, x, \cdot) \rangle$$

implies

$$\sum_{j=m}^n |\alpha_j u_j(x)|^2 \leq \sum_{j=0}^{\infty} |\langle u_j, G(\lambda, x, \cdot) \rangle|^2 = \int_0^1 |G(\lambda, x, t)|^2 dt \leq M(\lambda)^2,$$

where  $M(\lambda) := \max_{x,t \in [0,1]} |G(\lambda, x, t)|$ , again by (2.18).  $\square$

Moreover, it is even possible to weaken our assumptions for uniform convergence. To this end we consider the sequilinear form associated with  $L$ :

$$s_L(f, g) := \langle f, Lg \rangle = \int_0^1 (f'(x)^* g'(x) + q(x)f(x)^* g(x)) dx \quad (3.26)$$

for  $f, g \in \mathfrak{D}(L)$ , where we have used integration by parts as in (3.16). In fact, the above formula continues to hold for  $f$  in a slightly larger class of functions,

$$\mathfrak{Q}(L) := \{f \in C_p^1[0, 1] | f(0) = f(1) = 0\} \supseteq \mathfrak{D}(L), \quad (3.27)$$

which we call the **form domain** of  $L$ . Here  $C_p^1[a, b]$  denotes the set of piecewise continuously differentiable functions  $f$  in the sense that  $f$  is continuously differentiable except for a finite number of points at which it is continuous and the derivative has limits from the left and right. In fact, any class of functions for which the partial integration needed to obtain (3.26) can be justified would be good enough (e.g. the set of absolutely continuous functions to be discussed in Section 4.4 from [48]).

**Lemma 3.12.** *For a regular Sturm–Liouville problem (3.24) converges absolutely uniformly provided  $f \in \mathfrak{Q}(L)$ .*

**Proof.** By replacing  $L \rightarrow L - E_0 + 1$  (this will shift the eigenvalues  $E_n \rightarrow E_n - E_0 + 1$  and leave the eigenvectors unchanged) we can assume  $q_L(f) := s_L(f, f) > 0$  and  $E_j > 0$  without loss of generality.

Now let  $f \in \mathfrak{Q}(L)$  and consider (3.24). Then, observing that  $s_L(f, g)$  is a symmetric sesquilinear form (after our shift it is even a scalar product) as well as  $s_L(f, u_j) = E_j \langle f, u_j \rangle$  one obtains

$$\begin{aligned} 0 &\leq q_L\left(f - \sum_{j=m}^n \langle u_j, f \rangle u_j\right) = q_L(f) - \sum_{j=m}^n \langle u_j, f \rangle s_L(f, u_j) \\ &\quad - \sum_{j=m}^n \langle u_j, f \rangle^* s_L(u_j, f) + \sum_{j,k=m}^n \langle u_j, f \rangle^* \langle u_k, f \rangle s_L(u_j, u_k) \\ &= q_L(f) - \sum_{j=m}^n E_j |\langle u_j, f \rangle|^2 \end{aligned}$$

which implies

$$\sum_{j=m}^n E_j |\langle u_j, f \rangle|^2 \leq q_L(f).$$

In particular, note that this estimate applies to  $f(y) = G(\lambda, x, y)$ . Now from the proof of Theorem 3.11 (with  $\lambda = 0$  and  $\alpha_j = E_j^{-1}$ ) we have  $u_j(x) = E_j \langle u_j, G(0, x, \cdot) \rangle$  and hence

$$\begin{aligned} \sum_{j=m}^n |\langle u_j, f \rangle u_j(x)| &= \sum_{j=m}^n E_j |\langle u_j, f \rangle \langle u_j, G(0, x, \cdot) \rangle| \\ &\leq \left( \sum_{j=m}^n E_j |\langle u_j, f \rangle|^2 \sum_{j=m}^n E_j |\langle u_j, G(0, x, \cdot) \rangle|^2 \right)^{1/2} \\ &\leq \left( \sum_{j=m}^n E_j |\langle u_j, f \rangle|^2 \right)^{1/2} q_L(G(0, x, \cdot))^{1/2}, \end{aligned}$$

where we have used the Cauchy–Schwarz inequality for the weighted scalar product  $(f_j, g_j) \mapsto \sum_j f_j^* g_j E_j$ . Finally note that  $q_L(G(0, x, \cdot))$  is continuous with respect to  $x$  and hence can be estimated by its maximum over  $[0, 1]$ . This shows that the sum (3.24) is absolutely convergent, uniformly with respect to  $x$ .  $\square$

Another consequence of the computations in the previous proof is also worthwhile noting:

**Corollary 3.13.** *We have*

$$G(z, x, y) = \sum_{j=0}^{\infty} \frac{1}{E_j - z} u_j(x) u_j(y), \quad (3.28)$$

where the sum is uniformly convergent. Moreover, we have the following **trace formula**

$$\int_0^1 G(z, x, x) dx = \sum_{j=0}^{\infty} \frac{1}{E_j - z}. \quad (3.29)$$

**Proof.** Using the conventions from the proof of the previous lemma we have  $\langle u_j, G(0, x, \cdot) \rangle = E_j^{-1} u_j(x)$  and since  $G(0, x, \cdot) \in \mathfrak{Q}(L)$  for fixed  $x \in [a, b]$  we have

$$\sum_{j=0}^{\infty} \frac{1}{E_j} u_j(x) u_j(y) = G(0, x, y),$$

where the convergence is uniformly with respect to  $y$  (and  $x$  fixed). Moreover, for  $x = y$  Dini's theorem (cf. Problem B.38) shows that the convergence is uniform with respect to  $x = y$  and this also proves uniform convergence of our sum since

$$\sum_{j=0}^n \frac{1}{|E_j - z|} |u_j(x) u_j(y)| \leq C(z) \left( \sum_{j=0}^n \frac{1}{E_j} u_j(x)^2 \right)^{1/2} \left( \sum_{j=0}^n \frac{1}{E_j} u_j(y)^2 \right)^{1/2},$$

where  $C(z) := \sup_j \frac{E_j}{|E_j - z|}$ .

Finally, the last claim follows upon computing the integral using (3.28) and observing  $\|u_j\| = 1$ .  $\square$

**Example 3.8.** Let us look at the Sturm–Liouville problem with  $q = 0$ . Then the underlying differential equation is

$$-u''(x) = z u(x)$$

whose solution is given by  $u(x) = c_1 \sin(\sqrt{z}x) + c_2 \cos(\sqrt{z}x)$ . The solution satisfying the boundary condition at the left endpoint is  $u_-(z, x) = \sin(\sqrt{z}x)$  and it will be an eigenfunction if and only if  $u_-(z, 1) = \sin(\sqrt{z}) = 0$ . Hence the corresponding eigenvalues and normalized eigenfunctions are

$$E_n = \pi^2 n^2, \quad u_n(x) = \sqrt{2} \sin(n\pi x), \quad n \in \mathbb{N}.$$

Moreover, every function  $f \in \mathcal{L}_{cont}^2(0, 1)$  can be expanded into a **Fourier sine series**

$$f(x) = \sum_{n=1}^{\infty} f_n u_n(x), \quad f_n := \int_0^1 u_n(x) f(x) dx,$$

which is convergent with respect to our scalar product. If  $f \in C_p^1[0, 1]$  with  $f(0) = f(1) = 0$  the series will converge uniformly. For an application of the trace formula see Problem 3.14.  $\diamond$

**Example 3.9.** We could also look at the same equation as in the previous problem but with different boundary conditions

$$u'(0) = u'(1) = 0.$$

Then

$$E_n = \pi^2 n^2, \quad u_n(x) = \begin{cases} 1, & n = 0, \\ \sqrt{2} \cos(n\pi x), & n \in \mathbb{N}. \end{cases}$$

Moreover, every function  $f \in \mathcal{L}_{cont}^2(0, 1)$  can be expanded into a **Fourier cosine series**

$$f(x) = \sum_{n=1}^{\infty} f_n u_n(x), \quad f_n := \int_0^1 u_n(x) f(x) dx,$$

which is convergent with respect to our scalar product.  $\diamond$

**Example 3.10.** Combining the last two examples we see that every symmetric function on  $[-1, 1]$  can be expanded into a Fourier cosine series and every anti-symmetric function into a Fourier sine series. Moreover, since every function  $f(x)$  can be written as the sum of a symmetric function  $\frac{f(x)+f(-x)}{2}$  and an anti-symmetric function  $\frac{f(x)-f(-x)}{2}$ , it can be expanded into a Fourier series. Hence we recover Theorem 2.18.  $\diamond$

**Problem\* 3.11.** Show that for our Sturm–Liouville operator  $u_{\pm}(z, x)^* = u_{\pm}(z^*, x)$ . Conclude  $R_L(z)^* = R_L(z^*)$ . (Hint: Problem 3.7.)

**Problem 3.12.** Show that the resolvent  $R_A(z) = (A - z)^{-1}$  (provided it exists and is densely defined) of a symmetric operator  $A$  is again symmetric for  $z \in \mathbb{R}$ . (Hint:  $g \in \mathfrak{D}(R_A(z))$  if and only if  $g = (A - z)f$  for some  $f \in \mathfrak{D}(A)$ .)

**Problem 3.13.** Suppose  $E_0 > 0$  and equip  $\mathfrak{Q}(L)$  with the scalar product  $s_L$ . Show that

$$f(x) = s_L(G(0, x, \cdot), f).$$

In other words, point evaluations are continuous functionals associated with the vectors  $G(0, x, \cdot) \in \mathfrak{Q}(L)$ . In this context,  $G(0, x, y)$  is called a **reproducing kernel**.

**Problem 3.14.** Show that

$$\sum_{n=1}^{\infty} \frac{1}{n^2 - z} = \frac{1 - \pi\sqrt{z} \cot(\pi\sqrt{z})}{2z}, \quad z \in \mathbb{C} \setminus \mathbb{N}.$$

In particular, for  $z = 0$  this gives Euler's solution of the **Basel problem**:

$$\sum_{n=1}^{\infty} \frac{1}{n^2} = \frac{\pi^2}{6}.$$

In fact, comparing the power series of both sides at  $z = 0$  gives

$$\sum_{n=1}^{\infty} \frac{1}{n^{2k}} = \frac{(-1)^{k+1} (2\pi)^{2k} B_{2k}}{2(2k)!}, \quad k \in \mathbb{N},$$

where  $B_k$  are the **Bernoulli numbers** defined via  $\frac{x}{e^x - 1} = \sum_{k=0}^{\infty} \frac{B_k}{k!} z^k$ . (Hint: Use the trace formula (3.29).)

**Problem 3.15.** Consider the Sturm–Liouville problem on a compact interval  $[a, b]$  with domain

$$\mathfrak{D}(L) = \{f \in C^2[a, b] \mid f'(a) - \alpha f(a) = f'(b) - \beta f(b) = 0\}$$

for some real constants  $\alpha, \beta \in \mathbb{R}$ . Show that Theorem 3.11 continues to hold except for the lower bound on the eigenvalues.

### 3.4. Estimating eigenvalues

In general, there is no way of computing eigenvalues and their corresponding eigenfunctions explicitly. Hence it is important to be able to determine the eigenvalues at least approximately.

Let  $A$  be a symmetric operator which has a lowest eigenvalue  $\alpha_1$  (e.g.,  $A$  is a Sturm–Liouville operator). Suppose we have a vector  $f$  which is an approximation for the eigenvector  $u_1$  of this lowest eigenvalue  $\alpha_1$ . Moreover, suppose we can write

$$A := \sum_{j=1}^{\infty} \alpha_j \langle u_j, \cdot \rangle u_j, \quad \mathfrak{D}(A) := \{f \in \mathfrak{H} \mid \sum_{j=1}^{\infty} |\alpha_j \langle u_j, f \rangle|^2 < \infty\}, \quad (3.30)$$

where  $\{u_j\}_{j \in \mathbb{N}}$  is an orthonormal basis of eigenvectors. Since  $\alpha_1$  is supposed to be the lowest eigenvalue we have  $\alpha_j \geq \alpha_1$  for all  $j \in \mathbb{N}$ .

Writing  $f = \sum_j \gamma_j u_j$ ,  $\gamma_j = \langle u_j, f \rangle$ , one computes

$$\langle f, Af \rangle = \langle f, \sum_{j=1}^{\infty} \alpha_j \gamma_j u_j \rangle = \sum_{j=1}^{\infty} \alpha_j |\gamma_j|^2, \quad f \in \mathfrak{D}(A), \quad (3.31)$$

and we clearly have

$$\alpha_1 \leq \frac{\langle f, Af \rangle}{\|f\|^2}, \quad f \in \mathfrak{D}(A), \quad (3.32)$$

with equality for  $f = u_1$ . In particular, any  $f$  will provide an upper bound and if we add some free parameters to  $f$ , one can optimize them and obtain quite good upper bounds for the first eigenvalue. For example we could

take some orthogonal basis, take a finite number of coefficients and optimize them. This is known as the **Rayleigh–Ritz method**.

**Example 3.11.** Consider the Sturm–Liouville operator  $L$  with potential  $q(x) = x$  and Dirichlet boundary conditions  $f(0) = f(1) = 0$  on the interval  $[0, 1]$ . Our starting point is the quadratic form

$$q_L(f) := \langle f, Lf \rangle = \int_0^1 (|f'(x)|^2 + q(x)|f(x)|^2) dx$$

which gives us the lower bound

$$\langle f, Lf \rangle \geq \min_{0 \leq x \leq 1} q(x) = 0.$$

While the corresponding differential equation can in principle be solved in terms of Airy functions, there is no closed form for the eigenvalues.

First of all we can improve the above bound upon observing  $0 \leq q(x) \leq 1$  which implies

$$\langle f, L_0 f \rangle \leq \langle f, Lf \rangle \leq \langle f, (L_0 + 1)f \rangle, \quad f \in \mathfrak{D}(L) = \mathfrak{D}(L_0),$$

where  $L_0$  is the Sturm–Liouville operator corresponding to  $q(x) = 0$ . Since the lowest eigenvalue of  $L_0$  is  $\pi^2$  we obtain

$$\pi^2 \leq E_1 \leq \pi^2 + 1$$

for the lowest eigenvalue  $E_1$  of  $L$ .

Moreover, using the lowest eigenfunction  $f_1(x) = \sqrt{2} \sin(\pi x)$  of  $L_0$  one obtains the improved upper bound

$$E_1 \leq \langle f_1, Lf_1 \rangle = \pi^2 + \frac{1}{2} \approx 10.3696.$$

Taking the second eigenfunction  $f_2(x) = \sqrt{2} \sin(2\pi x)$  of  $L_0$  we can make the ansatz  $f(x) = (1 + \gamma^2)^{-1/2} (f_1(x) + \gamma f_2(x))$  which gives

$$\langle f, Lf \rangle = \pi^2 + \frac{1}{2} + \frac{\gamma}{1 + \gamma^2} (3\pi^2 \gamma - \frac{32}{9\pi^2}).$$

The right-hand side has a unique minimum at  $\gamma = \frac{32}{27\pi^4 + \sqrt{1024 + 729\pi^8}}$  giving the bound

$$E_1 \leq \frac{5}{2}\pi^2 + \frac{1}{2} - \frac{\sqrt{1024 + 729\pi^8}}{18\pi^2} \approx 10.3685$$

which coincides with the exact eigenvalue up to five digits.  $\diamond$

But is there also something one can say about the next eigenvalues? Suppose we know the first eigenfunction  $u_1$ . Then we can restrict  $A$  to the orthogonal complement of  $u_1$  and proceed as before:  $E_2$  will be the minimum of  $\langle f, Af \rangle$  over all  $f$  restricted to this subspace. If we restrict to the orthogonal complement of an approximating eigenfunction  $f_1$ , there will still be a component in the direction of  $u_1$  left and hence the infimum of the



expectations will be lower than  $E_2$ . Thus the optimal choice  $f_1 = u_1$  will give the maximal value  $E_2$ .

**Theorem 3.14** (Max-min). *Let  $A$  be a symmetric operator and let  $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_N$  be eigenvalues of  $A$  with corresponding orthonormal eigenvectors  $u_1, u_2, \dots, u_N$ . Suppose*

$$A = \sum_{j=1}^N \alpha_j \langle u_j, \cdot \rangle u_j + \tilde{A} \quad (3.33)$$

with  $\langle f, \tilde{A}f \rangle \geq \alpha_N \|f\|^2$  for all  $f \in \mathfrak{D}(A)$  and  $u_1, \dots, u_N \in \text{Ker}(\tilde{A})$ . Then

$$\alpha_j = \sup_{f_1, \dots, f_{j-1}} \inf_{f \in U(f_1, \dots, f_{j-1})} \langle f, Af \rangle, \quad 1 \leq j \leq N, \quad (3.34)$$

where

$$U(f_1, \dots, f_j) := \{f \in \mathfrak{D}(A) \mid \|f\| = 1, f \in \text{span}\{f_1, \dots, f_j\}^\perp\}. \quad (3.35)$$

**Proof.** We have

$$\inf_{f \in U(f_1, \dots, f_{j-1})} \langle f, Af \rangle \leq \alpha_j.$$

In fact, set  $f = \sum_{k=1}^j \gamma_k u_k$  and choose  $\gamma_k$  such that  $f \in U(f_1, \dots, f_{j-1})$ . Then

$$\langle f, Af \rangle = \sum_{k=1}^j |\gamma_k|^2 \alpha_k \leq \alpha_j$$

and the claim follows.

Conversely, let  $\gamma_k = \langle u_k, f \rangle$  and write  $f = \sum_{k=1}^j \gamma_k u_k + \tilde{f}$ . Then

$$\inf_{f \in U(u_1, \dots, u_{j-1})} \langle f, Af \rangle = \inf_{f \in U(u_1, \dots, u_{j-1})} \left( \sum_{k=j}^N |\gamma_k|^2 \alpha_k + \langle \tilde{f}, \tilde{A}\tilde{f} \rangle \right) = \alpha_j. \quad \square$$

Of course if we are interested in the largest eigenvalues all we have to do is consider  $-A$ .

Note that this immediately gives an estimate for eigenvalues if we have a corresponding estimate for the operators. To this end we will write

$$A \leq B \quad \Leftrightarrow \quad \langle f, Af \rangle \leq \langle f, Bf \rangle, \quad f \in \mathfrak{D}(A) \cap \mathfrak{D}(B). \quad (3.36)$$

**Corollary 3.15.** *Suppose  $A$  and  $B$  are symmetric operators with corresponding eigenvalues  $\alpha_j$  and  $\beta_j$  as in the previous theorem. If  $A \leq B$  and  $\mathfrak{D}(B) \subseteq \mathfrak{D}(A)$  then  $\alpha_j \leq \beta_j$ .*

**Proof.** By assumption we have  $\langle f, Af \rangle \leq \langle f, Bf \rangle$  for  $f \in \mathfrak{D}(B)$  implying

$$\inf_{f \in U_A(f_1, \dots, f_{j-1})} \langle f, Af \rangle \leq \inf_{f \in U_B(f_1, \dots, f_{j-1})} \langle f, Af \rangle \leq \inf_{f \in U_B(f_1, \dots, f_{j-1})} \langle f, Bf \rangle,$$

where we have indicated the dependence of  $U$  on the operator via a subscript. Taking the sup on both sides the claim follows.  $\square$

**Example 3.12.** Let  $L$  be again our Sturm–Liouville operator and  $L_0$  the corresponding operator with  $q(x) = 0$ . Set  $q_- = \min_{0 \leq x \leq 1} q(x)$  and  $q_+ = \max_{0 \leq x \leq 1} q(x)$ . Then  $L_0 + q_- \leq L \leq L_0 + q_+$  implies

$$\pi^2 n^2 + q_- \leq E_n \leq \pi^2 n^2 + q_+.$$

In particular, we have proven the famous **Weyl asymptotic**

$$E_n = \pi^2 n^2 + O(1)$$

for the eigenvalues.  $\diamond$

There is also an alternative version which can be proven similar (Problem 3.16):

**Theorem 3.16** (Min-max). *Let  $A$  be as in the previous theorem. Then*

$$\alpha_j = \inf_{V_j \subset \mathfrak{D}(A), \dim(V_j)=j} \sup_{f \in V_j, \|f\|=1} \langle f, Af \rangle, \quad (3.37)$$

where the inf is taken over subspaces with the indicated properties.

**Problem\* 3.16.** *Prove Theorem 3.16.*

**Problem 3.17.** *Suppose  $A, A_n$  are self-adjoint, bounded and  $A_n \rightarrow A$ . Then  $\alpha_k(A_n) \rightarrow \alpha_k(A)$ . (Hint: For  $B$  self-adjoint  $\|B\| \leq \varepsilon$  is equivalent to  $-\varepsilon \leq B \leq \varepsilon$ .)*

### 3.5. Singular value decomposition of compact operators

Our first aim is to find a generalization of Corollary 3.8 for general compact operators between Hilbert spaces. The key observation is that if  $K \in \mathcal{K}(\mathfrak{H}_1, \mathfrak{H}_2)$  is compact, then  $K^*K \in \mathcal{K}(\mathfrak{H}_1)$  is compact and symmetric and thus, by Corollary 3.8, there is a countable orthonormal set  $\{u_j\} \subset \mathfrak{H}_1$  and nonzero real numbers  $s_j^2 \neq 0$  such that

$$K^*Kf = \sum_j s_j^2 \langle u_j, f \rangle u_j. \quad (3.38)$$

Moreover,  $\|Ku_j\|^2 = \langle u_j, K^*Ku_j \rangle = \langle u_j, s_j^2 u_j \rangle = s_j^2$  shows that we can set

$$s_j := \|Ku_j\| > 0. \quad (3.39)$$

The numbers  $s_j = s_j(K)$  are called **singular values** of  $K$ . There are either finitely many singular values or they converge to zero.

**Theorem 3.17** (Singular value decomposition of compact operators). *Let  $K \in \mathcal{K}(\mathfrak{H}_1, \mathfrak{H}_2)$  be compact and let  $s_j$  be the singular values of  $K$  and  $\{u_j\} \subset \mathfrak{H}_1$  corresponding orthonormal eigenvectors of  $K^*K$ . Then*

$$K = \sum_j s_j \langle u_j, \cdot \rangle v_j, \quad (3.40)$$

where  $v_j = s_j^{-1} K u_j$ . The norm of  $K$  is given by the largest singular value

$$\|K\| = \max_j s_j(K). \quad (3.41)$$

Moreover, the vectors  $\{v_j\} \subset \mathfrak{H}_2$  are again orthonormal and satisfy  $K^*v_j = s_j u_j$ . In particular,  $v_j$  are eigenvectors of  $KK^*$  corresponding to the eigenvalues  $s_j^2$ .

**Proof.** For any  $f \in \mathfrak{H}_1$  we can write

$$f = \sum_j \langle u_j, f \rangle u_j + f_\perp$$

with  $f_\perp \in \text{Ker}(K^*K) = \text{Ker}(K)$  (Problem 3.18). Then

$$Kf = \sum_j \langle u_j, f \rangle K u_j = \sum_j s_j \langle u_j, f \rangle v_j$$

as required. Furthermore,

$$\langle v_j, v_k \rangle = (s_j s_k)^{-1} \langle K u_j, K u_k \rangle = (s_j s_k)^{-1} \langle K^* K u_j, u_k \rangle = s_j s_k^{-1} \langle u_j, u_k \rangle$$

shows that  $\{v_j\}$  are orthonormal. By definition  $K^*v_j = s_j^{-1} K^* K u_j = s_j u_j$  which also shows  $KK^*v_j = s_j K u_j = s_j^2 v_j$ .

Finally, (3.41) follows using Bessel's inequality

$$\|Kf\|^2 = \left\| \sum_j s_j \langle u_j, f \rangle v_j \right\|^2 = \sum_j s_j^2 |\langle u_j, f \rangle|^2 \leq \left( \max_j s_j(K)^2 \right) \|f\|^2,$$

where equality holds for  $f = u_{j_0}$  if  $s_{j_0} = \max_j s_j(K)$ .  $\square$

If  $K \in \mathcal{K}(\mathfrak{H})$  is self-adjoint, then  $u_j = \sigma_j v_j$ ,  $\sigma_j^2 = 1$ , are the eigenvectors of  $K$  and  $\sigma_j s_j$  are the corresponding eigenvalues. In particular, for a self-adjoint operators the singular values are the absolute values of the nonzero eigenvalues.

The above theorem also gives rise to the **polar decomposition**

$$K = U|K| = |K^*|U, \quad (3.42)$$

where

$$|K| := \sqrt{K^*K} = \sum_j s_j \langle u_j, \cdot \rangle u_j, \quad |K^*| = \sqrt{KK^*} = \sum_j s_j \langle v_j, \cdot \rangle v_j \quad (3.43)$$

are self-adjoint (in fact nonnegative) and

$$U := \sum_j \langle u_j, \cdot \rangle v_j \quad (3.44)$$

is an isometry from  $\overline{\text{Ran}(K^*)} = \overline{\text{span}\{u_j\}}$  onto  $\overline{\text{Ran}(K)} = \overline{\text{span}\{v_j\}}$ .

From the min-max theorem (Theorem 3.16) we obtain:

**Lemma 3.18.** *Let  $K \in \mathcal{K}(\mathfrak{H}_1, \mathfrak{H}_2)$  be compact; then*

$$s_j(K) = \min_{f_1, \dots, f_{j-1}} \max_{f \in U(f_1, \dots, f_{j-1})} \|Kf\|, \quad (3.45)$$

where  $U(f_1, \dots, f_j) := \{f \in \mathfrak{H}_1 \mid \|f\| = 1, f \in \text{span}\{f_1, \dots, f_j\}^\perp\}$ .

In particular, note

$$s_j(AK) \leq \|A\|s_j(K), \quad s_j(KA) \leq \|A\|s_j(K) \quad (3.46)$$

whenever  $K$  is compact and  $A$  is bounded (the second estimate follows from the first by taking adjoints).

An operator  $K \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  is called a **finite rank operator** if its range is finite dimensional. The dimension

$$\text{rank}(K) := \dim \text{Ran}(K)$$

is called the **rank** of  $K$ . Since for a compact operator

$$\overline{\text{Ran}(K)} = \overline{\text{span}\{v_j\}} \quad (3.47)$$

we see that a compact operator is finite rank if and only if the sum in (3.40) is finite. Note that the finite rank operators form an ideal in  $\mathcal{L}(\mathfrak{H})$  just as the compact operators do. Moreover, every finite rank operator is compact by the Heine–Borel theorem (Theorem B.22).

Now truncating the sum in the canonical form gives us a simple way to approximate compact operators by finite rank ones. Moreover, this is in fact the best approximation within the class of finite rank operators:

**Lemma 3.19.** *Let  $K \in \mathcal{K}(\mathfrak{H}_1, \mathfrak{H}_2)$  be compact and let its singular values be ordered. Then*

$$s_j(K) = \min_{\text{rank}(F) < j} \|K - F\|, \quad (3.48)$$

where the minimum is attained for

$$F_{j-1} := \sum_{k=1}^{j-1} s_k \langle u_k, \cdot \rangle v_k. \quad (3.49)$$

In particular, the closure of the ideal of finite rank operators in  $\mathcal{L}(\mathfrak{H})$  is the ideal of compact operators.

**Proof.** That there is equality for  $F = F_{j-1}$  follows from (3.41). In general, the restriction of  $F$  to  $\text{span}\{u_1, \dots, u_j\}$  will have a nontrivial kernel. Let  $f = \sum_{k=1}^j \alpha_k u_k$  be a normalized element of this kernel, then  $\|(K - F)f\|^2 = \|Kf\|^2 = \sum_{k=1}^j |\alpha_k s_k|^2 \geq s_j^2$ .

In particular, every compact operator can be approximated by finite rank ones and since the limit of compact operators is compact, we cannot get more than the compact operators.  $\square$

Two more consequences are worthwhile noting.

**Corollary 3.20.** *An operator  $K \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  is compact if and only if  $K^*K$  is.*

**Proof.** Just observe that  $K^*K$  compact is all that was used to show Theorem 3.17.  $\square$

**Corollary 3.21.** *An operator  $K \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  is compact (finite rank) if and only if  $K^* \in \mathcal{L}(\mathfrak{H}_2, \mathfrak{H}_1)$  is. In fact,  $s_j(K) = s_j(K^*)$  and*

$$K^* = \sum_j s_j \langle v_j, \cdot \rangle u_j. \quad (3.50)$$

**Proof.** First of all note that (3.50) follows from (3.40) since taking adjoints is continuous and  $(\langle u_j, \cdot \rangle v_j)^* = \langle v_j, \cdot \rangle u_j$  (cf. Problem 2.7). The rest is straightforward.  $\square$

From this last lemma one easily gets a number of useful inequalities for the singular values:

**Corollary 3.22.** *Let  $K_1$  and  $K_2$  be compact and let  $s_j(K_1)$  and  $s_j(K_2)$  be ordered. Then*

- (i)  $s_{j+k-1}(K_1 + K_2) \leq s_j(K_1) + s_k(K_2)$ ,
- (ii)  $s_{j+k-1}(K_1 K_2) \leq s_j(K_1) s_k(K_2)$ ,
- (iii)  $|s_j(K_1) - s_j(K_2)| \leq \|K_1 - K_2\|$ .

**Proof.** Let  $F_1$  be of rank  $j - 1$  and  $F_2$  of rank  $k - 1$  such that  $\|K_1 - F_1\| = s_j(K_1)$  and  $\|K_2 - F_2\| = s_k(K_2)$ . Then  $s_{j+k-1}(K_1 + K_2) \leq \|(K_1 + K_2) - (F_1 + F_2)\| = \|K_1 - F_1\| + \|K_2 - F_2\| = s_j(K_1) + s_k(K_2)$  since  $F_1 + F_2$  is of rank at most  $j + k - 2$ .

Similarly  $F = F_1(K_2 - F_2) + K_1 F_2$  is of rank at most  $j + k - 2$  and hence  $s_{j+k-1}(K_1 K_2) \leq \|K_1 K_2 - F\| = \|(K_1 - F_1)(K_2 - F_2)\| \leq \|K_1 - F_1\| \|K_2 - F_2\| = s_j(K_1) s_k(K_2)$ .

Next, choosing  $k = 1$  and replacing  $K_2 \rightarrow K_2 - K_1$  in (i) gives  $s_j(K_2) \leq s_j(K_1) + \|K_2 - K_1\|$ . Reversing the roles gives  $s_j(K_1) \leq s_j(K_2) + \|K_1 - K_2\|$  and proves (iii).  $\square$

**Example 3.13.** One might hope that item (i) from the previous corollary can be improved to  $s_j(K_1 + K_2) \leq s_j(K_1) + s_j(K_2)$ . However, this is not the case as the following example shows:

$$K_1 := \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad K_2 := \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

Then  $1 = s_2(K_1 + K_2) \not\leq s_2(K_1) + s_2(K_2) = 0$ .  $\diamond$

**Problem\* 3.18.** Show that  $\text{Ker}(A^*A) = \text{Ker}(A)$  for any  $A \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$ .

**Problem 3.19.** Let  $K$  be multiplication by a sequence  $k \in c_0(\mathbb{N})$  in the Hilbert space  $\ell^2(\mathbb{N})$ . What are the singular values of  $K$ ?

**Problem 3.20.** Let  $K$  be multiplication by a sequence  $k \in c_0(\mathbb{N})$  in the Hilbert space  $\ell^2(\mathbb{N})$  and consider  $L = KS^-$ . What are the singular values of  $L$ ? Does  $L$  have any eigenvalues?

**Problem 3.21.** Let  $K \in \mathcal{K}(\mathfrak{H}_1, \mathfrak{H}_2)$  be compact and let its singular values be ordered. Let  $M \subseteq \mathfrak{H}_1$ ,  $N \subseteq \mathfrak{H}_1$  be subspaces with corresponding orthogonal projections  $P_M$ ,  $P_N$ , respectively. Then

$$s_j(K) = \min_{\dim(M) < j} \|K - KP_M\| = \min_{\dim(N) < j} \|K - P_N K\|,$$

where the minimum is taken over all subspaces with the indicated dimension. Moreover, the minimum is attained for

$$M = \text{span}\{u_k\}_{k=1}^{j-1}, \quad N = \text{span}\{v_k\}_{k=1}^{j-1}.$$

### 3.6. Hilbert–Schmidt and trace class operators

We can further subdivide the class of compact operators  $\mathcal{K}(\mathfrak{H})$  according to the decay of their singular values. We define

$$\|K\|_p := \left( \sum_j s_j(K)^p \right)^{1/p} \quad (3.51)$$

plus corresponding spaces

$$\mathcal{J}_p(\mathfrak{H}) = \{K \in \mathcal{K}(\mathfrak{H}) \mid \|K\|_p < \infty\}, \quad (3.52)$$

which are known as **Schatten  $p$ -classes**. Even though our notation hints at the fact that  $\|\cdot\|_p$  is a norm, we will only prove this here for  $p = 1, 2$  (the only nontrivial part is the triangle inequality). Note that by (3.41)

$$\|K\| \leq \|K\|_p \quad (3.53)$$

and that by  $s_j(K) = s_j(K^*)$  we have

$$\|K\|_p = \|K^*\|_p. \quad (3.54)$$

The two most important cases are  $p = 1$  and  $p = 2$ :  $\mathcal{J}_2(\mathfrak{H})$  is the space of **Hilbert–Schmidt operators** and  $\mathcal{J}_1(\mathfrak{H})$  is the space of **trace class** operators.

**Example 3.14.** Any multiplication operator by a sequence from  $\ell^p(\mathbb{N})$  is in the Schatten  $p$ -class of  $\mathfrak{H} = \ell^2(\mathbb{N})$ .  $\diamond$

**Example 3.15.** By virtue of the Weyl asymptotics (see Example 3.12) the resolvent of our Sturm–Liouville operator is trace class.  $\diamond$

**Example 3.16.** Let  $k$  be a periodic function which is square integrable over  $[-\pi, \pi]$ . Then the integral operator

$$(Kf)(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} k(y-x)f(y)dy$$

has the eigenfunctions  $u_j(x) = (2\pi)^{-1/2}e^{-ijx}$  with corresponding eigenvalues  $\hat{k}_j$ ,  $j \in \mathbb{Z}$ , where  $\hat{k}_j$  are the Fourier coefficients of  $k$ . Since  $\{u_j\}_{j \in \mathbb{Z}}$  is an ONB we have found all eigenvalues. In particular, the Fourier transform maps  $K$  to the multiplication operator with the sequence of its eigenvalues  $\hat{k}_j$ . Hence the singular values are the absolute values of the nonzero eigenvalues and (3.40) reads

$$K = \sum_{j \in \mathbb{Z}} \hat{k}_j \langle u_j, \cdot \rangle u_j.$$

Moreover, since the eigenvalues are in  $\ell^2(\mathbb{Z})$  we see that  $K$  is a Hilbert–Schmidt operator. If  $k$  is continuous with summable Fourier coefficients (e.g.  $k \in C_{per}^2[-\pi, \pi]$ ), then  $K$  is trace class.  $\diamond$

We first prove an alternate definition for the Hilbert–Schmidt norm.

**Lemma 3.23.** *A bounded operator  $K$  is Hilbert–Schmidt if and only if*

$$\sum_{j \in J} \|Kw_j\|^2 < \infty \tag{3.55}$$

*for some orthonormal basis and*

$$\|K\|_2 = \left( \sum_{j \in J} \|Kw_j\|^2 \right)^{1/2}, \tag{3.56}$$

*for every orthonormal basis in this case.*

**Proof.** First of all note that (3.55) implies that  $K$  is compact. To see this, let  $P_n$  be the projection onto the space spanned by the first  $n$  elements of the orthonormal basis  $\{w_j\}$ . Then  $K_n = KP_n$  is finite rank and converges to  $K$  since

$$\|(K - K_n)f\| = \left\| \sum_{j > n} c_j Kw_j \right\| \leq \sum_{j > n} |c_j| \|Kw_j\| \leq \left( \sum_{j > n} \|Kw_j\|^2 \right)^{1/2} \|f\|,$$

where  $f = \sum_j c_j w_j$ .

The rest follows from (3.40) and

$$\begin{aligned} \sum_j \|Kw_j\|^2 &= \sum_{k,j} |\langle v_k, Kw_j \rangle|^2 = \sum_{k,j} |\langle K^*v_k, w_j \rangle|^2 = \sum_k \|K^*v_k\|^2 \\ &= \sum_k s_k(K)^2 = \|K\|_2^2. \end{aligned}$$

Here we have used  $\overline{\text{span}\{v_k\}} = \text{Ker}(K^*)^\perp = \overline{\text{Ran}(K)}$  in the first step.  $\square$

**Corollary 3.24.** *The Hilbert–Schmidt norm satisfies the triangle inequality and hence is indeed a norm.*

**Proof.** This follows from (3.56) upon using the triangle inequality for  $\mathfrak{H}$  and for  $\ell^2(J)$ .  $\square$

Now we can show

**Lemma 3.25.** *The set of Hilbert–Schmidt operators forms an ideal in  $\mathcal{L}(\mathfrak{H})$  and*

$$\|KA\|_2 \leq \|A\| \|K\|_2, \quad \text{respectively,} \quad \|AK\|_2 \leq \|A\| \|K\|_2. \quad (3.57)$$

**Proof.** If  $K_1$  and  $K_2$  are Hilbert–Schmidt operators, then so is their sum since

$$\begin{aligned} \|K_1 + K_2\|_2 &= \left( \sum_{j \in J} \|(K_1 + K_2)w_j\|^2 \right)^{1/2} \leq \left( \sum_{j \in J} (\|K_1w_j\| + \|K_2w_j\|)^2 \right)^{1/2} \\ &\leq \|K_1\|_2 + \|K_2\|_2, \end{aligned}$$

where we have used the triangle inequality for  $\ell^2(J)$ .

Let  $K$  be Hilbert–Schmidt and  $A$  bounded. Then  $AK$  is compact and

$$\|AK\|_2^2 = \sum_j \|AKw_j\|^2 \leq \|A\|^2 \sum_j \|Kw_j\|^2 = \|A\|^2 \|K\|_2^2.$$

For  $KA$  just consider adjoints.  $\square$

**Example 3.17.** Consider  $\ell^2(\mathbb{N})$  and let  $K$  be some compact operator. Let  $K_{jk} = \langle \delta^j, K\delta^k \rangle = (K\delta^j)_k$  be its matrix elements such that

$$(Ka)_j = \sum_{k=1}^{\infty} K_{jk} a_k.$$

Then, choosing  $w_j = \delta^j$  in (3.56) we get

$$\|K\|_2 = \left( \sum_{j=1}^{\infty} \|K\delta^j\|^2 \right)^{1/2} = \left( \sum_{j=1}^{\infty} \sum_{k=1}^{\infty} |K_{jk}|^2 \right)^{1/2}.$$



Hence  $K$  is Hilbert–Schmidt if and only if its matrix elements are in  $\ell^2(\mathbb{N} \times \mathbb{N})$  and the Hilbert–Schmidt norm coincides with the  $\ell^2(\mathbb{N} \times \mathbb{N})$  norm of the matrix elements. Especially in the finite dimensional case the Hilbert–Schmidt norm is also known as **Frobenius norm**.

Of course the same calculation shows that a bounded operator is Hilbert–Schmidt if and only if its matrix elements  $\langle w_j, Kw_k \rangle$  with respect to some orthonormal basis  $\{w_j\}_{j \in J}$  are in  $\ell^2(J \times J)$  and the Hilbert–Schmidt norm coincides with the  $\ell^2(J \times J)$  norm of the matrix elements.  $\diamond$

**Example 3.18.** Let  $I = [a, b]$  be a compact interval. Suppose  $K : L^2(I) \rightarrow C(I)$  is continuous, then  $K : L^2(I) \rightarrow L^2(I)$  is Hilbert–Schmidt with Hilbert–Schmidt norm  $\|K\|_2 \leq \sqrt{b-a}M$ , where  $M := \|K\|_{L^2(I) \rightarrow C(I)}$ .

To see this start by observing that point evaluations are continuous functionals on  $C(I)$  and hence  $f \mapsto (Kf)(x)$  is a continuous linear functional on  $L^2(I)$  satisfying  $|(Kf)(x)| \leq M\|f\|$ . By the Riesz lemma there is some  $K_x \in L^2(I)$  with  $\|K_x\| \leq M$  such that

$$(Kf)(x) = \langle K_x, f \rangle$$

and hence for any orthonormal basis  $\{w_j\}_{j \in \mathbb{N}}$  we have

$$\sum_{j \in \mathbb{N}} |(Kw_j)(x)|^2 = \sum_{j \in \mathbb{N}} |\langle K_x, w_j \rangle|^2 = \|K_x\|^2 \leq M^2.$$

But then

$$\begin{aligned} \sum_{j \in \mathbb{N}} \|Kw_j\|^2 &= \sum_{j \in \mathbb{N}} \int_a^b |(Kw_j)(x)|^2 dx = \int_a^b \left( \sum_{j \in \mathbb{N}} |(Kw_j)(x)|^2 \right) dx \\ &\leq (b-a)M^2 \end{aligned}$$

as claimed.  $\diamond$

Since Hilbert–Schmidt operators turn out easy to identify (cf. also Section 3.5 from [48]), it is important to relate  $\mathcal{J}_1(\mathfrak{H})$  with  $\mathcal{J}_2(\mathfrak{H})$ :

**Lemma 3.26.** *An operator is trace class if and only if it can be written as the product of two Hilbert–Schmidt operators,  $K = K_1 K_2$ , and in this case we have*

$$\|K\|_1 \leq \|K_1\|_2 \|K_2\|_2. \quad (3.58)$$

*In fact,  $K_1, K_2$  can be chosen such that  $\|K\|_1 = \|K_1\|_2 \|K_2\|_2$ .*

**Proof.** Using (3.40) (where we can extend  $u_n$  and  $v_n$  to orthonormal bases if necessary) and Cauchy–Schwarz we have

$$\begin{aligned}\|K\|_1 &= \sum_n \langle v_n, Ku_n \rangle = \sum_n |\langle K_1^* v_n, K_2 u_n \rangle| \\ &\leq \left( \sum_n \|K_1^* v_n\|^2 \sum_n \|K_2 u_n\|^2 \right)^{1/2} = \|K_1\|_2 \|K_2\|_2\end{aligned}$$

and hence  $K = K_1 K_2$  is trace class if both  $K_1$  and  $K_2$  are Hilbert–Schmidt operators. To see the converse, let  $K$  be given by (3.40) and choose  $K_1 = \sum_j \sqrt{s_j(K)} \langle u_j, \cdot \rangle v_j$ , respectively,  $K_2 = \sum_j \sqrt{s_j(K)} \langle u_j, \cdot \rangle u_j$ . Note that in this case  $\|K\|_1 = \|K_1\|_2^2 = \|K_2\|_2^2$ .  $\square$

Now we can also explain the name trace class:

**Lemma 3.27.** *If  $K$  is trace class, then for every orthonormal basis  $\{w_n\}$  the **trace***

$$\operatorname{tr}(K) = \sum_n \langle w_n, Kw_n \rangle \quad (3.59)$$

*is finite,*

$$|\operatorname{tr}(K)| \leq \|K\|_1, \quad (3.60)$$

*and independent of the orthonormal basis.*

**Proof.** If we write  $K = K_1 K_2$  with  $K_1, K_2$  Hilbert–Schmidt such that  $\|K\|_1 = \|K_1\|_2 \|K_2\|_2$ , then the Cauchy–Schwarz inequality implies  $|\operatorname{tr}(K)| \leq \|K_1^*\|_2 \|K_2\|_2 = \|K\|_1$ . Moreover, if  $\{\tilde{w}_n\}$  is another orthonormal basis, we have

$$\begin{aligned}\sum_n \langle w_n, K_1 K_2 w_n \rangle &= \sum_n \langle K_1^* w_n, K_2 w_n \rangle = \sum_{n,m} \langle K_1^* w_n, \tilde{w}_m \rangle \langle \tilde{w}_m, K_2 w_n \rangle \\ &= \sum_{m,n} \langle K_2^* \tilde{w}_m, w_n \rangle \langle w_n, K_1 \tilde{w}_m \rangle = \sum_m \langle K_2^* \tilde{w}_m, K_1 \tilde{w}_m \rangle \\ &= \sum_m \langle \tilde{w}_m, K_2 K_1 \tilde{w}_m \rangle.\end{aligned}$$

In the special case  $w = \tilde{w}$  we see  $\operatorname{tr}(K_1 K_2) = \operatorname{tr}(K_2 K_1)$  and the general case now shows that the trace is independent of the orthonormal basis.  $\square$

Clearly for self-adjoint trace class operators, the trace is the sum over all eigenvalues (counted with their multiplicity). To see this, one just has to choose the orthonormal basis to consist of eigenfunctions. This is even true for all trace class operators and is known as Lidskij trace theorem (see [37] for an easy to read introduction).

**Example 3.19.** We already mentioned that the resolvent of our Sturm–Liouville operator is trace class. Choosing a basis of eigenfunctions we see that the trace of the resolvent is the sum over its eigenvalues and combining this with our trace formula (3.29) gives

$$\operatorname{tr}(R_L(z)) = \sum_{j=0}^{\infty} \frac{1}{E_j - z} = \int_0^1 G(z, x, x) dx$$

for  $z \in \mathbb{C}$  no eigenvalue.  $\diamond$

**Example 3.20.** For our integral operator  $K$  from Example 3.16 we have in the trace class case

$$\operatorname{tr}(K) = \sum_{j \in \mathbb{Z}} \hat{k}_j = k(0).$$

Note that this can again be interpreted as the integral over the diagonal  $(2\pi)^{-1}k(x-x) = (2\pi)^{-1}k(0)$  of the kernel.  $\diamond$

We also note the following elementary properties of the trace:

**Lemma 3.28.** *Suppose  $K, K_1, K_2$  are trace class and  $A$  is bounded.*

- (i) *The trace is linear.*
- (ii)  $\operatorname{tr}(K^*) = \operatorname{tr}(K)^*$ .
- (iii) *If  $K_1 \leq K_2$ , then  $\operatorname{tr}(K_1) \leq \operatorname{tr}(K_2)$ .*
- (iv)  $\operatorname{tr}(AK) = \operatorname{tr}(KA)$ .

**Proof.** (i) and (ii) are straightforward. (iii) follows from  $K_1 \leq K_2$  if and only if  $\langle f, K_1 f \rangle \leq \langle f, K_2 f \rangle$  for every  $f \in \mathfrak{H}$ . (iv) By Problem 2.12 and (i), it is no restriction to assume that  $A$  is unitary. Let  $\{w_n\}$  be some ONB and note that  $\{\tilde{w}_n = Aw_n\}$  is also an ONB. Then

$$\begin{aligned} \operatorname{tr}(AK) &= \sum_n \langle \tilde{w}_n, AK \tilde{w}_n \rangle = \sum_n \langle Aw_n, AKAw_n \rangle \\ &= \sum_n \langle w_n, KAw_n \rangle = \operatorname{tr}(KA) \end{aligned}$$

and the claim follows.  $\square$

We also mention a useful criterion for  $K$  to be trace class.

**Lemma 3.29.** *An operator  $K$  is trace class if and only if it can be written as*

$$K = \sum_j \langle f_j, \cdot \rangle g_j \tag{3.61}$$

for some sequences  $f_j, g_j$  satisfying

$$\sum_j \|f_j\| \|g_j\| < \infty. \tag{3.62}$$

Moreover, in this case

$$\operatorname{tr}(K) = \sum_j \langle f_j, g_j \rangle \quad (3.63)$$

and

$$\|K\|_1 = \min \sum_j \|f_j\| \|g_j\|, \quad (3.64)$$

where the minimum is taken over all representations as in (3.61).

**Proof.** To see that a trace class operator (3.40) can be written in such a way choose  $f_j = u_j$ ,  $g_j = s_j v_j$ . This also shows that the minimum in (3.64) is attained. Conversely note that the sum converges in the operator norm and hence  $K$  is compact. Moreover, for every finite  $N$  we have

$$\begin{aligned} \sum_{k=1}^N s_k &= \sum_{k=1}^N \langle v_k, K u_k \rangle = \sum_{k=1}^N \sum_j \langle v_k, g_j \rangle \langle f_j, u_k \rangle = \sum_j \sum_{k=1}^N \langle v_k, g_j \rangle \langle f_j, u_k \rangle \\ &\leq \sum_j \left( \sum_{k=1}^N |\langle v_k, g_j \rangle|^2 \right)^{1/2} \left( \sum_{k=1}^N |\langle f_j, u_k \rangle|^2 \right)^{1/2} \leq \sum_j \|f_j\| \|g_j\|. \end{aligned}$$

This also shows that the right-hand side in (3.64) cannot exceed  $\|K\|_1$ . To see the last claim we choose an ONB  $\{w_k\}$  to compute the trace

$$\begin{aligned} \operatorname{tr}(K) &= \sum_k \langle w_k, K w_k \rangle = \sum_k \sum_j \langle w_k, \langle f_j, w_k \rangle g_j \rangle = \sum_j \sum_k \langle \langle w_k, f_j \rangle w_k, g_j \rangle \\ &= \sum_j \langle f_j, g_j \rangle. \end{aligned} \quad \square$$

An immediate consequence of (3.64) is:

**Corollary 3.30.** *The trace norm satisfies the triangle inequality and hence is indeed a norm.*

Finally, note that

$$\|K\|_2 = (\operatorname{tr}(K^* K))^{1/2} \quad (3.65)$$

which shows that  $\mathcal{J}_2(\mathfrak{H})$  is in fact a Hilbert space with scalar product given by

$$\langle K_1, K_2 \rangle = \operatorname{tr}(K_1^* K_2). \quad (3.66)$$

**Problem 3.22.** Let  $\mathfrak{H} := \ell^2(\mathbb{N})$  and let  $A$  be multiplication by a sequence  $a = (a_j)_{j=1}^\infty$ . Show that  $A$  is Hilbert–Schmidt if and only if  $a \in \ell^2(\mathbb{N})$ . Furthermore, show that  $\|A\|_2 = \|a\|$  in this case.

**Problem 3.23.** An operator of the form  $K : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$ ,  $f_n \mapsto \sum_{j \in \mathbb{N}} k_{n+j} f_j$  is called **Hankel operator**.

- Show that  $K$  is Hilbert–Schmidt if and only if  $\sum_{j \in \mathbb{N}} j |k_{j+1}|^2 < \infty$  and this number equals  $\|K\|_2$ .
- Show that  $K$  is Hilbert–Schmidt with  $\|K\|_2 \leq \|c\|_1$  if  $|k_j| \leq c_j$ , where  $c_j$  is decreasing and summable.

(Hint: For the first item use summation by parts.)

# The main theorems about Banach spaces

Despite the many advantages of Hilbert spaces, there are also situations where a non-Hilbert space is better suited (in fact the choice of the *right* space is typically crucial for many problems). Hence we will devote our attention to Banach spaces next.

## 4.1. The Baire theorem and its consequences

Recall that the interior of a set is the largest open subset (that is, the union of all open subsets). A set is called **nowhere dense** if its closure has empty interior. The key to several important theorems about Banach spaces is the observation that a Banach space cannot be the countable union of nowhere dense sets.

**Theorem 4.1** (Baire category theorem). *Let  $X$  be a (nonempty) complete metric space. Then  $X$  cannot be the countable union of nowhere dense sets.*

**Proof.** Suppose  $X = \bigcup_{n=1}^{\infty} X_n$ . We can assume that the sets  $X_n$  are closed and none of them contains a ball; that is,  $X \setminus X_n$  is open and nonempty for every  $n$ . We will construct a Cauchy sequence  $x_n$  which stays away from all  $X_n$ .

Since  $X \setminus X_1$  is open and nonempty, there is a ball  $B_{r_1}(x_1) \subseteq X \setminus X_1$ . Reducing  $r_1$  a little, we can even assume  $\overline{B_{r_1}(x_1)} \subseteq X \setminus X_1$ . Moreover, since  $X_2$  cannot contain  $\overline{B_{r_1}(x_1)}$ , there is some  $x_2 \in \overline{B_{r_1}(x_1)}$  that is not in  $X_2$ . Since  $\overline{B_{r_1}(x_1)} \cap (X \setminus X_2)$  is open, there is a closed ball  $\overline{B_{r_2}(x_2)} \subseteq \overline{B_{r_1}(x_1)} \cap (X \setminus X_2)$ . Proceeding recursively, we obtain a sequence (here we

use the axiom of choice) of balls such that

$$\overline{B_{r_n}(x_n)} \subseteq B_{r_{n-1}}(x_{n-1}) \cap (X \setminus X_n).$$

Now observe that in every step we can choose  $r_n$  as small as we please; hence without loss of generality  $r_n \rightarrow 0$ . Since by construction  $x_n \in \overline{B_{r_n}(x_n)}$  for  $n \geq N$ , we conclude that  $x_n$  is Cauchy and converges to some point  $x \in X$ . But  $x \in \overline{B_{r_n}(x_n)} \subseteq X \setminus X_n$  for every  $n$ , contradicting our assumption that the  $X_n$  cover  $X$ .  $\square$

In other words, if  $X_n \subseteq X$  is a sequence of closed subsets which cover  $X$ , at least one  $X_n$  contains a ball of radius  $\varepsilon > 0$ .

**Example 4.1.** The set of rational numbers  $\mathbb{Q}$  can be written as a countable union of its elements. This shows that the completeness assumption is crucial.  $\diamond$

Remark: Sets which can be written as the countable union of nowhere dense sets are said to be of **first category** or **meager** (also meagre). All other sets are **second category** or **fat** (also **residual**). Hence explaining the name category theorem.

Since a closed set is nowhere dense if and only if its complement is open and dense (cf. Problem B.7), there is a reformulation which is also worthwhile noting:

**Corollary 4.2.** *Let  $X$  be a complete metric space. Then any countable intersection of open dense sets is again dense.*

**Proof.** Let  $\{O_n\}$  be a family of open dense sets whose intersection is not dense. Then this intersection must be missing some closed ball  $\overline{B_\varepsilon}$ . This ball will lie in  $\bigcup_n X_n$ , where  $X_n := X \setminus O_n$  are closed and nowhere dense. Now note that  $\tilde{X}_n := X_n \cap \overline{B_\varepsilon}$  are closed nowhere dense sets in  $\overline{B_\varepsilon}$ . But  $\overline{B_\varepsilon}$  is a complete metric space, a contradiction.  $\square$

Countable intersections of open sets are in some sense the next general sets after open sets and are called  $G_\delta$  sets (here  $G$  and  $\delta$  stand for the German words *Gebiet* and *Durchschnitt*, respectively). The complement of a  $G_\delta$  set is a countable union of closed sets also known as an  $F_\sigma$  set (here  $F$  and  $\sigma$  stand for the French words *fermé* and *somme*, respectively). The complement of a dense  $G_\delta$  set will be a countable union of nowhere dense sets and hence by definition meager. Consequently properties which hold on a dense  $G_\delta$  are considered *generic* in this context.

**Example 4.2.** The irrational numbers are a dense  $G_\delta$  set in  $\mathbb{R}$ . To see this, let  $x_n$  be an enumeration of the rational numbers and consider the intersection of the open sets  $O_n := \mathbb{R} \setminus \{x_n\}$ . The rational numbers are hence an  $F_\sigma$  set.  $\diamond$

Now we are ready for the first important consequence:

**Theorem 4.3** (Banach–Steinhaus). *Let  $X$  be a Banach space and  $Y$  some normed vector space. Let  $\{A_\alpha\} \subseteq \mathcal{L}(X, Y)$  be a family of bounded operators. Then*

- *either  $\{A_\alpha\}$  is uniformly bounded,  $\|A_\alpha\| \leq C$ ,*
- *or the set  $\{x \in X \mid \sup_\alpha \|A_\alpha x\| = \infty\}$  is a dense  $G_\delta$ .*

**Proof.** Consider the sets

$$O_n := \{x \mid \|A_\alpha x\| > n \text{ for some } \alpha\} = \bigcup_\alpha \{x \mid \|A_\alpha x\| > n\}, \quad n \in \mathbb{N}.$$

By continuity of  $A_\alpha$  and the norm, each  $O_n$  is a union of open sets and hence open. Now either all of these sets are dense and hence their intersection

$$\bigcap_{n \in \mathbb{N}} O_n = \{x \mid \sup_\alpha \|A_\alpha x\| = \infty\}$$

is a dense  $G_\delta$  by Corollary 4.2. Otherwise,  $X \setminus \overline{O_n}$  is nonempty and open for one  $n$  and we can find a ball of positive radius  $\overline{B_\varepsilon(x_0)} \subset X \setminus O_n$ . Now observe

$$\|A_\alpha y\| = \|A_\alpha(y + x_0 - x_0)\| \leq \|A_\alpha(y + x_0)\| + \|A_\alpha x_0\| \leq 2n$$

for  $\|y\| \leq \varepsilon$ . Setting  $y = \varepsilon \frac{x}{\|x\|}$ , we obtain

$$\|A_\alpha x\| \leq \frac{2n}{\varepsilon} \|x\|$$

for every  $x$ . □

Note that there is also a variant of the Banach–Steinhaus theorem for pointwise limits of bounded operators which will be discussed in Lemma 4.30.

Hence there are two ways to use this theorem by excluding one of the two possible options. Showing that the pointwise bound holds on a sufficiently large set (e.g. a ball), thereby ruling out the second option, implies a uniform bound and is known as the **uniform boundedness principle**.

**Corollary 4.4.** *Let  $X$  be a Banach space and  $Y$  some normed vector space. Let  $\{A_\alpha\} \subseteq \mathcal{L}(X, Y)$  be a family of bounded operators. Suppose  $\|A_\alpha x\| \leq C(x)$  is bounded for every fixed  $x \in X$ . Then  $\{A_\alpha\}$  is uniformly bounded,  $\|A_\alpha\| \leq C$ .*

Conversely, if there is no uniform bound, the pointwise bound must fail on a dense  $G_\delta$ . This is illustrated in the following example.



**Example 4.3.** Consider the Fourier series (2.45) of a continuous periodic function  $f \in C_{per}[-\pi, \pi] = \{f \in C[-\pi, \pi] | f(-\pi) = f(\pi)\}$ . (Note that this is a closed subspace of  $C[-\pi, \pi]$  and hence a Banach space — it is the kernel of the linear functional  $\ell(f) = f(-\pi) - f(\pi)$ .) We want to show that for every fixed  $x \in [-\pi, \pi]$  there is a dense  $G_\delta$  set of functions in  $C_{per}[-\pi, \pi]$  for which the Fourier series will diverge at  $x$  (it will even be unbounded).

Without loss of generality we fix  $x = 0$  as our point. Then the  $n$ 'th partial sum gives rise to the linear functional

$$\ell_n(f) := S_n(f)(0) = \frac{1}{2\pi} \int_{-\pi}^{\pi} D_n(x) f(x) dx$$

and it suffices to show that the family  $\{\ell_n\}_{n \in \mathbb{N}}$  is not uniformly bounded.

By Example 1.22 (adapted to our present periodic setting) we have

$$\|\ell_n\| = \frac{1}{2\pi} \|D_n\|_1.$$

Now we estimate

$$\begin{aligned} \|D_n\|_1 &= 2 \int_0^\pi |D_n(x)| dx \geq 2 \int_0^\pi \frac{|\sin((n+1/2)x)|}{x/2} dx \\ &= 4 \int_0^{(n+1/2)\pi} |\sin(y)| \frac{dy}{y} \geq 4 \sum_{k=1}^n \int_{(k-1)\pi}^{k\pi} |\sin(y)| \frac{dy}{k\pi} = \frac{8}{\pi} \sum_{k=1}^n \frac{1}{k} \end{aligned}$$

and note that the harmonic series diverges.

In fact, we can even do better. Let  $G(x) \subset C_{per}[-\pi, \pi]$  be the dense  $G_\delta$  of functions whose Fourier series diverges at  $x$ . Then, given countably many points  $\{x_j\}_{j \in \mathbb{N}} \subset [-\pi, \pi]$ , the set  $G = \bigcap_{j \in \mathbb{N}} G(x_j)$  is still a dense  $G_\delta$  by Corollary 4.2. Hence there is a dense  $G_\delta$  of functions whose Fourier series diverges on a given countable set of points.  $\diamond$

**Example 4.4.** Recall that the Fourier coefficients of an absolutely continuous function satisfy the estimate

$$|\hat{f}_k| \leq \begin{cases} \|f\|_\infty, & k = 0, \\ \frac{\|f'\|_\infty}{|k|}, & k \neq 0. \end{cases}$$

This raises the question if a similar estimate can be true for continuous functions. More precisely, can we find a sequence  $c_k > 0$  such that

$$|\hat{f}_k| \leq C_f c_k,$$

where  $C_f$  is some constant depending on  $f$ . If this were true, the linear functionals

$$\ell_k(f) := \frac{\hat{f}_k}{c_k}, \quad k \in \mathbb{Z},$$

satisfy the assumptions of the uniform boundedness principle implying  $\|\ell_k\| \leq C$ . In other words, we must have an estimate of the type

$$|\hat{f}_k| \leq C \|f\|_\infty c_k$$

which implies  $1 \leq C c_k$  upon choosing  $f(x) = e^{ikx}$ . Hence our assumption cannot hold for any sequence  $c_k$  converging to zero and there is no universal decay rate for the Fourier coefficients of continuous functions beyond the fact that they must converge to zero by the Riemann–Lebesgue lemma.  $\diamond$

The next application is

**Theorem 4.5** (Open mapping). *Let  $A \in \mathcal{L}(X, Y)$  be a continuous linear operator between Banach spaces. Then  $A$  is open (i.e., maps open sets to open sets) if and only if it is onto.*

**Proof.** Set  $B_r^X := B_r^X(0)$  and similarly for  $B_r^Y(0)$ . By translating balls (using linearity of  $A$ ), it suffices to prove that for every  $\varepsilon > 0$  there is a  $\delta > 0$  such that  $B_\delta^Y \subseteq A(B_\varepsilon^X)$ .

So let  $\varepsilon > 0$  be given. Since  $A$  is surjective we have

$$Y = AX = A \bigcup_{n=1}^{\infty} nB_\varepsilon^X = \bigcup_{n=1}^{\infty} A(nB_\varepsilon^X) = \bigcup_{n=1}^{\infty} nA(B_\varepsilon^X)$$

and the Baire theorem implies that for some  $n$ ,  $\overline{nA(B_\varepsilon^X)}$  contains a ball. Since multiplication by  $n$  is a homeomorphism, the same must be true for  $n = 1$ , that is,  $B_\delta^Y(y) \subset \overline{A(B_\varepsilon^X)}$ . Consequently

$$B_\delta^Y \subseteq -y + \overline{A(B_\varepsilon^X)} \subset \overline{A(B_\varepsilon^X)} + \overline{A(B_\varepsilon^X)} \subseteq \overline{A(B_\varepsilon^X) + A(B_\varepsilon^X)} \subseteq \overline{A(2B_\varepsilon^X)}.$$

So it remains to get rid of the closure. To this end choose  $\varepsilon_n > 0$  such that  $\sum_{n=1}^{\infty} \varepsilon_n < \varepsilon$  and corresponding  $\delta_n \rightarrow 0$  such that  $B_{\delta_n}^Y \subset \overline{A(B_{\varepsilon_n}^X)}$ . Now for  $y \in B_{\delta_1}^Y \subset \overline{A(B_{\varepsilon_1}^X)}$  we have  $x_1 \in B_{\varepsilon_1}^X$  such that  $Ax_1$  is arbitrarily close to  $y$ , say  $y - Ax_1 \in B_{\delta_2}^Y \subset \overline{A(B_{\varepsilon_2}^X)}$ . Hence we can find  $x_2 \in B_{\varepsilon_2}^X$  such that  $(y - Ax_1) - Ax_2 \in B_{\delta_3}^Y \subset \overline{A(B_{\varepsilon_3}^X)}$  and proceeding like this a sequence  $x_n \in B_{\varepsilon_n}^X$  such that

$$y - \sum_{k=1}^n Ax_k \in B_{\delta_{n+1}}^Y.$$

By construction the limit  $x := \sum_{k=1}^{\infty} x_k$  exists and satisfies  $x \in B_\varepsilon^X$  as well as  $y = Ax \in A(B_\varepsilon^X)$ . That is,  $B_{\delta_1}^Y \subseteq A(B_\varepsilon^X)$  as desired.

Conversely, if  $A$  is open, then the image of the unit ball contains again some ball  $B_\varepsilon^Y \subseteq A(B_1^X)$ . Hence by scaling  $B_{r\varepsilon}^Y \subseteq A(B_r^X)$  and letting  $r \rightarrow \infty$  we see that  $A$  is onto:  $Y = A(X)$ .  $\square$

**Example 4.5.** Let  $X$  be a Banach space and  $M$  a closed subspace. Then the quotient map  $\pi : X \rightarrow X/M$  is open.  $\diamond$

**Example 4.6.** However, note that, under the assumptions of the open mapping theorem, the image of a closed set might not be closed. For example, consider the bounded linear operator  $A : \ell^2(\mathbb{N}) \rightarrow \ell^2(\mathbb{N})$ ,  $x \mapsto (x_2, x_4, \dots)$  which is clearly surjective. Then the image of the closed set  $U = \{x \in \ell^2(\mathbb{N}) \mid x_{2n} = x_{2n-1}/n\}$  is dense (it contains all sequences with finite support) but not all of  $\ell^2(\mathbb{N})$  (e.g.  $y_n = \frac{1}{n}$  is missing since this would imply  $x_{2n} = 1$ ).  $\diamond$

As an immediate consequence we get the inverse mapping theorem:

**Theorem 4.6** (Inverse mapping). *Let  $A \in \mathcal{L}(X, Y)$  be a continuous linear bijection between Banach spaces. Then  $A^{-1}$  is continuous.*

**Example 4.7.** Consider the operator  $(Aa)_{j=1}^n = (\frac{1}{j}a_j)_{j=1}^n$  in  $\ell^2(\mathbb{N})$ . Then its inverse  $(A^{-1}a)_{j=1}^n = (ja_j)_{j=1}^n$  is unbounded (show this!). This is in agreement with our theorem since its range is dense (why?) but not all of  $\ell^2(\mathbb{N})$ : For example,  $(b_j = \frac{1}{j})_{j=1}^\infty \notin \text{Ran}(A)$  since  $b = Aa$  gives the contradiction

$$\infty = \sum_{j=1}^{\infty} 1 = \sum_{j=1}^{\infty} |jb_j|^2 = \sum_{j=1}^{\infty} |a_j|^2 < \infty.$$

This should also be compared with Corollary 8.2 below.  $\diamond$

**Example 4.8.** Consider the Fourier series (2.45) of an integrable function. Using the inverse mapping theorem we can show that not every sequence tending to 0 (which is a necessary condition according to the Riemann–Lebesgue lemma) arises as the Fourier coefficients of an integrable function:

By the elementary estimate

$$\|\hat{f}\|_\infty \leq \frac{1}{2\pi} \|f\|_1$$

we see that that the mapping  $F(f) := \hat{f}$  continuously maps  $F : L^1(-\pi, \pi) \rightarrow c_0(\mathbb{Z})$  (the Banach space of sequences converging to 0). In fact, this estimate holds for continuous functions and hence there is a unique continuous extension of  $F$  to all of  $L^1(-\pi, \pi)$  by Theorem 1.16. Moreover, it can be shown that  $F$  is injective (for  $f \in L^2$  this follows from Theorem 2.18, for the general case  $f \in L^1$  see Example 3.8 from [48]). Now if  $F$  were onto, the inverse mapping theorem would show that the inverse is also continuous, that is, we would have an estimate  $\|\hat{f}\|_\infty \geq C\|f\|_1$  for some  $C > 0$ . However, considering the Dirichlet kernel  $D_n$  we have  $\|\hat{D}_n\|_\infty = 1$  but  $\|D_n\|_1 \rightarrow \infty$  as shown in Example 4.3.  $\diamond$

Another important consequence is the closed graph theorem. The **graph** of an operator  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  between Banach spaces is

$$\Gamma(A) := \{(x, Ax) | x \in \mathfrak{D}(A)\}. \quad (4.1)$$

If  $A$  is linear, the graph is a subspace of the Banach space  $X \oplus Y$ , which is just the Cartesian product together with the norm

$$\|(x, y)\|_{X \oplus Y} := \|x\|_X + \|y\|_Y. \quad (4.2)$$

Note that  $(x_n, y_n) \rightarrow (x, y)$  if and only if  $x_n \rightarrow x$  and  $y_n \rightarrow y$ . We say that  $A$  has a closed graph if  $\Gamma(A)$  is a closed subset of  $X \oplus Y$ .

**Theorem 4.7** (Closed graph). *Let  $A : X \rightarrow Y$  be a linear map from a Banach space  $X$  to another Banach space  $Y$ . Then  $A$  is continuous if and only if its graph is closed.*

**Proof.** If  $\Gamma(A)$  is closed, then it is again a Banach space. Now the projection  $\pi_1(x, Ax) = x$  onto the first component is a continuous bijection onto  $X$ . So by the inverse mapping theorem its inverse  $\pi_1^{-1}$  is again continuous. Moreover, the projection  $\pi_2(x, Ax) = Ax$  onto the second component is also continuous and consequently so is  $A = \pi_2 \circ \pi_1^{-1}$ . The converse is easy.  $\square$

Remark: The crucial fact here is that  $A$  is defined on *all* of  $X$ !

Operators whose graphs are closed are called **closed operators**. Warning: By Example 4.6 a closed operator will not map closed sets to closed sets in general. In particular, the concept of a closed operator should not be confused with the concept of a closed map in topology!

Being closed is the next option you have once an operator turns out to be unbounded. These operators play an important role and we will have a closer look at them in Section 8.1. For now we only point out that the closed graph theorem tells us that closed linear operators can be defined on all of  $X$  if and only if they are bounded. So if we have an unbounded operator we cannot have both! That is, if we want our operator to be at least closed, we have to live with domains. This is the reason why in quantum mechanics most operators are defined on domains. In fact, there is another important property which does not allow unbounded operators to be defined on the entire space:

**Theorem 4.8** (Hellinger–Toeplitz). *Let  $A : \mathfrak{H} \rightarrow \mathfrak{H}$  be a linear operator on some Hilbert space  $\mathfrak{H}$ . If  $A$  is symmetric, that is  $\langle g, Af \rangle = \langle Ag, f \rangle$ ,  $f, g \in \mathfrak{H}$ , then  $A$  is bounded.*

**Proof.** It suffices to prove that  $A$  is closed. In fact,  $f_n \rightarrow f$  and  $Af_n \rightarrow g$  implies

$$\langle h, g \rangle = \lim_{n \rightarrow \infty} \langle h, Af_n \rangle = \lim_{n \rightarrow \infty} \langle Ah, f_n \rangle = \langle Ah, f \rangle = \langle h, Af \rangle$$

for every  $h \in \mathfrak{H}$ . Hence  $Af = g$ .  $\square$

**Problem 4.1.** Every subset of a meager set is again meager.

**Problem 4.2.** An infinite dimensional Banach space cannot have a countable Hamel basis (see Problem 1.7). (Hint: Apply Baire's theorem to  $X_n := \text{span}\{u_j\}_{j=1}^n$ .)

**Problem 4.3.** Let  $X := C[0, 1]$ . Show that the set of functions which are nowhere differentiable contains a dense  $G_\delta$ . (Hint: Consider  $F_k := \{f \in X \mid \exists x \in [0, 1] : |f(x) - f(y)| \leq k|x - y|, \forall y \in [0, 1]\}$ . Show that this set is closed and nowhere dense. For the first property Bolzano–Weierstraß might be useful, for the latter property show that the set of piecewise linear functions whose slopes are bounded below by some fixed number in absolute value are dense.)

**Problem 4.4.** Let  $X$  be a complete metric space without isolated points. Show that a dense  $G_\delta$  set cannot be countable. (Hint: A single point is nowhere dense.)

**Problem 4.5.** Let  $X$  be the space of sequences with finitely many nonzero terms together with the sup norm. Consider the family of operators  $\{A_n\}_{n \in \mathbb{N}}$  given by  $(A_n a)_j := ja_j$ ,  $j \leq n$  and  $(A_n a)_j := 0$ ,  $j > n$ . Then this family is pointwise bounded but not uniformly bounded. Does this contradict the Banach–Steinhaus theorem?

**Problem 4.6.** Show that a bilinear map  $B : X \times Y \rightarrow Z$  is bounded,  $\|B(x, y)\| \leq C\|x\|\|y\|$ , if and only if it is separately continuous with respect to both arguments. (Hint: Uniform boundedness principle.)

**Problem 4.7.** Consider a Schauder basis as in (1.31). Show that the coordinate functionals  $\alpha_n$  are continuous. (Hint: Denote the set of all possible sequences of Schauder coefficients by  $\mathcal{A}$  and equip it with the norm  $\|\alpha\| := \sup_n \|\sum_{k=1}^n \alpha_k u_k\|$ ; note that  $\mathcal{A}$  is precisely the set of sequences for which this norm is finite. By construction the operator  $A : \mathcal{A} \rightarrow X$ ,  $\alpha \mapsto \sum_k \alpha_k u_k$  has norm one. Now show that  $\mathcal{A}$  is complete and apply the inverse mapping theorem.)

**Problem 4.8.** Show that a compact symmetric operator in an infinite-dimensional Hilbert space cannot be surjective.

## 4.2. The Hahn–Banach theorem and its consequences

Let  $X$  be a Banach space. Recall that we have called the set of all bounded linear functionals the dual space  $X^*$  (which is again a Banach space by Theorem 1.17).

**Example 4.9.** Consider the Banach space  $\ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ . Taking the Kronecker deltas  $\delta^n$  as a Schauder basis the  $n$ 'th term  $x_n$  of a sequence  $x \in \ell^p(\mathbb{N})$  can also be considered as the  $n$ 'th coordinate of  $x$  with respect to this basis. Moreover, the map  $l_n(x) = x_n$  is a bounded linear functional, that is,  $l_n \in \ell^p(\mathbb{N})^*$ , since  $|l_n(x)| = |x_n| \leq \|x\|_p$ . It is a special case of the following more general example (in fact, we have  $l_n = l_{\delta^n}$ ). Since the coordinates of a vector carry all the information this explains why understanding linear functionals is of key importance.  $\diamond$

**Example 4.10.** Consider the Banach space  $\ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ . We have already seen that by Hölder's inequality (1.25) every  $y \in \ell^q(\mathbb{N})$  gives rise to a bounded linear functional

$$l_y(x) := \sum_{n \in \mathbb{N}} y_n x_n \quad (4.3)$$

whose norm is  $\|l_y\| = \|y\|_q$  (Problem 4.14). But can every element of  $\ell^p(\mathbb{N})^*$  be written in this form?

Suppose  $p := 1$  and choose  $l \in \ell^1(\mathbb{N})^*$ . Now define

$$y_n := l(\delta^n).$$

Then

$$|y_n| = |l(\delta^n)| \leq \|l\| \|\delta^n\|_1 = \|l\|$$

shows  $\|y\|_\infty \leq \|l\|$ , that is,  $y \in \ell^\infty(\mathbb{N})$ . By construction  $l(x) = l_y(x)$  for every  $x \in \text{span}\{\delta^n\}$ . By continuity of  $l$  it even holds for  $x \in \overline{\text{span}\{\delta^n\}} = \ell^1(\mathbb{N})$ . Hence the map  $y \mapsto l_y$  is an isomorphism, that is,  $\ell^1(\mathbb{N})^* \cong \ell^\infty(\mathbb{N})$ . A similar argument shows  $\ell^p(\mathbb{N})^* \cong \ell^q(\mathbb{N})$ ,  $1 \leq p < \infty$  (Problem 4.15). One usually identifies  $\ell^p(\mathbb{N})^*$  with  $\ell^q(\mathbb{N})$  using this canonical isomorphism and simply writes  $\ell^p(\mathbb{N})^* = \ell^q(\mathbb{N})$ . In the case  $p = \infty$  this is not true, as we will see soon.  $\diamond$

It turns out that many questions are easier to handle after applying a linear functional  $\ell \in X^*$ . For example, suppose  $x(t)$  is a function  $\mathbb{R} \rightarrow X$  (or  $\mathbb{C} \rightarrow X$ ), then  $\ell(x(t))$  is a function  $\mathbb{R} \rightarrow \mathbb{C}$  (respectively  $\mathbb{C} \rightarrow \mathbb{C}$ ) for any  $\ell \in X^*$ . So to investigate  $\ell(x(t))$  we have all tools from real/complex analysis at our disposal. But how do we translate this information back to  $x(t)$ ? Suppose we have  $\ell(x(t)) = \ell(y(t))$  for all  $\ell \in X^*$ . Can we conclude  $x(t) = y(t)$ ? The answer is yes and will follow from the Hahn–Banach theorem.

We first prove the real version from which the complex one then follows easily.

**Theorem 4.9** (Hahn–Banach, real version). *Let  $X$  be a real vector space and  $\varphi : X \rightarrow \mathbb{R}$  a convex function (i.e.,  $\varphi(\lambda x + (1 - \lambda)y) \leq \lambda\varphi(x) + (1 - \lambda)\varphi(y)$  for  $\lambda \in (0, 1)$ ).*

If  $\ell$  is a linear functional defined on some subspace  $Y \subset X$  which satisfies  $\ell(y) \leq \varphi(y)$ ,  $y \in Y$ , then there is an extension  $\bar{\ell}$  to all of  $X$  satisfying  $\bar{\ell}(x) \leq \varphi(x)$ ,  $x \in X$ .

**Proof.** Let us first try to extend  $\ell$  in just one direction: Take  $x \notin Y$  and set  $\tilde{Y} = \text{span}\{x, Y\}$ . If there is an extension  $\tilde{\ell}$  to  $\tilde{Y}$  it must clearly satisfy

$$\tilde{\ell}(y + \alpha x) = \ell(y) + \alpha \tilde{\ell}(x).$$

So all we need to do is to choose  $\tilde{\ell}(x)$  such that  $\tilde{\ell}(y + \alpha x) \leq \varphi(y + \alpha x)$ . But this is equivalent to

$$\sup_{\alpha > 0, y \in Y} \frac{\varphi(y - \alpha x) - \ell(y)}{-\alpha} \leq \tilde{\ell}(x) \leq \inf_{\alpha > 0, y \in Y} \frac{\varphi(y + \alpha x) - \ell(y)}{\alpha}$$

and is hence only possible if

$$\frac{\varphi(y_1 - \alpha_1 x) - \ell(y_1)}{-\alpha_1} \leq \frac{\varphi(y_2 + \alpha_2 x) - \ell(y_2)}{\alpha_2}$$

for every  $\alpha_1, \alpha_2 > 0$  and  $y_1, y_2 \in Y$ . Rearranging this last equations we see that we need to show

$$\alpha_2 \ell(y_1) + \alpha_1 \ell(y_2) \leq \alpha_2 \varphi(y_1 - \alpha_1 x) + \alpha_1 \varphi(y_2 + \alpha_2 x).$$

Starting with the left-hand side we have

$$\begin{aligned} \alpha_2 \ell(y_1) + \alpha_1 \ell(y_2) &= (\alpha_1 + \alpha_2) \ell(\lambda y_1 + (1 - \lambda) y_2) \\ &\leq (\alpha_1 + \alpha_2) \varphi(\lambda y_1 + (1 - \lambda) y_2) \\ &= (\alpha_1 + \alpha_2) \varphi(\lambda(y_1 - \alpha_1 x) + (1 - \lambda)(y_2 + \alpha_2 x)) \\ &\leq \alpha_2 \varphi(y_1 - \alpha_1 x) + \alpha_1 \varphi(y_2 + \alpha_2 x), \end{aligned}$$

where  $\lambda = \frac{\alpha_2}{\alpha_1 + \alpha_2}$ . Hence one dimension works.

To finish the proof we appeal to Zorn's lemma (see Appendix A): Let  $E$  be the collection of all extensions  $\tilde{\ell}$  satisfying  $\tilde{\ell}(x) \leq \varphi(x)$ . Then  $E$  can be partially ordered by inclusion (with respect to the domain) and every linear chain has an upper bound (defined on the union of all domains). Hence there is a maximal element  $\bar{\ell}$  by Zorn's lemma. This element is defined on  $X$ , since if it were not, we could extend it as before contradicting maximality.  $\square$

Note that linearity gives us a corresponding lower bound  $-\varphi(-x) \leq \bar{\ell}(x)$ ,  $x \in X$ , for free. In particular, if  $\varphi(x) = \varphi(-x)$  then  $|\bar{\ell}(x)| \leq \varphi(x)$ .

**Theorem 4.10** (Hahn–Banach, complex version). *Let  $X$  be a complex vector space and  $\varphi : X \rightarrow \mathbb{R}$  a convex function satisfying  $\varphi(\alpha x) \leq \varphi(x)$  if  $|\alpha| = 1$ .*

*If  $\ell$  is a linear functional defined on some subspace  $Y \subset X$  which satisfies  $|\ell(y)| \leq \varphi(y)$ ,  $y \in Y$ , then there is an extension  $\bar{\ell}$  to all of  $X$  satisfying  $|\bar{\ell}(x)| \leq \varphi(x)$ ,  $x \in X$ .*

**Proof.** Set  $\ell_r = \operatorname{Re}(\ell)$  and observe

$$\ell(x) = \ell_r(x) - i\ell_r(ix).$$

By our previous theorem, there is a real linear extension  $\bar{\ell}_r$  satisfying  $\bar{\ell}_r(x) \leq \varphi(x)$ . Now set  $\bar{\ell}(x) = \bar{\ell}_r(x) - i\bar{\ell}_r(ix)$ . Then  $\bar{\ell}(x)$  is real linear and by  $\bar{\ell}(ix) = \bar{\ell}_r(ix) + i\bar{\ell}_r(x) = i\bar{\ell}(x)$  also complex linear. To show  $|\bar{\ell}(x)| \leq \varphi(x)$  we abbreviate  $\alpha = \frac{\bar{\ell}(x)^*}{|\bar{\ell}(x)|}$  and use

$$|\bar{\ell}(x)| = \alpha \bar{\ell}(x) = \bar{\ell}(\alpha x) = \bar{\ell}_r(\alpha x) \leq \varphi(\alpha x) \leq \varphi(x),$$

which finishes the proof.  $\square$

Note that  $\varphi(\alpha x) \leq \varphi(x)$ ,  $|\alpha| = 1$  is in fact equivalent to  $\varphi(\alpha x) = \varphi(x)$ ,  $|\alpha| = 1$ .

If  $\ell$  is a bounded linear functional defined on some subspace, the choice  $\varphi(x) = \|\ell\| \|x\|$  implies:

**Corollary 4.11.** *Let  $X$  be a normed space and let  $\ell$  be a bounded linear functional defined on some subspace  $Y \subseteq X$ . Then there is an extension  $\bar{\ell} \in X^*$  preserving the norm.*

**Example 4.11.** Note that in a Hilbert space this result is trivial: First of all there is a unique extension to  $\bar{Y}$  by Theorem 1.16. Now set  $\bar{\ell} = 0$  on  $Y^\perp$ . Moreover, any other extension is of the form  $\bar{\ell} + \ell_1$ , where  $\ell_1$  vanishes on  $Y$ . Then  $\|\bar{\ell} + \ell_1\|^2 = \|\ell\|^2 + \|\ell_1\|^2$  and the norm will increase unless  $\ell_1 = 0$ .  $\diamond$

**Example 4.12.** In a Banach space this extension will in general not be unique: Consider  $X := \ell^1(\mathbb{N})$  and  $\ell(x) := x_1$  on  $Y := \operatorname{span}\{\delta^1\}$ . Then by Example 4.10 any extension is of the form  $\bar{\ell} = l_y$  with  $y \in \ell^\infty(\mathbb{N})$  and  $y_1 = 1$ ,  $\|y\|_\infty \leq 1$ . (Sometimes it still might be unique: Problems 4.9 and 4.10).  $\diamond$

Moreover, we can now easily prove our anticipated result

**Corollary 4.12.** *Let  $X$  be a normed space and  $x \in X$  fixed. Suppose  $\ell(x) = 0$  for all  $\ell$  in some total subset  $Y \subseteq X^*$ . Then  $x = 0$ .*

**Proof.** Clearly, if  $\ell(x) = 0$  holds for all  $\ell$  in some total subset, this holds for all  $\ell \in X^*$ . If  $x \neq 0$  we can construct a bounded linear functional on  $\operatorname{span}\{x\}$  by setting  $\ell(\alpha x) = \alpha$  and extending it to  $X^*$  using the previous corollary. But this contradicts our assumption.  $\square$

**Example 4.13.** Let us return to our example  $\ell^\infty(\mathbb{N})$ . Let  $c(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$  be the subspace of convergent sequences. Set

$$l(x) = \lim_{n \rightarrow \infty} x_n, \quad x \in c(\mathbb{N}), \quad (4.4)$$

then  $l$  is bounded since

$$|l(x)| = \lim_{n \rightarrow \infty} |x_n| \leq \|x\|_\infty. \quad (4.5)$$



Hence we can extend it to  $\ell^\infty(\mathbb{N})$  by Corollary 4.11. Then  $l(x)$  cannot be written as  $l(x) = l_y(x)$  for some  $y \in \ell^1(\mathbb{N})$  (as in (4.3)) since  $y_n = l(\delta^n) = 0$  shows  $y = 0$  and hence  $\ell_y = 0$ . The problem is that  $\overline{\text{span}\{\delta^n\}} = c_0(\mathbb{N}) \neq \ell^\infty(\mathbb{N})$ , where  $c_0(\mathbb{N})$  is the subspace of sequences converging to 0.

Moreover, there is also no other way to identify  $\ell^\infty(\mathbb{N})^*$  with  $\ell^1(\mathbb{N})$ , since  $\ell^1(\mathbb{N})$  is separable whereas  $\ell^\infty(\mathbb{N})$  is not. This will follow from Lemma 4.17 (iii) below.  $\diamond$

Another useful consequence is

**Corollary 4.13.** *Let  $Y \subseteq X$  be a subspace of a normed vector space and let  $x_0 \in X \setminus \bar{Y}$ . Then there exists an  $\ell \in X^*$  such that (i)  $\ell(y) = 0$ ,  $y \in Y$ , (ii)  $\ell(x_0) = \text{dist}(x_0, Y)$ , and (iii)  $\|\ell\| = 1$ .*

**Proof.** Replacing  $Y$  by  $\bar{Y}$  we see that it is no restriction to assume that  $Y$  is closed. (Note that  $x_0 \in X \setminus \bar{Y}$  if and only if  $\text{dist}(x_0, Y) > 0$ .) Let  $\tilde{Y} = \text{span}\{x_0, Y\}$ . Since every element of  $\tilde{Y}$  can be uniquely written as  $y + \alpha x_0$  we can define

$$\ell(y + \alpha x_0) = \alpha \text{dist}(x_0, Y).$$

By construction  $\ell$  is linear on  $\tilde{Y}$  and satisfies (i) and (ii). Moreover, by  $\text{dist}(x_0, Y) \leq \|x_0 - \frac{-y}{\alpha}\|$  for every  $y \in Y$  we have

$$|\ell(y + \alpha x_0)| = |\alpha| \text{dist}(x_0, Y) \leq \|y + \alpha x_0\|, \quad y \in Y.$$

Hence  $\|\ell\| \leq 1$  and there is an extension to  $X$  by Corollary 4.11. To see that the norm is in fact equal to one, take a sequence  $y_n \in Y$  such that  $\text{dist}(x_0, Y) \geq (1 - \frac{1}{n})\|x_0 + y_n\|$ . Then

$$|\ell(y_n + x_0)| = \text{dist}(x_0, Y) \geq (1 - \frac{1}{n})\|y_n + x_0\|$$

establishing (iii).  $\square$

Two more straightforward consequences of the last corollary are also worthwhile noting:

**Corollary 4.14.** *Let  $Y \subseteq X$  be a subspace of a normed vector space. Then  $x \in \bar{Y}$  if and only if  $\ell(x) = 0$  for every  $\ell \in X^*$  which vanishes on  $Y$ .*

**Corollary 4.15.** *Let  $Y$  be a closed subspace and let  $\{x_j\}_{j=1}^n$  be a linearly independent subset of  $X$ . If  $Y \cap \text{span}\{x_j\}_{j=1}^n = \{0\}$ , then there exists a **biorthogonal system**  $\{\ell_j\}_{j=1}^n \subset X^*$  such that  $\ell_j(x_k) = 0$  for  $j \neq k$ ,  $\ell_j(x_j) = 1$  and  $\ell(y) = 0$  for  $y \in Y$ .*

**Proof.** Fix  $j_0$ . Since  $Y_{j_0} = Y + \text{span}\{x_j\}_{1 \leq j \leq n; j \neq j_0}$  is closed (Corollary 1.19),  $x_{j_0} \notin Y_{j_0}$  implies  $\text{dist}(x_{j_0}, Y_{j_0}) > 0$  and existence of  $\ell_{j_0}$  follows from Corollary 4.13.  $\square$

If we take the **bidual** (or **double dual**)  $X^{**}$  of a normed space  $X$ , then the Hahn–Banach theorem tells us, that  $X$  can be identified with a subspace of  $X^{**}$ . In fact, consider the linear map  $J : X \rightarrow X^{**}$  defined by  $J(x)(\ell) = \ell(x)$  (i.e.,  $J(x)$  is evaluation at  $x$ ). Then

**Theorem 4.16.** *Let  $X$  be a normed space. Then  $J : X \rightarrow X^{**}$  is isometric (norm preserving).*

**Proof.** Fix  $x_0 \in X$ . By  $|J(x_0)(\ell)| = |\ell(x_0)| \leq \|\ell\|_* \|x_0\|$  we have at least  $\|J(x_0)\|_{**} \leq \|x_0\|$ . Next, by Hahn–Banach there is a linear functional  $\ell_0$  with norm  $\|\ell_0\|_* = 1$  such that  $\ell_0(x_0) = \|x_0\|$ . Hence  $|J(x_0)(\ell_0)| = |\ell_0(x_0)| = \|x_0\|$  shows  $\|J(x_0)\|_{**} = \|x_0\|$ .  $\square$

**Example 4.14.** This gives another quick way of showing that a normed space has a completion: Take  $\bar{X} := \overline{J(X)} \subseteq X^{**}$  and recall that a dual space is always complete (Theorem 1.17).  $\diamond$

Thus  $J : X \rightarrow X^{**}$  is an isometric embedding. In many cases we even have  $J(X) = X^{**}$  and  $X$  is called **reflexive** in this case.

**Example 4.15.** The Banach spaces  $\ell^p(\mathbb{N})$  with  $1 < p < \infty$  are reflexive: Identify  $\ell^p(\mathbb{N})^*$  with  $\ell^q(\mathbb{N})$  (cf. Problem 4.15) and choose  $z \in \ell^p(\mathbb{N})^{**}$ . Then there is some  $x \in \ell^p(\mathbb{N})$  such that

$$z(y) = \sum_{j \in \mathbb{N}} y_j x_j, \quad y \in \ell^q(\mathbb{N}) \cong \ell^p(\mathbb{N})^*.$$

But this implies  $z(y) = y(x)$ , that is,  $z = J(x)$ , and thus  $J$  is surjective. (Warning: It does not suffice to just argue  $\ell^p(\mathbb{N})^{**} \cong \ell^q(\mathbb{N})^* \cong \ell^p(\mathbb{N})$ .)

However,  $\ell^1$  is not reflexive since  $\ell^1(\mathbb{N})^* \cong \ell^\infty(\mathbb{N})$  but  $\ell^\infty(\mathbb{N})^* \not\cong \ell^1(\mathbb{N})$  as noted earlier. Things get even a bit more explicit if we look at  $c_0(\mathbb{N})$ , where we can identify (cf. Problem 4.16)  $c_0(\mathbb{N})^*$  with  $\ell^1(\mathbb{N})$  and  $c_0(\mathbb{N})^{**}$  with  $\ell^\infty(\mathbb{N})$ . Under this identification  $J(c_0(\mathbb{N})) = c_0(\mathbb{N}) \subseteq \ell^\infty(\mathbb{N})$ .  $\diamond$

**Example 4.16.** By the same argument, every Hilbert space is reflexive. In fact, by the Riesz lemma we can identify  $\mathfrak{H}^*$  with  $\mathfrak{H}$  via the (conjugate linear) map  $x \mapsto \langle x, \cdot \rangle$ . Taking  $z \in \mathfrak{H}^{**}$  we have, again by the Riesz lemma, that  $z(y) = \langle \langle x, \cdot \rangle, \langle y, \cdot \rangle \rangle_{\mathfrak{H}^*} = \langle x, y \rangle^* = \langle y, x \rangle = J(x)(y)$ .  $\diamond$

**Example 4.17.** The sum of reflexive spaces is reflexive. Indeed, recall  $(X \oplus Y)^* \cong X^* \oplus Y^*$  with  $(x', y')(x, y) := x'(x) + y'(y)$  and hence also  $(X \oplus Y)^{**} \cong X^{**} \oplus Y^{**}$  with  $(x'', y'')(x', y') := x''(x') + y''(y')$ . Hence for  $(x'', y'') \in (X \oplus Y)^{**}$  there is  $x \in X$  and  $y \in Y$  such that  $(x'', y'') = (J_X(x), J_Y(y)) = J(x, y)$  and hence  $J$  is surjective. This even extends to countable sums — Problem 4.18.  $\diamond$

**Lemma 4.17.** *Let  $X$  be a Banach space.*

- (i) *If  $X$  is reflexive, so is every closed subspace.*

- (ii)  $X$  is reflexive if and only if  $X^*$  is.
- (iii) If  $X^*$  is separable, so is  $X$ .

**Proof.** (i) Let  $Y$  be a closed subspace. Denote by  $j : Y \hookrightarrow X$  the natural inclusion and define  $j_{**} : Y^{**} \rightarrow X^{**}$  via  $(j_{**}(y''))(\ell) = y''(\ell|_Y)$  for  $y'' \in Y^{**}$  and  $\ell \in X^*$ . Note that  $j_{**}$  is isometric by Corollary 4.11. Then

$$\begin{array}{ccc} X & \xrightarrow{J_X} & X^{**} \\ j \uparrow & & \uparrow j_{**} \\ Y & \xrightarrow{J_Y} & Y^{**} \end{array}$$

commutes. In fact, we have  $j_{**}(J_Y(y))(\ell) = J_Y(y)(\ell|_Y) = \ell(y) = J_X(y)(\ell)$ . Moreover, since  $J_X$  is surjective, for every  $y'' \in Y^{**}$  there is an  $x \in X$  such that  $j_{**}(y'') = J_X(x)$ . Since  $j_{**}(y'')(\ell) = y''(\ell|_Y)$  vanishes on all  $\ell \in X^*$  which vanish on  $Y$ , so does  $\ell(x) = J_X(x)(\ell) = j_{**}(y'')(\ell)$  and thus  $x \in Y$  by Corollary 4.14. That is,  $j_{**}(Y^{**}) = J_X(Y)$  and  $J_Y = j \circ J_X \circ j_{**}^{-1}$  is surjective.

(ii) Suppose  $X$  is reflexive. Then the two maps

$$\begin{array}{ccc} (J_X)_* : X^* & \rightarrow & X^{***} \\ x' & \mapsto & x' \circ J_X^{-1} \end{array} \quad \begin{array}{ccc} (J_X)^* : X^{***} & \rightarrow & X^* \\ x''' & \mapsto & x''' \circ J_X \end{array}$$

are inverse of each other. Moreover, fix  $x'' \in X^{**}$  and let  $x = J_X^{-1}(x'')$ . Then  $J_X^*(x')(x'') = x''(x') = J(x)(x') = x'(x) = x'(J_X^{-1}(x''))$ , that is  $J_X^* = (J_X)_*$  respectively  $(J_X^*)^{-1} = (J_X)_*$ , which shows  $X^*$  reflexive if  $X$  reflexive. To see the converse, observe that  $X^*$  reflexive implies  $X^{**}$  reflexive and hence  $J_X(X) \cong X$  is reflexive by (i).

(iii) Let  $\{\ell_n\}_{n=1}^\infty$  be a dense set in  $X^*$ . Then we can choose  $x_n \in X$  such that  $\|x_n\| = 1$  and  $\ell_n(x_n) \geq \|\ell_n\|/2$ . We will show that  $\{x_n\}_{n=1}^\infty$  is total in  $X$ . If it were not, we could find some  $x \in X \setminus \overline{\text{span}\{x_n\}_{n=1}^\infty}$  and hence there is a functional  $\ell \in X^*$  as in Corollary 4.13. Choose a subsequence  $\ell_{n_k} \rightarrow \ell$ . Then

$$\|\ell - \ell_{n_k}\| \geq |(\ell - \ell_{n_k})(x_{n_k})| = |\ell_{n_k}(x_{n_k})| \geq \|\ell_{n_k}\|/2,$$

which implies  $\ell_{n_k} \rightarrow 0$  and contradicts  $\|\ell\| = 1$ .  $\square$

If  $X$  is reflexive, then the converse of (iii) is also true (since  $X \cong X^{**}$  separable implies  $X^*$  separable), but in general this fails as the example  $\ell^1(\mathbb{N})^* \cong \ell^\infty(\mathbb{N})$  shows. In fact, this can be used to show that a separable space is not reflexive, by showing that its dual is not separable.

**Example 4.18.** The space  $C(I)$  is not reflexive. To see this observe that the dual space contains point evaluations  $\ell_{x_0}(f) := f(x_0)$ ,  $x_0 \in I$ . Moreover, for  $x_0 \neq x_1$  we have  $\|\ell_{x_0} - \ell_{x_1}\| = 2$  and hence  $C(I)^*$  is not separable. You should appreciate the fact that it was not necessary to know the full dual space which is quite intricate (see Theorem 6.5 from [48]).  $\diamond$

Finally we discuss the analog of the orthogonal complement of a set. Given subsets  $M \subseteq X$  and  $N \subseteq X^*$  we define their **annihilator** as

$$\begin{aligned} M^\perp &:= \{\ell \in X^* \mid \ell(x) = 0 \ \forall x \in M\} = \{\ell \in X^* \mid M \subseteq \text{Ker}(\ell)\} \\ &= \bigcap_{x \in M} \{\ell \in X^* \mid \ell(x) = 0\} = \bigcap_{x \in M} \{x\}^\perp, \\ N_\perp &:= \{x \in X \mid \ell(x) = 0 \ \forall \ell \in N\} = \bigcap_{\ell \in N} \text{Ker}(\ell) = \bigcap_{\ell \in N} \{\ell\}_\perp. \end{aligned} \quad (4.6)$$

In particular,  $\{\ell\}_\perp = \text{Ker}(\ell)$  while  $\{x\}^\perp = \text{Ker}(J(x))$  (with  $J : X \hookrightarrow X^{**}$  the canonical embedding). Note  $\{0\}^\perp = X^*$  and  $\{0\}_\perp = X$ .

**Example 4.19.** In a Hilbert space the annihilator is simply the orthogonal complement.  $\diamond$

The following properties are immediate from the definition (by linearity and continuity):

- $M^\perp$  is a closed subspace of  $X^*$  and  $M^\perp = \overline{\text{span}(M)}^\perp$ .
- $N_\perp$  is a closed subspace of  $X$  and  $N_\perp = \overline{\text{span}(N)}_\perp$ .

Note that we can also consider  $N^\perp \subseteq X^{**}$  and that we have  $J(N_\perp) \subseteq N^\perp$  with equality if  $X$  is reflexive. Similarly we have  $J(M)_\perp = M^\perp$ .

Note also that

$$\begin{aligned} \overline{\text{span}(M)} = X &\Leftrightarrow M^\perp = \{0\}, \\ \overline{\text{span}(N)} = X^* &\Rightarrow N_\perp = \{0\} \end{aligned} \quad (4.7)$$

by Corollary 4.13 and Corollary 4.12, respectively. The converse of the last statement is wrong in general (unless  $X$  is reflexive, see the following lemma).

**Example 4.20.** Consider  $X := \ell^1(\mathbb{N})$  and  $N := \{\delta^n\}_{n \in \mathbb{N}} \subset \ell^\infty(\mathbb{N}) \cong X^*$ . Then  $\overline{\text{span}(N)} = c_0(\mathbb{N})$  but  $N_\perp = \{0\}$ .  $\diamond$

**Lemma 4.18.** We have  $(M^\perp)_\perp = \overline{\text{span}(M)}$  and  $(N_\perp)^\perp \supseteq \overline{\text{span}(N)}$  with equality if  $X$  is reflexive.

**Proof.** By the preceding remarks we can assume  $M, N$  to be closed subspaces. The first part

$$(M^\perp)_\perp = \{x \in X \mid \ell(x) = 0 \ \forall \ell \in X^* \text{ with } M \subseteq \text{Ker}(\ell)\} = \overline{\text{span}(M)}$$

is Corollary 4.14 and for the second part one just has to spell out the definition:

$$(N_\perp)^\perp = \{\ell \in X^* \mid \bigcap_{\tilde{\ell} \in N} \text{Ker}(\tilde{\ell}) \subseteq \text{Ker}(\ell)\} \supseteq \overline{\text{span}(N)}.$$

If  $X$  is reflexive we can use the first part to conclude

$$(N_\perp)^\perp = (J(N_\perp))_\perp = (N^\perp)_\perp = \overline{\text{span}(N)}. \quad \square$$

Note that we also have equality in the preceding lemma if  $N$  is finite dimensional (Problem 4.22). For non-reflexive spaces the inclusion can be strict as the previous example shows. Moreover, with a little more machinery one can identify  $(N_\perp)^\perp$  as the weak-\* closure of  $\overline{\text{span}(N)}$  (Problem 6.20).

**Warning:** Some authors call a set  $N \subseteq X^*$  total if  $N_\perp = \{0\}$ . By the preceding discussion this is equivalent to our definition if  $X$  is reflexive, but otherwise might differ.

With the help of annihilators we can also describe the dual spaces of subspaces.

**Theorem 4.19.** *Let  $M$  be a closed subspace of a normed space  $X$ . Then there are canonical isometries*

$$(X/M)^* \cong M^\perp, \quad M^* \cong X^*/M^\perp. \quad (4.8)$$

**Proof.** In the first case the isometry is given by  $\ell \mapsto \ell \circ j$ , where  $j : X \rightarrow X/M$  is the quotient map. In the second case  $x' + M^\perp \mapsto x'|_M$ . The details are easy to check.  $\square$

**Problem 4.9.** Let  $X := \mathbb{C}^3$  equipped with the norm  $|(x, y, z)|_1 := |x| + |y| + |z|$  and  $Y := \{(x, y, z) | x + y = 0, z = 0\}$ . Find at least two extensions of  $\ell(x, y, z) := x$  from  $Y$  to  $X$  which preserve the norm. What if we take  $Y := \{(x, y, z) | x + y = 0\}$ ?

**Problem 4.10.** Show that the extension from Corollary 4.11 is unique if  $X^*$  is strictly convex. (Hint: Problem 1.13.)

**Problem\* 4.11.** Let  $X$  be some normed space. Show that

$$\|x\| = \sup_{\ell \in V, \|\ell\|=1} |\ell(x)|, \quad (4.9)$$

where  $V \subset X^*$  is some dense subspace. Show that equality is attained if  $V = X^*$ .

**Problem 4.12.** Let  $X$  be some normed space. By definition we have

$$\|\ell\| = \sup_{x \in X, \|x\|=1} |\ell(x)|$$

for every  $\ell \in X^*$ . One calls  $\ell \in X^*$  **norm-attaining**, if the supremum is attained, that is, there is some  $x \in X$  such that  $\|\ell\| = |\ell(x)|$ .

Show that in a reflexive Banach space every linear functional is norm-attaining. Give an example of a linear functional which is not norm-attaining. For uniqueness see Problem 6.39. (Hint: For the first part apply the previous problem to  $X^*$ . For the second part consider Problem 4.16 below.)

**Problem 4.13.** Let  $X, Y$  be some normed spaces and  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$ . Show

$$\|A\| = \sup_{x \in X, \|x\|=1; \ell \in V, \|\ell\|=1} |\ell(Ax)|, \quad (4.10)$$

where  $V \subset Y^*$  is a dense subspace.

**Problem\* 4.14.** Show that  $\|l_y\| = \|y\|_q$ , where  $l_y \in \ell^p(\mathbb{N})^*$  as defined in (4.3). (Hint: Choose  $x \in \ell^p$  such that  $x_n y_n = |y_n|^q$ .)

**Problem\* 4.15.** Show that every  $l \in \ell^p(\mathbb{N})^*$ ,  $1 \leq p < \infty$ , can be written as

$$l(x) = \sum_{n \in \mathbb{N}} y_n x_n$$

with some  $y \in \ell^q(\mathbb{N})$ . (Hint: To see  $y \in \ell^q(\mathbb{N})$  consider  $x^N$  defined such that  $x_n^N = |y_n|^q / y_n$  for  $n \leq N$  with  $y_n \neq 0$  and  $x_n^N = 0$  else. Now look at  $|l(x^N)| \leq \|l\| \|x^N\|_p$ .)

**Problem\* 4.16.** Let  $c_0(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$  be the subspace of sequences which converge to 0, and  $c(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$  the subspace of convergent sequences.

- (i) Show that  $c_0(\mathbb{N})$ ,  $c(\mathbb{N})$  are both Banach spaces and that  $c(\mathbb{N}) = \text{span}\{c_0(\mathbb{N}), e\}$ , where  $e = (1, 1, 1, \dots) \in c(\mathbb{N})$ .
- (ii) Show that every  $l \in c_0(\mathbb{N})^*$  can be written as

$$l(a) = \sum_{n \in \mathbb{N}} b_n a_n$$

with some  $b \in \ell^1(\mathbb{N})$  which satisfies  $\|b\|_1 = \|\ell\|$ .

- (iii) Show that every  $l \in c(\mathbb{N})^*$  can be written as

$$l(a) = \sum_{n \in \mathbb{N}} b_n a_n + b_0 \lim_{n \rightarrow \infty} a_n$$

with some  $b \in \ell^1(\mathbb{N})$  which satisfies  $|b_0| + \|b\|_1 = \|\ell\|$ .

**Problem 4.17.** Let  $u_n \in X$  be a Schauder basis and suppose the complex numbers  $c_n$  satisfy  $|c_n| \leq c \|u_n\|$ . Is there a bounded linear functional  $\ell \in X^*$  with  $\ell(u_n) = c_n$ ? (Hint: Consider e.g.  $X = \ell^2(\mathbb{Z})$ .)

**Problem\* 4.18.** Let  $X := \bigoplus_{p,j \in \mathbb{N}} X_j$  be defined as in Problem 1.42 and let  $\frac{1}{p} + \frac{1}{q} = 1$ . Show that for  $1 \leq p < \infty$  we have  $X^* \cong \bigoplus_{q,j \in \mathbb{N}} X_j^*$ , where the identification is given by

$$y(x) = \sum_{j \in \mathbb{N}} y_j(x_j), \quad x = (x_j)_{j \in \mathbb{N}} \in \bigoplus_{p,j \in \mathbb{N}} X_j, \quad y = (y_j)_{j \in \mathbb{N}} \in \bigoplus_{q,j \in \mathbb{N}} X_j^*.$$

Moreover, if all  $X_j$  are reflexive, so is  $X$  for  $1 < p < \infty$ .

**Problem 4.19** (Banach limit). Let  $\mathfrak{c}(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$  be the subspace of all bounded sequences for which the limit of the Cesàro means

$$L(x) := \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n x_k$$

exists. Note that  $c(\mathbb{N}) \subseteq \mathfrak{c}(\mathbb{N})$  and  $L(x) = \lim_{n \rightarrow \infty} x_n$  for  $x \in c(\mathbb{N})$ .

Show that  $L$  can be extended to all of  $\ell^\infty(\mathbb{N})$  such that

- (i)  $L$  is linear,
- (ii)  $|L(x)| \leq \|x\|_\infty$ ,
- (iii)  $L(Sx) = L(x)$  where  $(Sx)_n = x_{n+1}$  is the shift operator,
- (iv)  $L(x) \geq 0$  when  $x_n \geq 0$  for all  $n$ ,
- (v)  $\liminf_n x_n \leq L(x) \leq \limsup_n x_n$  for all real-valued sequences.

(Hint: Of course existence follows from Hahn–Banach and (i), (ii) will come for free. Also (iii) will be inherited from the construction. For (iv) note that the extension can be assumed to be real-valued and investigate  $L(e - x)$  for  $x \geq 0$  with  $\|x\|_\infty = 1$  where  $e = (1, 1, 1, \dots)$ . (v) then follows from (iv).)

**Problem\* 4.20.** Show that a finite dimensional subspace  $M$  of a Banach space  $X$  can be complemented. (Hint: Start with a basis  $\{x_j\}$  for  $M$  and choose a corresponding dual basis  $\{\ell_k\}$  with  $\ell_k(x_j) = \delta_{j,k}$  which can be extended to  $X^*$ .)

**Problem 4.21.** Suppose  $X$  is separable. Show that there exists a countable set  $N \subset X^*$  with  $N_\perp = \{0\}$ .

**Problem\* 4.22.** Suppose  $X$  is a vector space and  $\ell, \ell_1, \dots, \ell_n$  are linear functionals such that  $\bigcap_{j=1}^n \text{Ker}(\ell_j) \subseteq \text{Ker}(\ell)$ . Then  $\ell = \sum_{j=1}^n \alpha_j \ell_j$  for some constants  $\alpha_j \in \mathbb{C}$ . (Hint: Find a dual basis  $x_k \in X$  such that  $\ell_j(x_k) = \delta_{j,k}$  and look at  $x - \sum_{j=1}^n \ell_j(x)x_j$ .)

**Problem 4.23.** Suppose  $M_1, M_2$  are closed subspaces of  $X$ . Show

$$M_1 \cap M_2 = (M_1^\perp + M_2^\perp)^\perp, \quad M_1^\perp \cap M_2^\perp = (M_1 + M_2)^\perp$$

and

$$(M_1 \cap M_2)^\perp \supseteq \overline{(M_1^\perp + M_2^\perp)}, \quad (M_1^\perp \cap M_2^\perp)^\perp = \overline{(M_1 + M_2)}.$$

### 4.3. The adjoint operator

Given two normed spaces  $X$  and  $Y$  and a bounded operator  $A \in \mathcal{L}(X, Y)$  we can define its **adjoint**  $A' : Y^* \rightarrow X^*$  via  $A'y' = y' \circ A$ ,  $y' \in Y^*$ . It is

immediate that  $A'$  is linear and boundedness follows from

$$\begin{aligned}\|A'\| &= \sup_{y' \in Y^*: \|y'\|=1} \|A'y'\| = \sup_{y' \in Y^*: \|y'\|=1} \left( \sup_{x \in X: \|x\|=1} |(A'y')(x)| \right) \\ &= \sup_{y' \in Y^*: \|y'\|=1} \left( \sup_{x \in X: \|x\|=1} |y'(Ax)| \right) = \sup_{x \in X: \|x\|=1} \|Ax\| = \|A\|,\end{aligned}$$

where we have used Problem 4.11 to obtain the fourth equality. In summary,

**Theorem 4.20.** *Suppose  $X, Y$  are Banach spaces. Let  $A \in \mathcal{L}(X, Y)$ , then  $A' \in \mathcal{L}(Y^*, X^*)$  with  $\|A\| = \|A'\|$ .*

Note that for  $A, B \in \mathcal{L}(X, Y)$  and  $\alpha, \beta \in \mathbb{C}$  we have

$$(\alpha A + \beta B)' = \alpha A' + \beta B' \quad (4.11)$$

and for  $A \in \mathcal{L}(X, Y)$  and  $B \in \mathcal{L}(Y, Z)$  we have

$$(BA)' = A'B' \quad (4.12)$$

which is immediate from the definition.

**Example 4.21.** Given a Hilbert space  $\mathfrak{H}$  we have the conjugate linear isometry  $C : \mathfrak{H} \rightarrow \mathfrak{H}^*$ ,  $f \mapsto \langle f, \cdot \rangle$ . Hence for given  $A \in \mathcal{L}(\mathfrak{H}_1, \mathfrak{H}_2)$  we have  $A'C_2f = \langle f, A \cdot \rangle = \langle A^*f, \cdot \rangle$  which shows  $A' = C_1A^*C_2^{-1}$ .  $\diamond$

**Example 4.22.** Let  $X := Y := \ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ , such that  $X^* = \ell^q(\mathbb{N})$ ,  $\frac{1}{p} + \frac{1}{q} = 1$ . Consider the right shift  $R \in \mathcal{L}(\ell^p(\mathbb{N}))$  given by

$$Rx := (0, x_1, x_2, \dots).$$

Then for  $y' \in \ell^q(\mathbb{N})$

$$y'(Rx) = \sum_{j=1}^{\infty} y'_j(Rx)_j = \sum_{j=2}^{\infty} y'_j x_{j-1} = \sum_{j=1}^{\infty} y'_{j+1} x_j$$

which shows  $(R'y')_k = y'_{k+1}$  upon choosing  $x = \delta^k$ . Hence  $R' = L$  is the left shift:  $Ly := (y_2, y_3, \dots)$ . Similarly,  $L' = R$ .  $\diamond$

**Example 4.23.** Consider the multiplication operator  $A$  in  $\ell^p(\mathbb{N})$  with  $1 \leq p < \infty$  defined by  $Aa_j := \frac{1}{j}a_j$ . Then  $X^* \cong \ell^q(\mathbb{N})$  with  $\frac{1}{q} + \frac{1}{p} = 1$ . Given  $b \in \mathfrak{D}(A')$  we need to find the  $c \in \ell^q(\mathbb{N})$  such that  $\ell_c(a) = \ell_b(Aa)$  for all  $a \in \ell^p(\mathbb{N})$ . That is,

$$\sum_{j=1}^{\infty} c_j a_j = \sum_{j=1}^{\infty} b_j \left( \frac{1}{j} a_j \right), \quad a \in \ell^p(\mathbb{N}).$$

Choosing  $a = \delta^j$  shows  $c_j = \frac{1}{j}b_j$  and hence  $A'b_j = \frac{1}{j}b_j$ .  $\diamond$



**Example 4.24.** Let us compute the adjoint of  $B$  from Example 8.4. Proceeding as in the previous example we get  $c_j = b_1$  which is in  $\ell^2$  if and only if  $b_1 = 0$ . Thus  $\mathfrak{D}(B') = \{b \in \ell^2(\mathbb{N}) | b_1 = 0\}$  and  $B'b = 0$ . So don't expect the adjoint of a noncloseable operator to contain much information about the operator.  $\diamond$

**Example 4.25.** Recall that an operator  $K \in \mathcal{L}(X, Y)$  is called a **finite rank operator** if its range is finite dimensional. The dimension of its range  $\text{rank}(K) := \dim \text{Ran}(K)$  is called the **rank** of  $K$ . Choosing a basis  $\{y_j = Kx_j\}_{j=1}^n$  for  $\text{Ran}(K)$  and a corresponding dual basis  $\{y'_j\}_{j=1}^n$  (cf. Problem 4.20), then  $x'_j := K'y'_j$  is a dual basis for  $x_j$  and

$$Kx = \sum_{j=1}^n y'_j(Kx)y_j = \sum_{j=1}^n x'_j(x)y_j, \quad K'y' = \sum_{j=1}^n y'(y_j)x'_j.$$

In particular,  $\text{rank}(K) = \text{rank}(K')$ .  $\diamond$

Of course we can also consider the doubly adjoint operator  $A''$ . Then a simple computation

$$A''(J_X(x))(y') = J_X(x)(A'y') = (A'y')(x) = y'(Ax) = J_Y(Ax)(y') \quad (4.13)$$

shows that the following diagram commutes

$$\begin{array}{ccc} X & \xrightarrow{A} & Y \\ J_X \downarrow & & \downarrow J_Y \\ X^{**} & \xrightarrow{A''} & Y^{**} \end{array}$$

Consequently

$$A'' \upharpoonright_{\text{Ran}(J_X)} = J_Y A J_X^{-1}, \quad A = J_Y^{-1} A'' J_X. \quad (4.14)$$

Hence, regarding  $X$  as a subspace  $J_X(X) \subseteq X^{**}$  and  $Y$  as a subspace  $J_Y(Y) \subseteq Y^{**}$ , then  $A''$  is an extension of  $A$  to  $X^{**}$  but with values in  $Y^{**}$ . In particular, note that  $B \in \mathcal{L}(Y^*, X^*)$  is the adjoint of some other operator  $B = A'$  if and only if  $B'(J_X(X)) = A''(J_X(X)) \subseteq J_Y(Y)$  (for the converse note that  $A := J_Y^{-1} B' J_X$  will do the trick). This can be used to show that not every operator is an adjoint (Problem 4.24).

**Theorem 4.21** (Schauder). *Suppose  $X, Y$  are Banach spaces and  $A \in \mathcal{L}(X, Y)$ . Then  $A$  is compact if and only if  $A'$  is.*

**Proof.** If  $A$  is compact, then  $A(B_1^X(0))$  is relatively compact and hence  $K = \overline{A(B_1^X(0))}$  is a compact metric space. Let  $y'_n \in Y^*$  be a bounded sequence and consider the family of functions  $f_n := y'_n|_K \in C(K)$ . Then this family is bounded and equicontinuous since

$$|f_n(y_1) - f_n(y_2)| \leq \|y'_n\| \|y_1 - y_2\| \leq C \|y_1 - y_2\|.$$

Hence the Arzelà–Ascoli theorem (Theorem B.40) implies existence of a uniformly converging subsequence  $f_{n_j}$ . For this subsequence we have

$$\|A'y'_{n_j} - A'y'_{n_k}\| \leq \sup_{x \in B_1^X(0)} |y'_{n_j}(Ax) - y'_{n_k}(Ax)| = \|f_{n_j} - f_{n_k}\|_\infty$$

since  $A(B_1^X(0)) \subseteq K$  is dense. Thus  $y'_{n_j}$  is the required subsequence and  $A'$  is compact.

To see the converse note that if  $A'$  is compact then so is  $A''$  by the first part and hence also  $A = J_Y^{-1}A''J_X$ .  $\square$

**Theorem 4.22.** *Suppose  $X, Y$  are Banach spaces. If  $A \in \mathcal{L}(X, Y)$ , then  $A^{-1}$  exists and is in  $\mathcal{L}(Y, X)$  if and only if  $(A')^{-1}$  exists and is in  $\mathcal{L}(X^*, Y^*)$ . Moreover, in this case we have*

$$(A')^{-1} = (A^{-1})'. \quad (4.15)$$

**Proof.** If  $A$  is invertible, then  $A'(A^{-1})' = (A^{-1}A)' = (\mathbb{I}_X)' = \mathbb{I}_{X^*}$  and  $(A^{-1})'A' = (A^{-1}A)' = (\mathbb{I}_Y)' = \mathbb{I}_{Y^*}$  shows that  $A'$  is invertible with  $(A')^{-1} = (A^{-1})'$ .

Conversely, let  $(A')^{-1} \in \mathcal{L}(X^*, Y^*)$ . Then by the first part  $(A'')^{-1}$  exists and is in  $\mathcal{L}(X^{**}, Y^{**})$ . Moreover,  $A^{-1} = J_X^{-1}(A'')^{-1}J_Y \in \mathcal{L}(Y, X)$ .  $\square$

Finally we discuss the relation between solvability of  $Ax = y$  and the corresponding adjoint equation  $A'y' = x'$ . We begin with the following analog of (2.28) (Problem 4.26):

**Lemma 4.23.** *If  $A \in \mathcal{L}(X, Y)$ , then  $\text{Ran}(A)^\perp = \text{Ker}(A')$  and  $\text{Ran}(A')_\perp = \text{Ker}(A)$ .*

Taking annihilators in these formulas we obtain

$$\text{Ker}(A')_\perp = (\text{Ran}(A)^\perp)_\perp = \overline{\text{Ran}(A)} \quad (4.16)$$

and

$$\text{Ker}(A)^\perp = (\text{Ran}(A')_\perp)^\perp \supseteq \overline{\text{Ran}(A')} \quad (4.17)$$

which raises the question of equality in the latter.

Note that the first identity tells us that, for an operator  $A$  with closed range, a necessary and sufficient solvability criterion for the equation  $Ax = y$  is  $y \in \text{Ker}(A')_\perp$  (that is,  $\ell(y) = 0$  for all  $\ell \in \text{Ker}(A')$ ). The second identity would imply an analogous criterion for the adjoint equation.

**Example 4.26.** Consider  $A$  from Example 4.23 in the case  $p = 1$ . Then  $\text{Ran}(A) = \{b \in \ell^1(\mathbb{N}) \mid (jb_j)_{j=1}^\infty \in \ell^1(\mathbb{N})\}$  is dense while  $\text{Ran}(A') = \{b \in \ell^\infty(\mathbb{N}) \mid (jb_j)_{j=1}^\infty \in \ell^\infty(\mathbb{N})\}$  is not dense since  $\overline{\text{Ran}(A')} = c_0(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$ . In particular,  $\text{Ker}(A)^\perp = \{0\}^\perp = \ell^\infty(\mathbb{N}) \supset \overline{\text{Ran}(A')} = c_0(\mathbb{N})$ . Note that in this example neither range is closed.  $\diamond$

If the range of  $A$  or  $A'$  is closed, these problems do not occur:

**Theorem 4.24** (Banach; Closed range). *Suppose  $X, Y$  are Banach spaces and  $A \in \mathcal{L}(X, Y)$ . Then the following items are equivalent:*

- (i)  $\text{Ran}(A)$  is closed.
- (ii)  $\text{Ker}(A)^\perp = \text{Ran}(A')$ .
- (iii)  $\text{Ran}(A')$  is closed.
- (iv)  $\text{Ker}(A')^\perp = \text{Ran}(A)$ .

**Proof.** Consider  $\tilde{X} = X/\text{Ker}(A)$  and  $\tilde{Y} = \overline{\text{Ran}(A)}$  and the corresponding operator  $\tilde{A}$  as in Problem 1.46. Then  $\tilde{A}$  is a bounded injective operator whose range is dense. In particular,  $\text{Ran}(A)$  is closed if and only if  $\tilde{A}^{-1} \in \mathcal{L}(\tilde{Y}, \tilde{X})$ . Moreover,  $\text{Ran}(\tilde{A}') \subseteq (X/\text{Ker}(A))^* \cong \text{Ker}(A)^\perp$  (cf. Theorem 4.19) and the canonical isometry (i.e. composition with the quotient map) will map  $\text{Ran}(A') \subseteq \text{Ker}(A)^\perp$  to  $\text{Ran}(\tilde{A}')$ . In particular,  $\text{Ran}(A')$  is closed if and only if  $(\tilde{A}')^{-1} \in \mathcal{L}(\tilde{X}^*, \tilde{Y}^*)$ .

Hence (i)  $\Leftrightarrow$  (iii) follows from Theorem 4.22 applied to  $\tilde{A}$ . (ii)  $\Rightarrow$  (iii) is clear since annihilators are closed. (i)  $\Leftrightarrow$  (iv) is immediate from (4.16).

(i)  $\Rightarrow$  (ii): Note that if  $\ell \in \text{Ran}(A')$  then  $\ell = A'(\tilde{\ell}) = \tilde{\ell} \circ A$  vanishes on  $\text{Ker}(A)$  and hence  $\ell \in \text{Ker}(A)^\perp$ . Conversely, if  $\ell \in \text{Ker}(A)^\perp$  it is well-defined on  $\tilde{X}$  and we can set  $\tilde{\ell}(y) = \ell(\tilde{A}^{-1}y)$  for  $y \in \text{Ran}(A)$  and extend it to all of  $\tilde{Y}$  using Corollary 4.11. By construction  $\ell = A'(\tilde{\ell}) \in \text{Ran}(A')$ .  $\square$

**Problem\* 4.24.** Let  $X := Y := c_0(\mathbb{N})$  and recall that  $X^* \cong \ell^1(\mathbb{N})$  and  $X^{**} \cong \ell^\infty(\mathbb{N})$ . Consider the operator  $A \in \mathcal{L}(\ell^1(\mathbb{N}))$  given by

$$Ax := \left( \sum_{n \in \mathbb{N}} x_n, 0, \dots \right).$$

Show that

$$A'x' = (x'_1, x'_1, \dots).$$

Conclude that  $A$  is not the adjoint of an operator from  $\mathcal{L}(c_0(\mathbb{N}))$ .

**Problem 4.25.** Show that for  $A \in \mathcal{L}(X, Y)$  we have

$$\text{rank}(A) = \text{rank}(A').$$

**Problem 4.26.** Show Lemma 4.23.

**Problem 4.27.** Let us write  $\ell_n \xrightarrow{*} \ell$  provided the sequence converges point-wise, that is,  $\ell_n(x) \rightarrow \ell(x)$  for all  $x \in X$ . Let  $N \subseteq X^*$  and suppose  $\ell_n \xrightarrow{*} \ell$  with  $\ell_n \in N$ . Show that  $\ell \in (N_\perp)^\perp$ .

#### 4.4. Weak convergence

In Section 4.2 we have seen that  $\ell(x) = 0$  for all  $\ell \in X^*$  implies  $x = 0$ . Now what about convergence? Does  $\ell(x_n) \rightarrow \ell(x)$  for every  $\ell \in X^*$  imply  $x_n \rightarrow x$ ? In fact, in a finite dimensional space component-wise convergence is equivalent to convergence. Unfortunately in the infinite dimensional this is no longer true in general:

**Example 4.27.** Let  $u_n$  be an infinite orthonormal set in some Hilbert space. Then  $\langle g, u_n \rangle \rightarrow 0$  for every  $g$  since these are just the expansion coefficients of  $g$ , which are in  $\ell^2(\mathbb{N})$  by Bessel's inequality. Since by the Riesz lemma (Theorem 2.10), every bounded linear functional is of this form, we have  $\ell(u_n) \rightarrow 0$  for every bounded linear functional. (Clearly  $u_n$  does not converge to 0, since  $\|u_n\| = 1$ .)  $\diamond$

If  $\ell(x_n) \rightarrow \ell(x)$  for every  $\ell \in X^*$  we say that  $x_n$  **converges weakly** to  $x$  and write

$$\text{w-lim}_{n \rightarrow \infty} x_n = x \quad \text{or} \quad x_n \rightharpoonup x. \quad (4.18)$$

Clearly,  $x_n \rightarrow x$  implies  $x_n \rightharpoonup x$  and hence this notion of convergence is indeed weaker. Moreover, the weak limit is unique, since  $\ell(x_n) \rightarrow \ell(x)$  and  $\ell(x_n) \rightarrow \ell(\tilde{x})$  imply  $\ell(x - \tilde{x}) = 0$ . A sequence  $x_n$  is called a **weak Cauchy sequence** if  $\ell(x_n)$  is Cauchy (i.e. converges) for every  $\ell \in X^*$ .

**Lemma 4.25.** *Let  $X$  be a Banach space.*

- (i)  $x_n \rightharpoonup x$ ,  $y_n \rightharpoonup y$  and  $\alpha_n \rightarrow \alpha$  implies  $x_n + y_n \rightharpoonup x + y$  and  $\alpha_n x_n \rightharpoonup \alpha x$ .
- (ii)  $x_n \rightharpoonup x$  implies  $\|x\| \leq \liminf \|x_n\|$ .
- (iii) Every weak Cauchy sequence  $x_n$  is bounded:  $\|x_n\| \leq C$ .
- (iv) If  $X$  is reflexive, then every weak Cauchy sequence converges weakly.
- (v) A sequence  $x_n$  is Cauchy if and only if  $\ell(x_n)$  is Cauchy, uniformly for  $\ell \in X^*$  with  $\|\ell\| = 1$ .

**Proof.** (i) Follows from  $\ell(\alpha_n x_n + y_n) = \alpha_n \ell(x_n) + \ell(y_n) \rightarrow \alpha \ell(x) + \ell(y)$ . (ii) Choose  $\ell \in X^*$  such that  $\ell(x) = \|x\|$  (for the limit  $x$ ) and  $\|\ell\| = 1$ . Then

$$\|x\| = \ell(x) = \liminf \ell(x_n) \leq \liminf \|x_n\|.$$

(iii) For every  $\ell$  we have that  $|J(x_n)(\ell)| = |\ell(x_n)| \leq C(\ell)$  is bounded. Hence by the uniform boundedness principle we have  $\|x_n\| = \|J(x_n)\| \leq C$ .

(iv) If  $x_n$  is a weak Cauchy sequence, then  $\ell(x_n)$  converges and we can define  $j(\ell) = \lim \ell(x_n)$ . By construction  $j$  is a linear functional on  $X^*$ . Moreover, by (iii) we have  $|j(\ell)| \leq \sup |\ell(x_n)| \leq \|\ell\| \sup \|x_n\| \leq C\|\ell\|$  which shows  $j \in X^{**}$ . Since  $X$  is reflexive,  $j = J(x)$  for some  $x \in X$  and by construction

$\ell(x_n) \rightarrow J(x)(\ell) = \ell(x)$ , that is,  $x_n \rightharpoonup x$ .

(v) This follows from

$$\|x_n - x_m\| = \sup_{\|\ell\|=1} |\ell(x_n - x_m)|$$

(cf. Problem 4.11). □

Item (ii) says that the norm is sequentially weakly lower semicontinuous (cf. Problem B.19) while the previous example shows that it is not sequentially weakly continuous (this will in fact be true for any convex function as we will see later). However, bounded linear operators turn out to be sequentially weakly continuous (Problem 4.29). Nonlinear operations are more tricky as the next example shows:

**Example 4.28.** Consider  $L^2(0, 1)$  and recall (see Example 3.8) that

$$u_n(x) = \sqrt{2} \sin(n\pi x), \quad n \in \mathbb{N},$$

form an ONB and hence  $u_n \rightharpoonup 0$ . However,  $v_n = u_n^2 \rightharpoonup 1$ . In fact, one easily computes

$$\langle u_m, v_n \rangle = \frac{\sqrt{2}(1 - (-1)^m)}{m\pi} \frac{4n^2}{(4n^2 - m^2)} \rightarrow \frac{\sqrt{2}(1 - (-1)^m)}{m\pi} = \langle u_m, 1 \rangle$$

and the claim follows from Problem 4.33 since  $\|v_n\| = \sqrt{\frac{3}{2}}$ . ◇

**Example 4.29.** Let  $X := c_0(\mathbb{N})$  and hence  $X^* \cong \ell^1(\mathbb{N})$ . Let  $a_j^n := 1$  for  $1 \leq j \leq n$  and  $a_j^n := 0$  for  $j > n$ . Then for every  $b \in \ell^1(\mathbb{N})$  we have

$$\lim_{n \rightarrow \infty} l_b(a^n) = \lim_{n \rightarrow \infty} \sum_{j=1}^{\infty} b_j a_j^n = \lim_{n \rightarrow \infty} \sum_{j=1}^n b_j = \sum_{j=1}^{\infty} b_j$$

and hence  $a^n$  is a weak Cauchy sequence which, does not converge. Indeed,  $a^n \rightharpoonup a$  would imply  $a_j = 1$  for all  $j \in \mathbb{N}$  (upon choosing  $b = \delta^j$ ) which is clearly not in  $X$ . The limit is however in  $X^{**} \cong \ell^\infty(\mathbb{N})$ . ◇

**Remark:** One can equip  $X$  with the weakest topology for which all  $\ell \in X^*$  remain continuous. This topology is called the **weak topology** and it is given by taking all finite intersections of inverse images of open sets as a base. By construction, a sequence will converge in the weak topology if and only if it converges weakly. By Corollary 4.13 the weak topology is Hausdorff, but it will not be metrizable in general. In particular, sequences do not suffice to describe this topology. Nevertheless we will stick with sequences for now and come back to this more general point of view in Section 6.3.

In a Hilbert space there is also a simple criterion for a weakly convergent sequence to converge in norm (see Theorem 6.19 for a generalization).

**Lemma 4.26.** *Let  $\mathfrak{H}$  be a Hilbert space and let  $f_n \rightharpoonup f$ . Then  $f_n \rightarrow f$  if and only if  $\limsup \|f_n\| \leq \|f\|$ .*

**Proof.** By (ii) of the previous lemma we have  $\lim \|f_n\| = \|f\|$  and hence

$$\|f - f_n\|^2 = \|f\|^2 - 2\operatorname{Re}(\langle f, f_n \rangle) + \|f_n\|^2 \rightarrow 0.$$

The converse is straightforward.  $\square$

Now we come to the main reason why weakly convergent sequences are of interest: A typical approach for solving a given equation in a Banach space is as follows:

- (i) Construct a (bounded) sequence  $x_n$  of approximating solutions (e.g. by solving the equation restricted to a finite dimensional subspace and increasing this subspace).
- (ii) Use a compactness argument to extract a convergent subsequence.
- (iii) Show that the limit solves the equation.

Our aim here is to provide some results for the step (ii). In a finite dimensional vector space the most important compactness criterion is boundedness (Heine–Borel theorem, Theorem B.22). In infinite dimensions this breaks down as we have already seen in Section 1.5. We even have

**Theorem 4.27** (F. Riesz). *The closed unit ball in a Banach space  $X$  is compact if and only if  $X$  is finite dimensional.*

**Proof.** If  $X$  is finite dimensional, then by Theorem 1.8 we can assume  $X = \mathbb{C}^n$  and the closed unit ball is compact by the Heine–Borel theorem.

Conversely, suppose  $S = \{x \in X \mid \|x\| = 1\}$  is compact. Then  $\{X \setminus \operatorname{Ker}(\ell)\}_{\ell \in X^*}$  is an open cover since for every  $x \in S$  there is some  $\ell \in X^*$  with  $\ell(x) \neq 0$  by Corollary 4.11. This cover has a finite subcover,  $S \subset \bigcup_{j=1}^n (X \setminus \operatorname{Ker}(\ell_j)) = X \setminus \bigcap_{j=1}^n \operatorname{Ker}(\ell_j)$ . Hence  $\bigcap_{j=1}^n \operatorname{Ker}(\ell_j) = \{0\}$  and the map  $X \rightarrow \mathbb{C}^n$ ,  $x \mapsto (\ell_1(x), \dots, \ell_n(x))$  is injective, that is,  $\dim(X) \leq n$ .  $\square$

However, if we are willing to treat convergence for weak convergence, the situation looks much brighter!

**Theorem 4.28** (Šmulian). *Let  $X$  be a reflexive Banach space. Then every bounded sequence has a weakly convergent subsequence.*

**Proof.** Let  $x_n$  be some bounded sequence and consider  $Y = \overline{\operatorname{span}\{x_n\}}$ . Then  $Y$  is reflexive by Lemma 4.17 (i). Moreover, by construction  $Y$  is separable and so is  $Y^*$  by the remark after Lemma 4.17.

Let  $\ell_k$  be a dense set in  $Y^*$ . Then by the usual diagonal sequence argument we can find a subsequence  $x_{n_m}$  such that  $\ell_k(x_{n_m})$  converges for every  $k$ . Denote this subsequence again by  $x_n$  for notational simplicity. Then,

$$\begin{aligned} |\ell(x_n) - \ell(x_m)| &\leq |\ell(x_n) - \ell_k(x_n)| + |\ell_k(x_n) - \ell_k(x_m)| \\ &\quad + |\ell_k(x_m) - \ell(x_m)| \\ &\leq 2C\|\ell - \ell_k\| + |\ell_k(x_n) - \ell_k(x_m)| \end{aligned}$$

shows that  $\ell(x_n)$  converges for every  $\ell \in \overline{\text{span}\{\ell_k\}} = Y^*$ . Thus there is a limit by Lemma 4.25 (iv).  $\square$

Note that this theorem breaks down if  $X$  is not reflexive.

**Example 4.30.** Consider the sequence of vectors  $\delta^n$  (with  $\delta_n^n = 1$  and  $\delta_m^n = 0$ ,  $n \neq m$ ) in  $\ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$ . Then  $\delta^n \rightharpoonup 0$  for  $1 < p < \infty$ . In fact, since every  $l \in \ell^p(\mathbb{N})^*$  is of the form  $l = l_y$  for some  $y \in \ell^q(\mathbb{N})$  we have  $l_y(\delta^n) = y_n \rightarrow 0$ .

If we consider the same sequence in  $\ell^1(\mathbb{N})$  there is no weakly convergent subsequence. In fact, since  $l_y(\delta^n) \rightarrow 0$  for every sequence  $y \in \ell^\infty(\mathbb{N})$  with finitely many nonzero entries, the only possible weak limit is zero. On the other hand choosing the constant sequence  $y = (1)_{j=1}^\infty$  we see  $l_y(\delta^n) = 1 \not\rightarrow 0$ , a contradiction.  $\diamond$

**Example 4.31.** Let  $X := L^1(-1, 1)$ . Every bounded integrable  $\varphi$  gives rise to a linear functional

$$\ell_\varphi(f) := \int f(x)\varphi(x) dx$$

in  $L^1[-1, 1]^*$ . Take some nonnegative  $u_1$  with compact support,  $\|u_1\|_1 = 1$ , and set  $u_k(x) = ku_1(kx)$  (implying  $\|u_k\|_1 = 1$ ). Then we have

$$\int u_k(x)\varphi(x) dx \rightarrow \varphi(0)$$

(see Problem 3.30 from [48]) for every continuous  $\varphi$ . Furthermore, if  $u_{k_j} \rightharpoonup u$  we conclude

$$\int u(x)\varphi(x) dx = \varphi(0).$$

In particular, choosing  $\varphi_k(x) = \max(0, 1 - k|x|)$  we infer from the dominated convergence theorem

$$1 = \int u(x)\varphi_k(x) dx \rightarrow \int u(x)\chi_{\{0\}}(x) dx = 0,$$

a contradiction.

In fact,  $u_k$  converges to the Dirac measure centered at 0, which is not in  $L^1(-1, 1)$ .  $\diamond$

Note that the above theorem also shows that in an infinite dimensional reflexive Banach space weak convergence is always weaker than strong convergence since otherwise every bounded sequence had a weakly, and thus by assumption also norm, convergent subsequence contradicting Theorem 4.27. In a non-reflexive space this situation can however occur.

**Example 4.32.** In  $\ell^1(\mathbb{N})$  every weakly convergent sequence is in fact (norm) convergent (such Banach spaces are said to have the **Schur property**). First of all recall that  $\ell^1(\mathbb{N})^* \cong \ell^\infty(\mathbb{N})$  and  $a^n \rightharpoonup 0$  implies

$$l_b(a^n) = \sum_{k=1}^{\infty} b_k a_k^n \rightarrow 0, \quad \forall b \in \ell^\infty(\mathbb{N}).$$

Now suppose we could find a sequence  $a^n \rightharpoonup 0$  for which  $\limsup_n \|a^n\|_1 \geq \varepsilon > 0$ . After passing to a subsequence we can assume  $\|a^n\|_1 \geq \varepsilon/2$  and after rescaling the vector even  $\|a^n\|_1 = 1$ . Now weak convergence  $a^n \rightharpoonup 0$  implies  $a_j^n = l_{\delta^j}(a^n) \rightarrow 0$  for every fixed  $j \in \mathbb{N}$ . Hence the main contribution to the norm of  $a^n$  must move towards  $\infty$  and we can find a subsequence  $n_j$  and a corresponding increasing sequence of integers  $k_j$  such that  $\sum_{k_j \leq k < k_{j+1}} |a_k^{n_j}| \geq \frac{2}{3}$ . Now set

$$b_k = \text{sign}(a_k^{n_j}), \quad k_j \leq k < k_{j+1}.$$

Then

$$|l_b(a^{n_j})| \geq \sum_{k_j \leq k < k_{j+1}} |a_k^{n_j}| - \left| \sum_{1 \leq k < k_j; k_{j+1} \leq k} b_k a_k^{n_j} \right| \geq \frac{2}{3} - \frac{1}{3} = \frac{1}{3},$$

contradicting  $a^{n_j} \rightharpoonup 0$ .  $\diamond$

It is also useful to observe that compact operators will turn weakly convergent into (norm) convergent sequences.

**Theorem 4.29.** *Let  $A \in \mathcal{K}(X, Y)$  be compact. Then  $x_n \rightharpoonup x$  implies  $Ax_n \rightarrow Ax$ . If  $X$  is reflexive the converse is also true.*

**Proof.** If  $x_n \rightharpoonup x$  we have  $\sup_n \|x_n\| \leq C$  by Lemma 4.25 (ii). Consequently  $Ax_n$  is bounded and we can pass to a subsequence such that  $Ax_{n_k} \rightarrow y$ . Moreover, by Problem 4.29 we even have  $y = Ax$  and Lemma B.5 shows  $Ax_n \rightarrow Ax$ .

Conversely, if  $X$  is reflexive, then by Theorem 4.28 every bounded sequence  $x_n$  has a subsequence  $x_{n_k} \rightharpoonup x$  and by assumption  $Ax_{n_k} \rightarrow x$ . Hence  $A$  is compact.  $\square$

Operators which map weakly convergent sequences to convergent sequences are also called **completely continuous**. However, be warned that some authors use completely continuous for compact operators. By the above



theorem every compact operator is completely continuous and the converse also holds in reflexive spaces. However, the last example shows that the identity map in  $\ell^1(\mathbb{N})$  is completely continuous but it is clearly not compact by Theorem 4.27.

Similar concepts can be introduced for operators. This is of particular importance for the case of unbounded operators, where convergence in the operator norm makes no sense at all.

A sequence of operators  $A_n$  is said to **converge strongly** to  $A$ ,

$$\text{s-lim}_{n \rightarrow \infty} A_n = A \quad :\Leftrightarrow \quad A_n x \rightarrow Ax \quad \forall x \in \mathfrak{D}(A) \subseteq \mathfrak{D}(A_n). \quad (4.19)$$

It is said to **converge weakly** to  $A$ ,

$$\text{w-lim}_{n \rightarrow \infty} A_n = A \quad :\Leftrightarrow \quad A_n x \rightharpoonup Ax \quad \forall x \in \mathfrak{D}(A) \subseteq \mathfrak{D}(A_n). \quad (4.20)$$

Clearly norm convergence implies strong convergence and strong convergence implies weak convergence. If  $Y$  is finite dimensional strong and weak convergence will be the same and this is in particular the case for  $Y = \mathbb{C}$ .

**Example 4.33.** Consider the operator  $S_n \in \mathcal{L}(\ell^2(\mathbb{N}))$  which shifts a sequence  $n$  places to the left, that is,

$$S_n(x_1, x_2, \dots) = (x_{n+1}, x_{n+2}, \dots) \quad (4.21)$$

and the operator  $S_n^* \in \mathcal{L}(\ell^2(\mathbb{N}))$  which shifts a sequence  $n$  places to the right and fills up the first  $n$  places with zeros, that is,

$$S_n^*(x_1, x_2, \dots) = (\underbrace{0, \dots, 0}_{n \text{ places}}, x_1, x_2, \dots). \quad (4.22)$$

Then  $S_n$  converges to zero strongly but not in norm (since  $\|S_n\| = 1$ ) and  $S_n^*$  converges weakly to zero (since  $\langle x, S_n^* y \rangle = \langle S_n x, y \rangle$ ) but not strongly (since  $\|S_n^* x\| = \|x\|$ ).  $\diamond$

**Lemma 4.30.** Suppose  $A_n, B_n \in \mathcal{L}(X, Y)$  are sequences of bounded operators.

- (i)  $\text{s-lim}_{n \rightarrow \infty} A_n = A$ ,  $\text{s-lim}_{n \rightarrow \infty} B_n = B$ , and  $\alpha_n \rightarrow \alpha$  implies  $\text{s-lim}_{n \rightarrow \infty} (A_n + B_n) = A + B$  and  $\text{s-lim}_{n \rightarrow \infty} \alpha_n A_n = \alpha A$ .
- (ii)  $\text{s-lim}_{n \rightarrow \infty} A_n = A$  implies  $\|A\| \leq \liminf_{n \rightarrow \infty} \|A_n\|$ .
- (iii) If  $A_n x$  converges for all  $x \in X$  then  $\|A_n\| \leq C$  and there is an operator  $A \in \mathcal{L}(X, Y)$  such that  $\text{s-lim}_{n \rightarrow \infty} A_n = A$ .
- (iv) If  $A_n y$  converges for  $y$  in a total set and  $\|A_n\| \leq C$ , then there is an operator  $A \in \mathcal{L}(X, Y)$  such that  $\text{s-lim}_{n \rightarrow \infty} A_n = A$ .

The same result holds if strong convergence is replaced by weak convergence.

**Proof.** (i)  $\lim_{n \rightarrow \infty} (\alpha_n A_n + B_n)x = \lim_{n \rightarrow \infty} (\alpha_n A_n x + B_n x) = \alpha A x + B x$ .  
(ii) follows from

$$\|Ax\| = \lim_{n \rightarrow \infty} \|A_n x\| \leq \liminf_{n \rightarrow \infty} \|A_n\|$$

for every  $x \in X$  with  $\|x\| = 1$ .

(iii) By linearity of the limit,  $Ax := \lim_{n \rightarrow \infty} A_n x$  is a linear operator. Moreover, since convergent sequences are bounded,  $\|A_n x\| \leq C(x)$ , the uniform boundedness principle implies  $\|A_n\| \leq C$ . Hence  $\|Ax\| = \lim_{n \rightarrow \infty} \|A_n x\| \leq C\|x\|$ .

(iv) By taking linear combinations we can replace the total set by a dense one. Moreover, we can define a linear operator  $A$  on this dense set via  $Ay := \lim_{n \rightarrow \infty} A_n y$ . By  $\|A_n\| \leq C$  we see  $\|A\| \leq C$  and there is a unique extension to all of  $X$ . Now just use

$$\begin{aligned} \|A_n x - Ax\| &\leq \|A_n x - A_n y\| + \|A_n y - Ay\| + \|Ay - Ax\| \\ &\leq 2C\|x - y\| + \|A_n y - Ay\| \end{aligned}$$

and choose  $y$  in the dense subspace such that  $\|x - y\| \leq \frac{\varepsilon}{4C}$  and  $n$  large such that  $\|A_n y - Ay\| \leq \frac{\varepsilon}{2}$ .

The case of weak convergence is left as an exercise (Problem 4.31).  $\square$

Item (iii) of this lemma is sometimes also known as Banach–Steinhaus theorem. For an application of this lemma see Lemma 3.21 from [48].

**Example 4.34.** Let  $X$  be a Banach space of functions  $f : [-\pi, \pi] \rightarrow \mathbb{C}$  such that the functions  $\{e_k(x) := e^{ikx}\}_{k \in \mathbb{Z}}$  are total. E.g.  $X = C_{per}[-\pi, \pi]$  or  $X = L^p[-\pi, \pi]$  for  $1 \leq p < \infty$ . Then the Fourier series (2.45) converges on a total set and hence it will converge on all of  $X$  if and only if  $\|S_n\| \leq C$ . For example, if  $X = C_{per}[-\pi, \pi]$  then

$$\|S_n\| = \sup_{\|f\|_\infty=1} \|S_n(f)\| = \sup_{\|f\|_\infty=1} |S_n(f)(0)| = \frac{1}{2\pi} \|D_n\|_1$$

which is unbounded as we have seen in Example 4.3. In fact, in this example we have even shown failure of pointwise convergence and hence this is nothing new. However, if we consider  $X = L^1[-\pi, \pi]$  we have (recall the Fejér kernel which satisfies  $\|F_n\|_1 = 1$  and use (2.53) together with  $S_n(D_m) = D_{\min(m,n)}$ )

$$\|S_n\| = \sup_{\|f\|_1=1} \|S_n(f)\| \geq \lim_{m \rightarrow \infty} \|S_n(F_m)\|_1 = \|D_n\|_1$$

and we get that the Fourier series does not converge for some  $L^1$  function.  $\diamond$

**Lemma 4.31.** Suppose  $A_n \in \mathcal{L}(Y, Z)$ ,  $B_n \in \mathcal{L}(X, Y)$  are two sequences of bounded operators.

- (i)  $\text{s-lim}_{n \rightarrow \infty} A_n = A$  and  $\text{s-lim}_{n \rightarrow \infty} B_n = B$  implies  $\text{s-lim}_{n \rightarrow \infty} A_n B_n = AB$ .

- (ii)  $\text{w-lim}_{n \rightarrow \infty} A_n = A$  and  $\text{s-lim}_{n \rightarrow \infty} B_n = B$  implies  $\text{w-lim}_{n \rightarrow \infty} A_n B_n = AB$ .  
 (iii)  $\lim_{n \rightarrow \infty} A_n = A$  and  $\text{w-lim}_{n \rightarrow \infty} B_n = B$  implies  $\text{w-lim}_{n \rightarrow \infty} A_n B_n = AB$ .

**Proof.** For the first case just observe

$$\|(A_n B_n - AB)x\| \leq \|(A_n - A)Bx\| + \|A_n\| \|(B_n - B)x\| \rightarrow 0.$$

The remaining cases are similar and again left as an exercise.  $\square$

**Example 4.35.** Consider again the last example. Then

$$S_n^* S_n(x_1, x_2, \dots) = (\underbrace{0, \dots, 0}_{n \text{ places}}, x_{n+1}, x_{n+2}, \dots)$$

converges to 0 weakly (in fact even strongly) but

$$S_n S_n^*(x_1, x_2, \dots) = (x_1, x_2, \dots)$$

does not! Hence the order in the second claim is important.  $\diamond$

For a sequence of linear functionals  $\ell_n$ , strong convergence is also called **weak-\*** convergence. That is, the weak-\* limit of  $\ell_n$  is  $\ell$  if  $\ell_n(x) \rightarrow \ell(x)$  for all  $x \in X$  and we will write

$$\text{w}^*\text{-lim}_{n \rightarrow \infty} \ell_n = \ell \quad \text{or} \quad \ell_n \xrightarrow{*} \ell \quad (4.23)$$

in this case. Note that this is not the same as weak convergence on  $X^*$  unless  $X$  is reflexive:  $\ell$  is the weak limit of  $\ell_n$  if

$$j(\ell_n) \rightarrow j(\ell) \quad \forall j \in X^{**}, \quad (4.24)$$

whereas for the weak-\* limit this is only required for  $j \in J(X) \subseteq X^{**}$  (recall  $J(x)(\ell) = \ell(x)$ ).

**Example 4.36.** In a Hilbert space weak-\* convergence of the linear functionals  $\langle x_n, \cdot \rangle$  is the same as weak convergence of the vectors  $x_n$ .  $\diamond$

**Example 4.37.** Consider  $X := c_0(\mathbb{N})$ ,  $X^* \cong \ell^1(\mathbb{N})$ , and  $X^{**} \cong \ell^\infty(\mathbb{N})$  with  $J$  corresponding to the inclusion  $c_0(\mathbb{N}) \hookrightarrow \ell^\infty(\mathbb{N})$ . Then weak convergence on  $X^*$  implies

$$l_b(a^n - a) = \sum_{k=1}^{\infty} b_k(a_k^n - a_k) \rightarrow 0$$

for all  $b \in \ell^\infty(\mathbb{N})$  and weak-\* convergence implies that this holds for all  $b \in c_0(\mathbb{N})$ . Whereas we already have seen that weak convergence is equivalent to norm convergence, it is not hard to see that weak-\* convergence is equivalent to the fact that the sequence is bounded and each component converges (cf. Problem 4.34).  $\diamond$

With this notation the proof of Lemma 4.25 (iv) shows that (without assuming  $X$  to be reflexive) every weak Cauchy sequence converges weak-\*. Similarly, it is also possible to slightly generalize Theorem 4.28 (Problem 4.35):

**Lemma 4.32** (Helly). *Suppose  $X$  is a separable Banach space. Then every bounded sequence  $\ell_n \in X^*$  has a weak-\* convergent subsequence.*

**Example 4.38.** Let us return to the example after Theorem 4.28. Consider the Banach space of continuous functions  $X := C[-1, 1]$ . Using  $\ell_f(\varphi) := \int \varphi f dx$  we can regard  $L^1[-1, 1]$  as a subspace of  $X^*$ . Then the Dirac measure centered at 0 is also in  $X^*$  and it is the weak-\* limit of the sequence  $u_k$ .  $\diamond$

**Example 4.39.** Consider  $X := \ell^\infty(\mathbb{N})$ . Then the sequence of projections  $l_k \in X^*$  given by  $l_k(x) := x_k$  has no weak-\* convergent subsequence (if there were such a subsequence  $k_j$ , choose  $x \in X$  such that  $x_{k_j}$  does not converge to get a contradiction). Hence the assumption that  $X$  is separable cannot be dropped in Lemma 4.32.  $\diamond$

**Problem 4.28.** *Suppose  $\ell_n \rightarrow \ell$  in  $X^*$  and  $x_n \rightarrow x$  in  $X$ . Then  $\ell_n(x_n) \rightarrow \ell(x)$ . Similarly, suppose  $\text{s-lim } \ell_n \rightarrow \ell$  and  $x_n \rightarrow x$ . Then  $\ell_n(x_n) \rightarrow \ell(x)$ . Does this still hold if  $\text{s-lim } \ell_n \rightarrow \ell$  and  $x_n \rightarrow x$ ?*

**Problem\* 4.29.** *Show that  $x_n \rightarrow x$  implies  $Ax_n \rightarrow Ax$  for  $A \in \mathcal{L}(X, Y)$ . Conversely, show that if  $x_n \rightarrow 0$  implies  $Ax_n \rightarrow 0$  then  $A \in \mathcal{L}(X, Y)$ .*

**Problem 4.30.** *Let  $X := X_1 \oplus X_2$  show that  $(x_{1,n}, x_{2,n}) \rightarrow (x_1, x_2)$  if and only if  $x_{j,n} \rightarrow x_j$  for  $j = 1, 2$ .*

**Problem 4.31.** *Establish Lemma 4.30 in the case of weak convergence. (Hint: Problem 4.13 might be useful.)*

**Problem 4.32.** *Suppose  $A_n, A \in \mathcal{L}(X, Y)$ . Show that  $\text{s-lim } A_n = A$  and  $\lim x_n = x$  implies  $\lim A_n x_n = Ax$ .*

**Problem\* 4.33.** *Show that if  $\{\ell_j\} \subseteq X^*$  is some total set, then  $x_n \rightarrow x$  if and only if  $x_n$  is bounded and  $\ell_j(x_n) \rightarrow \ell_j(x)$  for all  $j$ . Show that this is wrong without the boundedness assumption (Hint: Take e.g.  $X = \ell^2(\mathbb{N})$ ).*

**Problem\* 4.34.** *Show that if  $\{x_j\} \subseteq X$  is some total set, then  $\ell_n \xrightarrow{*} \ell$  if and only if  $\ell_n \in X^*$  is bounded and  $\ell_n(x_j) \rightarrow \ell(x_j)$  for all  $j$ .*

**Problem\* 4.35.** *Prove Lemma 4.32.*



# Bounded linear operators

We have started out our study by looking at eigenvalue problems which, from a historic view point, were one of the key problems driving the development of functional analysis. In Chapter 3 we have investigated compact operators in Hilbert space and we have seen that they allow a treatment similar to what is known from matrices. However, more sophisticated problems will lead to operators whose spectra consist of more than just eigenvalues. Hence we want to go one step further and look at spectral theory for bounded operators. Here one of the driving forces was the development of quantum mechanics (there even the boundedness assumption is too much — but first things first). A crucial role is played by the algebraic structure, namely recall from Section 1.6 that the bounded linear operators on  $X$  form a Banach space which has a (non-commutative) multiplication given by composition. In order to emphasize that it is only this algebraic structure which matters, we will develop the theory from this abstract point of view. While the reader should always remember that bounded operators on a Hilbert space is what we have in mind as the prime application, examples will apply these ideas also to other cases thereby justifying the abstract approach.

To begin with, the operators could be on a Banach space (note that even if  $X$  is a Hilbert space,  $\mathcal{L}(X)$  will only be a Banach space) but eventually again self-adjointness will be needed. Hence we will need the additional operation of taking adjoints.

## 5.1. Banach algebras

A Banach space  $X$  together with a multiplication satisfying

$$(x + y)z = xz + yz, \quad x(y + z) = xy + xz, \quad x, y, z \in X, \quad (5.1)$$

and

$$(xy)z = x(yz), \quad \alpha(xy) = (\alpha x)y = x(\alpha y), \quad \alpha \in \mathbb{C}, \quad (5.2)$$

and

$$\|xy\| \leq \|x\|\|y\|. \quad (5.3)$$

is called a **Banach algebra**. In particular, note that (5.3) ensures that multiplication is continuous (Problem 5.1). In fact, one can show that (separate) continuity of multiplication implies existence of an equivalent norm satisfying (5.3) (Problem 5.2).

An element  $e \in X$  satisfying

$$ex = xe = x, \quad \forall x \in X \quad (5.4)$$

is called **identity** (show that  $e$  is unique) and we will assume  $\|e\| = 1$  in this case (by Problem 5.2 this can be done without loss of generality).

**Example 5.1.** The continuous functions  $C(I)$  over some compact interval form a commutative Banach algebra with identity 1.  $\diamond$

**Example 5.2.** The differentiable functions  $C^n(I)$  over some compact interval do not form a commutative Banach algebra since (5.3) fails for  $n \geq 1$ . However, the equivalent norm

$$\|f\|_{\infty, n} := \sum_{k=0}^n \frac{\|f^{(k)}\|_{\infty}}{k!}$$

remedies this problem.  $\diamond$

**Example 5.3.** The bounded linear operators  $\mathcal{L}(X)$  form a Banach algebra with identity  $\mathbb{I}$ .  $\diamond$

**Example 5.4.** The bounded sequences  $\ell^\infty(\mathbb{N})$  together with the component-wise product form a commutative Banach algebra with identity 1.  $\diamond$

**Example 5.5.** The space of all periodic continuous functions which have an absolutely convergent Fourier series  $\mathcal{A}$  together with the norm

$$\|f\|_{\mathcal{A}} := \sum_{k \in \mathbb{Z}} |\hat{f}_k|$$

and the usual product is known as the **Wiener algebra**. Of course as a Banach space it is isomorphic to  $\ell^1(\mathbb{Z})$  via the Fourier transform. To see that it is a Banach algebra note that

$$\begin{aligned} f(x)g(x) &= \sum_{k \in \mathbb{Z}} \hat{f}_k e^{ikx} \sum_{j \in \mathbb{Z}} \hat{g}_j e^{ijx} = \sum_{k, j \in \mathbb{Z}} \hat{f}_k \hat{g}_j e^{i(k+j)x} \\ &= \sum_{k \in \mathbb{Z}} \left( \sum_{j \in \mathbb{Z}} \hat{f}_j \hat{g}_{k-j} \right) e^{ikx}. \end{aligned}$$

Moreover, interchanging the order of summation

$$\|fg\|_{\mathcal{A}} = \sum_{k \in \mathbb{Z}} \left| \sum_{j \in \mathbb{Z}} \hat{f}_j \hat{g}_{k-j} \right| \leq \sum_{j \in \mathbb{Z}} \sum_{k \in \mathbb{Z}} |\hat{f}_j| |\hat{g}_{k-j}| = \|f\|_{\mathcal{A}} \|g\|_{\mathcal{A}}$$

shows that  $\mathcal{A}$  is a Banach algebra. The identity is of course given by  $e(x) \equiv 1$ . Moreover, note that  $\mathcal{A} \subseteq C_{per}[-\pi, \pi]$  and  $\|f\|_{\infty} \leq \|f\|_{\mathcal{A}}$ .  $\diamond$

**Example 5.6.** The space  $L^1(\mathbb{R}^n)$  together with the convolution

$$(g * f)(x) := \int_{\mathbb{R}^n} g(x-y)f(y)dy = \int_{\mathbb{R}^n} g(y)f(x-y)dy \quad (5.5)$$

is a commutative Banach algebra (Problem 5.14) without identity.  $\diamond$

A Banach algebra with identity is also known as **unital** and we will assume  $X$  to be a Banach algebra with identity  $e$  throughout the rest of this section. Note that an identity can always be added if needed (Problem 5.3).

An element  $x \in X$  is called **invertible** if there is some  $y \in X$  such that

$$xy = yx = e. \quad (5.6)$$

In this case  $y$  is called the inverse of  $x$  and is denoted by  $x^{-1}$ . It is straightforward to show that the inverse is unique (if one exists at all) and that

$$(xy)^{-1} = y^{-1}x^{-1}, \quad (x^{-1})^{-1} = x. \quad (5.7)$$

In particular, the set of invertible elements  $\mathcal{G}(X)$  forms a group under multiplication.

**Example 5.7.** If  $X = \mathcal{L}(\mathbb{C}^n)$  is the set of  $n$  by  $n$  matrices, then  $\mathcal{G}(X) = \text{GL}(n)$  is the general linear group.  $\diamond$

**Example 5.8.** Let  $X = \mathcal{L}(\ell^p(\mathbb{N}))$  and recall the shift operators  $S^{\pm}$  defined via  $(S^{\pm}a)_j = a_{j \pm 1}$  with the convention that  $a_0 = 0$ . Then  $S^+S^- = \mathbb{I}$  but  $S^-S^+ \neq \mathbb{I}$ . Moreover, note that  $S^+S^-$  is invertible while  $S^-S^+$  is not. So you really need to check both  $xy = e$  and  $yx = e$  in general.  $\diamond$

If  $x$  is invertible, then the same will be true all elements in a neighborhood. This will be a consequence from the following straightforward generalization of the geometric series to our abstract setting.

**Lemma 5.1.** *Let  $X$  be a Banach algebra with identity  $e$ . Suppose  $\|x\| < 1$ . Then  $e - x$  is invertible and*

$$(e - x)^{-1} = \sum_{n=0}^{\infty} x^n. \quad (5.8)$$

**Proof.** Since  $\|x\| < 1$  the series converges and

$$(e - x) \sum_{n=0}^{\infty} x^n = \sum_{n=0}^{\infty} x^n - \sum_{n=1}^{\infty} x^n = e$$



respectively

$$\left(\sum_{n=0}^{\infty} x^n\right)(e - x) = \sum_{n=0}^{\infty} x^n - \sum_{n=1}^{\infty} x^n = e. \quad \square$$

**Corollary 5.2.** *Suppose  $x$  is invertible and  $\|x^{-1}y\| < 1$  or  $\|yx^{-1}\| < 1$ . Then  $(x - y)$  is invertible as well and*

$$(x - y)^{-1} = \sum_{n=0}^{\infty} (x^{-1}y)^n x^{-1} \quad \text{or} \quad (x - y)^{-1} = \sum_{n=0}^{\infty} x^{-1} (yx^{-1})^n. \quad (5.9)$$

*In particular, both conditions are satisfied if  $\|y\| < \|x^{-1}\|^{-1}$  and the set of invertible elements  $\mathcal{G}(X)$  is open and taking the inverse is continuous:*

$$\|(x - y)^{-1} - x^{-1}\| \leq \frac{\|y\| \|x^{-1}\|^2}{1 - \|x^{-1}y\|}. \quad (5.10)$$

**Proof.** Just observe  $x - y = x(e - x^{-1}y) = (e - yx^{-1})x$ .  $\square$

The **resolvent set** is defined as

$$\rho(x) := \{\alpha \in \mathbb{C} \mid (x - \alpha) \text{ is invertible in } X\} \subseteq \mathbb{C}, \quad (5.11)$$

where we have used the shorthand notation  $x - \alpha := x - \alpha e$ . Its complement is called the **spectrum**

$$\sigma(x) := \mathbb{C} \setminus \rho(x). \quad (5.12)$$

It is important to observe that the inverse has to exist as an element of  $X$ . That is, if the elements of  $X$  are bounded linear operators, it does not suffice that  $x - \alpha$  is injective, as it might not be surjective. If it is bijective, boundedness of the inverse will come for free from the inverse mapping theorem.

**Example 5.9.** If  $X := \mathcal{L}(\mathbb{C}^n)$  is the space of  $n$  by  $n$  matrices, then the spectrum is just the set of eigenvalues. More general, if  $X$  are the bounded linear operators on an infinite-dimensional Hilbert or Banach space, then every eigenvalue will be in the spectrum but the converse is not true in general as an injective operator might not be surjective. In fact, this already can happen for compact operators where 0 could be in the spectrum without being an eigenvalue.  $\diamond$

**Example 5.10.** If  $X := C(I)$ , then the spectrum of a function  $x \in C(I)$  is just its range,  $\sigma(x) = x(I)$ . Indeed, if  $\alpha \notin \text{Ran}(x)$  then  $t \mapsto (x(t) - \alpha)^{-1}$  is the inverse of  $x - \alpha$  (note that  $\text{Ran}(x)$  is compact). Conversely, if  $\alpha \in \text{Ran}(x)$  and  $y$  were an inverse, then  $y(t_0)(x(t_0) - \alpha) = 1$  gives a contradiction for any  $t_0 \in I$  with  $f(t_0) = \alpha$ .  $\diamond$

**Example 5.11.** If  $X = \mathcal{A}$  is the Wiener algebra, then, as in the previous example, every function which vanishes at some point cannot be inverted. If it does not vanish anywhere, it can be inverted and the inverse will be a

continuous function. But will it again have a convergent Fourier series, that is, will it be in the Wiener Algebra? The affirmative answer of this question is a famous theorem of Wiener, which will be given later in Theorem 7.22.  $\diamond$

The map  $\alpha \mapsto (x - \alpha)^{-1}$  is called the **resolvent** of  $x \in X$ . If  $\alpha_0 \in \rho(x)$  we can choose  $x \rightarrow x - \alpha_0$  and  $y \rightarrow \alpha - \alpha_0$  in (5.9) which implies

$$(x - \alpha)^{-1} = \sum_{n=0}^{\infty} (\alpha - \alpha_0)^n (x - \alpha_0)^{-n-1}, \quad |\alpha - \alpha_0| < \|(x - \alpha_0)^{-1}\|^{-1}. \quad (5.13)$$

In particular, since the radius of convergence cannot exceed the distance to the spectrum (since everything within the radius of convergent must belong to the resolvent set), we see that the norm of the resolvent must diverge

$$\|(x - \alpha)^{-1}\| \geq \frac{1}{\text{dist}(\alpha, \sigma(x))} \quad (5.14)$$

as  $\alpha$  approaches the spectrum. Moreover, this shows that  $(x - \alpha)^{-1}$  has a convergent power series with coefficients in  $X$  around every point  $\alpha_0 \in \rho(x)$ . As in the case of coefficients in  $\mathbb{C}$ , such functions will be called **analytic**.

**Example 5.12.** If  $A \in \mathcal{L}(\mathbb{C}^n)$  is an  $n$  by  $n$  matrices, then the resolvent is given by

$$(A - \alpha)^{-1} = \frac{1}{\det(A - \alpha)} (A - \alpha)^{\text{adj}},$$

where  $A^{\text{adj}}$  denotes the **adjugate** (transpose of the cofactor matrix) of  $A$ . Since  $(A - \alpha)^{\text{adj}}$  is a polynomial in  $\alpha$ , the resolvent has a pole at each eigenvalue whose order is at most the algebraic multiplicity of the eigenvalue. In fact the order of the pole equals the algebraic multiplicity which can be seen using (e.g.) the Jordan canonical form.  $\diamond$

In particular,  $\ell((x - \alpha)^{-1})$  is a complex-valued analytic function for every  $\ell \in X^*$  and we can apply well-known results from complex analysis:

**Theorem 5.3.** *For every  $x \in X$ , the spectrum  $\sigma(x)$  is compact, nonempty and satisfies*

$$\sigma(x) \subseteq \{\alpha \mid |\alpha| \leq \|x\|\}. \quad (5.15)$$

**Proof.** Equation (5.13) already shows that  $\rho(x)$  is open. Hence  $\sigma(x)$  is closed. Moreover,  $x - \alpha = -\alpha(e - \frac{1}{\alpha}x)$  together with Lemma 5.1 shows

$$(x - \alpha)^{-1} = -\frac{1}{\alpha} \sum_{n=0}^{\infty} \left(\frac{x}{\alpha}\right)^n, \quad |\alpha| > \|x\|,$$

which implies  $\sigma(x) \subseteq \{\alpha \mid |\alpha| \leq \|x\|\}$  is bounded and thus compact. Moreover, taking norms shows

$$\|(x - \alpha)^{-1}\| \leq \frac{1}{|\alpha|} \sum_{n=0}^{\infty} \frac{\|x\|^n}{|\alpha|^n} = \frac{1}{|\alpha| - \|x\|}, \quad |\alpha| > \|x\|,$$

which implies  $(x - \alpha)^{-1} \rightarrow 0$  as  $\alpha \rightarrow \infty$ . In particular, if  $\sigma(x)$  is empty, then  $\ell((x - \alpha)^{-1})$  is an entire analytic function which vanishes at infinity. By Liouville's theorem we must have  $\ell((x - \alpha)^{-1}) = 0$  for all  $\ell \in X^*$  in this case, and so  $(x - \alpha)^{-1} = 0$ , which is impossible.  $\square$

**Example 5.13.** The spectrum of the matrix

$$A := \begin{pmatrix} 0 & 1 & & & \\ & 0 & 1 & & \\ & & \ddots & \ddots & \\ & & & 0 & 1 \\ -c_0 & -c_1 & \cdots & \cdots & -c_{n-1} \end{pmatrix}$$

is given by the zeros of the polynomial (show this)

$$\det(z\mathbb{I} - A) = z^n + c_{n-1}z^{n-1} + \cdots + c_1z + c_0.$$

Hence the fact that  $\sigma(A)$  is nonempty implies the **fundamental theorem of algebra**, that every non-constant polynomial has at least one zero.  $\diamond$

As another simple consequence we obtain:

**Theorem 5.4** (Gelfand–Mazur). *Suppose  $X$  is a Banach algebra in which every element except 0 is invertible. Then  $X$  is isomorphic to  $\mathbb{C}$ .*

**Proof.** Pick  $x \in X$  and  $\alpha \in \sigma(x)$ . Then  $x - \alpha$  is not invertible and hence  $x - \alpha = 0$ , that is  $x = \alpha$ . Thus every element is a multiple of the identity.  $\square$

Now we look at functions of  $x$ . Given a polynomial  $p(\alpha) = \sum_{j=0}^n p_j \alpha^j$  we of course set

$$p(x) := \sum_{j=0}^n p_j x^j. \quad (5.16)$$

In fact, we could easily extend this definition to arbitrary convergent power series whose radius of convergence is larger than  $\|x\|$  (cf. Problem 1.39). While this will give a nice functional calculus sufficient for many applications our aim is the spectral theorem which will allow us to handle arbitrary continuous functions. Since continuous functions can be approximated by polynomials by the Weierstraß theorem, polynomials will be sufficient for now. Moreover, the following result will be one of the two key ingredients for the proof of the spectral theorem.

**Theorem 5.5** (Spectral mapping). *For every polynomial  $p$  and  $x \in X$  we have*

$$\sigma(p(x)) = p(\sigma(x)), \quad (5.17)$$

where  $p(\sigma(x)) := \{p(\alpha) | \alpha \in \sigma(x)\}$ .

**Proof.** Let  $\alpha \in \sigma(x)$  and observe

$$p(x) - p(\alpha) = (x - \alpha)q(x).$$

But since  $(x - \alpha)$  is not invertible, the same is true for  $(x - \alpha)q(x) = q(x)(x - \alpha)$  by Problem 5.10 and hence  $p(\alpha) \in p(\sigma(x))$ .

Conversely, let  $\beta \in \sigma(p(x))$ . Then

$$p(x) - \beta = a(x - \lambda_1) \cdots (x - \lambda_n)$$

and at least one  $\lambda_j \in \sigma(x)$  since otherwise the right-hand side would be invertible. But then  $\beta = p(\lambda_j) \in p(\sigma(x))$ .  $\square$

The second key ingredient for the proof of the spectral theorem is the **spectral radius**

$$r(x) := \sup_{\alpha \in \sigma(x)} |\alpha| \quad (5.18)$$

of  $x$ . Note that by (5.15) we have

$$r(x) \leq \|x\|. \quad (5.19)$$

As our next theorem shows, it is related to the radius of convergence of the **Neumann series** for the resolvent

$$(x - \alpha)^{-1} = -\frac{1}{\alpha} \sum_{n=0}^{\infty} \left(\frac{x}{\alpha}\right)^n \quad (5.20)$$

encountered in the proof of Theorem 5.3 (which is just the Laurent expansion around infinity).

**Theorem 5.6** (Beurling–Gelfand). *The spectral radius satisfies*

$$r(x) = \inf_{n \in \mathbb{N}} \|x^n\|^{1/n} = \lim_{n \rightarrow \infty} \|x^n\|^{1/n}. \quad (5.21)$$

**Proof.** By spectral mapping we have  $r(x)^n = r(x^n) \leq \|x^n\|$  and hence

$$r(x) \leq \inf \|x^n\|^{1/n}.$$

Conversely, fix  $\ell \in X^*$ , and consider

$$\ell((x - \alpha)^{-1}) = -\frac{1}{\alpha} \sum_{n=0}^{\infty} \frac{1}{\alpha^n} \ell(x^n). \quad (5.22)$$

Then  $\ell((x - \alpha)^{-1})$  is analytic in  $|\alpha| > r(x)$  and hence (5.22) converges absolutely for  $|\alpha| > r(x)$  by Cauchy's integral formula for derivatives. Hence for fixed  $\alpha$  with  $|\alpha| > r(x)$ ,  $\ell(x^n/\alpha^n)$  converges to zero for every  $\ell \in X^*$ . Since every weakly convergent sequence is bounded we have

$$\frac{\|x^n\|}{|\alpha|^n} \leq C(\alpha)$$

and thus

$$\limsup_{n \rightarrow \infty} \|x^n\|^{1/n} \leq \limsup_{n \rightarrow \infty} C(\alpha)^{1/n} |\alpha| = |\alpha|.$$

Since this holds for every  $|\alpha| > r(x)$  we have

$$r(x) \leq \inf \|x^n\|^{1/n} \leq \liminf_{n \rightarrow \infty} \|x^n\|^{1/n} \leq \limsup_{n \rightarrow \infty} \|x^n\|^{1/n} \leq r(x),$$

which finishes the proof.  $\square$

Note that it might be tempting to conjecture that the sequence  $\|x^n\|^{1/n}$  is monotone, however this is false in general – see Problem 5.11. By the ratio test, the Neumann series (5.20) converges for  $|\alpha| > r(x)$ .

Next let us look at some examples illustrating these ideas.

**Example 5.14.** In  $X := C(I)$  we have  $\sigma(x) = x(I)$  and hence  $r(x) = \|x\|_\infty$  for all  $x$ .  $\diamond$

**Example 5.15.** If  $X := \mathcal{L}(\mathbb{C}^2)$  and  $x := \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$  such that  $x^2 = 0$  and consequently  $r(x) = 0$ . This is not surprising, since  $x$  has the only eigenvalue 0. In particular, the spectral radius can be strictly smaller than the norm (note that  $\|x\| = 1$  in our example). The same is true for any nilpotent matrix. In general,  $x$  will be called **nilpotent** if  $x^n = 0$  for some  $n \in \mathbb{N}$  and any nilpotent element will satisfy  $r(x) = 0$ . Note that in this case the Neumann series terminates after  $n$  terms,

$$(x - \alpha)^{-1} = -\frac{1}{\alpha} \sum_{j=0}^{n-1} \left(\frac{x}{\alpha}\right)^j, \quad \alpha \neq 0,$$

and the resolvent has a pole of order  $n$  at 0.  $\diamond$

**Example 5.16.** Consider the linear **Volterra integral operator**

$$K(x)(t) := \int_0^t k(t, s)x(s)ds, \quad x \in C[0, 1]. \quad (5.23)$$

Then, using induction, it is not hard to verify (Problem 5.13)

$$|K^n(x)(t)| \leq \frac{\|k\|_\infty^n t^n}{n!} \|x\|_\infty. \quad (5.24)$$

Consequently

$$\|K^n x\|_\infty \leq \frac{\|k\|_\infty^n}{n!} \|x\|_\infty,$$

that is  $\|K^n\| \leq \frac{\|k\|_\infty^n}{n!}$ , which shows

$$r(K) \leq \lim_{n \rightarrow \infty} \frac{\|k\|_\infty}{(n!)^{1/n}} = 0.$$

Hence  $r(K) = 0$  and for every  $\lambda \in \mathbb{C}$  and every  $y \in C[0, 1]$  the equation

$$x - \lambda K x = y \quad (5.25)$$

has a unique solution given by

$$x = (\mathbb{I} - \lambda K)^{-1}y = \sum_{n=0}^{\infty} \lambda^n K^n y. \quad (5.26)$$

Note that  $\sigma(K) = \{0\}$  but 0 is in general not an eigenvalue (consider e.g.  $k(t, s) = 1$ ). Elements of a Banach algebra with  $r(x) = 0$  are called **quasinilpotent**. Since the Neumann series (5.20) converges for  $|\alpha| > 0$  in this case, the resolvent has an essential singularity at 0 if  $x$  is quasinilpotent (but not nilpotent).  $\diamond$

In the last two examples we have seen a strict inequality in (5.19). If we regard  $r(x)$  as a *spectral norm* for  $x$ , then the spectral norm does not control the *algebraic* norm in such a situation. On the other hand, if we had equality for some  $x$ , and moreover, this were also true for any polynomial  $p(x)$ , then spectral mapping would imply that the spectral norm  $\sup_{\alpha \in \sigma(x)} |p(\alpha)|$  equals the algebraic norm  $\|p(x)\|$  and convergence on one side would imply convergence on the other side. So by taking limits we could get an isometric identification of elements of the form  $f(x)$  with functions  $f \in C(\sigma(x))$ . But this is nothing but the content of the spectral theorem and self-adjointness will be the property which will make all this work.

We end this section with the remark that neither the spectrum nor the spectral radius is continuous. All one can say is

**Lemma 5.7.** *Let  $x_n \in X$  be a convergent sequence and  $x := \lim_{n \rightarrow \infty} x_n$ . Then whenever  $\alpha \in \rho(x)$  we have  $\alpha \in \rho(x_n)$  eventually and*

$$(x_n - \alpha)^{-1} \rightarrow (x - \alpha)^{-1}. \quad (5.27)$$

Moreover,

$$\lim_{n \rightarrow \infty} \sigma(x_n) \subseteq \sigma(x), \quad (5.28)$$

where  $\lim_{n \rightarrow \infty} \sigma(x_n) := \{\alpha \in \mathbb{C} \mid \exists \alpha_n \in \sigma(x_n) \rightarrow \alpha\}$ , and

$$r(x) \geq \limsup_{n \rightarrow \infty} r(x_n). \quad (5.29)$$

**Proof.** The first claim is immediate since taking the inverse is continuous by Corollary 5.2. Furthermore, Corollary 5.2 also shows that for  $\alpha \in \rho(A)$  and  $\|x - x_n\| + |\alpha - \alpha_n| < \|(x - \alpha)^{-1}\|^{-1}$  we have  $\alpha_n \in \rho(x_n)$ , which implies the second claim.

Concerning the last claim, observe that  $r(x_k) \leq \|x_k^n\|^{1/n}$  implies that  $\limsup_{k \rightarrow \infty} r(x_k) \leq \|x^n\|^{1/n}$ .  $\square$

**Example 5.17.** That the spectrum can expand is shown by the following example due to Kakutani. We consider the bounded linear operators on

$\ell^2(\mathbb{N})$  and look at shift-type operators of the form

$$(Aa)_j := q_j a_{j+1},$$

where  $q \in \ell^\infty(\mathbb{N})$ . Then we have  $\|A\| = \sup_{j \in \mathbb{N}} |q_j|$  and

$$(A^n a)_j = (q_j q_{j+1} \cdots q_{j+n-1}) a_{j+n}$$

with  $\|A^n\| = \sup_{j \in \mathbb{N}} |q_j q_{j+1} \cdots q_{j+n-1}|$ .

Now note that every integer can be written as  $j = 2^k(2l+1)$  and write  $k(j) := k$  in this case. Choose

$$q_j := e^{-k(j)}.$$

To compute the above products we group integers into blocks  $2^{m-1}, \dots, 2^m - 1$  of  $2^{m-1}$  elements and observe that  $K_m := \sum_{j=2^{m-1}}^{2^m-1} k(j) = 2^{m-1} - 1$ . Indeed, note that since odd numbers do not contribute to this sum, we can drop them and divide the remaining even ones by 2 to get the previous block. This shows  $K_m = 2^{m-2} + K_{m-1}$  and establishes the claim. Summing over all blocks we have  $\sum_{m=1}^n K_m = 2^n - n - 1$  implying

$$\|A^{2^n}\|^{1/2^n} = q_1 q_2 \cdots q_{2^n-1} = \exp(-1 + (n+1)2^{-n}).$$

Taking the limit  $n \rightarrow \infty$  shows  $r(A) = \frac{1}{e}$ .

Next define

$$(A_k a)_j := \begin{cases} 0, & k(j) = k, \\ q_j a_{j+1}, & \text{else,} \end{cases}$$

and observe that  $A_k$  is nilpotent since  $A_k^{2^{k+1}} = 0$ . Indeed note that  $(A_k a)_j = 0$  for  $j = 2^k, 2^k 3, 2^k 5, \dots$  which are a distance  $2^{k+1} - 1$  apart. Hence applying  $A$  once more the result will vanish at the previous points as well, etc. Moreover,

$$((A_k - A)a)_j = \begin{cases} q_j a_{j+1}, & k(j) = k, \\ 0, & \text{else,} \end{cases}$$

implying  $\|A_k - A\| = e^{-k}$ . Hence we have  $A_k \rightarrow A$  with  $r(A_k) = 0 \rightarrow 0 < e^{-1} = r(A)$  and  $\sigma(A_k) = \{0\} \rightarrow \{0\} \subsetneq \sigma(A)$ .  $\diamond$

**Problem\* 5.1.** Show that the multiplication in a Banach algebra  $X$  is continuous:  $x_n \rightarrow x$  and  $y_n \rightarrow y$  imply  $x_n y_n \rightarrow xy$ .

**Problem\* 5.2.** Suppose that  $X$  satisfies all requirements for a Banach algebra except that (5.3) is replaced by

$$\|xy\| \leq C\|x\|\|y\|, \quad C > 0.$$

Of course one can rescale the norm to reduce it to the case  $C = 1$ . However, this might have undesirable side effects in case there is a unit. Show that if  $X$  has a unit  $e$ , then  $\|e\| \geq C^{-1}$  and there is an equivalent norm  $\|\cdot\|_0$  which satisfies (5.3) and  $\|e\|_0 = 1$ .

Finally, note that for this construction to work it suffices to assume that multiplication is separately continuous by Problem 4.6.

(Hint: Identify  $x \in X$  with the operator  $L_x : X \rightarrow X$ ,  $y \mapsto xy$  in  $\mathcal{L}(X)$ . For the last part use the uniform boundedness principle.)

**Problem\* 5.3** (Unitization). Show that if  $X$  is a Banach algebra then  $\mathbb{C} \oplus X$  is a unital Banach algebra, where we set  $\|(\alpha, x)\| = |\alpha| + \|x\|$  and  $(\alpha, x)(\beta, y) = (\alpha\beta, \alpha y + \beta x + xy)$ .

**Problem 5.4.** Show  $\sigma(x^{-1}) = \sigma(x)^{-1}$  if  $x$  is invertible.

**Problem 5.5.** An element  $x \in X$  satisfying  $x^2 = x$  is called a **projection**. Compute the spectrum of a projection.

**Problem 5.6.** If  $X := \mathcal{L}(L^p(I))$ , then every  $x \in C(I)$  gives rise to a multiplication operator  $M_x \in X$  defined as  $M_x f := x f$ . Show  $r(M_x) = \|M_x\| = \|x\|_\infty$  and  $\sigma(M_x) = \text{Ran}(x)$ .

**Problem 5.7.** If  $X := \mathcal{L}(\ell^p(\mathbb{N}))$ , then every  $m \in \ell^\infty(\mathbb{N})$  gives rise to a multiplication operator  $M \in X$  defined as  $(Ma)_n := m_n a_n$ . Show  $r(M) = \|M\| = \|m\|_\infty$  and  $\sigma(M) = \overline{\text{Ran}(m)}$ .

**Problem 5.8.** Can every compact set  $K \subset \mathbb{C}$  arise as the spectrum of an element of some Banach algebra?

**Problem\* 5.9.** Suppose  $x$  has both a right inverse  $y$  (i.e.,  $xy = e$ ) and a left inverse  $z$  (i.e.,  $zx = e$ ). Show that  $y = z = x^{-1}$ .

**Problem\* 5.10.** Suppose  $xy$  and  $yx$  are both invertible, then so are  $x$  and  $y$ :

$$y^{-1} = (xy)^{-1}x = x(yx)^{-1}, \quad x^{-1} = (yx)^{-1}y = y(xy)^{-1}.$$

(Hint: Previous problem.)

**Problem\* 5.11.** Let  $X := \mathcal{L}(\mathbb{C}^2)$  and compute  $\|x^n\|^{1/n}$  for  $x := \begin{pmatrix} 0 & \alpha \\ \beta & 0 \end{pmatrix}$ . Conclude that this sequence is not monotone in general.

**Problem 5.12.** Let  $X := \ell^\infty(\mathbb{N})$ . Show  $\sigma(x) = \overline{\{x_n\}_{n \in \mathbb{N}}}$ . Also show that  $r(x) = \|x\|$  for all  $x \in X$ .

**Problem\* 5.13.** Show (5.24).

**Problem 5.14.** Show that  $L^1(\mathbb{R}^n)$  with convolution as multiplication is a commutative Banach algebra without identity (Hint: Lemma 3.20 from [48]).

**Problem 5.15.** Show the **first resolvent identity**

$$\begin{aligned} (x - \alpha)^{-1} - (x - \beta)^{-1} &= (\alpha - \beta)(x - \alpha)^{-1}(x - \beta)^{-1} \\ &= (\alpha - \beta)(x - \beta)^{-1}(x - \alpha)^{-1}, \end{aligned} \quad (5.30)$$

for  $\alpha, \beta \in \rho(x)$ .



**Problem 5.16.** Show  $\sigma(xy) \setminus \{0\} = \sigma(yx) \setminus \{0\}$ . (Hint: Find a relation between  $(xy - \alpha)^{-1}$  and  $(yx - \alpha)^{-1}$ .)

## 5.2. The $C^*$ algebra of operators and the spectral theorem

We begin by recalling that if  $\mathfrak{H}$  is some Hilbert space, then for every  $A \in \mathcal{L}(\mathfrak{H})$  we can define its adjoint  $A^* \in \mathcal{L}(\mathfrak{H})$ . Hence the Banach algebra  $\mathcal{L}(\mathfrak{H})$  has an additional operation in this case which will also give us self-adjointness, a property which has already turned out crucial for the spectral theorem in the case of compact operators. Even though this is not immediately evident, in some sense this additional structure adds the convenient geometric properties of Hilbert spaces to the picture.

A Banach algebra  $X$  together with an **involution** satisfying

$$(x + y)^* = x^* + y^*, \quad (\alpha x)^* = \alpha^* x^*, \quad x^{**} = x, \quad (xy)^* = y^* x^*, \quad (5.31)$$

and

$$\|x\|^2 = \|x^* x\| \quad (5.32)$$

is called a  $C^*$  **algebra**. Any subalgebra (we do not require a subalgebra to contain the identity) which is also closed under involution, is called a  $*$ -subalgebra.

The condition (5.32) might look a bit artificial at this point. Maybe a requirement like  $\|x^*\| = \|x\|$  might seem more natural. In fact, at this point the only justification is that it holds for our guiding example  $\mathcal{L}(\mathfrak{H})$  (cf. Lemma 2.14). Furthermore, it is important to emphasize that (5.32) is a rather strong condition as it implies that the norm is already uniquely determined by the algebraic structure. More precisely, Lemma 5.8 below implies that the norm of  $x$  can be computed from the spectral radius of  $x^*x$  via  $\|x\| = r(x^*x)^{1/2}$ . So while there might be several norms which turn  $X$  into a Banach algebra, there is at most one which will give a  $C^*$  algebra.

Note that (5.32) implies  $\|x\|^2 = \|x^* x\| \leq \|x\| \|x^*\|$  and hence  $\|x\| \leq \|x^*\|$ . By  $x^{**} = x$  this also implies  $\|x^*\| \leq \|x^{**}\| = \|x\|$  and hence

$$\|x\| = \|x^*\|, \quad \|x\|^2 = \|x^* x\| = \|x x^*\|. \quad (5.33)$$

**Example 5.18.** The continuous functions  $C(I)$  together with complex conjugation form a commutative  $C^*$  algebra.  $\diamond$

**Example 5.19.** The Banach algebra  $\mathcal{L}(\mathfrak{H})$  is a  $C^*$  algebra by Lemma 2.14. The compact operators  $\mathcal{K}(\mathfrak{H})$  are a  $*$ -subalgebra.  $\diamond$

**Example 5.20.** The bounded sequences  $\ell^\infty(\mathbb{N})$  together with complex conjugation form a commutative  $C^*$  algebra. The set  $c_0(\mathbb{N})$  of sequences converging to 0 are a  $*$ -subalgebra.  $\diamond$

If  $X$  has an identity  $e$ , we clearly have  $e^* = e$ ,  $\|e\| = 1$ ,  $(x^{-1})^* = (x^*)^{-1}$  (show this), and

$$\sigma(x^*) = \sigma(x)^*. \quad (5.34)$$

We will always assume that we have an identity and we note that it is always possible to add an identity (Problem 5.17).

If  $X$  is a  $C^*$  algebra, then  $x \in X$  is called **normal** if  $x^*x = xx^*$ , **self-adjoint** if  $x^* = x$ , and **unitary** if  $x^* = x^{-1}$ . Moreover,  $x$  is called **positive** if  $x = y^2$  for some  $y = y^* \in X$ . Clearly both self-adjoint and unitary elements are normal and positive elements are self-adjoint. If  $x$  is normal, then so is any polynomial  $p(x)$  (it will be self-adjoint if  $x$  is and  $p$  is real-valued).

As already pointed out in the previous section, it is crucial to identify elements for which the spectral radius equals the norm. The key ingredient will be (5.32) which implies  $\|x^2\| = \|x\|^2$  if  $x$  is self-adjoint. For unitary elements we have  $\|x\| = \sqrt{\|x^*x\|} = \sqrt{\|e\|} = 1$ . Moreover, for normal elements we get

**Lemma 5.8.** *If  $x \in X$  is normal, then  $\|x^2\| = \|x\|^2$  and  $r(x) = \|x\|$ .*

**Proof.** Using (5.32) three times we have

$$\|x^2\| = \|(x^2)^*(x^2)\|^{1/2} = \|(x^*x)^*(x^*x)\|^{1/2} = \|x^*x\| = \|x\|^2$$

and hence  $r(x) = \lim_{k \rightarrow \infty} \|x^{2^k}\|^{1/2^k} = \|x\|$ .  $\square$

The next result generalizes the fact that self-adjoint operators have only real eigenvalues.

**Lemma 5.9.** *If  $x$  is self-adjoint, then  $\sigma(x) \subseteq \mathbb{R}$ . If  $x$  is positive, then  $\sigma(x) \subseteq [0, \infty)$ .*

**Proof.** Suppose  $\alpha + i\beta \in \sigma(x)$ ,  $\lambda \in \mathbb{R}$ . Then  $\alpha + i(\beta + \lambda) \in \sigma(x + i\lambda)$  and

$$\alpha^2 + (\beta + \lambda)^2 \leq \|x + i\lambda\|^2 = \|(x + i\lambda)(x - i\lambda)\| = \|x^2 + \lambda^2\| \leq \|x\|^2 + \lambda^2.$$

Hence  $\alpha^2 + \beta^2 + 2\beta\lambda \leq \|x\|^2$  which gives a contradiction if we let  $|\lambda| \rightarrow \infty$  unless  $\beta = 0$ .

The second claim follows from the first using spectral mapping (Theorem 5.5).  $\square$

**Example 5.21.** If  $X := \mathcal{L}(\mathbb{C}^2)$  and  $x := \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$  then  $\sigma(x) = \{0\}$ . Hence the converse of the above lemma is not true in general.  $\diamond$

Given  $x \in X$  we can consider the  $C^*$  algebra  $C^*(x)$  (with identity) generated by  $x$  (i.e., the smallest closed  $*$ -subalgebra containing  $e$  and  $x$ ). If  $x$  is normal we explicitly have

$$C^*(x) = \overline{\{p(x, x^*) | p : \mathbb{C}^2 \rightarrow \mathbb{C} \text{ polynomial}\}}, \quad xx^* = x^*x, \quad (5.35)$$

and, in particular,  $C^*(x)$  is commutative (Problem 5.18). In the self-adjoint case this simplifies to

$$C^*(x) := \overline{\{p(x)|p : \mathbb{C} \rightarrow \mathbb{C} \text{ polynomial}\}}, \quad x = x^*. \quad (5.36)$$

Moreover, in this case  $C^*(x)$  is isomorphic to  $C(\sigma(x))$  (the continuous functions on the spectrum):

**Theorem 5.10** (Spectral theorem). *If  $X$  is a  $C^*$  algebra and  $x \in X$  is self-adjoint, then there is an isometric isomorphism  $\Phi : C(\sigma(x)) \rightarrow C^*(x)$  such that  $f(t) = t$  maps to  $\Phi(t) = x$  and  $f(t) = 1$  maps to  $\Phi(1) = e$ .*

*Moreover, for every  $f \in C(\sigma(x))$  we have*

$$\sigma(f(x)) = f(\sigma(x)), \quad (5.37)$$

*where  $f(x) = \Phi(f)$ .*

**Proof.** First of all,  $\Phi$  is well defined for polynomials  $p$  and given by  $\Phi(p) = p(x)$ . Moreover, since  $p(x)$  is normal, spectral mapping implies

$$\|p(x)\| = r(p(x)) = \sup_{\alpha \in \sigma(p(x))} |\alpha| = \sup_{\alpha \in \sigma(x)} |p(\alpha)| = \|p\|_\infty$$

for every polynomial  $p$ . Hence  $\Phi$  is isometric. Next we use that the polynomials are dense in  $C(\sigma(x))$ . In fact, to see this one can either consider a compact interval  $I$  containing  $\sigma(x)$  and use the Tietze extension theorem (Theorem B.30) to extend  $f$  to  $I$  and then approximate the extension using polynomials (Theorem 1.3) or use the Stone–Weierstraß theorem (Theorem B.42). Thus  $\Phi$  uniquely extends to a map on all of  $C(\sigma(x))$  by Theorem 1.16. By continuity of the norm this extension is again isometric. Similarly, we have  $\Phi(fg) = \Phi(f)\Phi(g)$  and  $\Phi(f)^* = \Phi(f^*)$  since both relations hold for polynomials.

To show  $\sigma(f(x)) = f(\sigma(x))$  fix some  $\alpha \in \mathbb{C}$ . If  $\alpha \notin f(\sigma(x))$ , then  $g(t) = \frac{1}{f(t) - \alpha} \in C(\sigma(x))$  and  $\Phi(g) = (f(x) - \alpha)^{-1} \in X$  shows  $\alpha \notin \sigma(f(x))$ . Conversely, if  $\alpha \notin \sigma(f(x))$  then  $g = \Phi^{-1}((f(x) - \alpha)^{-1}) = \frac{1}{f - \alpha}$  is continuous, which shows  $\alpha \notin f(\sigma(x))$ .  $\square$

In particular, this last theorem tells us that we have a functional calculus for self-adjoint operators, that is, if  $A \in \mathcal{L}(\mathfrak{H})$  is self-adjoint, then  $f(A)$  is well defined for every  $f \in C(\sigma(A))$ . Specifically, we can compute  $f(A)$  by choosing a sequence of polynomials  $p_n$  which converge to  $f$  uniformly on  $\sigma(A)$ , then we have  $p_n(A) \rightarrow f(A)$  in the operator norm. In particular, if  $f$  is given by a power series, then  $f(A)$  defined via  $\Phi$  coincides with  $f(A)$  defined via its power series (cf. Problem 1.39).

**Problem\* 5.17** (Unitization). *Show that if  $X$  is a non-unital  $C^*$  algebra then  $\mathbb{C} \oplus X$  is a unital  $C^*$  algebra, where we set  $\|(\alpha, x)\| := \sup\{\|\alpha y + xy\| | y \in X\}$ .*

$X, \|y\| \leq 1\}$ ,  $(\alpha, x)(\beta, y) = (\alpha\beta, \alpha y + \beta x + xy)$  and  $(\alpha, x)^* = (\alpha^*, x^*)$ . (Hint: It might be helpful to identify  $x \in X$  with the operator  $L_x : X \rightarrow X, y \mapsto xy$  in  $\mathcal{L}(X)$ . Moreover, note  $\|L_x\| = \|x\|$ .)

**Problem\* 5.18.** Let  $X$  be a  $C^*$  algebra and  $Y$  a  $*$ -subalgebra. Show that if  $Y$  is commutative, then so is  $\overline{Y}$ .

**Problem 5.19.** Show that the map  $\Phi$  from the spectral theorem is positivity preserving, that is,  $f \geq 0$  if and only if  $\Phi(f)$  is positive.

**Problem 5.20.** Show that the map  $\Phi$  from the spectral theorem is positivity preserving, that is,  $f \geq 0$  if and only if  $\Phi(f)$  is positive.

**Problem 5.21.** Let  $x$  be self-adjoint. Show that the following are equivalent:

- (i)  $\sigma(x) \subseteq [0, \infty)$ .
- (ii)  $x$  is positive.
- (iii)  $\|\lambda - x\| \leq \lambda$  for all  $\lambda \geq \|x\|$ .
- (iv)  $\|\lambda - x\| \leq \lambda$  for one  $\lambda \geq \|x\|$ .

**Problem 5.22.** Let  $A \in \mathcal{L}(\mathfrak{H})$ . Show that  $A$  is normal if and only if

$$\|Au\| = \|A^*u\|, \quad \forall u \in \mathfrak{H}.$$

In particular,  $\text{Ker}(A) = \text{Ker}(A^*)$ . (Hint: Problem 1.21.)

**Problem 5.23.** Show that the **Cayley transform** of a self-adjoint element  $x$ ,

$$y := (x - i)(x + i)^{-1}$$

is unitary. Show that  $1 \notin \sigma(y)$  and

$$x = i(1 + y)(1 - y)^{-1}.$$

**Problem 5.24.** Show if  $x$  is unitary then  $\sigma(x) \subseteq \{\alpha \in \mathbb{C} \mid |\alpha| = 1\}$ .

**Problem 5.25.** Suppose  $x$  is self-adjoint. Show that

$$\|(x - \alpha)^{-1}\| = \frac{1}{\text{dist}(\alpha, \sigma(x))}.$$

### 5.3. Spectral measures

The purpose of this section is to derive another formulation of the spectral theorem which is important in quantum mechanics. This reformulation requires familiarity with measure theory and can be skipped as the results will not be needed in the sequel.

Using the Riesz representation theorem we get a formulation in terms of spectral measures:

**Theorem 5.11.** *Let  $\mathfrak{H}$  be a Hilbert space, and let  $A \in \mathcal{L}(\mathfrak{H})$  be self-adjoint. For every  $u, v \in \mathfrak{H}$  there is a corresponding complex Borel measure  $\mu_{u,v}$  supported on  $\sigma(A)$  (the **spectral measure**) such that*

$$\langle u, f(A)v \rangle = \int_{\sigma(A)} f(t) d\mu_{u,v}(t), \quad f \in C(\sigma(A)). \quad (5.38)$$

We have

$$\mu_{u,v_1+v_2} = \mu_{u,v_1} + \mu_{u,v_2}, \quad \mu_{u,\alpha v} = \alpha \mu_{u,v}, \quad \mu_{v,u} = \mu_{u,v}^* \quad (5.39)$$

and  $|\mu_{u,v}|(\sigma(A)) \leq \|u\|\|v\|$ . Furthermore,  $\mu_u = \mu_{u,u}$  is a positive Borel measure with  $\mu_u(\sigma(A)) = \|u\|^2$ .

**Proof.** Consider the continuous functions on  $I = [-\|A\|, \|A\|]$  and note that every  $f \in C(I)$  gives rise to some  $f \in C(\sigma(A))$  by restricting its domain. Clearly  $\ell_{u,v}(f) = \langle u, f(A)v \rangle$  is a bounded linear functional and the existence of a corresponding measure  $\mu_{u,v}$  with  $|\mu_{u,v}|(I) = \|\ell_{u,v}\| \leq \|u\|\|v\|$  follows from the Riesz representation theorem (Theorem 6.5 from [48]). Since  $\ell_{u,v}(f)$  depends only on the value of  $f$  on  $\sigma(A) \subseteq I$ ,  $\mu_{u,v}$  is supported on  $\sigma(A)$ .

Moreover, if  $f \geq 0$  we have  $\ell_u(f) = \langle u, f(A)u \rangle = \langle f(A)^{1/2}u, f(A)^{1/2}u \rangle = \|f(A)^{1/2}u\|^2 \geq 0$  and hence  $\ell_u$  is positive and the corresponding measure  $\mu_u$  is positive. The rest follows from the properties of the scalar product.  $\square$

It is often convenient to regard  $\mu_{u,v}$  as a complex measure on  $\mathbb{R}$  by using  $\mu_{u,v}(\Omega) = \mu_{u,v}(\Omega \cap \sigma(A))$ . If we do this, we can also consider  $f$  as a function on  $\mathbb{R}$ . However, note that  $f(A)$  depends only on the values of  $f$  on  $\sigma(A)$ ! Moreover, it suffices to consider  $\mu_u$  since using the polarization identity (1.55) we have

$$\mu_{u,v}(\Omega) = \frac{1}{4}(\mu_{u+v}(\Omega) - \mu_{u-v}(\Omega) + i\mu_{u-iv}(\Omega) - i\mu_{u+iv}(\Omega)). \quad (5.40)$$

Now the last theorem can be used to define  $f(A)$  for every bounded measurable function  $f \in B(\sigma(A))$  via Lemma 2.12 and extend the functional calculus from continuous to measurable functions:

**Theorem 5.12** (Spectral theorem). *If  $\mathfrak{H}$  is a Hilbert space and  $A \in \mathcal{L}(\mathfrak{H})$  is self-adjoint, then there is an homomorphism  $\Phi : B(\sigma(A)) \rightarrow \mathcal{L}(\mathfrak{H})$  given by*

$$\langle u, f(A)v \rangle = \int_{\sigma(A)} f(t) d\mu_{u,v}(t), \quad f \in B(\sigma(A)). \quad (5.41)$$

Moreover, if  $f_n(t) \rightarrow f(t)$  pointwise and  $\sup_n \|f_n\|_\infty$  is bounded, then  $f_n(A)u \rightarrow f(A)u$  for every  $u \in \mathfrak{H}$ .

**Proof.** The map  $\Phi$  is a well-defined linear operator by Lemma 2.12 since we have

$$\left| \int_{\sigma(A)} f(t) d\mu_{u,v}(t) \right| \leq \|f\|_\infty |\mu_{u,v}|(\sigma(A)) \leq \|f\|_\infty \|u\| \|v\|$$

and (5.39). Next, observe that  $\Phi(f)^* = \Phi(f^*)$  and  $\Phi(fg) = \Phi(f)\Phi(g)$  holds at least for continuous functions. To obtain it for arbitrary bounded functions, choose a (bounded) sequence  $f_n$  converging to  $f$  in  $L^2(\sigma(A), d\mu_u)$  and observe

$$\|(f_n(A) - f(A))u\|^2 = \int |f_n(t) - f(t)|^2 d\mu_u(t)$$

(use  $\|h(A)u\|^2 = \langle h(A)u, h(A)u \rangle = \langle u, h(A)^* h(A)u \rangle$ ). Thus  $f_n(A)u \rightarrow f(A)u$  and for bounded  $g$  we also have that  $(gf_n)(A)u \rightarrow (gf)(A)u$  and  $g(A)f_n(A)u \rightarrow g(A)f(A)u$ . This establishes the case where  $f$  is bounded and  $g$  is continuous. Similarly, approximating  $g$  removes the continuity requirement from  $g$ .

The last claim follows since  $f_n \rightarrow f$  in  $L^2$  by dominated convergence in this case.  $\square$

Our final aim is to generalize Corollary 3.8 to bounded self-adjoint operators. Since the spectrum of an arbitrary self-adjoint might contain more than just eigenvalues we need to replace the sum by an integral. To this end we recall the family of Borel sets  $\mathfrak{B}(\mathbb{R})$  and begin by defining the **spectral projections**

$$P_A(\Omega) = \chi_\Omega(A), \quad \Omega \in \mathfrak{B}(\mathbb{R}), \quad (5.42)$$

such that

$$\mu_{u,v}(\Omega) = \langle u, P_A(\Omega)v \rangle. \quad (5.43)$$

By  $\chi_\Omega^2 = \chi_\Omega$  and  $\chi_\Omega^* = \chi_\Omega$  they are **orthogonal projections**, that is  $P^2 = P$  and  $P^* = P$ . Recall that any orthogonal projection  $P$  decomposes  $\mathfrak{H}$  into an orthogonal sum

$$\mathfrak{H} = \text{Ker}(P) \oplus \text{Ran}(P), \quad (5.44)$$

where  $\text{Ker}(P) = (\mathbb{I} - P)\mathfrak{H}$ ,  $\text{Ran}(P) = P\mathfrak{H}$ .

In addition, the spectral projections satisfy

$$P_A(\mathbb{R}) = \mathbb{I}, \quad P_A\left(\bigcup_{n=1}^{\infty} \Omega_n\right)u = \sum_{n=1}^{\infty} P_A(\Omega_n)u, \quad \Omega_n \cap \Omega_m = \emptyset, \quad m \neq n, \quad (5.45)$$

for every  $u \in \mathfrak{H}$ . Here the dot inside the union just emphasizes that the sets are mutually disjoint. Such a family of projections is called a **projection-valued measure**. Indeed the first claim follows since  $\chi_{\mathbb{R}} = 1$  and by  $\chi_{\Omega_1 \cup \Omega_2} = \chi_{\Omega_1} + \chi_{\Omega_2}$  if  $\Omega_1 \cap \Omega_2 = \emptyset$  the second claim follows at least for finite unions. The case of countable unions follows from the last part of the

previous theorem since  $\sum_{n=1}^N \chi_{\Omega_n} = \chi_{\bigcup_{n=1}^N \Omega_n} \rightarrow \chi_{\bigcup_{n=1}^{\infty} \Omega_n}$  pointwise (note that the limit will not be uniform unless the  $\Omega_n$  are eventually empty and hence there is no chance that this series will converge in the operator norm). Moreover, since all spectral measures are supported on  $\sigma(A)$  the same is true for  $P_A$  in the sense that

$$P_A(\sigma(A)) = \mathbb{I}. \quad (5.46)$$

I also remark that in this connection the corresponding distribution function

$$P_A(t) := P_A((-\infty, t]) \quad (5.47)$$

is called a **resolution of the identity**.

Using our projection-valued measure we can define an operator-valued integral as follows: For every simple function  $f = \sum_{j=1}^n \alpha_j \chi_{\Omega_j}$  (where  $\Omega_j = f^{-1}(\alpha_j)$ ), we set

$$\int_{\mathbb{R}} f(t) dP_A(t) := \sum_{j=1}^n \alpha_j P_A(\Omega_j). \quad (5.48)$$

By (5.43) we conclude that this definition agrees with  $f(A)$  from Theorem 5.12:

$$\int_{\mathbb{R}} f(t) dP_A(t) = f(A). \quad (5.49)$$

Extending this integral to functions from  $B(\sigma(A))$  by approximating such functions with simple functions we get an alternative way of defining  $f(A)$  for such functions. This can in fact be done by just using the definition of a projection-valued measure and hence there is a one-to-one correspondence between projection-valued measures (with bounded support) and (bounded) self-adjoint operators such that

$$A = \int t dP_A(t). \quad (5.50)$$

If  $P_A(\{\alpha\}) \neq 0$ , then  $\alpha$  is an eigenvalue and  $\text{Ran}(P_A(\{\alpha\}))$  is the corresponding eigenspace (Problem 5.27). The fact that eigenspaces to different eigenvalues are orthogonal now generalizes to

**Lemma 5.13.** *Suppose  $\Omega_1 \cap \Omega_2 = \emptyset$ . Then*

$$\text{Ran}(P_A(\Omega_1)) \perp \text{Ran}(P_A(\Omega_2)). \quad (5.51)$$

**Proof.** Clearly  $\chi_{\Omega_1} \chi_{\Omega_2} = \chi_{\Omega_1 \cap \Omega_2}$  and hence

$$P_A(\Omega_1)P_A(\Omega_2) = P_A(\Omega_1 \cap \Omega_2).$$

Now if  $\Omega_1 \cap \Omega_2 = \emptyset$ , then

$$\langle P_A(\Omega_1)u, P_A(\Omega_2)v \rangle = \langle u, P_A(\Omega_1)P_A(\Omega_2)v \rangle = \langle u, P_A(\emptyset)v \rangle = 0,$$

which shows that the ranges are orthogonal to each other.  $\square$

**Example 5.22.** Let  $A \in \mathcal{L}(\mathbb{C}^n)$  be some symmetric matrix and let  $\alpha_1, \dots, \alpha_m$  be its (distinct) eigenvalues. Then

$$A = \sum_{j=1}^m \alpha_j P_A(\{\alpha_j\}),$$

where  $P_A(\{\alpha_j\})$  is the projection onto the eigenspace  $\text{Ker}(A - \alpha_j)$  corresponding to the eigenvalue  $\alpha_j$  by Problem 5.27. In fact, using that  $P_A$  is supported on the spectrum,  $P_A(\sigma(A)) = \mathbb{I}$ , we see

$$P_A(\Omega) = P_A(\sigma(A))P_A(\Omega) = P_A(\sigma(A) \cap \Omega) = \sum_{\alpha_j \in \Omega} P_A(\{\alpha_j\}).$$

Hence using that any  $f \in B(\sigma(A))$  is given as a simple function  $f = \sum_{j=1}^m f(\alpha_j) \chi_{\{\alpha_j\}}$  we obtain

$$f(A) = \int f(t) dP_A(t) = \sum_{j=1}^m f(\alpha_j) P_A(\{\alpha_j\}).$$

In particular, for  $f(t) = t$  we recover the above representation for  $A$ .  $\diamond$

**Example 5.23.** Let  $A \in \mathcal{L}(\mathbb{C}^n)$  be self-adjoint and let  $\alpha$  be an eigenvalue. Let  $P = P_A(\{\alpha\})$  be the projection onto the corresponding eigenspace and consider the restriction  $\tilde{A} = A|_{\tilde{\mathfrak{H}}}$  onto the orthogonal complement of this eigenspace  $\tilde{\mathfrak{H}} = (1 - P)\mathfrak{H}$ . Then by Lemma 5.13 we have  $\mu_{u,v}(\{\alpha\}) = 0$  for  $u, v \in \tilde{\mathfrak{H}}$ . Hence the integral in (5.41) does not see the point  $\alpha$  in the sense that

$$\langle u, f(A)v \rangle = \int_{\sigma(A)} f(t) d\mu_{u,v}(t) = \int_{\sigma(A) \setminus \{\alpha\}} f(t) d\mu_{u,v}(t), \quad u, v \in \tilde{\mathfrak{H}}.$$

Hence  $\Phi$  extends to a homomorphism  $\tilde{\Phi} : B(\sigma(A) \setminus \{\alpha\}) \rightarrow \mathcal{L}(\tilde{\mathfrak{H}})$ . In particular, if  $\alpha$  is an isolated eigenvalue, that is  $(\alpha - \varepsilon, \alpha + \varepsilon) \cap \sigma(A) = \{\alpha\}$  for  $\varepsilon > 0$  sufficiently small, we have  $(\cdot - \alpha)^{-1} \in B(\sigma(A) \setminus \{\alpha\})$  and hence  $\alpha \in \rho(\tilde{A})$ .  $\diamond$

**Problem 5.26.** Suppose  $A$  is self-adjoint. Let  $\alpha$  be an eigenvalue and  $u$  a corresponding normalized eigenvector. Show  $\int f(t) d\mu_u(t) = f(\alpha)$ , that is,  $\mu_u$  is the Dirac delta measure (with mass one) centered at  $\alpha$ .

**Problem\* 5.27.** Suppose  $A$  is self-adjoint. Show

$$\text{Ran}(P_A(\{\alpha\})) = \text{Ker}(A - \alpha).$$

(Hint: Start by verifying  $\text{Ran}(P_A(\{\alpha\})) \subseteq \text{Ker}(A - \alpha)$ . To see the converse, let  $u \in \text{Ker}(A - \alpha)$  and use the previous example.)





---

*Part 2*

# Advanced Functional Analysis



# More on convexity

## 6.1. The geometric Hahn–Banach theorem

The Hahn–Banach theorem is about constructing (continuous) linear functionals with given properties. In our original version this was done by showing that a functional defined on a subspace can be extended to the entire space. In this section we will establish a geometric version which establishes existence of functionals separating given convex sets. The key ingredient will be an association between convex sets and convex functions such that we can apply our original version of the Hahn–Banach theorem.

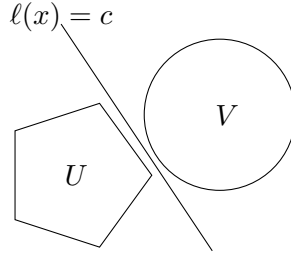
Let  $X$  be a vector space. For every subset  $U \subset X$  we define its **Minkowski functional** (or **gauge**)

$$p_U(x) := \inf\{t > 0 \mid x \in tU\}. \quad (6.1)$$

Here  $tU := \{tx \mid x \in U\}$ . Note that  $0 \in U$  implies  $p_U(0) = 0$  and  $p_U(x)$  will be finite for all  $x$  when  $U$  is **absorbing**, that is, for every  $x \in X$  there is some  $r$  such that  $x \in \alpha U$  for every  $|\alpha| \geq r$ . Note that every absorbing set contains 0 and every neighborhood of 0 in a Banach space is absorbing. Also observe that for  $U \subseteq V$  we have  $p_V(x) \leq p_U(x)$  for all  $x \in X$ .

**Example 6.1.** Let  $X$  be a Banach space and  $U := B_1(0)$ , then  $p_U(x) = \|x\|$ . If  $X := \mathbb{R}^2$  and  $U := (-1, 1) \times \mathbb{R}$  then  $p_U(x) = |x_1|$ . If  $X := \mathbb{R}^2$  and  $U := (-1, 1) \times \{0\}$  then  $p_U(x) = |x_1|$  if  $x_2 = 0$  and  $p_U(x) = \infty$  else.  $\diamond$

We will only need minimal requirements and it will suffice if  $X$  is a **topological vector space**, that is, a vector space which carries a topology such that both vector addition  $X \times X \rightarrow X$  and scalar multiplication  $\mathbb{C} \times X \rightarrow X$  are continuous mappings. Of course every normed vector space is a topological vector space with the usual topology generated by open balls.



**Figure 6.1.** Separation of convex sets via a hyperplane

As in the case of normed linear spaces,  $X^*$  will denote the vector space of all continuous linear functionals on  $X$ .

**Lemma 6.1.** *Let  $X$  be a vector space and  $U$  a convex subset containing 0. Then*

$$p_U(x + y) \leq p_U(x) + p_U(y), \quad p_U(\lambda x) = \lambda p_U(x), \quad \lambda \geq 0. \quad (6.2)$$

Moreover,  $\{x | p_U(x) < 1\} \subseteq U \subseteq \{x | p_U(x) \leq 1\}$ . If, in addition,  $X$  is a topological vector space and  $U$  is open, then  $U = \{x | p_U(x) < 1\}$ .

**Proof.** The homogeneity condition  $p(\lambda x) = \lambda p(x)$  for  $\lambda > 0$  is straightforward. To see the sublinearity, suppose  $p_U(x)$  and  $p_U(y)$  are both finite (otherwise there is nothing to do) and let  $t, s > 0$  with  $x \in tU$  and  $y \in sU$ , then

$$\frac{t}{t+s} \frac{x}{t} + \frac{s}{t+s} \frac{y}{s} = \frac{x+y}{t+s}$$

is in  $U$  by convexity. Moreover,  $p_U(x+y) \leq s+t$  and taking the infimum over all  $t$  and  $s$  we find  $p_U(x+y) \leq p_U(x) + p_U(y)$ .

Suppose  $p_U(x) < 1$ , then  $t^{-1}x \in U$  for some  $t < 1$  and thus  $x = t(t^{-1}x) + (1-t)0 \in U$  by convexity. Similarly, if  $x \in U$  then  $t^{-1}x = t^{-1}x + (1-t^{-1})0 \in U$  for  $t \geq 1$  by convexity and thus  $p_U(x) \leq 1$ . Finally, let  $U$  be open and  $x \in U$ , then  $(1+\varepsilon)x \in U$  for some  $\varepsilon > 0$  and thus  $p(x) \leq (1+\varepsilon)^{-1}$ .  $\square$

Note that (6.2) implies convexity

$$p_U(\lambda x + (1-\lambda)y) \leq \lambda p_U(x) + (1-\lambda)p_U(y), \quad \lambda \in [0, 1]. \quad (6.3)$$

**Theorem 6.2** (geometric Hahn–Banach, real version). *Let  $U, V$  be disjoint nonempty convex subsets of a real topological vector space  $X$  and let  $U$  be open. Then there is a linear functional  $\ell \in X^*$  and some  $c \in \mathbb{R}$  such that*

$$\ell(x) < c \leq \ell(y), \quad x \in U, y \in V. \quad (6.4)$$

If  $V$  is also open, then the second inequality is also strict.

**Proof.** Choose  $x_0 \in U$  and  $y_0 \in V$ , then

$$W := (U - x_0) - (V - y_0) = \{(x - x_0) - (y - y_0) | x \in U, y \in V\}$$

is open (since  $U$  is), convex (since  $U$  and  $V$  are) and contains 0. Moreover, since  $U$  and  $V$  are disjoint we have  $z_0 = y_0 - x_0 \notin W$ . By the previous lemma, the associated Minkowski functional  $p_W$  is convex and by the Hahn–Banach theorem (Theorem 4.9) there is a linear functional satisfying

$$\ell(tz_0) = t, \quad \ell(x) \leq p_W(x).$$

Note that since  $z_0 \notin W$  we have  $p_W(z_0) \geq 1$ . Moreover, given  $\varepsilon > 0$  we consider  $W_\varepsilon := (\varepsilon W) \cap (-\varepsilon W)$  which is an open neighborhood of 0. Then, for  $x \in W_\varepsilon$  we have  $p_W(\pm x) \leq \varepsilon$  implying  $|\ell(x)| \leq \varepsilon$ . This shows that  $\ell$  is continuous at 0 and hence continuous everywhere by linearity.

Finally we again use  $p_W(z) < 1$  for  $z \in W$  implying

$$\ell(x) - \ell(y) + 1 = \ell(x - y + z_0) \leq p_W(x - y + z_0) < 1$$

and hence  $\ell(x) < \ell(y)$  for  $x \in U$  and  $y \in V$ . Therefore  $\ell(U)$  and  $\ell(V)$  are disjoint convex subsets of  $\mathbb{R}$ . Finally, let us suppose that there is some  $x_1 \in U$  for which  $\ell(x_1) = \sup \ell(U)$ . Then, by continuity of the map  $t \mapsto x_1 + tz_0$  there is some  $\varepsilon > 0$  such that  $x_1 + \varepsilon z_0 \in U$ . But this gives a contradiction  $\ell(x_1) + \varepsilon = \ell(x_1 + \varepsilon z_0) \leq \ell(x_1)$ . Thus the claim holds with  $c = \sup \ell(U)$ . If  $V$  is also open, an analogous argument shows  $\inf \ell(V) < \ell(y)$  for all  $y \in V$ .  $\square$

Of course there is also a complex version.

**Theorem 6.3** (geometric Hahn–Banach, complex version). *Let  $U, V$  be disjoint nonempty convex subsets of a topological vector space  $X$  and let  $U$  be open. Then there is a linear functional  $\ell \in X^*$  and some  $c \in \mathbb{R}$  such that*

$$\operatorname{Re}(\ell(x)) < c \leq \operatorname{Re}(\ell(y)), \quad x \in U, y \in V. \quad (6.5)$$

*If in addition  $V$  is open, then the second inequality is also strict.*

**Proof.** Consider  $X$  as a real topological vector space. Then there is a continuous real-linear functional  $\ell_r : X \rightarrow \mathbb{R}$  by the real version of the geometric Hahn–Banach theorem. Then  $\ell(x) = \ell_r(x) - i\ell_r(ix)$  is the functional we are looking for (check this).  $\square$

**Example 6.2.** The assumption that one set is open is crucial. If you let  $X := \mathbb{R}^2$  and consider  $U_0 := \{x \in \mathbb{R}^2 | x_1 > 0\}$ ,  $V := \{0\}$ . Then both sets can be separated as in the theorem using  $\ell(x) := -x_1$  (and  $c = 0$ ). However, if we consider  $U := U_0 \cup \{(x_1, 0) | x_1 > 0\}$  this is no longer possible.  $\diamond$

Note that two disjoint closed convex sets can be separated strictly if one of them is compact. However, this will require that every point has a neighborhood base of convex open sets. Such topological vector spaces

are called **locally convex spaces** and they will be discussed further in Section 6.4. For now we just remark that every normed vector space is locally convex since balls are convex.

**Corollary 6.4.** *Let  $U, V$  be disjoint nonempty closed convex subsets of a locally convex space  $X$  and let  $U$  be compact. Then there is a linear functional  $\ell \in X^*$  and some  $c, d \in \mathbb{R}$  such that*

$$\operatorname{Re}(\ell(x)) \leq d < c \leq \operatorname{Re}(\ell(y)), \quad x \in U, y \in V. \quad (6.6)$$

**Proof.** Since  $V$  is closed, for every  $x \in U$  there is a convex open neighborhood  $N_x$  of 0 such that  $x + N_x$  does not intersect  $V$ . By compactness of  $U$  there are  $x_1, \dots, x_n$  such that the corresponding neighborhoods  $x_j + \frac{1}{2}N_{x_j}$  cover  $U$ . Set  $N := \bigcap_{j=1}^n N_{x_j}$  which is a convex open neighborhood of 0. Then

$$\tilde{U} := U + \frac{1}{2}N \subseteq \bigcup_{j=1}^n (x_j + \frac{1}{2}N_{x_j}) + \frac{1}{2}N \subseteq \bigcup_{j=1}^n (x_j + \frac{1}{2}N_{x_j} + \frac{1}{2}N_{x_j}) = \bigcup_{j=1}^n (x_j + N_{x_j})$$

is a convex open set which is disjoint from  $V$ . Hence by the previous theorem we can find some  $\ell$  such that  $\operatorname{Re}(\ell(x)) < c \leq \operatorname{Re}(\ell(y))$  for all  $x \in \tilde{U}$  and  $y \in V$ . Moreover, since  $\operatorname{Re}(\ell(U))$  is a compact interval  $[e, d]$ , the claim follows.  $\square$

Note that if  $U$  and  $V$  are **absolutely convex**, that is,  $\alpha U + \beta U \subseteq U$  for  $|\alpha| + |\beta| \leq 1$ , then we can write the previous condition equivalently as

$$|\ell(x)| \leq d < c \leq |\ell(y)|, \quad x \in U, y \in V, \quad (6.7)$$

since  $x \in U$  implies  $\theta x \in U$  for  $\theta = \operatorname{sign}(\ell(x))$  and thus  $|\ell(x)| = \theta \ell(x) = \ell(\theta x) = \operatorname{Re}(\ell(\theta x))$ .

From the last corollary we can also obtain versions of Corollaries 4.13 and 4.11 for locally convex vector spaces.

**Corollary 6.5.** *Let  $Y \subseteq X$  be a subspace of a locally convex space and let  $x_0 \in X \setminus \overline{Y}$ . Then there exists an  $\ell \in X^*$  such that (i)  $\ell(y) = 0$ ,  $y \in Y$  and (ii)  $\ell(x_0) = 1$ .*

**Proof.** Consider  $\ell$  from Corollary 6.4 applied to  $U = \{x_0\}$  and  $V = \overline{Y}$ . Now observe that  $\ell(Y)$  must be a subspace of  $\mathbb{C}$  and hence  $\ell(Y) = \{0\}$  implying  $\operatorname{Re}(\ell(x_0)) < 0$ . Finally  $\ell(x_0)^{-1}\ell$  is the required functional.  $\square$

**Corollary 6.6.** *Let  $Y \subseteq X$  be a subspace of a locally convex space and let  $\ell : Y \rightarrow \mathbb{C}$  be a continuous linear functional. Then there exists a continuous extension  $\bar{\ell} \in X^*$ .*

**Proof.** Without loss of generality we can assume that  $\ell$  is nonzero such that we can find  $x_0 \in Y$  with  $\ell(x_0) = 1$ . Since  $Y$  has the subset topology  $x_0 \notin Y_0 := \overline{\text{Ker}(\ell)}$ , where the closure is taken in  $X$ . Now Corollary 6.5 gives a functional  $\bar{\ell}$  with  $\bar{\ell}(x_0) = 1$  and  $Y_0 \subseteq \text{Ker}(\bar{\ell})$ . Moreover,

$$\bar{\ell}(x) - \ell(x) = \bar{\ell}(x) - \ell(x)\bar{\ell}(x_0) = \bar{\ell}(x - \ell(x)x_0) = 0, \quad x \in Y,$$

since  $x - \ell(x)x_0 \in \text{Ker}(\ell) \subseteq Y_0$ .  $\square$

**Problem 6.1.** Suppose  $\alpha \in \mathbb{C} \setminus \{0\}$ . Show that  $\alpha U = U$  implies  $p_U(\alpha x) = p_U(x)$ . Show that if  $U$  is symmetric,  $U = -U$ , and convex, then  $p_U$  satisfies the inverse triangle inequality  $|p_U(x) - p_U(y)| \leq p_U(x - y)$ .

**Problem 6.2.** Show that in a Banach space  $B_r(0) \subseteq U$  implies  $p_U(x) \leq \frac{1}{r}\|x\|$ .

**Problem\* 6.3.** Let  $X$  be a topological vector space. Show that  $U + V$  is open if one of the sets is open.

**Problem 6.4.** Show that Corollary 6.4 fails even in  $\mathbb{R}^2$  unless one set is compact.

**Problem 6.5.** Let  $X$  be a topological vector space and  $M \subseteq X$ ,  $N \subseteq X^*$ . Then the corresponding **polar**, **prepolar sets** are

$$M^\circ := \{\ell \in X^* \mid |\ell(x)| \leq 1 \forall x \in M\}, \quad N_\circ := \{x \in X \mid |\ell(x)| \leq 1 \forall \ell \in N\},$$

respectively. Show

- (i)  $M^\circ$  is closed and absolutely convex.
- (ii)  $M_1 \subseteq M_2$  implies  $M_2^\circ \subseteq M_1^\circ$ .
- (iii) For  $\alpha \neq 0$  we have  $(\alpha M)^\circ = |\alpha|^{-1}M^\circ$ .
- (iv) If  $M$  is a subspace we have  $M^\circ = M^\perp$ .

The same claims hold for prepolar sets.

**Problem 6.6** (Bipolar theorem). Let  $X$  be a locally convex space and suppose  $M \subseteq X$  is absolutely convex. Show  $(M^\circ)_\circ = \overline{M}$ . (Hint: Use Corollary 6.4 to show that for every  $y \notin \overline{M}$  there is some  $\ell \in X^*$  with  $\text{Re}(\ell(x)) \leq 1 < \ell(y)$ ,  $x \in \overline{M}$ .)

## 6.2. Convex sets and the Krein–Milman theorem

Let  $X$  be a locally convex vector space. Since the intersection of arbitrary convex sets is again convex we can define the convex hull of a set  $U$  as the smallest convex set containing  $U$ , that is, the intersection of all convex sets



containing  $U$ . It is straightforward to show (Problem 6.7) that the convex hull is given by

$$\text{conv}(U) := \left\{ \sum_{j=1}^n \lambda_j x_j \mid n \in \mathbb{N}, x_j \in U, \sum_{j=1}^n \lambda_j = 1, \lambda_j \geq 0 \right\}. \quad (6.8)$$

A line segment is convex and can be generated as the convex hull of its endpoints. Similarly, a full triangle is convex and can be generated as the convex hull of its vertices. However, if we look at a ball, then we need its entire boundary to recover it as the convex hull. So how can we characterize those points which determine a convex set via the convex hull?

Let  $K$  be a set and  $M \subseteq K$  a nonempty subset. Then  $M$  is called an **extremal subset** of  $K$  if no point of  $M$  can be written as a convex combination of two points unless both are in  $M$ : For given  $x, y \in K$  and  $\lambda \in (0, 1)$  we have that

$$\lambda x + (1 - \lambda)y \in M \Rightarrow x, y \in M. \quad (6.9)$$

If  $M = \{x\}$  is extremal, then  $x$  is called an **extremal point** of  $K$ . Hence an extremal point cannot be written as a convex combination of two other points from  $K$ .

Note that we did not require  $K$  to be convex. If  $K$  is convex and  $M$  is extremal, then  $K \setminus M$  is convex. Conversely, if  $K$  and  $K \setminus \{x\}$  are convex, then  $x$  is an extremal point. Note that the nonempty intersection of extremal sets is extremal. Moreover, if  $L \subseteq M$  is extremal and  $M \subseteq K$  is extremal, then  $L \subseteq K$  is extremal as well (Problem 6.8).

**Example 6.3.** Consider  $\mathbb{R}^2$  with the norms  $\|\cdot\|_p$ . Then the extremal points of the closed unit ball (cf. Figure 1.1) are the boundary points for  $1 < p < \infty$  and the vertices for  $p = 1, \infty$ . In any case the boundary is an extremal set. Slightly more general, in a strictly convex space, (ii) of Problem 1.13 says that the extremal points of the unit ball are precisely its boundary points.  $\diamond$

**Example 6.4.** Consider  $\mathbb{R}^3$  and let  $C := \{(x_1, x_2, 0) \in \mathbb{R}^3 \mid x_1^2 + x_2^2 = 1\}$ . Take two more points  $x_{\pm} = (0, 0, \pm 1)$  and consider the convex hull  $K$  of  $M := C \cup \{x_+, x_-\}$ . Then  $M$  is extremal in  $K$  and, moreover, every point from  $M$  is an extremal point. However, if we change the two extra points to be  $x_{\pm} = (1, 0, \pm 1)$ , then the point  $(1, 0, 0)$  is no longer extremal. Hence the extremal points are now  $M \setminus \{(1, 0, 0)\}$ . Note in particular that the set of extremal points is not closed in this case.  $\diamond$

Extremal sets arise naturally when minimizing linear functionals.

**Lemma 6.7.** Suppose  $K \subseteq X$  and  $\ell \in X^*$ . If

$$K_{\ell} := \{x \in K \mid \text{Re}(\ell(x)) = \inf_{y \in K} \text{Re}(\ell(y))\}$$

is nonempty (e.g. if  $K$  is compact), then it is extremal in  $K$ . If  $K$  is closed and convex, then  $K_\ell$  is closed and convex.

**Proof.** Set  $m = \inf_{y \in K} \operatorname{Re}(\ell(y))$ . Let  $x, y \in K$ ,  $\lambda \in (0, 1)$  and suppose  $\lambda x + (1 - \lambda)y \in K_\ell$ . Then

$$m = \operatorname{Re}(\ell(\lambda x + (1 - \lambda)y)) = \lambda \operatorname{Re}(\ell(x)) + (1 - \lambda) \operatorname{Re}(\ell(y)) \geq \lambda m + (1 - \lambda)m = m$$

with strict inequality if  $\operatorname{Re}(\ell(x)) > m$  or  $\operatorname{Re}(\ell(y)) > m$ . Hence we must have  $x, y \in K_\ell$ . Finally, by linearity  $K_\ell$  is convex and by continuity it is closed.  $\square$

If  $K$  is a closed convex set, then nonempty subsets of the type  $K_\ell$  are called **faces** of  $K$  and  $H_\ell := \{x \in X \mid \operatorname{Re}(\ell(x)) = \inf_{y \in K} \operatorname{Re}(\ell(y))\}$  is called a **support hyperplane** of  $K$ .

Conversely, if  $K$  is convex with nonempty interior, then every point  $x$  on the boundary has a supporting hyperplane (observe that the interior is convex and apply the geometric Hahn–Banach theorem with  $U = K^\circ$  and  $V = \{x\}$ ).

Next we want to look into existence of extremal points.

**Example 6.5.** Note that an interior point can never be extremal as it can be written as convex combination of some neighboring points. In particular, an open convex set will not have any extremal points (e.g.  $X$ , which is also closed, has no extremal points).  $\diamond$

**Example 6.6.** Suppose  $X$  is a strictly convex Banach space. Then every nonempty compact subset  $K$  has an extremal point. Indeed, let  $x \in K$  be such that  $\|x\| = \sup_{y \in K} \|y\|$ , then  $x$  is extremal: If  $x = \lambda y + (1 - \lambda)z$  then  $\|x\| \leq \lambda \|y\| + (1 - \lambda)\|z\| \leq \|x\|$  shows that we have equality in the triangle inequality and hence  $x = y = z$  by Problem 1.13 (i).  $\diamond$

**Example 6.7.** In a not strictly convex space the situation is quite different. For example, consider the closed unit ball in  $\ell^\infty(\mathbb{N})$ . Let  $a \in \ell^\infty(\mathbb{N})$  with  $\|a\|_\infty = 1$ . If there is some index  $j$  such that  $\lambda := |a_j| < 1$  then  $a = \frac{1}{2}b + \frac{1}{2}c$  where  $b = a + \varepsilon \delta^j$  and  $c = a - \varepsilon \delta^j$  with  $\varepsilon \leq 1 - |a_j|$ . Hence the only possible extremal points are those with  $|a_j| = 1$  for all  $j \in \mathbb{N}$ . If we have such an  $a$ , then if  $a = \lambda b + (1 - \lambda)c$  we must have  $1 = |\lambda b_n + (1 - \lambda)c_n| \leq \lambda |b_n| + (1 - \lambda)|c_n| \leq 1$  and hence  $a_n = b_n = c_n$  by strict convexity of the absolute value. Hence all such sequences are extremal.

However, if we consider  $c_0(\mathbb{N})$  the same argument shows that the closed unit ball contains no extremal points. In particular, the following lemma implies that there is no locally convex topology for which the closed unit ball in  $c_0(\mathbb{N})$  is compact. Together with the Banach–Alaoglu theorem (Theorem 6.10) this will show that  $c_0(\mathbb{N})$  is not the dual of any Banach space.  $\diamond$

**Lemma 6.8** (Krein–Milman). *Let  $X$  be a locally convex Hausdorff space. Suppose  $K \subseteq X$  is compact and nonempty. Then it contains at least one extremal point.*

**Proof.** We want to apply Zorn's lemma. To this end consider the family

$$\mathcal{M} = \{M \subseteq K \mid \text{compact and extremal in } K\}$$

with the partial order given by reversed inclusion. Since  $K \in \mathcal{M}$  this family is nonempty. Moreover, given a linear chain  $\mathcal{C} \subset \mathcal{M}$  we consider  $M := \bigcap \mathcal{C}$ . Then  $M \subseteq K$  is nonempty by the finite intersection property and since it is closed also compact. Moreover, as the nonempty intersection of extremal sets it is also extremal. Hence  $M \in \mathcal{M}$  and thus  $\mathcal{M}$  has a maximal element. Denote this maximal element by  $M$ .

We will show that  $M$  contains precisely one point (which is then extremal by construction). Indeed, suppose  $x, y \in M$ . If  $x \neq y$  then  $\{x\}, \{y\}$  are two disjoint and compact sets to which we can apply Corollary 6.4 to obtain a linear functional  $\ell \in X^*$  with  $\operatorname{Re}(\ell(x)) \neq \operatorname{Re}(\ell(y))$ . Then by Lemma 6.7  $M_\ell \subset M$  is extremal in  $M$  and hence also in  $K$ . But by  $\operatorname{Re}(\ell(x)) \neq \operatorname{Re}(\ell(y))$  it cannot contain both  $x$  and  $y$  contradicting maximality of  $M$ .  $\square$

Finally, we want to recover a convex set as the convex hull of its extremal points. In our infinite dimensional setting an additional closure will be necessary in general.

Since the intersection of arbitrary closed convex sets is again closed and convex we can define the closed convex hull of a set  $U$  as the smallest closed convex set containing  $U$ , that is, the intersection of all closed convex sets containing  $U$ . Since the closure of a convex set is again convex (Problem 6.11) the closed convex hull is simply the closure of the convex hull.

**Theorem 6.9** (Krein–Milman). *Let  $X$  be a locally convex Hausdorff space. Suppose  $K \subseteq X$  is convex and compact. Then it is the closed convex hull of its extremal points.*

**Proof.** Let  $E$  be the extremal points and  $M := \overline{\operatorname{conv}(E)} \subseteq K$  be its closed convex hull. Suppose  $x \in K \setminus M$  and use Corollary 6.4 to choose a linear functional  $\ell \in X^*$  with

$$\min_{y \in M} \operatorname{Re}(\ell(y)) > \operatorname{Re}(\ell(x)) \geq \min_{y \in K} \operatorname{Re}(\ell(y)).$$

Now consider  $K_\ell$  from Lemma 6.7 which is nonempty and hence contains an extremal point  $y \in E$ . But  $y \notin M$ , a contradiction.  $\square$

While in the finite dimensional case the closure is not necessary (Problem 6.13), it is important in general as the following example shows.

**Example 6.8.** Consider the closed unit ball in  $\ell^1(\mathbb{N})$ . Then the extremal points are  $\{e^{i\theta}\delta^n | n \in \mathbb{N}, \theta \in \mathbb{R}\}$ . Indeed, suppose  $\|a\|_1 = 1$  with  $\lambda := |a_j| \in (0, 1)$  for some  $j \in \mathbb{N}$ . Then  $a = \lambda b + (1 - \lambda)c$  where  $b := \lambda^{-1}a_j\delta^j$  and  $c := (1 - \lambda)^{-1}(a - a_j\delta^j)$ . Hence the only possible extremal points are of the form  $e^{i\theta}\delta^n$ . Moreover, if  $e^{i\theta}\delta^n = \lambda b + (1 - \lambda)c$  we must have  $1 = |\lambda b_n + (1 - \lambda)c_n| \leq \lambda|b_n| + (1 - \lambda)|c_n| \leq 1$  and hence  $a_n = b_n = c_n$  by strict convexity of the absolute value. Thus the convex hull of the extremal points are the sequences from the unit ball which have finitely many terms nonzero. While the closed unit ball is not compact in the norm topology it will be in the weak-\* topology by the Banach–Alaoglu theorem (Theorem 6.10). To this end note that  $\ell^1(\mathbb{N}) \cong c_0(\mathbb{N})^*$ .  $\diamond$

Also note that in the infinite dimensional case the extremal points can be dense.

**Example 6.9.** Let  $X = C([0, 1], \mathbb{R})$  and consider the convex set  $K = \{f \in C^1([0, 1], \mathbb{R}) | f(0) = 0, \|f'\|_\infty \leq 1\}$ . Note that the functions  $f_\pm(x) = \pm x$  are extremal. For example, assume

$$x = \lambda f(x) + (1 - \lambda)g(x)$$

then

$$1 = \lambda f'(x) + (1 - \lambda)g'(x)$$

which implies  $f'(x) = g'(x) = 1$  and hence  $f(x) = g(x) = x$ .

To see that there are no other extremal functions, suppose  $|f'(x)| \leq 1 - \varepsilon$  on some interval  $I$ . Choose a nontrivial continuous function  $g$  which is 0 outside  $I$  and has integral 0 over  $I$  and  $\|g\|_\infty \leq \varepsilon$ . Let  $G = \int_0^x g(t)dt$ . Then  $f = \frac{1}{2}(f + G) + \frac{1}{2}(f - G)$  and hence  $f$  is not extremal. Thus  $f_\pm$  are the only extremal points and their (closed) convex hull is given by  $f_\lambda(x) = \lambda x$  for  $\lambda \in [-1, 1]$ .

Of course the problem is that  $K$  is not closed. Hence we consider the Lipschitz continuous functions  $\bar{K} := \{f \in C^{0,1}([0, 1], \mathbb{R}) | f(0) = 0, [f]_1 \leq 1\}$  (this is in fact the closure of  $K$ , but this is a bit tricky to see and we won't need this here). By the Arzelà–Ascoli theorem (Theorem 1.13)  $\bar{K}$  is relatively compact and since the Lipschitz estimate clearly is preserved under uniform limits it is even compact.

Now note that piecewise linear functions with  $f'(x) \in \{\pm 1\}$  away from the kinks are extremal in  $\bar{K}$ . Moreover, these functions are dense: Split  $[0, 1]$  into  $n$  pieces of equal length using  $x_j = \frac{j}{n}$ . Set  $f_n(x_0) = 0$  and  $f_n(x) = f_n(x_j) \pm (x - x_j)$  for  $x \in [x_j, x_{j+1}]$  where the sign is chosen such that  $|f(x_{j+1}) - f_n(x_{j+1})|$  gets minimal. Then  $\|f - f_n\|_\infty \leq \frac{1}{n}$ .  $\diamond$

**Problem\* 6.7.** Show that the convex hull is given by (6.8).

**Problem\* 6.8.** Show that the nonempty intersection of extremal sets is extremal. Show that if  $L \subseteq M$  is extremal and  $M \subseteq K$  is extremal, then  $L \subseteq K$  is extremal as well.

**Problem 6.9.** Show that the closed unit ball in  $L^1(0, 1)$  has no extremal points.

**Problem 6.10.** Find the extremal points of the closed unit ball in  $C([0, 1], \mathbb{R})$ .

**Problem 6.11.** Let  $X$  be a topological vector space. Show that the closure and the interior of a convex set is convex. (Hint: One way of showing the first claim is to consider the continuous map  $f : X \times X \rightarrow X$  given by  $(x, y) \mapsto \lambda x + (1 - \lambda)y$  and use Problem B.15.)

**Problem 6.12.** Let  $X$  be a separable Banach space. Suppose  $K \subseteq X$  is convex and  $x$  is a boundary point of  $K$ . Then there is a supporting hyperplane at  $x$ . That is, there is some  $\ell \in X^*$  such that  $\ell(x) = 0$  and  $K$  is contained in the closed half-plane  $\{y | \operatorname{Re}(\ell(y - x)) \leq 0\}$ . (Hint: Lemma 4.32.)

**Problem 6.13** (Carathéodory). Show that for a compact convex set  $K \subseteq \mathbb{R}^n$  every point can be written as convex combination of  $n + 1$  extremal points. (Hint: Induction on  $n$ . Without loss assume that  $0$  is an extremal point. If  $K$  is contained in an  $n - 1$  dimensional subspace we are done. Otherwise  $K$  has an open interior. Now for a given point the line through this point and  $0$  intersects the boundary where we have a corresponding face.)

### 6.3. Weak topologies

In Section 4.4 we have defined weak convergence for sequences and this raises the question about a natural topology associated with this convergence. To this end we define the **weak topology** on a Banach space  $X$  as the weakest topology for which all  $\ell \in X^*$  remain continuous. Recall that a base for this topology is given by sets of the form

$$x + \bigcap_{j=1}^n |\ell_j|^{-1}([0, \varepsilon_j)) = \{\tilde{x} \in X | |\ell_j(x - \tilde{x})| < \varepsilon_j, 1 \leq j \leq n\},$$

$$x \in X, \ell_j \in X^*, \varepsilon_j > 0. \quad (6.10)$$

In particular, it is straightforward to check that a sequence converges with respect to this topology if and only if it converges weakly. Since the linear functionals separate points (cf. Corollary 4.12) the weak topology is Hausdorff. Moreover, the sets in the above base are clearly convex and hence we even have a locally convex space.

Note that, if  $X^*$  is separable, given a total set  $\{\ell_n\}_{n \in \mathbb{N}} \subset X^*$  of (w.l.o.g.) normalized linear functionals

$$d(x, \tilde{x}) := \sum_{n=1}^{\infty} \frac{1}{2^n} |\ell_n(x - \tilde{x})| \quad (6.11)$$

defines a metric on the unit ball  $B_1(0) \subset X$  which can be shown to generate the weak topology (Problem 6.17). However, on all of  $X$  the weak topology is not first countable unless  $X$  is finite dimensional (Problem 6.18).

Similarly, we define the **weak-\* topology** on  $X^*$  as the weakest topology for which all  $j \in J(X) \subseteq X^{**}$  remain continuous. In particular, the weak-\* topology is weaker than the weak topology on  $X^*$  and both are equal if  $X$  is reflexive. Like the weak topology it is Hausdorff (since different linear functionals must differ at least at one point) and not first countable unless  $X$  is finite dimensional (Problem 6.18). It also turns out that the continuous linear functionals with respect to the weak-\* topology are precisely  $J(X)$  (Problem 6.19).

A base for the weak-\* topology is given by sets of the form

$$\ell + \bigcap_{j=1}^n |J(x_j)|^{-1}([0, \varepsilon_j]) = \{\tilde{\ell} \in X^* \mid |(\ell - \tilde{\ell})(x_j)| < \varepsilon_j, 1 \leq j \leq n\},$$

$$\ell \in X^*, x_j \in X, \varepsilon_j > 0. \quad (6.12)$$

Again these sets are convex and we have a locally convex space. Note that, if  $X$  is separable, given a total set  $\{x_n\}_{n \in \mathbb{N}} \subset X$  of (w.l.o.g.) normalized vectors

$$d(\ell, \tilde{\ell}) := \sum_{n=1}^{\infty} \frac{1}{2^n} |(\ell - \tilde{\ell})(x_n)| \quad (6.13)$$

defines a metric on the unit ball  $B_1^*(0) \subset X^*$  which can be shown to generate the weak-\* topology (Problem 6.17). Hence Lemma 4.32 could also be stated as  $\bar{B}_1^*(0) \subset X^*$  being weak-\* compact. This is in fact true without assuming  $X$  to be separable and is known as Banach–Alaoglu theorem.

**Theorem 6.10** (Banach–Alaoglu). *Let  $X$  be a Banach space. Then  $\bar{B}_1^*(0) \subset X^*$  is compact in the weak-\* topology.*

**Proof.** Abbreviate  $B := \bar{B}_1^X(0)$ ,  $B^* := \bar{B}_1^{X^*}(0)$ , and  $B_x := \bar{B}_{\|x\|}^{\mathbb{C}}(0)$ . Consider the (injective) map  $\Phi : X^* \rightarrow \mathbb{C}^X$  given by  $\Phi(\ell)(x) := \ell(x)$  and identify  $X^*$  with  $\Phi(X^*)$ . Then the weak-\* topology on  $X^*$  coincides with the relative topology on  $\Phi(X^*) \subseteq \mathbb{C}^X$  (recall that the product topology on  $\mathbb{C}^X$  is the weakest topology which makes all point evaluations continuous). Moreover,  $|\Phi(\ell)| \leq \|\ell\| \|x\|$  implies  $\Phi(B^*) \subset \times_{x \in X} B_x$ , where the last product is compact by Tychonoff's theorem (Theorem B.18). Hence it suffices to show that

$\Phi(B^*)$  is closed. To this end let  $l \in \overline{\Phi(B^*)}$ . We need to show that  $l \in \Phi(B^*)$ . Fix  $x_1, x_2 \in X$ ,  $\alpha \in \mathbb{C}$ , and consider the open neighborhood

$$U(l) = \left\{ h \in \bigcap_{x \in X} B_x \mid \begin{array}{l} |h(x_1 + \alpha x_2) - l(x_1 + \alpha x_2)| < \varepsilon, \\ |h(x_1) - l(x_1)| < \varepsilon, \quad |\alpha| |h(x_2) - l(x_2)| < \varepsilon \end{array} \right\}$$

of  $l$ . Since  $U(l) \cap \Phi(B^*)$  is nonempty we can choose an element  $h$  from this intersection to show  $|l(x_1 + \alpha x_2) - l(x_1) - \alpha l(x_2)| < 3\varepsilon$ . Since  $\varepsilon > 0$  is arbitrary we conclude  $l(x_1 + \alpha x_2) = l(x_1) + \alpha l(x_2)$ . Moreover,  $|l(x_1)| \leq |h(x_1)| + \varepsilon \leq \|x_1\| + \varepsilon$  shows  $\|l\| \leq 1$  and thus  $l \in \Phi(B^*)$ .  $\square$

Please note that Example 4.39 shows that  $\bar{B}_1^*(0) \subset X^*$  might not be sequentially compact in the weak-\* topology unless  $X$  is separable.

If  $X$  is a reflexive space and we apply this to  $X^*$ , we get that the closed unit ball is compact in the weak topology. In fact, the converse is also true.

**Theorem 6.11** (Kakutani). *A Banach space  $X$  is reflexive if and only if the closed unit ball  $\bar{B}_1(0)$  is weakly compact.*

**Proof.** Suppose  $X$  is not reflexive and choose  $x'' \in \bar{B}_1^{**}(0) \setminus J(\bar{B}_1(0))$  with  $\|x''\| = 1$ . Then, if  $\bar{B}_1(0)$  is weakly compact,  $J(\bar{B}_1(0))$  is weak-\* compact (note that  $J$  is a homeomorphism if we equip  $X$  with the weak and  $X^{**}$  with the weak-\* topology). Moreover, if we equip  $X^{**}$  with the weak-\* topology, then its dual is isomorphic to  $X^*$  (Problem 6.19) and by Corollary 6.4 we can find some  $\ell \in X^*$  with  $\|\ell\| = 1$  and

$$\operatorname{Re}(x''(\ell)) < \inf_{y'' \in J(\bar{B}_1(0))} \operatorname{Re}(y''(\ell)) = \inf_{y \in \bar{B}_1(0)} \operatorname{Re}(\ell(y)) = -1.$$

But this contradicts  $|x''(\ell)| \leq 1$ .  $\square$

Note that in this context Theorem 4.28 says that in a reflexive Banach space  $X$  the closed unit ball  $\bar{B}_1(0)$  is weakly sequentially compact. The converse is also true and known as the Eberlein theorem.

Since the weak topology is weaker than the norm topology, every weakly closed set is also (norm) closed. Moreover, the weak closure of a set will in general be larger than the norm closure. However, for convex sets both will coincide. In fact, we have the following characterization in terms of closed (affine) **half-spaces**, that is, sets of the form  $\{x \in X \mid \operatorname{Re}(\ell(x)) \leq \alpha\}$  for some  $\ell \in X^*$  and some  $\alpha \in \mathbb{R}$ .

**Theorem 6.12** (Mazur). *The weak as well as the norm closure of a convex set  $K$  is the intersection of all half-spaces containing  $K$ . In particular, for a convex set  $K \subseteq X$  the following are equivalent:*

- $K$  is weakly closed,
- $K$  is weakly sequentially closed,

- $K$  is (norm) closed.

**Proof.** Since the intersection of closed-half spaces is (weakly) closed, it suffices to show that for every  $x$  not in the (weak) closure there is a closed half-plane not containing  $x$ . Moreover, if  $x$  is not in the weak closure it is also not in the norm closure (the norm closure is contained in the weak closure) and by Theorem 6.3 with  $U := B_{\text{dist}(x,K)}(x)$  and  $V := K$  there is a functional  $\ell \in X^*$  such that  $K \subseteq \text{Re}(\ell)^{-1}([c, \infty))$  and  $x \notin \text{Re}(\ell)^{-1}([c, \infty))$ .

For the last claim note that a weakly closed set is weakly sequentially closed (Example B.13). Moreover, a weakly sequentially closed set is also sequentially closed and hence closed. Finally, a closed convex set is weakly closed by the first part.  $\square$

**Example 6.10.** Suppose  $X$  is infinite dimensional. The weak closure  $\bar{S}^w$  of  $S := \{x \in X \mid \|x\| = 1\}$  is the closed unit ball  $\bar{B}_1(0)$ . Indeed, since  $\bar{B}_1(0)$  is convex the previous theorem shows  $\bar{S}^w \subseteq \bar{B}_1(0)$ . Conversely, if  $x \in \bar{B}_1(0)$  is not in the weak closure, then there must be an open neighborhood  $x + \bigcap_{j=1}^n |\ell_j|^{-1}([0, \varepsilon))$  not contained in the weak closure. Since  $X$  is infinite dimensional we can find a nonzero element  $x_0 \in \bigcap_{j=1}^n \text{Ker}(\ell_j)$  (by Problem 4.22 the  $\ell_j$  would otherwise be a basis) such that the affine line  $x + tx_0$  is in this neighborhood and hence also avoids  $\bar{S}^w$ . But this is impossible since by the intermediate value theorem there is some  $t_0 > 0$  such that  $\|x + t_0 x_0\| = 1$ . Hence  $\bar{B}_1(0) \subseteq \bar{S}^w$ .  $\diamond$

Note that this example also shows that in an infinite dimensional space the weak and norm topologies are always different! In a finite dimensional space both topologies of course agree.

**Corollary 6.13** (Mazur lemma). *Suppose  $x_k \rightharpoonup x$ , then there are convex combinations  $y_k := \sum_{j=1}^{n_k} \lambda_{k,j} x_j$  (with  $\sum_{j=1}^{n_k} \lambda_{k,j} = 1$  and  $\lambda_{k,j} \geq 0$ ) such that  $y_k \rightarrow x$ .*

**Proof.** Let  $K := \{\sum_{j=1}^n \lambda_j x_j \mid n \in \mathbb{N}, \sum_{j=1}^n \lambda_j = 1, \lambda_j \geq 0\}$  be the convex hull of the points  $\{x_k\}$ . Then by the previous result  $x \in \bar{K}$ .  $\square$

**Example 6.11.** Let  $\mathfrak{H}$  be a Hilbert space and  $\{\varphi_j\}$  some infinite ONS. Then we already know  $\varphi_j \rightharpoonup 0$ . Moreover, the convex combination  $\psi_j := \frac{1}{j} \sum_{k=1}^j \varphi_k \rightarrow 0$  since  $\|\psi_j\| = j^{-1/2}$ .  $\diamond$

For the last result note that since  $X^{**}$  is the dual of  $X^*$  it has a corresponding weak-\* topology and by the Banach–Alaoglu theorem  $\bar{B}_1^{**}(0)$  is weak-\* compact and hence weak-\* closed.

**Theorem 6.14** (Goldstine). *The image of the closed unit ball  $\bar{B}_1(0)$  under the canonical embedding  $J$  into the closed unit ball  $\bar{B}_1^{**}(0)$  is weak-\* dense.*



**Proof.** Let  $j \in \bar{B}_1^{**}(0)$  be given. Since sets of the form  $\{\tilde{j} \in X^{**} \mid |j(\ell_k) - \tilde{j}(\ell_k)| < \varepsilon, 1 \leq k \leq n\}$  provide a neighborhood base (where we can assume the  $\ell_k \in X^*$  to be linearly independent without loss of generality), it suffices to find some  $x \in \bar{B}_{1+\varepsilon}(0)$  with  $\ell_k(x) = j(\ell_k)$  for  $1 \leq k \leq n$  since then  $(1 + \varepsilon)^{-1}J(x)$  will be in the above neighborhood. Without the requirement  $\|x\| \leq 1 + \varepsilon$  this follows from surjectivity of the map  $F : X \rightarrow \mathbb{C}^n$ ,  $x \mapsto (\ell_1(x), \dots, \ell_n(x))$ . Moreover, given one such  $x$  the same is true for every element from  $x + Y$ , where  $Y := \bigcap_k \text{Ker}(\ell_k)$ . But if  $(x + Y) \cap \bar{B}_{1+\varepsilon}(0)$  were empty, we would have  $\text{dist}(x, Y) \geq 1 + \varepsilon$  and by Corollary 4.13 we could find some normalized  $\ell \in X^*$  which vanishes on  $Y$  and satisfies  $\ell(x) \geq 1 + \varepsilon$ . By Problem 4.22 we have  $\ell \in \text{span}(\ell_1, \dots, \ell_n)$  implying

$$1 + \varepsilon \leq \ell(x) = j(\ell) \leq \|j\| \|\ell\| \leq 1$$

a contradiction.  $\square$

Note that if  $\bar{B}_1(0) \subset X$  is weakly compact, then  $J(\bar{B}_1(0))$  is compact (and thus closed) in the weak-\* topology on  $X^{**}$ . Hence Goldstine's theorem implies  $J(\bar{B}_1(0)) = \bar{B}_1^{**}(0)$  and we get an alternative proof of Kakutani's theorem.

**Example 6.12.** Consider  $X := c_0(\mathbb{N})$ ,  $X^* \cong \ell^1(\mathbb{N})$ , and  $X^{**} \cong \ell^\infty(\mathbb{N})$  with  $J$  corresponding to the inclusion  $c_0(\mathbb{N}) \hookrightarrow \ell^\infty(\mathbb{N})$ . Then we can consider the linear functionals  $\ell_j(x) = x_j$  which are total in  $X^*$  and a sequence in  $X^{**}$  will be weak-\* convergent if and only if it is bounded and converges when composed with any of the  $\ell_j$  (in other words, when the sequence converges componentwise — cf. Problem 4.34). So for example, cutting off a sequence in  $\bar{B}_1^{**}(0)$  after  $n$  terms (setting the remaining terms equal to 0) we get a sequence from  $\bar{B}_1(0) \hookrightarrow \bar{B}_1^{**}(0)$  which is weak-\* convergent (but of course not norm convergent).

Also observe that  $c_0(\mathbb{N}) \subseteq \ell^\infty(\mathbb{N})$  is closed but not weak-\* closed and hence Mazur's theorem does not hold if we replace the weak by the weak-\* topology.  $\diamond$

**Problem 6.14.** Show that in an infinite dimensional space, a weakly open neighborhood of 0 contains a nontrivial subspace. Show the analogue statement for weak-\* open neighborhoods of 0.

**Problem 6.15.** Show that a weakly sequentially compact set is bounded. Similarly, show that a weakly compact set is bounded.

**Problem 6.16** (von Neumann). Consider  $X := \ell^2(\mathbb{N})$  with the canonical basis  $\delta^n$ . Let  $A := \{\delta^n + n\delta^m\}_{m,n \in \mathbb{N}}$ . Show that the weak sequential closure of  $A$  contains  $\{\delta^n\}_{n \in \mathbb{N}}$  but not 0. Hence the weak sequential closure of  $A$  is not weakly sequentially closed.

**Problem\* 6.17.** Show that (6.11) generates the weak topology on  $B_1(0) \subset X$ . Show that (6.13) generates the weak-\* topology on  $B_1^*(0) \subset X^*$ .

**Problem 6.18.** Show that neither the weak nor the weak-\* topology is first countable if  $X$  is infinite dimensional. (Hint: If there is a countable neighborhood base, you can find, using Problem 6.14, an unbounded sequence of vectors which converge weakly to zero.)

**Problem\* 6.19.** The dual of  $X^*$  with respect to the weak-\* topology is  $J(X)$ . (Hint: Use that a continuous linear functional is bounded and Problem 4.22.)

**Problem\* 6.20.** Show that the annihilator  $M^\perp$  of a set  $M \subseteq X$  is weak-\* closed. Moreover show that  $(N_\perp)^\perp = \overline{\text{span}(N)}^{\text{weak-*}}$ . (Hint: The first part and hence one inclusion of the second part are straightforward. For the other inclusion use Problem 6.19 and apply Corollary 6.5.)

## 6.4. Beyond Banach spaces: Locally convex spaces

We have already seen that it is often important to weaken the notion of convergence (i.e., to weaken the underlying topology) to get a larger class of converging sequences. It turns out that all cases considered so far fit within a general framework which we want to discuss in this section. We start with an alternate definition of a locally convex vector space which we already briefly encountered in Corollary 6.4 (equivalence of both definitions will be established below).

A vector space  $X$  together with a topology is called a **locally convex vector space** if there exists a family of seminorms  $\{q_\alpha\}_{\alpha \in A}$  which generates the topology in the sense that the topology is the weakest topology for which the family of functions  $\{q_\alpha(\cdot - x)\}_{\alpha \in A, x \in X}$  is continuous. Hence the topology is generated by sets of the form  $x + q_\alpha^{-1}(I)$ , where  $I \subseteq [0, \infty)$  is open (in the relative topology). Moreover, sets of the form

$$x + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j)) \quad (6.14)$$

are a neighborhood base at  $x$  and hence it is straightforward to check that a locally convex vector space is a topological vector space, that is, both vector addition and scalar multiplication are continuous. For example, if  $z = x + y$  then the preimage of the open neighborhood  $z + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j))$  contains the open neighborhood  $(x + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j/2)), y + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j/2))$  by virtue of the triangle inequality. Similarly, if  $z = \gamma x$  then the preimage of the open neighborhood  $z + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j))$  contains the open neighborhood  $(B_\varepsilon(\gamma), x + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \frac{\varepsilon_j}{2(|\gamma| + \varepsilon)})))$  with  $\varepsilon < \frac{\varepsilon_j}{2q_{\alpha_j}(x)}$ .

Moreover, note that a sequence  $x_n$  will converge to  $x$  in this topology if and only if  $q_\alpha(x_n - x) \rightarrow 0$  for all  $\alpha$ .

**Example 6.13.** Of course every Banach space equipped with the norm topology is a locally convex vector space if we choose the single seminorm  $q(x) = \|x\|$ .  $\diamond$

**Example 6.14.** A Banach space  $X$  equipped with the weak topology is a locally convex vector space. In this case we have used the continuous linear functionals  $\ell \in X^*$  to generate the topology. However, note that the corresponding seminorms  $q_\ell(x) := |\ell(x)|$  generate the same topology since  $x + q_\ell^{-1}([0, \varepsilon]) = \ell^{-1}(B_\varepsilon(x))$  in this case. The same is true for  $X^*$  equipped with the weak or the weak-\* topology.  $\diamond$

**Example 6.15.** The bounded linear operators  $\mathcal{L}(X, Y)$  together with the seminorms  $q_x(A) := \|Ax\|$  for all  $x \in X$  (strong convergence) or the seminorms  $q_{\ell, x}(A) := |\ell(Ax)|$  for all  $x \in X, \ell \in Y^*$  (weak convergence) are locally convex vector spaces.  $\diamond$

**Example 6.16.** The continuous functions  $C(I)$  together with the pointwise topology generated by the seminorms  $q_x(f) := |f(x)|$  for all  $x \in I$  is a locally convex vector space.  $\diamond$

In all these examples we have one additional property which is often required as part of the definition: The seminorms are called **separated** if for every  $x \in X \setminus \{0\}$  there is a seminorm with  $q_\alpha(x) \neq 0$ . In this case the corresponding locally convex space is Hausdorff, since for  $x \neq y$  the neighborhoods  $U(x) = x + q_\alpha^{-1}([0, \varepsilon])$  and  $U(y) = y + q_\alpha^{-1}([0, \varepsilon])$  will be disjoint for  $\varepsilon = \frac{1}{2}q_\alpha(x - y) > 0$  (the converse is also true; Problem 6.27).

It turns out crucial to understand when a seminorm is continuous.

**Lemma 6.15.** *Let  $X$  be a locally convex vector space with corresponding family of seminorms  $\{q_\alpha\}_{\alpha \in A}$ . Then a seminorm  $q$  is continuous if and only if there are seminorms  $q_{\alpha_j}$  and constants  $c_j > 0$ ,  $1 \leq j \leq n$ , such that  $q(x) \leq \sum_{j=1}^n c_j q_{\alpha_j}(x)$ .*

**Proof.** If  $q$  is continuous, then  $q^{-1}([0, 1])$  contains an open neighborhood of 0 of the form  $\bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j])$  and choosing  $c_j = \max_{1 \leq j \leq n} \varepsilon_j^{-1}$  we obtain that  $\sum_{j=1}^n c_j q_{\alpha_j}(x) < 1$  implies  $q(x) < 1$  and the claim follows from Problem 6.22. Conversely note that if  $q(x) = r$  then  $q^{-1}((r - \varepsilon, r + \varepsilon))$  contains the set  $U(x) := x + \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j])$  provided  $\sum_{j=1}^n c_j \varepsilon_j \leq \varepsilon$  since by the inverse triangle inequality  $|q(y) - q(x)| \leq q(y - x) \leq \sum_{j=1}^n c_j q_{\alpha_j}(y - x) < \varepsilon$  for  $y \in U(x)$ .  $\square$

**Example 6.17.** The weak topology on an infinite dimensional space cannot be generated by a norm. Indeed, let  $q$  be a continuous seminorm and  $q_{\alpha_j} = |\ell_{\alpha_j}|$  as in the lemma. Then  $\bigcap_{j=1}^n \text{Ker}(\ell_{\alpha_j})$  has codimension at most  $n$  and

hence contains some  $x \neq 0$  implying that  $q(x) \leq \sum_{j=1}^n c_j q_{\alpha_j}(x) = 0$ . Thus  $q$  is no norm. Similarly, the other examples cannot be generated by a norm except in finite dimensional cases.  $\diamond$

Moreover, note that the topology is translation invariant in the sense that  $U(x)$  is a neighborhood of  $x$  if and only if  $U(x) - x = \{y - x | y \in U(x)\}$  is a neighborhood of 0. Hence we can restrict our attention to neighborhoods of 0 (this is of course true for any topological vector space). Hence if  $X$  and  $Y$  are topological vector spaces, then a linear map  $A : X \rightarrow Y$  will be continuous if and only if it is continuous at 0. Moreover, if  $Y$  is a locally convex space with respect to some seminorms  $p_\beta$ , then  $A$  will be continuous if and only if  $p_\beta \circ A$  is continuous for every  $\beta$  (Lemma B.11). Finally, since  $p_\beta \circ A$  is a seminorm, the previous lemma implies:

**Corollary 6.16.** *Let  $(X, \{q_\alpha\})$  and  $(Y, \{p_\beta\})$  be locally convex vector spaces. Then a linear map  $A : X \rightarrow Y$  is continuous if and only if for every  $\beta$  there are some seminorms  $q_{\alpha_j}$  and constants  $c_j > 0$ ,  $1 \leq j \leq n$ , such that  $p_\beta(Ax) \leq \sum_{j=1}^n c_j q_{\alpha_j}(x)$ .*

It will shorten notation when sums of the type  $\sum_{j=1}^n c_j q_{\alpha_j}(x)$ , which appeared in the last two results, can be replaced by a single expression  $c q_\alpha$ . This can be done if the family of seminorms  $\{q_\alpha\}_{\alpha \in A}$  is **directed**, that is, for given  $\alpha, \beta \in A$  there is a  $\gamma \in A$  such that  $q_\alpha(x) + q_\beta(x) \leq C q_\gamma(x)$  for some  $C > 0$ . Moreover, if  $\mathcal{F}(A)$  is the set of all finite subsets of  $A$ , then  $\{\tilde{q}_F = \sum_{\alpha \in F} q_\alpha\}_{F \in \mathcal{F}(A)}$  is a directed family which generates the same topology (since every  $\tilde{q}_F$  is continuous with respect to the original family we do not get any new open sets).

While the family of seminorms is in most cases more convenient to work with, it is important to observe that different families can give rise to the same topology and it is only the topology which matters for us. In fact, it is possible to characterize locally convex vector spaces as topological vector spaces which have a neighborhood basis at 0 of absolutely convex sets. Here a set  $U$  is called **absolutely convex**, if for  $|\alpha| + |\beta| \leq 1$  we have  $\alpha U + \beta U \subseteq U$ .

**Example 6.18.** The absolutely convex sets in  $\mathbb{C}$  are precisely the (open and closed) balls.  $\diamond$

Since the sets  $q_\alpha^{-1}([0, \varepsilon))$  are absolutely convex we always have such a basis in our case. To see the converse note that such a neighborhood  $U$  of 0 is also absorbing (Problem 6.21) and hence the corresponding Minkowski functional (6.1) is a seminorm (Problem 6.26). By construction, these seminorms generate the topology since if  $U_0 = \bigcap_{j=1}^n q_{\alpha_j}^{-1}([0, \varepsilon_j)) \subseteq U$  we have for the corresponding Minkowski functionals  $p_U(x) \leq p_{U_0}(x) \leq \varepsilon^{-1} \sum_{j=1}^n q_{\alpha_j}(x)$ , where  $\varepsilon = \min \varepsilon_j$ . With a little more work (Problem 6.25), one can even

show that it suffices to assume to have a neighborhood basis at 0 of convex open sets.

Given a topological vector space  $X$  we can define its dual space  $X^*$  as the set of all continuous linear functionals. However, while it can happen in general that the dual space is trivial,  $X^*$  will always be nontrivial for a locally convex space since the Hahn–Banach theorem can be used to construct linear functionals (using a continuous seminorm for  $\varphi$  in Theorem 4.10) and also the geometric Hahn–Banach theorem (Theorem 6.3) holds; see also its corollaries. In this respect note that for every continuous linear functional  $\ell$  in a topological vector space,  $|\ell|^{-1}([0, \varepsilon))$  is an absolutely convex open neighborhoods of 0 and hence existence of such sets is necessary for the existence of nontrivial continuous functionals. As a natural topology on  $X^*$  we could use the weak-\* topology defined to be the weakest topology generated by the family of all point evaluations  $q_x(\ell) = |\ell(x)|$  for all  $x \in X$ . Since different linear functionals must differ at least at one point, the weak-\* topology is Hausdorff. Given a continuous linear operator  $A : X \rightarrow Y$  between locally convex spaces we can define its adjoint  $A' : Y^* \rightarrow X^*$  as before,

$$(A'y^*)(x) := y^*(Ax). \quad (6.15)$$

A brief calculation

$$q_x(A'y^*) = |(A'y^*)(x)| = |y^*(Ax)| = q_{Ax}(y^*) \quad (6.16)$$

verifies that  $A'$  is continuous in the weak-\* topology by virtue of Corollary 6.16.

The remaining theorems we have established for Banach spaces were consequences of the Baire theorem (which requires a complete metric space) and this leads us to the question when a locally convex space is a metric space. From our above analysis we see that a locally convex vector space will be first countable if and only if countably many seminorms suffice to determine the topology. In this case  $X$  turns out to be metrizable.

**Theorem 6.17.** *A locally convex Hausdorff space is metrizable if and only if it is first countable. In this case there is a countable family of separated seminorms  $\{q_n\}_{n \in \mathbb{N}}$  generating the topology and a metric is given by*

$$d(x, y) := \max_{n \in \mathbb{N}} \frac{1}{2^n} \frac{q_n(x - y)}{1 + q_n(x - y)}. \quad (6.17)$$

**Proof.** If  $X$  is first countable there is a countable neighborhood base at 0 and hence also a countable neighborhood base of absolutely convex sets. The Minkowski functionals corresponding to the latter base are seminorms of the required type.

Now in this case it is straightforward to check that (6.17) defines a metric (see also Problem B.3). Moreover, the metric balls  $B_r(x) := \bigcap_{n: 2^{-n} > r} \{y | q_n(y - x) < 1\}$

$x) < \frac{r}{2^{-n}-r}\}$  are clearly open and convex (note that the intersection is finite). Conversely, for every set of the form (6.14) we can choose  $\varepsilon := \min\{2^{-\alpha_j} \frac{\varepsilon_j}{1+\varepsilon_j} | 1 \leq j \leq n\}$  such that  $B_\varepsilon(x)$  will be contained in this set. Hence both topologies are equivalent (cf. Lemma B.2).  $\square$

In general, a locally convex vector space  $X$  which has a separated countable family of seminorms is called a **Fréchet space** if it is complete with respect to the metric (6.17). Note that the metric (6.17) is **translation invariant**

$$d(f, g) = d(f - h, g - h). \quad (6.18)$$

**Example 6.19.** The continuous functions  $C(\mathbb{R})$  together with local uniform convergence are a Fréchet space. A countable family of seminorms is for example

$$\|f\|_j := \sup_{|x| \leq j} |f(x)|, \quad j \in \mathbb{N}. \quad (6.19)$$

Then  $f_k \rightarrow f$  if and only if  $\|f_k - f\|_j \rightarrow 0$  for all  $j \in \mathbb{N}$  and it follows that  $C(\mathbb{R})$  is complete.  $\diamond$

**Example 6.20.** The space  $C^\infty(\mathbb{R}^m)$  together with the seminorms

$$\|f\|_{j,k} := \sum_{|\alpha| \leq k} \sup_{|x| \leq j} |\partial_\alpha f(x)|, \quad j \in \mathbb{N}, k \in \mathbb{N}_0, \quad (6.20)$$

is a Fréchet space.

Note that  $\partial_\alpha : C^\infty(\mathbb{R}^m) \rightarrow C^\infty(\mathbb{R}^m)$  is continuous. Indeed by Corollary 6.16 it suffices to observe that  $\|\partial_\alpha f\|_{j,k} \leq \|f\|_{j,k+|\alpha|}$ .  $\diamond$

**Example 6.21.** The **Schwartz space**

$$\mathcal{S}(\mathbb{R}^m) := \{f \in C^\infty(\mathbb{R}^m) | \sup_x |x^\alpha (\partial_\beta f)(x)| < \infty, \forall \alpha, \beta \in \mathbb{N}_0^m\} \quad (6.21)$$

together with the seminorms

$$q_{\alpha,\beta}(f) := \|x^\alpha (\partial_\beta f)(x)\|_\infty, \quad \alpha, \beta \in \mathbb{N}_0^m, \quad (6.22)$$

is a Fréchet space. To see completeness note that a Cauchy sequence  $f_n$  is in particular a Cauchy sequence in  $C^\infty(\mathbb{R}^m)$ . Hence there is a limit  $f \in C^\infty(\mathbb{R}^m)$  such that all derivatives converge uniformly. Moreover, since Cauchy sequences are bounded  $\|x^\alpha (\partial_\beta f_n)(x)\|_\infty \leq C_{\alpha,\beta}$  we obtain  $f \in \mathcal{S}(\mathbb{R}^m)$ .

Again  $\partial_\gamma : \mathcal{S}(\mathbb{R}^m) \rightarrow \mathcal{S}(\mathbb{R}^m)$  is continuous since  $q_{\alpha,\beta}(\partial_\gamma f) \leq q_{\alpha,\beta+\gamma}(f)$  and so is  $x^\gamma : \mathcal{S}(\mathbb{R}^m) \rightarrow \mathcal{S}(\mathbb{R}^m)$  since  $q_{\alpha,\beta}(x^\gamma f) \leq \sum_{\eta \leq \beta} \binom{\beta}{\eta} \frac{\gamma!}{(\gamma-\eta)!} q_{\alpha+\gamma-\eta,\beta-\eta}(f)$ .

The dual space  $\mathcal{S}^*(\mathbb{R}^m)$  is known as the space of **tempered distributions**.  $\diamond$

**Example 6.22.** The space of all entire functions  $f$  (i.e. functions which are holomorphic on all of  $\mathbb{C}$ ) together with the seminorms  $\|f\|_j := \sup_{|z| \leq j} |f(z)|$ ,  $j \in \mathbb{N}$ , is a Fréchet space. Completeness follows from the Weierstraß convergence theorem which states that a limit of holomorphic functions which is uniform on every compact subset is again holomorphic.  $\diamond$

**Example 6.23.** In all of the previous examples the topology cannot be generated by a norm. For example, if  $q$  is a norm for  $C(\mathbb{R})$ , then by Lemma 6.15 there is some index  $j$  such that  $q(f) \leq C\|f\|_j$ . Now choose a nonzero function which vanishes on  $[-j, j]$  to get a contradiction.  $\diamond$

There is another useful criterion when the topology can be described by a single norm. To this end we call a set  $B \subseteq X$  bounded if  $\sup_{x \in B} q_\alpha(x) < \infty$  for every  $\alpha$ . By Corollary 6.16 this will then be true for any continuous seminorm on  $X$ .

**Theorem 6.18** (Kolmogorov). *A locally convex vector space can be generated from a single seminorm if and only if it contains a bounded open set.*

**Proof.** In a normed space every open ball is bounded and hence only the converse direction is nontrivial. So let  $U$  be a bounded open set. By shifting and decreasing  $U$  if necessary we can assume  $U$  to be an absolutely convex open neighborhood of 0 and consider the associated Minkowski functional  $q = p_U$ . Then since  $U = \{x | q(x) < 1\}$  and  $\sup_{x \in U} q_\alpha(x) = C_\alpha < \infty$  we infer  $q_\alpha(x) \leq C_\alpha q(x)$  (Problem 6.22) and thus the single seminorm  $q$  generates the topology.  $\square$

Finally, we mention that, since the Baire category theorem holds for arbitrary complete metric spaces, the open mapping theorem (Theorem 4.5), the inverse mapping theorem (Theorem 4.6) and the closed graph theorem (Theorem 4.7) hold for Fréchet spaces without modifications. In fact, they are formulated such that it suffices to replace Banach by Fréchet in these theorems as well as their proofs (concerning the proof of Theorem 4.5 take into account Problems 6.21 and 6.28).

**Problem\* 6.21.** *In a topological vector space every neighborhood  $U$  of 0 is absorbing.*

**Problem\* 6.22.** *Let  $p, q$  be two seminorms. Then  $p(x) \leq Cq(x)$  if and only if  $q(x) < 1$  implies  $p(x) < C$ .*

**Problem 6.23.** *Let  $X$  be a vector space. We call a set  $U$  **balanced** if  $\alpha U \subseteq U$  for every  $|\alpha| \leq 1$ . Show that a set is balanced and convex if and only if it is absolutely convex.*

**Problem\* 6.24.** *The intersection of arbitrary absolutely convex/balanced sets is again absolutely convex/balanced convex. Hence we can define the*

absolutely convex/balanced hull of a set  $U$  as the smallest absolutely convex/balanced set containing  $U$ , that is, the intersection of all absolutely convex/balanced sets containing  $U$ . Show that the absolutely convex hull is given by

$$\text{ahull}(U) := \left\{ \sum_{j=1}^n \lambda_j x_j \mid n \in \mathbb{N}, x_j \in U, \sum_{j=1}^n |\lambda_j| \leq 1 \right\}$$

and the balanced hull by

$$\text{bhull}(U) := \{ \alpha x \mid x \in U, |\alpha| \leq 1 \}.$$

Show that  $\text{ahull}(U) = \text{conv}(\text{bhull}(U))$ .

**Problem\* 6.25.** In a topological vector space every convex open neighborhood  $U$  of zero contains a balanced open neighborhood. Moreover, every convex open neighborhood  $U$  of zero contains an absolutely convex open neighborhood of zero. (Hint: By continuity of the scalar multiplication  $U$  contains a set of the form  $B_\varepsilon^\mathbb{C}(0) \cdot V$ , where  $V$  is an open neighborhood of zero.)

**Problem\* 6.26.** Let  $X$  be a vector space. Show that the Minkowski functional of a balanced, convex, absorbing set is a seminorm.

**Problem\* 6.27.** If a locally convex space is Hausdorff then any corresponding family of seminorms is separated.

**Problem\* 6.28.** Suppose  $X$  is a complete vector space with a translation invariant metric  $d$ . Show that  $\sum_{j=1}^\infty d(0, x_j) < \infty$  implies that

$$\sum_{j=1}^\infty x_j = \lim_{n \rightarrow \infty} \sum_{j=1}^n x_j$$

exists and

$$d(0, \sum_{j=1}^\infty x_j) \leq \sum_{j=1}^\infty d(0, x_j)$$

in this case (compare also Problem 1.5).

**Problem 6.29.** Let  $X$  be a locally convex Hausdorff space. Then for a non-zero linear functional  $f$  the following are equivalent:

- (i)  $f$  is continuous.
- (ii)  $\text{Ker}(f)$  is closed.
- (iii)  $\text{Ker}(f)$  is not dense in  $X$ .
- (iv) There exists a neighborhood  $U$  of zero such that  $f(U)$  is bounded.

**Problem 6.30.** Instead of (6.17) one frequently uses

$$\tilde{d}(x, y) := \sum_{n \in \mathbb{N}} \frac{1}{2^n} \frac{q_n(x - y)}{1 + q_n(x - y)}.$$



Show that this metric generates the same topology.

Consider the Fréchet space  $C(\mathbb{R})$  with  $q_n(f) = \sup_{[-n,n]} |f|$ . Show that the metric balls with respect to  $\tilde{d}$  are not convex.

**Problem 6.31.** Suppose  $X$  is a metric vector space. Then balls are convex if and only if the metric is quasiconvex:

$$d(\lambda x + (1 - \lambda)y, z) \leq \max\{d(x, z), d(y, z)\}, \quad \lambda \in (0, 1).$$

(See also Problem 9.15.)

**Problem 6.32.** Consider  $\ell^p(\mathbb{N})$  for  $p \in (0, 1)$  — compare Problem 1.15. Show that  $\|\cdot\|_p$  is not convex. Show that every convex open set is unbounded. Conclude that it is not a locally convex vector space. (Hint: Consider  $B_R(0)$ . Then for  $r < R$  all vectors which have one entry equal to  $r$  and all other entries zero are in this ball. By taking convex combinations all vectors which have  $n$  entries equal to  $r/n$  are in the convex hull. The quasinorm of such a vector is  $n^{1/p-1}r$ .)

**Problem 6.33.** Show that  $C_c^\infty(\mathbb{R}^m)$  is dense in  $\mathcal{S}(\mathbb{R}^m)$ .

**Problem 6.34.** Let  $X$  be a topological vector space and  $M$  a closed subspace. Show that the quotient space  $X/M$  is again a topological vector space and that  $\pi : X \rightarrow X/M$  is linear, continuous, and open. Show that points in  $X/M$  are closed.

## 6.5. Uniformly convex spaces

In a Banach space  $X$ , the unit ball is convex by the triangle inequality. Moreover,  $X$  is called **strictly convex** if the unit ball is a strictly convex set, that is, if for any two points on the unit sphere their average is inside the unit ball. See Problem 1.13 for some equivalent definitions. This is illustrated in Figure 1.1 which shows that in  $\mathbb{R}^2$  this is only true for  $1 < p < \infty$ .

**Example 6.24.** By Problem 1.13 it follows that  $\ell^p(\mathbb{N})$  is strictly convex for  $1 < p < \infty$  but not for  $p = 1, \infty$ .  $\diamond$

A more qualitative notion is to require that if two unit vectors  $x, y$  satisfy  $\|x - y\| \geq \varepsilon$  for some  $\varepsilon > 0$ , then there is some  $\delta > 0$  such that  $\|\frac{x+y}{2}\| \leq 1 - \delta$ . In this case one calls  $X$  **uniformly convex** and

$$\delta(\varepsilon) := \inf \left\{ 1 - \left\| \frac{x+y}{2} \right\| \mid \|x\| = \|y\| = 1, \|x - y\| \geq \varepsilon \right\}, \quad 0 \leq \varepsilon \leq 2, \quad (6.23)$$

is called the modulus of convexity. Of course every uniformly convex space is strictly convex. In finite dimensions the converse is also true (Problem 6.38).

Note that  $\delta$  is nondecreasing and

$$1 - \left\| \frac{x+y}{2} \right\| = 1 - \left\| x - \frac{x-y}{2} \right\| \leq \frac{\|x - y\|}{2}$$

for  $\|x\| = \|y\| = 1$  shows  $0 \leq \delta(\varepsilon) \leq \frac{\varepsilon}{2}$ . Moreover,  $\delta(2) = 1$  implies  $X$  strictly convex. In fact in this case  $1 = \delta(2) \leq 1 - \|\frac{x+y}{2}\| \leq 1$  for  $2 \leq \|x - y\| \leq 2$ . That is,  $x = -y$  whenever  $\|x - y\| = 2 = \|x\| + \|y\|$ .

**Example 6.25.** Every Hilbert space is uniformly convex with modulus of convexity  $\delta(\varepsilon) = 1 - \sqrt{1 - \frac{\varepsilon^2}{4}}$  (Problem 6.36).  $\diamond$

**Example 6.26.** Consider  $C[0, 1]$  with the norm

$$\|x\| := \|x\|_\infty + \|x\|_2 = \max_{t \in [0, 1]} |x(t)| + \left( \int_0^1 |x(t)|^2 dt \right)^{1/2}.$$

Note that by  $\|x\|_2 \leq \|x\|_\infty$  this norm is equivalent to the usual one:  $\|x\|_\infty \leq \|x\| \leq 2\|x\|_\infty$ . While with the usual norm  $\|\cdot\|_\infty$  this space is not strictly convex, it is with the new one. To see this we use (i) from Problem 1.13. Then if  $\|x + y\| = \|x\| + \|y\|$  we must have both  $\|x + y\|_\infty = \|x\|_\infty + \|y\|_\infty$  and  $\|x + y\|_2 = \|x\|_2 + \|y\|_2$ . Hence strict convexity of  $\|\cdot\|_2$  implies strict convexity of  $\|\cdot\|$ .

Note however, that  $\|\cdot\|$  is not uniformly convex. In fact, since by the Milman–Pettis theorem below, every uniformly convex space is reflexive, there cannot be an equivalent norm on  $C[0, 1]$  which is uniformly convex (cf. Example 4.18).  $\diamond$

**Example 6.27.** It can be shown that  $\ell^p(\mathbb{N})$  is uniformly convex for  $1 < p < \infty$  (see Theorem 3.11 from [48]).  $\diamond$

Equivalently, uniform convexity implies that if the average of two unit vectors is close to the boundary, then they must be close to each other. Specifically, if  $\|x\| = \|y\| = 1$  and  $\|\frac{x+y}{2}\| > 1 - \delta(\varepsilon)$  then  $\|x - y\| < \varepsilon$ . The following result (which generalizes Lemma 4.26) uses this observation:

**Theorem 6.19** (Radon–Riesz theorem). *Let  $X$  be a uniformly convex Banach space and let  $x_n \rightharpoonup x$ . Then  $x_n \rightarrow x$  if and only if  $\limsup \|x_n\| \leq \|x\|$ .*

**Proof.** By Lemma 4.25 (ii) we have in fact  $\lim \|x_n\| = \|x\|$ . If  $x = 0$  there is nothing to prove. Hence we can assume  $x_n \neq 0$  for all  $n$  and consider  $y_n := \frac{x_n}{\|x_n\|}$ . Then  $y_n \rightharpoonup y := \frac{x}{\|x\|}$  and it suffices to show  $y_n \rightarrow y$ . Next choose a linear functional  $\ell \in X^*$  with  $\|\ell\| = 1$  and  $\ell(y) = 1$ . Then

$$\ell\left(\frac{y_n + y}{2}\right) \leq \left\| \frac{y_n + y}{2} \right\| \leq 1$$

and letting  $n \rightarrow \infty$  shows  $\|\frac{y_n + y}{2}\| \rightarrow 1$ . Finally uniform convexity shows  $y_n \rightarrow y$ .  $\square$

For the proof of the next result we need the following equivalent condition.

**Lemma 6.20.** *Let  $X$  be a Banach space. Then*

$$\delta(\varepsilon) = \inf \left\{ 1 - \left\| \frac{x+y}{2} \right\| \mid \|x\| \leq 1, \|y\| \leq 1, \|x - y\| \geq \varepsilon \right\} \quad (6.24)$$

for  $0 \leq \varepsilon \leq 2$ .

**Proof.** It suffices to show that for given  $x$  and  $y$  which are not both on the unit sphere there is an equivalent pair on the unit sphere within the real subspace spanned by these vectors. By scaling we could get a better pair if both were strictly inside the unit ball and hence we can assume at least one vector to have norm one, say  $\|x\| = 1$ . Moreover, consider

$$u(t) := \frac{\cos(t)x + \sin(t)y}{\|\cos(t)x + \sin(t)y\|}, \quad v(t) := u(t) + (y - x).$$

Then  $\|v(0)\| = \|y\| < 1$ . Moreover, let  $t_0 \in (\frac{\pi}{2}, \frac{3\pi}{4})$  be the value such that the line from  $x$  to  $u(t_0)$  passes through  $y$ . Then we must have  $\|v(t_0)\| > 1$  and by the intermediate value theorem there is some  $0 < t_1 < t_0$  with  $\|v(t_1)\| = 1$ . Let  $u := u(t_1)$ ,  $v := v(t_1)$ . The line through  $u$  and  $x$  is not parallel to the line through  $0$  and  $x + y$  and hence there are  $\alpha, \lambda \geq 0$  such that

$$\frac{\alpha}{2}(x + y) = \lambda u + (1 - \lambda)x.$$

Moreover, since the line from  $x$  to  $u$  is above the line from  $x$  to  $y$  (since  $t_1 < t_0$ ) we have  $\alpha \geq 1$ . Rearranging this equation we get

$$\frac{\alpha}{2}(u + v) = (\alpha + \lambda)u + (1 - \alpha - \lambda)x.$$

Now, consider the convex function  $f(t) := \|tu + (1 - t)x\|$  which satisfies  $f(0) = f(1) = 1$ . Then for  $0 \leq \lambda \leq 1$ ,  $\alpha \geq 1$  we have  $f(\lambda) \leq 1 \leq f(\lambda + \alpha)$  and for  $\lambda \geq 1$ ,  $\alpha \geq 1$  we have  $f(\lambda) \leq f(\lambda + \alpha)$ . Hence we always have  $f(\lambda) \leq f(\lambda + \alpha)$  or equivalently  $\|\frac{1}{2}(x + y)\| \leq \|\frac{1}{2}(u + v)\|$  and  $u, v$  is as required.  $\square$

Now we can prove:

**Theorem 6.21** (Milman–Pettis). *A uniformly convex Banach space is reflexive.*

**Proof.** Pick some  $x'' \in X^{**}$  with  $\|x''\| = 1$ . It suffices to find some  $x \in \bar{B}_1(0)$  with  $\|x'' - J(x)\| \leq \varepsilon$ . So fix  $\varepsilon > 0$  and  $\delta := \delta(\varepsilon)$ , where  $\delta(\varepsilon)$  is the modulus of convexity. Then  $\|x''\| = 1$  implies that we can find some  $\ell \in X^*$  with  $\|\ell\| = 1$  and  $|x''(\ell)| > 1 - \frac{\delta}{2}$ . Consider the weak-\* neighborhood

$$U := \{y'' \in X^{**} \mid |(y'' - x'')(\ell)| < \frac{\delta}{2}\}$$

of  $x''$ . By Goldstine's theorem (Theorem 6.14) there is some  $x \in \bar{B}_1(0)$  with  $J(x) \in U$  and this is the  $x$  we are looking for. In fact, suppose this were not the case. Then the set  $V := X^{**} \setminus \bar{B}_\varepsilon^{**}(J(x))$  is another weak-\* neighborhood

of  $x''$  (since  $\bar{B}_\varepsilon^{**}(J(x))$  is weak-\* compact by the Banach-Alaoglu theorem) and appealing again to Goldstine's theorem there is some  $y \in \bar{B}_1(0)$  with  $J(y) \in U \cap V$ . Since  $x, y \in U$  we obtain

$$1 - \frac{\delta}{2} < |x''(\ell)| \leq |\ell(\frac{x+y}{2})| + \frac{\delta}{2} \Rightarrow 1 - \delta < |\ell(\frac{x+y}{2})| \leq \|\frac{x+y}{2}\|,$$

a contradiction to uniform convexity since  $\|x - y\| \geq \varepsilon$ .  $\square$

**Problem 6.35.** Find an equivalent norm for  $\ell^1(\mathbb{N})$  such that it becomes strictly convex (cf. Problems 1.13 and 1.18).

**Problem\* 6.36.** Show that a Hilbert space is uniformly convex. (Hint: Use the parallelogram law.)

**Problem 6.37.** A Banach space  $X$  is uniformly convex if and only if  $\|x_n\| = \|y_n\| = 1$  and  $\|\frac{x_n + y_n}{2}\| \rightarrow 1$  implies  $\|x_n - y_n\| \rightarrow 0$ .

**Problem\* 6.38.** Show that a finite dimensional space is uniformly convex if and only if it is strictly convex.

**Problem 6.39.** Let  $X$  be strictly convex. Show that every nonzero linear functional attains its norm for at most one unit vector (cf. Problem 4.12).



# Advanced Spectral theory

## 7.1. Spectral theory for compact operators

So far we have developed spectral theory on an algebraic level based on the fact that bounded operators form a Banach algebra. In this section we want to take a more operator centered view and consider bounded linear operators  $\mathcal{L}(X)$ , where  $X$  is some Banach space. Now we can make a finer subdivision of the spectrum based on why our operator fails to have a bounded inverse. Since in the bijective case boundedness of the inverse comes for free from the inverse mapping theorem (Theorem 4.6), there are basically two things which can go wrong: Either our map is not injective or it is not surjective. Moreover, in the latter case one can also ask how far it is away from being surjective, that is, if the range is dense or not. Accordingly one defines the **point spectrum**

$$\sigma_p(A) := \{\alpha \in \sigma(A) \mid \text{Ker}(A - \alpha) \neq \{0\}\} \quad (7.1)$$

as the set of all eigenvalues, the **continuous spectrum**

$$\sigma_c(A) := \{\alpha \in \sigma(A) \setminus \sigma_p(A) \mid \overline{\text{Ran}(A - \alpha)} = X\} \quad (7.2)$$

and finally the **residual spectrum**

$$\sigma_r(A) := \{\alpha \in \sigma(A) \setminus \sigma_p(A) \mid \overline{\text{Ran}(A - \alpha)} \neq X\}. \quad (7.3)$$

Clearly we have

$$\sigma(A) = \sigma_p(A) \dot{\cup} \sigma_c(A) \dot{\cup} \sigma_r(A). \quad (7.4)$$

Here the dot indicates that the union is disjoint. Note that in a Hilbert space  $\sigma_x(A^*) = \sigma_x(A')^*$  for  $x \in \{p, c, r\}$ .

**Example 7.1.** Suppose  $\mathfrak{H}$  is a Hilbert space and  $A = A^*$  is bounded and self-adjoint. Then by (2.28),  $\sigma_r(A) = \emptyset$ .  $\diamond$

**Example 7.2.** Suppose  $X := \ell^p(\mathbb{N})$  and  $L$  is the left shift. Then  $\sigma(L) = \bar{B}_1(0)$ . Indeed, a simple calculation shows that  $\text{Ker}(L - \alpha) = \text{span}\{(\alpha^j)_{j \in \mathbb{N}}\}$  for  $|\alpha| < 1$  if  $1 \leq p < \infty$  and for  $|\alpha| \leq 1$  if  $p = \infty$ . Hence  $\sigma_p(L) = B_1(0)$  for  $1 \leq p < \infty$  and  $\sigma_p(L) = \bar{B}_1(0)$  if  $p = \infty$ . In particular, since the spectrum is closed and  $\|L\| = 1$  we have  $\sigma(L) = \bar{B}_1(0)$ . Moreover, for  $y \in \ell_c(\mathbb{N})$  we set  $x_j := -\sum_{k=j}^{\infty} \alpha^{j-k-1} y_k$  such that  $(L - \alpha)x = y$ . In particular,  $\ell_c(\mathbb{N}) \subset \text{Ran}(L - \alpha)$  and hence  $\text{Ran}(L - \alpha)$  is dense for  $1 \leq p < \infty$ . Thus  $\sigma_c(L) = \partial B_1(0)$  for  $1 \leq p < \infty$ . Consequently,  $\sigma_r(L) = \emptyset$ .  $\diamond$

Since  $A$  is invertible if and only if  $A'$  is by Theorem 4.22 we obtain:

**Lemma 7.1.** *Suppose  $A \in \mathcal{L}(X)$ . Then*

$$\sigma(A) = \sigma(A'). \quad (7.5)$$

Moreover,

$$\begin{aligned} \sigma_p(A') &\subseteq \sigma_p(A) \cup \sigma_r(A), & \sigma_p(A) &\subseteq \sigma_p(A') \cup \sigma_r(A'), \\ \sigma_r(A') &\subseteq \sigma_p(A) \cup \sigma_c(A), & \sigma_r(A) &\subseteq \sigma_p(A'), \\ \sigma_c(A') &\subseteq \sigma_c(A), & \sigma_c(A) &\subseteq \sigma_r(A') \cup \sigma_c(A'). \end{aligned} \quad (7.6)$$

If in addition,  $X$  is reflexive we have  $\sigma_r(A') \subseteq \sigma_p(A)$  as well as  $\sigma_c(A') = \sigma_c(A)$ .

**Proof.** As already indicated, the first claim follows from Theorem 4.22. The remaining items follow from Lemma 4.23 and (4.7). For example, if  $\alpha \in \sigma_p(A')$ , then  $\text{Ran}(A)^\perp = \text{Ker}(A' - \alpha) \neq \{0\}$  and hence  $\overline{\text{Ran}(A)} \neq X$ , so  $\alpha \notin \sigma_c(X)$ . If  $\alpha \in \sigma_r(A')$ , then  $\text{Ker}(A') = \{0\}$  and hence  $\overline{\text{Ran}(A)} = X$ , so  $\alpha \notin \sigma_r(X)$ . Etc. In the reflexive case use  $A \cong A''$  by (4.14).  $\square$

**Example 7.3.** Consider  $L'$  from the previous example, which is just the right shift in  $\ell^q(\mathbb{N})$  if  $1 \leq p < \infty$ . Then  $\sigma(L') = \sigma(L) = \bar{B}_1(0)$ . Moreover, it is easy to see that  $\sigma_p(L') = \emptyset$ . Thus in the reflexive case  $1 < p < \infty$  we have  $\sigma_c(L') = \sigma_c(L) = \partial B_1(0)$  as well as  $\sigma_r(L') = \sigma(L') \setminus \sigma_c(L') = B_1(0)$ . Otherwise, if  $p = 1$ , we only get  $B_1(0) \subseteq \sigma_r(L')$  and  $\sigma_c(L') \subseteq \sigma_c(L) = \partial B_1(0)$ . Hence it remains to investigate  $\text{Ran}(L' - \alpha)$  for  $|\alpha| = 1$ : If we have  $(L' - \alpha)x = y$  with some  $y \in \ell^\infty(\mathbb{N})$ , we must have  $x_j := -\alpha^{-j-1} \sum_{k=1}^j \alpha^k y_k$ . Thus  $y = ((\alpha^*)^n)_{n \in \mathbb{N}}$  is clearly not in  $\text{Ran}(L' - \alpha)$ . Moreover, if  $\|y - \tilde{y}\|_\infty \leq \varepsilon$  we have  $|\tilde{x}_j| = |\sum_{k=1}^j \alpha^k \tilde{y}_k| \geq (1 - \varepsilon)j$  and hence  $\tilde{y} \notin \text{Ran}(L' - \alpha)$ , which shows that the range is not dense and hence  $\sigma_r(L') = \bar{B}_1(0)$ ,  $\sigma_c(L') = \emptyset$ .  $\diamond$

**Example 7.4.** Consider the bilateral left shift  $(Sx)_j = x_{j+1}$  on  $X := \ell^p(\mathbb{Z})$ . By  $\|S\| = 1$  we conclude  $\sigma(S) \subseteq \bar{B}_1(0)$ . Moreover, since its inverse is the corresponding right shift  $(S^{-1}x)_j = x_{j-1}$ , we also have  $\|S^{-1}\| = 1$  and thus

$\sigma(S^{-1}) \subseteq \bar{B}_1(0)$ . Since we also have  $\sigma(S^{-1}) = \sigma(S)^{-1}$  (cf. Problem 5.4), we arrive at  $\sigma(S) \subseteq \partial B_1(0)$  as well as  $\sigma(S^{-1}) \subseteq \partial B_1(0)$ .

Now for  $|\alpha| = 1$  there are two cases: If  $p = \infty$ , then clearly  $x_j := \alpha^j$  is in  $\text{Ker}(S - \alpha)$  and hence  $\sigma(S) = \sigma_p(S) = \partial B_1(0)$  as well as  $\sigma(S^{-1}) = \sigma_p(S^{-1}) = \partial B_1(0)$ . If  $p < \infty$ , then one concludes that if  $y = (S - \alpha)x$  has compact support, so has  $x$  (as outside the support of  $y$  the sequence  $x$  must equal a multiple of  $\alpha^j$  and this multiple must be 0 since  $x \in \ell^p(\mathbb{Z})$ ). Moreover, in this case we further infer  $\sum_{j \in \mathbb{Z}} \alpha^{-j} y_j = \sum_{j \in \mathbb{Z}} \alpha^{-j} (x_{j+1} - \alpha x_j) = 0$ . Hence any sequence  $y$  with compact support violating this condition cannot be in the range of  $S - \alpha$ . Hence  $\sigma(S) = \partial B_1(0)$ . Moreover, since  $S' = S^{-1}$  we see from the fact that  $\sigma_p(S) = \sigma_p(S^{-1}) = \emptyset$ , that  $\sigma_r(S) = \sigma_r(S^{-1}) = \emptyset$ . Hence  $\sigma_c(S) = \sigma_c(S^{-1}) = \partial B_1(0)$ .  $\diamond$

Moreover, for compact operators the spectrum is particularly simple (cf. also Theorem 3.7). We start with the following observation:

**Lemma 7.2.** *Suppose that  $K \in \mathcal{K}(X)$  and  $\alpha \in \mathbb{C} \setminus \{0\}$ . Then  $\text{Ker}(K - \alpha)$  is finite dimensional and the range  $\text{Ran}(K - \alpha)$  is closed.*

**Proof.** For  $\alpha \neq 0$  we can consider  $\mathbb{I} - \alpha^{-1}K$  and assume  $\alpha = 1$  without loss of generality. First of all note that  $K$  restricted to  $\text{Ker}(\mathbb{I} - K)$  is the identity and since the identity is compact the corresponding space must be finite dimensional by Theorem 4.27. In particular, it is complemented (Problem 4.20), that is, there exists a closed subspace  $X_0 \subseteq X$  such that  $X = \text{Ker}(\mathbb{I} - K) \dot{+} X_0$ .

To see that  $\text{Ran}(\mathbb{I} - K)$  is closed we consider  $\mathbb{I} - K$  restricted to  $X_0$  which is injective and has the same range. Hence if  $\text{Ran}(\mathbb{I} - K)$  were not closed, Corollary 8.4 would imply that there is a sequence  $x_n \in X_0$  with  $\|x_n\| = 1$  and  $x_n - Kx_n \rightarrow 0$ . By compactness of  $K$  we can pass to a subsequence such that  $Kx_n \rightarrow y$  implying  $x_n \rightarrow y \in X_0$  and hence  $y \in \text{Ker}(\mathbb{I} - K)$  contradicting  $y \in X_0$  with  $\|y\| = 1$ .  $\square$

Next, we want to have a closer look at eigenvalues. Note that eigenvectors corresponding to different eigenvalues are always linearly independent (Problem 7.5). In Theorem 3.7 we have seen that for a symmetric compact operator in a Hilbert space we can choose an orthonormal basis of eigenfunctions. Without the symmetry assumption we know that even in the finite dimensional case we can in general no longer find a basis of eigenfunctions and that the Jordan canonical form is the best one can do. There the generalized eigenspaces  $\text{Ker}((A - \alpha)^k)$  play an important role. In this respect one looks at the following ascending and descending chains of invariant subspaces associated to  $A \in \mathcal{L}(X)$  (where we have assumed  $\alpha = 0$  without loss



of generality):

$$\{0\} \subseteq \text{Ker}(A) \subseteq \text{Ker}(A^2) \subseteq \text{Ker}(A^3) \subseteq \dots \quad (7.7)$$

and

$$X \supseteq \text{Ran}(A) \supseteq \text{Ran}(A^2) \supseteq \text{Ran}(A^3) \supseteq \dots \quad (7.8)$$

We will say that the kernel chain **stabilizes** at  $n \in \mathbb{N}_0$  if  $\text{Ker}(A^{n+1}) = \text{Ker}(A^n)$ . In this case the number  $n$  is also called the **ascent** of  $A$ . Substituting  $x = Ay$  in the equivalence  $A^n x = 0 \Leftrightarrow A^{n+1} x = 0$  gives  $A^{n+1} y = 0 \Leftrightarrow A^{n+2} y = 0$  and hence by induction we have  $\text{Ker}(A^{n+k}) = \text{Ker}(A^n)$  for all  $k \in \mathbb{N}_0$  in this case. Note that if  $\text{Ker}(A) = \{0\}$ , then the kernel chain stabilizes at 0. Similarly, will say that the range chain **stabilizes** at  $m \in \mathbb{N}_0$  if  $\text{Ran}(A^{m+1}) = \text{Ran}(A^m)$  and call  $m$  the **descent** of  $A$ . Again, if  $x = A^{m+1} y \in \text{Ran}(A^{m+1})$  we can write  $A^m y = A^{m+1} z$  for some  $z$  which shows  $x = A^{m+2} z \in \text{Ran}(A^{m+2})$  and thus  $\text{Ran}(A^{m+k}) = \text{Ran}(A^m)$  for all  $k \in \mathbb{N}_0$  in this case. While in a finite dimensional space both chains eventually have to stabilize, there is no reason why the same should happen in an infinite dimensional space.

**Example 7.5.** For the left shift operator  $L$  we have  $\text{Ran}(L^n) = \ell^p(\mathbb{N})$  for all  $n \in \mathbb{N}$  while the kernel chain does not stabilize as  $\text{Ker}(L^n) = \{a \in \ell^p(\mathbb{N}) | a_j = 0, j > n\}$ . Similarly, for the right shift operator  $R$  we have  $\text{Ker}(R^n) = \{0\}$  while the range chain does not stabilize as  $\text{Ran}(R^n) = \{a \in \ell^p(\mathbb{N}) | a_j = 0, 1 \leq j \leq n\}$ .  $\diamond$

**Lemma 7.3.** Suppose  $A : X \rightarrow X$  is a linear operator.

- (i) The kernel chain stabilizes at  $n$  if  $\text{Ran}(A^n) \cap \text{Ker}(A) = \{0\}$ . Conversely, if the kernel chain stabilizes at  $n$ , then  $\text{Ran}(A^n) \cap \text{Ker}(A^n) = \{0\}$  and  $A$  restricted to  $\text{Ran}(A^n) \rightarrow \text{Ran}(A^n)$  is injective.
- (ii) The range chain stabilizes at  $m$  if  $\text{Ker}(A^m) + \text{Ran}(A) = X$ . Conversely, if the range chain stabilizes at  $m$ , then  $\text{Ker}(A^m) + \text{Ran}(A^m) = X$  and  $A$  restricted to  $\text{Ran}(A^m) \rightarrow \text{Ran}(A^m)$  is surjective.
- (iii) If both chains stabilize, then  $m = n$  and  $\text{Ker}(A^m) \dot{+} \text{Ran}(A^m) = X$  and  $A$  restricted to  $\text{Ran}(A^m) \rightarrow \text{Ran}(A^m)$  is bijective.

**Proof.** (i). If  $\text{Ran}(A^n) \cap \text{Ker}(A) = \{0\}$  then  $x \in \text{Ker}(A^{n+1})$  implies  $A^n x \in \text{Ran}(A^n) \cap \text{Ker}(A) = \{0\}$  and the kernel chain stabilizes at  $n$ . Conversely, let  $x \in \text{Ran}(A^n) \cap \text{Ker}(A^n)$ , then  $x = A^n y$  and  $A^n x = A^{2n} y = 0$  implying  $y \in \text{Ker}(A^{2n}) = \text{Ker}(A^n)$ , that is,  $x = A^n y = 0$ . Moreover, if  $Ax = 0$  for some  $x = A^n y \in \text{Ran}(A^n)$ , then  $y \in \text{Ker}(A^{n+1}) = \text{Ker}(A^n)$  implying  $x = 0$ .

(ii). If  $\text{Ker}(A^m) + \text{Ran}(A) = X$ , then for any  $x = z + Ay$  we have  $A^m x = A^{m+1} y$  and hence  $\text{Ran}(A^m) = \text{Ran}(A^{m+1})$ . Conversely, if the range chain stabilizes at  $m$ , then  $A^m x = A^{2m} y$  and  $x = A^m y + (x - A^m y)$ . Moreover,  $A \text{Ran}(A^m) = \text{Ran}(A^{m+1}) = \text{Ran}(A^m)$  shows surjectivity.

(iii). Suppose  $\text{Ran}(A^{m+1}) = \text{Ran}(A^m)$  but  $\text{Ker}(A^m) \subsetneq \text{Ker}(A^{m+1})$ . Let  $x \in \text{Ker}(A^{m+1}) \setminus \text{Ker}(A^m)$  and observe that by  $0 \neq A^m x = A^{m+1} y$  there is an  $y \in \text{Ker}(A^{m+2}) \setminus \text{Ker}(A^{m+1})$ . Iterating this argument would show that the kernel chain does not stabilize, contradiction our assumption. Hence  $n \leq m$ .

Conversely, suppose  $\text{Ker}(A^{n+1}) = \text{Ker}(A^n)$  and  $\text{Ran}(A^{m+1}) = \text{Ran}(A^m)$  for  $m \geq n$ . Then for every  $x$  there is some  $y$  such that

$$A^m x = A^{m+1} y \Rightarrow x - Ay \in \text{Ker}(A^m) = \text{Ker}(A^n) \Rightarrow A^n x = A^{n+1} y$$

shows  $\text{Ran}(A^{n+1}) = \text{Ran}(A^n)$ , that is,  $m \leq n$ . The rest follows by combining (i) and (ii).  $\square$

Of course the desired case is (iii), where we have a splitting of  $X$  into a direct sum of invariant subspaces such that  $A$  restricted to  $\text{Ker}(A^m) \rightarrow \text{Ker}(A^m)$  is nilpotent and restricted to  $\text{Ran}(A^m) \rightarrow \text{Ran}(A^m)$  is bijective.

**Example 7.6.** In a finite dimensional space we are of course always in case (iii).  $\diamond$

**Example 7.7.** Let  $A \in \mathcal{L}(\mathfrak{H})$  be a self-adjoint operator in a Hilbert space. Then by (2.28) the kernel chain always stabilizes at  $n = 1$  and the range chain stabilizes at  $n = 1$  if  $\text{Ran}(A)$  is closed.  $\diamond$

As a further preparation we establish the following result which should be thought of a replacement of an orthogonal vector.

**Lemma 7.4** (Riesz lemma). *Let  $X$  be a normed vector space and  $Y \subset X$  some subspace. which is not dense  $\overline{Y} \neq X$ . Then for every  $\varepsilon \in (0, 1)$  there exists an  $x_\varepsilon$  with  $\|x_\varepsilon\| = 1$  and*

$$\inf_{y \in Y} \|x_\varepsilon - y\| \geq 1 - \varepsilon. \quad (7.9)$$

**Proof.** Pick  $x \in X \setminus \overline{Y}$  and abbreviate  $d := \text{dist}(x, Y) > 0$ . Choose  $y_\varepsilon \in Y$  such that  $\|x - y_\varepsilon\| \leq \frac{d}{1-\varepsilon}$ . Then  $x_\varepsilon := \frac{x - y_\varepsilon}{\|x - y_\varepsilon\|}$  is the vector we are looking for since

$$\|x_\varepsilon - y\| = \frac{1}{\|x - y_\varepsilon\|} \|x - (y_\varepsilon + \|x - y_\varepsilon\|y)\| \geq \frac{d}{\|x - y_\varepsilon\|} \geq 1 - \varepsilon$$

as required.  $\square$

As already indicated, in a Hilbert space the claim holds with  $\varepsilon = 0$  for any normalized  $x$  in the orthogonal complement of  $Y$ . Slightly more general, if for every  $x \in X \setminus Y$  we can find some  $y_0 \in \overline{Y}$  with  $\|x - y_0\| = \text{dist}(x, Y)$ , then the lemma holds with  $\varepsilon = 0$ . This is the case in reflexive spaces — Example 9.22.

Now we can apply this to our situation.

**Lemma 7.5.** *Suppose that  $K \in \mathcal{K}(X)$  and  $\alpha \in \mathbb{C} \setminus \{0\}$ . Then the space  $\text{Ker}(K - \alpha)^m$  is finite dimensional and the space  $\text{Ran}(K - \alpha)^m$  is closed for every  $m \in \mathbb{N}$ . Moreover, there is some  $n = n(\alpha) \in \mathbb{N}_0$  such that the kernel and range chain of both  $K - \alpha$  and  $K' - \alpha$  stabilize and hence*

$$\begin{aligned} X &= \text{Ker}(K - \alpha)^n \dot{+} \text{Ran}(K - \alpha)^n, \\ X^* &= \text{Ker}(K' - \alpha)^n \dot{+} \text{Ran}(K' - \alpha)^n. \end{aligned} \quad (7.10)$$

**Proof.** Since  $\alpha \neq 0$  we can consider  $\mathbb{I} - \alpha^{-1}K$  and assume  $\alpha = 1$  without loss of generality. Moreover, since  $(\mathbb{I} - K)^n - \mathbb{I} \in \mathcal{K}(X)$  we see that  $\text{Ker}(\mathbb{I} - K)^n$  is finite dimensional and  $\text{Ran}(\mathbb{I} - K)^n$  is closed for every  $n \in \mathbb{N}$ . Next suppose the kernel chain does not stabilize. Abbreviate  $K_n := \text{Ker}(\mathbb{I} - K)^n$ . Then, by Lemma 7.4, we can choose  $x_n \in K_{n+1} \setminus K_n$  such that  $\|x_n\| = 1$  and  $\text{dist}(x_n, K_n) \geq \frac{1}{2}$ . But since  $(\mathbb{I} - K)x_n \in K_n$  and  $Kx_n \in K_{n+1}$ , we see that

$$\|Kx_n - Kx_m\| = \|x_n - (\mathbb{I} - K)x_n - Kx_m\| \geq \text{dist}(x_n, K_n) \geq \frac{1}{2}$$

for  $n > m$  and hence the image of the bounded sequence  $Kx_n$  has no convergent subsequence, a contradiction.

Consequently the kernel sequence for  $\mathbb{I} - K'$  also stabilizes. Moreover, by the closed range theorem (Theorem 4.24) we have

$$\text{Ker}((\mathbb{I} - K')^n)^\perp = \text{Ran}(\mathbb{I} - K)^n, \quad \text{Ker}((\mathbb{I} - K)^n)^\perp = \text{Ran}(\mathbb{I} - K')^n$$

and hence the kernel chain of  $\mathbb{I} - K'$  stabilizes at  $n$  when the range chain of  $\mathbb{I} - K$  stabilizes at  $n$  and the kernel chain of  $\mathbb{I} - K$  stabilizes at  $n$  when the range chain of  $\mathbb{I} - K'$  stabilizes at  $n$ . The rest follows from the previous lemma.  $\square$

As an immediate consequence we get the famous Fredholm alternative:

**Theorem 7.6** (Fredholm alternative). *Suppose that  $K \in \mathcal{K}(X)$  and  $\alpha \in \mathbb{C} \setminus \{0\}$ . Then the following are equivalent:*

- $\text{Ker}(K - \alpha) = \{0\}$ .
- $\text{Ran}(K - \alpha) = X$ .
- $\text{Ker}(K' - \alpha) = \{0\}$ .
- $\text{Ran}(K' - \alpha) = X^*$ .

Of course in terms of equations this can equivalently be phrased as, either the inhomogeneous equation

$$(K - \alpha)x = y \quad (7.11)$$

has a unique solution for every  $y \in X$  or the corresponding homogeneous equation

$$(K - \alpha)x = 0 \quad (7.12)$$

has a nontrivial solution. Moreover, by the closed range theorem (Theorem 4.24), there will be a solution if and only if  $\ell(y) = 0$  for all  $\ell \in \text{Ker}(K' - \alpha)$ . Of course the analogous statement holds for  $K'$ .

In particular, this applies to the case where  $K$  is a compact integral operator (cf. Lemma 3.4), which was the case originally studied by Fredholm.

Note that this also implies

$$\sigma_p(K) \setminus \{0\} = \sigma_p(K') \setminus \{0\}. \quad (7.13)$$

For an eigenvalue  $\alpha \in \mathbb{C}$  the dimension  $\dim(\text{Ker}(A - \alpha))$  is called the **geometric multiplicity** of  $\alpha$  and if the kernel chain stabilizes at  $n$ , then  $n$  is called the **index** of  $\alpha$  and  $\dim(\text{Ker}(A - \alpha)^n)$  is called the **algebraic multiplicity** of  $\alpha$ . Otherwise, if the kernel chain does not stabilize, both the index and the algebraic multiplicity are set equal to infinity. The **order** of a generalized eigenvector  $u$  corresponding to an eigenvalue  $\alpha$  is the smallest  $n$  such that  $(A - \alpha)^n u = 0$ .

**Example 7.8.** Consider  $X := \ell^p(\mathbb{N})$  and  $K \in \mathcal{K}(X)$  given by  $(Ka)_n := \frac{1}{n} a_{n+1}$ . Then  $Ka = \alpha a$  implies  $a_n = \alpha^{n-1}(n-1)!a_1$  and hence  $a_1 = 0$  for  $\alpha \neq 0$  and  $a = a_1 \delta^1$  for  $\alpha = 0$ . Hence  $\sigma(K) = \{0\}$  with 0 being an eigenvalue of geometric multiplicity one. Since  $K^n a = 0$  implies  $a_j = 0$  for  $j > n$  we see that its index as well as its algebraic multiplicity is  $\infty$ . Moreover,  $\delta^n$  is a generalized eigenvalue of order  $n$ .  $\diamond$

**Theorem 7.7** (Spectral theorem for compact operators; F. Riesz). *Suppose that  $K \in \mathcal{K}(X)$ . Then every  $\alpha \in \sigma(K) \setminus \{0\}$  is an eigenvalue of finite algebraic multiplicity and  $X$  can be decomposed into invariant closed subspaces according to (7.10), where  $n$  is the index of  $\alpha$ . Furthermore, there are at most countably many eigenvalues which can only accumulate at 0. If  $X$  is infinite dimensional, we have  $0 \in \sigma(K)$ . In this case either  $0 \in \sigma_p(K)$  with  $\dim(\text{Ker}(K)) = \infty$  or  $\text{Ran}(K)$  is not closed.*

**Proof.** That every eigenvalue  $\alpha \neq 0$  has finite algebraic multiplicity follows from the previous two lemmas. Moreover if  $\text{Ker}(K - \alpha) = \{0\}$ , then the kernel chain stabilizes as  $n = 0$  and hence  $\text{Ran}(K - \alpha) = X$ , that is  $\alpha \notin \sigma(K)$ .

Let  $\alpha_n$  be a sequence of different eigenvalues with  $|\alpha_n| \geq \varepsilon$ . Let  $x_n$  be corresponding normalized eigenvectors and let  $X_n := \text{span}\{x_j\}_{j=1}^n$ . The sequence of spaces  $X_n$  is increasing and by Lemma 7.4 we can choose normalized vectors  $\tilde{x}_n \in X_n$  such that  $\text{dist}(\tilde{x}_n, X_{n-1}) \geq \frac{1}{2}$ . Now, since  $(K - \alpha_n)X_n \subset X_{n-1}$ ,  $\|K\tilde{x}_n - K\tilde{x}_m\| = \|\alpha_n \tilde{x}_n + ((K - \alpha_n)\tilde{x}_n - K\tilde{x}_m)\| \geq \frac{|\alpha_n|}{2} \geq \frac{\varepsilon}{2}$  for  $m < n$  and hence there is no convergent subsequence, a contradiction. Moreover, if  $0 \in \rho(K)$  then  $K^{-1}$  is bounded and hence  $\mathbb{I} = K^{-1}K$  is compact, implying that  $X$  is finite dimensional.

Finally, if  $\text{Ran}(K)$  is closed we can consider the bijective operator  $\tilde{K} : X/\text{Ker}(K) \rightarrow \text{Ran}(K)$  (cf. Problem 1.46) which is again compact. Hence  $\mathbb{I}_{X/\text{Ker}(K)} = \tilde{K}^{-1}\tilde{K}$  is compact and thus  $X/\text{Ker}(K)$  is finite dimensional.  $\square$

**Example 7.9.** Note that in contradistinction to the symmetric case, there might be no eigenvalues at all, as the Volterra integral operator from Example 5.16 shows.  $\diamond$

This result says in particular, that for  $\alpha \in \sigma_p(K) \setminus \{0\}$  we can split  $X = X_1 \oplus X_2$  where both  $X_1 := \text{Ker}(K - \alpha)^n$  and  $X_2 := \text{Ran}(K - \alpha)^n$  are invariant subspaces for  $K$ . Consequently we can split  $K = K_1 \oplus K_2$ , where  $K_1$  is the restriction of  $K$  to the finite dimensional subspace  $X_1$  with  $K_1 - \alpha$  a nilpotent matrix and  $K_2$  is the restriction of  $K$  to  $X_2$  with  $K_2 - \alpha$  bijective and hence  $\alpha \in \rho(K_2)$ .

**Problem 7.1.** Discuss the spectrum of the right shift  $R$  on  $\ell^1(\mathbb{N})$ . Show  $\sigma(R) = \sigma_r(R) = \bar{B}_1(0)$  and  $\sigma_p(R) = \sigma_c(R) = \emptyset$ .

**Problem 7.2.** Compute the point, continuous, and residual spectrum of the multiplication operator  $A : C[0, 1] \rightarrow C[0, 1]$ ,  $x(t) \mapsto a(t)x(t)$ , where  $a(t) := |3t - 2| - |3t - 1|$ .

**Problem 7.3.** Compute the point, continuous, and residual spectrum of the operator  $A : C[0, 1] \rightarrow C[0, 1]$ ,  $x(t) \mapsto x(0) + tx(1)$ . Compute the kernel and the range for each point in the spectrum. Is the range closed or dense?

**Problem 7.4.** Compute the point, continuous, and residual spectrum of the operator  $A : C[0, 1] \rightarrow C[0, 1]$ ,  $x(t) \mapsto \int_0^t x(s)ds$ .

**Problem\* 7.5.** Suppose  $A \in \mathcal{L}(X)$ . Show that generalized eigenvectors corresponding to different eigenvalues or with different order are linearly independent.

**Problem 7.6.** Suppose  $\mathfrak{H}$  is a Hilbert space and  $A \in \mathcal{L}(\mathfrak{H})$  is normal. Then  $\sigma_p(A) = \sigma_p(A^*)^*$ ,  $\sigma_c(A) = \sigma_c(A^*)^*$ , and  $\sigma_r(A) = \sigma_r(A^*) = \emptyset$ . (Hint: Problem 5.22.)

**Problem 7.7.** Suppose  $A_j \in \mathcal{L}(X_j)$ ,  $j = 1, 2$ . Then  $A_1 \oplus A_2 \in \mathcal{L}(X_1 \oplus X_2)$  and  $\sigma(A_1 \oplus A_2) = \sigma(A_1) \cup \sigma(A_2)$ .

**Problem 7.8.** Let  $A : X \rightarrow Y$ ,  $B : Y \rightarrow Z$ . Show  $\dim(\text{Ker}(BA)) \leq \dim(\text{Ker}(A)) + \dim(\text{Ker}(B))$  and hence  $\dim(\text{Ker}(A^n)) \leq n \dim(\text{Ker}(A))$  if  $A : X \rightarrow X$ .

## 7.2. Fredholm operators

In this section we want to investigate solvability of the equation

$$Ax = y \tag{7.14}$$

for  $A \in \mathcal{L}(X, Y)$  given  $y \in Y$ . Clearly there exists a solution if  $y \in \text{Ran}(A)$  and this solution is unique if  $\text{Ker}(A) = \{0\}$ . Hence these subspaces play a crucial role. Moreover, if the underlying Banach spaces are finite dimensional, the kernel has a complement  $X = \text{Ker}(A) \dot{+} X_0$  and after factoring out the kernel this complement is isomorphic to the range of  $A$ . As a consequence, the dimensions of these spaces are connected by the famous **rank-nullity theorem**

$$\dim \text{Ker}(A) + \dim \text{Ran}(A) = \dim X \quad (7.15)$$

from linear algebra. In our infinite dimensional setting (apart from the technical difficulties that the kernel might not be complemented and the range might not be closed) this formula does not contain much information, but if we rewrite it in terms of the index,

$$\text{ind}(A) := \dim \text{Ker}(A) - \dim \text{Coker}(A) = \dim(X) - \dim(Y), \quad (7.16)$$

at least the left-hand side will be finite if we assume both  $\text{Ker}(A)$  and  $\text{Coker}(A) := Y/\text{Ran}(A)$  to be finite dimensional. One of the most useful consequences of the rank-nullity theorem is that in the case  $X = Y$  the index will vanish and hence uniqueness of solutions for  $Ax = y$  will automatically give you existence for free (and vice versa). Indeed, for equations of the form  $x + Kx = y$  with  $K$  compact originally studied by Fredholm this is still true by the famous Fredholm alternative (Theorem 7.6). It took a while until Fritz Noether found an example of singular integral equations which have a nonzero index and started investigating the general case.

We first note that in this case  $\text{Ran}(A)$  will be automatically closed.

**Lemma 7.8.** *Suppose  $A \in \mathcal{L}(X, Y)$  with finite dimensional cokernel. Then  $\text{Ran}(A)$  is closed.*

**Proof.** First of all note that the induced map  $\tilde{A} : X/\text{Ker}(A) \rightarrow Y$  is injective (Problem 1.46). Moreover, the assumption that the cokernel is finite says that there is a finite subspace  $Y_0 \subset Y$  such that  $Y = Y_0 \dot{+} \text{Ran}(A)$ . Then

$$\hat{A} : X/\text{Ker}(A) \oplus Y_0 \rightarrow Y, \quad \hat{A}(x, y) = \tilde{A}x + y$$

is bijective and hence a homeomorphism by Theorem 4.6. Since  $\tilde{X} := X/\text{Ker}(A) \oplus \{0\}$  is a closed subspace of  $X/\text{Ker}(A) \oplus Y_0$  we see that  $\text{Ran}(A) = \hat{A}(\tilde{X})$  is closed in  $Y$ .  $\square$

Hence we call an operator  $A \in \mathcal{L}(X, Y)$  a **Fredholm operator** (also **Noether operator**) if both its kernel and cokernel are finite dimensional. In this case we define its **index** as

$$\text{ind}(A) := \dim \text{Ker}(A) - \dim \text{Coker}(A). \quad (7.17)$$

The set of Fredholm operators will be denoted by  $\Phi(X, Y)$  and the set of Fredholm operators with index zero will be denoted by  $\Phi_0(X, Y)$ .

**Example 7.10.** We have  $\mathbb{I} \in \Phi_0(X)$  but clearly  $0 \notin \Phi(X)$  unless  $X$  is finite dimensional. Moreover,  $A \in \Phi(X, Y)$  implies  $\alpha A \in \Phi(X, Y)$  for  $\alpha \in \mathbb{C} \setminus \{0\}$  with  $\text{ind}(\alpha A) = \text{ind}(A)$  but the sum of two Fredholm operators is in general not Fredholm (e.g.  $A - A = 0$ ). In particular,  $\Phi(X, Y)$  is *not* a linear subspace of  $\mathcal{L}(X, Y)$  unless both  $X$  and  $Y$  are finite dimensional.  $\diamond$

**Lemma 7.9.** *Suppose  $A \in \mathcal{L}(X, Y)$  with  $\text{Ran}(A)$  closed. Then*

$$\text{Ker}(A') \cong \text{Coker}(A)^*, \quad \text{Coker}(A') \cong \text{Ker}(A)^*. \quad (7.18)$$

*In particular, we have*

$$\dim \text{Ker}(A') = \dim \text{Coker}(A), \quad \dim \text{Ker}(A) = \dim \text{Coker}(A'). \quad (7.19)$$

**Proof.** Using Lemma 4.23 and Theorems 4.24, 4.19 we have  $\text{Ker}(A') = \text{Ran}(A)^\perp \cong (Y/\text{Ran}(A))^* = \text{Coker}(A)^*$  and  $\text{Coker}(A') = X^*/\text{Ran}(A') = X^*/\text{Ker}(A)^\perp \cong \text{Ker}(A)^*$ . The second claim follows since for a finite dimensional space the dual space has the same dimension.  $\square$

An immediate consequence is:

**Theorem 7.10 (Riesz).** *A bounded operator  $A$  is Fredholm if and only if  $A'$  is and in this case*

$$\text{ind}(A') = -\text{ind}(A). \quad (7.20)$$

**Example 7.11.** The left shift operator  $L$  in  $X = Y := \ell^p(\mathbb{N})$ ,  $1 \leq p < \infty$  is Fredholm. In fact, we have  $\text{Ker}(L) = \text{span}\{\delta^1\}$  and  $\text{Ran}(L) = X$  implying  $\text{ind}(L) = 1$ . Consequently  $L'$ , which is just the right shift, is also Fredholm with  $\text{ind}(L') = -1$ . Of course the last fact can also be checked directly.  $\diamond$

In the case of Hilbert spaces  $\text{Ran}(A)$  closed implies  $\mathfrak{H} = \text{Ran}(A) \oplus \text{Ran}(A)^\perp$  and thus  $\text{Coker}(A) \cong \text{Ran}(A)^\perp$ . Hence an operator is Fredholm if  $\text{Ran}(A)$  is closed and  $\text{Ker}(A)$  and  $\text{Ran}(A)^\perp$  are both finite dimensional. In this case

$$\text{ind}(A) = \dim \text{Ker}(A) - \dim \text{Ran}(A)^\perp \quad (7.21)$$

and  $\text{ind}(A^*) = -\text{ind}(A)$  as is immediate from (2.28).

**Example 7.12.** Suppose  $\mathfrak{H}$  is a Hilbert space and  $A = A^*$  is a self-adjoint Fredholm operator, then (2.28) shows that  $\text{ind}(A) = 0$ . In particular, a self-adjoint operator is Fredholm if  $\dim \text{Ker}(A) < \infty$  and  $\text{Ran}(A)$  is closed. For example, according to Example 5.23,  $A - \lambda$  is Fredholm if  $\lambda$  is an eigenvalue of finite multiplicity (in fact, inspecting this example shows that the converse is also true).

It is however important to notice that  $\text{Ran}(A)^\perp$  finite dimensional does *not* imply  $\text{Ran}(A)$  closed! For example consider  $(Ax)_n = \frac{1}{n}x_n$  in  $\ell^2(\mathbb{N})$  whose range is dense but not closed.  $\diamond$

Another useful formula concerns the product of two Fredholm operators. For its proof it will be convenient to use the notion of an exact sequence: Let  $X_j$  be Banach spaces. A sequence of operators  $A_j \in \mathcal{L}(X_j, X_{j+1})$

$$X_1 \xrightarrow{A_1} X_2 \xrightarrow{A_2} X_3 \cdots X_n \xrightarrow{A_n} X_{n+1}$$

is said to be **exact** if  $\text{Ran}(A_j) = \text{Ker}(A_{j+1})$  for  $0 \leq j < n$ . We will also need the following two easily checked facts: If  $X_{j-1}$  and  $X_{j+1}$  are finite dimensional, so is  $X_j$  (Problem 7.9) and if the sequence of finite dimensional spaces starts with  $X_0 = \{0\}$  and ends with  $X_{n+1} = \{0\}$ , then the alternating sum over the dimensions vanishes (Problem 7.10).

**Lemma 7.11** (Atkinson). *Suppose  $A \in \mathcal{L}(X, Y)$ ,  $B \in \mathcal{L}(Y, Z)$ . If two of the operators  $A$ ,  $B$ ,  $BA$  are Fredholm, so is the third and we have*

$$\text{ind}(BA) = \text{ind}(A) + \text{ind}(B). \quad (7.22)$$

**Proof.** It is straightforward to check that the sequence

$$\begin{aligned} 0 \longrightarrow \text{Ker}(A) \longrightarrow \text{Ker}(BA) \xrightarrow{A} \text{Ker}(B) \longrightarrow \text{Coker}(A) \\ \xrightarrow{\tilde{B}} \text{Coker}(BA) \longrightarrow \text{Coker}(B) \longrightarrow 0 \end{aligned}$$

is exact. Here the maps which are not explicitly stated are canonical inclusions/quotient maps. Hence by Problem 7.9, if two operators are Fredholm, so is the third. Moreover, the formula for the index follows from Problem 7.10.  $\square$

Next we want to look a bit further into the structure of Fredholm operators. First of all, since  $\text{Ker}(A)$  is finite dimensional, it is complemented (Problem 4.20), that is, there exists a closed subspace  $X_0 \subseteq X$  such that  $X = \text{Ker}(A) \dot{+} X_0$  and a corresponding projection  $P \in \mathcal{L}(X)$  with  $\text{Ran}(P) = \text{Ker}(A)$ . Similarly,  $\text{Ran}(A)$  is complemented (Problem 1.47) and there exists a closed subspace  $Y_0 \subseteq Y$  such that  $Y = Y_0 \dot{+} \text{Ran}(A)$  and a corresponding projection  $Q \in \mathcal{L}(Y)$  with  $\text{Ran}(Q) = Y_0$ . With respect to the decomposition  $\text{Ker}(A) \oplus X_0 \rightarrow Y_0 \oplus \text{Ran}(A)$  our Fredholm operator is given by

$$A = \begin{pmatrix} 0 & 0 \\ 0 & A_0 \end{pmatrix}, \quad (7.23)$$

where  $A_0$  is the restriction of  $A$  to  $X_0 \rightarrow \text{Ran}(A)$ . By construction  $A_0$  is bijective and hence a homeomorphism (Theorem 4.6).



**Example 7.13.** In case of an Hilbert space we can choose  $X_0 = \text{Ker}(A)^\perp = \text{Ran}(A^*)$  and  $Y_0 = \text{Ran}(A)^\perp = \text{Ker}(A^*)$  by (2.28). In particular,  $A : \text{Ker}(A)^\perp \rightarrow \text{Ker}(A^*)^\perp$  has a bounded inverse.  $\diamond$

Defining

$$B := \begin{pmatrix} 0 & 0 \\ 0 & A_0^{-1} \end{pmatrix} \quad (7.24)$$

we get

$$AB = \mathbb{I} - Q, \quad BA = \mathbb{I} - P \quad (7.25)$$

and hence  $A$  is invertible up to finite rank operators. Now we are ready for showing that the index is stable under small perturbations.

**Theorem 7.12** (Dieudonné). *The set of Fredholm operators  $\Phi(X, Y)$  is open in  $\mathcal{L}(X, Y)$  and the index is locally constant.*

**Proof.** Let  $C \in \mathcal{L}(X, Y)$  and write it as

$$C = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix}$$

with respect to the above splitting. Then if  $\|C_{22}\| < \|A_0^{-1}\|^{-1}$  we have that  $A_0 + C_{22}$  is still invertible (Problem 1.38). Now introduce

$$D_1 = \begin{pmatrix} \mathbb{I} & -C_{12}(A_0 + C_{22})^{-1} \\ 0 & \mathbb{I} \end{pmatrix}, \quad D_2 = \begin{pmatrix} \mathbb{I} & 0 \\ -(A_0 + C_{22})^{-1}C_{21} & \mathbb{I} \end{pmatrix}$$

and observe

$$D_1(A + C)D_2 = \begin{pmatrix} C_{11} - C_{12}(A_0 + C_{22})^{-1}C_{21} & 0 \\ 0 & A_0 + C_{22} \end{pmatrix}.$$

Since  $D_1, D_2$  are homeomorphisms we see that  $A + C$  is Fredholm since the right-hand side obviously is. Moreover,  $\text{ind}(A + C) = \text{ind}(C_{11} - C_{12}(A_0 + C_{22})^{-1}C_{21}) = \dim(\text{Ker}(A)) - \dim(Y_0) = \text{ind}(A)$  since the second operator is between finite dimensional spaces and the index can be evaluated using (7.16).  $\square$

Since the index is locally constant, it is constant on every connected component of  $\Phi(X, Y)$  which often is useful for computing the index. The next result identifies an important class of Fredholm operators and uses this observation for computing the index.

**Theorem 7.13** (Riesz). *For every  $K \in \mathcal{K}(X)$  we have  $\mathbb{I} - K \in \Phi_0(X)$ .*

**Proof.** That  $\mathbb{I} - K$  is Fredholm follows from Lemma 7.2 since  $K'$  is compact as well and  $\text{Coker}(\mathbb{I} - K)^* \cong \text{Ker}(\mathbb{I} - K')$  by Lemma 7.9. Furthermore, the index is constant along  $[0, 1] \rightarrow \Phi(X)$ ,  $\alpha \mapsto \mathbb{I} - \alpha K$  and hence  $\text{ind}(\mathbb{I} - K) = \text{ind}(\mathbb{I}) = 0$ .  $\square$

Next we show that an operator is Fredholm if and only if it has a left/right inverse up to compact operators.

**Theorem 7.14** (Atkinson). *An operator  $A \in \mathcal{L}(X, Y)$  is Fredholm if and only if there exist  $B_1, B_2 \in \mathcal{L}(Y, X)$  such that  $B_1A - \mathbb{I} \in \mathcal{K}(X)$  and  $AB_2 - \mathbb{I} \in \mathcal{K}(Y)$ .*

**Proof.** If  $A$  is Fredholm we have already given an operator  $B$  in (7.24) such that  $BA - \mathbb{I}$  and  $AB - \mathbb{I}$  are finite rank. Conversely, according to Theorem 7.13  $B_1A$  and  $AB_2$  are Fredholm. Since  $\text{Ker}(A) \subseteq \text{Ker}(B_1A)$  and  $\text{Ran}(AB_2) \subseteq \text{Ran}(A)$  this shows that  $A$  is Fredholm.  $\square$

Operators  $B_1$  and  $B_2$  as in the previous theorem are also known as a left and right **parametrix**, respectively. As a consequence we can now strengthen Theorem 7.13:

**Corollary 7.15** (Yood). *For every  $A \in \Phi(X, Y)$  and  $K \in \mathcal{K}(X, Y)$  we have  $A + K \in \Phi(X, Y)$  with  $\text{ind}(A + K) = \text{ind}(A)$ .*

**Proof.** Using (7.24) we see that  $B(A + K) - \mathbb{I} = -P + BK \in \mathcal{K}(X)$  and  $(A + K)B - \mathbb{I} = -Q + KB \in \mathcal{K}(Y)$  and hence  $A + K$  is Fredholm. In fact,  $A + \alpha K$  is Fredholm for  $\alpha \in [0, 1]$  and hence  $\text{ind}(A + K) = \text{ind}(A)$  since the index is locally constant.  $\square$

Fredholm operators are also used to split the spectrum. For  $A \in \mathcal{L}(X)$  one defines the **essential spectrum**

$$\sigma_{\text{ess}}(A) := \{\alpha \in \mathbb{C} \mid A - \alpha \notin \Phi_0(X)\} \subseteq \sigma(A) \quad (7.26)$$

and the **Fredholm spectrum**

$$\sigma_{\Phi}(A) := \{\alpha \in \mathbb{C} \mid A - \alpha \notin \Phi(X)\} \subseteq \sigma_{\text{ess}}(A). \quad (7.27)$$

By Dieudonné's theorem both  $\sigma_{\text{ess}}(A)$  and  $\sigma_{\Phi}(A)$  are closed. Also note that we have  $\sigma_c(A) \subseteq \sigma_{\Phi}(A)$ . Warning: These definitions are not universally accepted and several variants can be found in the literature.

**Example 7.14.** Let  $X$  be infinite dimensional and  $K \in \mathcal{K}(X)$ . Then  $\sigma_{\text{ess}}(K) = \sigma_{\Phi}(K) = \{0\}$ .  $\diamond$

**Example 7.15.** If  $X$  is a Hilbert space and  $A$  is self-adjoint, then  $\sigma(A) \subseteq \mathbb{R}$  and for  $\alpha \in \mathbb{R} \setminus \sigma_{\Phi}(A)$  the identity  $\text{ind}(A - \alpha) = -\text{ind}((A - \alpha)^*) = -\text{ind}(A - \alpha)$  shows that the index is always zero. Thus  $\sigma_{\text{ess}}(A) = \sigma_{\Phi}(A)$  for self-adjoint operators.  $\diamond$

By Corollary 7.15 both the Fredholm spectrum and the essential spectrum are invariant under compact perturbations:

**Theorem 7.16** (Weyl). *Let  $A \in \mathcal{L}(X)$ , then*

$$\sigma_{\Phi}(A + K) = \sigma_{\Phi}(A), \quad \sigma_{\text{ess}}(A + K) = \sigma_{\text{ess}}(A), \quad K \in \mathcal{K}(X). \quad (7.28)$$

Moreover, if  $\alpha \notin \sigma_{ess}(A)$ , using the notation from (7.23) for  $A - \alpha$ , we can find an operator  $K_0 : \text{Ker}(A - \alpha) \rightarrow Y_0$  such that  $K_0 - \alpha$  is invertible and extend it to a finite rank operator  $K : X \rightarrow X$  by setting it equal to 0 on  $X_0$ . Then  $(A + K) - \alpha$  has a bounded inverse implying  $\alpha \in \rho(A + K)$ . Thus the essential spectrum is precisely the part which is stable under compact perturbations

$$\sigma_{ess}(A) = \bigcap_{K \in \mathcal{K}(X)} \sigma(A + K). \quad (7.29)$$

The complement is called the **discrete spectrum**

$$\sigma_d(A) := \sigma(A) \setminus \sigma_{ess}(A) = \{\alpha \in \sigma(A) \mid A - \alpha \in \Phi_0(X)\} \subseteq \sigma_p(A). \quad (7.30)$$

Clearly points in the discrete spectrum are eigenvalues with finite geometric multiplicity. Moreover, if the algebraic multiplicity of an eigenvalue in  $\sigma_d(A)$  is finite, then by definition the kernel chain stabilizes and so does the range chain (Problem 7.12). Hence by Lemma 7.3 we have  $X = \text{Ker}((A - \alpha)^n) \dot{+} \text{Ran}((A - \alpha)^n)$ , where  $n$  is the index of  $\alpha$ . These spaces are invariant and  $\text{Ker}((A - \alpha)^n)$  is still finite dimensional. With respect to this decomposition  $A$  has a simple block structure with the first block  $A_0 : \text{Ker}((A - \alpha)^n) \rightarrow \text{Ker}((A - \alpha)^n)$  such that  $A_0 - \alpha$  is nilpotent and the second block  $A_1 : \text{Ran}((A - \alpha)^n) \rightarrow \text{Ran}((A - \alpha)^n)$  such that  $\alpha \in \rho(A_1)$ . Hence for sufficiently small  $\varepsilon > 0$  we will have  $\alpha + \varepsilon \in \rho(A_0)$  and  $\alpha + \varepsilon \in \rho(A_1)$  implying  $\alpha + \varepsilon \in \rho(A)$  such that  $\alpha$  is an isolated point of the spectrum. This happens for example if  $A$  is self-adjoint (in which case  $n = 1$ ). However, in the general case,  $\sigma_d(A)$  could contain much more than just isolated eigenvalues with finite algebraic multiplicity as the following example shows.

**Example 7.16.** Let  $X = \ell^2(\mathbb{N}) \oplus \ell^2(\mathbb{N})$  with  $A = L \oplus R$ , where  $L, R$  are the left, right shift, respectively. Explicitly,  $A(x, y) = ((x_2, x_3, \dots), (0, y_1, y_2, \dots))$ . Then  $\sigma(A) = \sigma(L) \cup \sigma(R) = \bar{B}_1(0)$  and  $\sigma_p(A) = \sigma_p(L) \cup \sigma_p(R) = B_1(0)$ . Moreover, note that  $A \in \Phi_0$  and that 0 is an eigenvalue of infinite algebraic multiplicity.

Now consider the rank-one operator  $K(x, y) := ((0, 0, \dots), (x_1, 0, \dots))$  such that  $(A + K)(x, y) = ((x_2, x_3, \dots), (x_1, y_1, y_2, \dots))$ . Then  $A + K$  is unitary (note that this is essentially a two-sided shift) and hence  $\sigma(A + K) \subseteq \partial B_1(0)$ . Consequently  $\sigma_{ess}(A) \subseteq \partial B_1(0)$  and  $\sigma_d(A) \subseteq B_1(0)$  which shows  $\sigma_d(A) = B_1(0)$  and  $\sigma_{ess}(A) = \partial B_1(0)$ .  $\diamond$

It is important to emphasize, that Weyl's theorem makes it possible to determine these spectra even in nontrivial situations. This makes the splitting into essential and discrete spectrum much more versatile than the splitting into point, continuous, and residual spectrum.

**Problem\* 7.9.** Suppose  $X \xrightarrow{A} Y \xrightarrow{B} Z$  is exact. Show that if  $X$  and  $Z$  are finite dimensional, so is  $Y$ .

**Problem\* 7.10.** Let  $X_j$  be finite dimensional vector spaces and suppose

$$0 \longrightarrow X_1 \xrightarrow{A_1} X_2 \xrightarrow{A_2} X_3 \cdots X_{n-1} \xrightarrow{A_{n-1}} X_n \longrightarrow 0$$

is exact. Show that

$$\sum_{j=1}^n (-1)^j \dim(X_j) = 0.$$

(Hint: Rank-nullity theorem.)

**Problem 7.11.** Let  $A \in \Phi(X, Y)$  with a corresponding parametrix  $B_1, B_2 \in \mathcal{L}(Y, X)$ . Set  $K_1 := \mathbb{I} - B_1 A \in \mathcal{K}(X)$ ,  $K_2 := \mathbb{I} - A B_2 \in \mathcal{K}(Y)$  and show

$$B_1 - B = B_1 Q - K_1 B \in \mathcal{K}(Y, X), \quad B_2 - B = P B_2 - B K_2 \in \mathcal{K}(Y, X).$$

Hence a parametrix is unique up to compact operators. Moreover,  $B_1, B_2 \in \Phi(Y, X)$ .

**Problem\* 7.12.** Suppose  $A \in \Phi(X)$ . If the kernel chain stabilizes then  $\text{ind}(A) \leq 0$ . If the range chain stabilizes then  $\text{ind}(A) \geq 0$ . Moreover, if  $A \in \Phi_0(X)$ , then the kernel chain stabilizes if and only if the range chain stabilizes.

### 7.3. The Gelfand representation theorem

In this section we look at an alternative approach to the spectral theorem by trying to find a canonical representation for a Banach algebra. The idea is as follows: Given the Banach algebra  $C[a, b]$  we have a one-to-one correspondence between points  $x_0 \in [a, b]$  and point evaluations  $m_{x_0}(f) = f(x_0)$ . These point evaluations are linear functionals which at the same time preserve multiplication. In fact, we will see that these are the only (non-trivial) multiplicative functionals and hence we also have a one-to-one correspondence between points in  $[a, b]$  and multiplicative functionals. Now  $m_{x_0}(f) = f(x_0)$  says that the action of a multiplicative functional on a function is the same as the action of the function on a point. So for a general algebra  $X$  we can try to go the other way: Consider the multiplicative functionals  $m$  as points and the elements  $x \in X$  as functions acting on these points (the value of this function being  $m(x)$ ). This gives a map, the Gelfand representation, from  $X$  into an algebra of functions.

A nonzero algebra homeomorphism  $m : X \rightarrow \mathbb{C}$  will be called a **multiplicative linear functional** or **character**:

$$m(xy) = m(x)m(y), \quad m(e) = 1. \quad (7.31)$$

Note that the last equation comes for free from multiplicativity since  $m$  is nontrivial. Moreover, there is no need to require that  $m$  is continuous as this will also follow automatically (cf. Lemma 7.18 below).

As we will see, they are closely related to **ideals**, that is linear subspaces  $I$  of  $X$  for which  $a \in I$ ,  $x \in X$  implies  $ax \in I$  and  $xa \in I$ . An ideal is called **proper** if it is not equal to  $X$  and it is called **maximal** if it is not contained in any other proper ideal.

**Example 7.17.** Let  $X := C[a, b]$  be the continuous functions over some compact interval. Then for fixed  $x_0 \in [a, b]$ , the linear functional  $m_{x_0}(f) := f(x_0)$  is multiplicative. Moreover, its kernel  $\text{Ker}(m_{x_0}) = \{f \in C[a, b] \mid f(x_0) = 0\}$  is a maximal ideal (we will prove this in more generality in Lemma 7.18 below).  $\diamond$

**Example 7.18.** Let  $X$  be a Banach space. Then the compact operators are a closed ideal in  $\mathcal{L}(X)$  (cf. Theorem 3.1).  $\diamond$

We first collect a few elementary properties of ideals.

**Lemma 7.17.** *Let  $X$  be a unital Banach algebra.*

- (i) *A proper ideal can never contain an invertible element.*
- (ii) *If  $X$  is commutative, every non-invertible element is contained in a proper ideal.*
- (iii) *The closure of a (proper) ideal is again a (proper) ideal.*
- (iv) *Maximal ideals are closed.*
- (v) *Every proper ideal is contained in a maximal ideal.*

**Proof.** (i). Let  $I$  be a proper ideal. If  $x \in I$  is invertible then  $y = x(x^{-1}y) \in I$  shows  $I = X$ . (ii). Consider the ideal  $xX = \{xy \mid y \in X\}$ . Then  $xX = X$  if and only if there is some  $y \in X$  with  $xy = e$ , that is,  $y = x^{-1}$ . (iii) and (iv). That the closure of an ideal is again an ideal follows from continuity of the product. Indeed, for  $a \in \bar{I}$  choose a sequence  $a_n \in I$  converging to  $a$ . Then  $xa_n \in I \rightarrow xa \in \bar{I}$  as well as  $a_nx \in I \rightarrow ax \in \bar{I}$ . Moreover, note that by Lemma 5.1 all elements in the ball  $B_1(e)$  are invertible and hence every proper ideal must be contained in the closed set  $X \setminus B_1(e)$ . So the closure of a proper ideal is proper and any maximal ideal must be closed. (v). To see that every ideal  $I$  is contained in a maximal ideal, consider the family of proper ideals containing  $I$  ordered by inclusion. Then, since any union of a chain of proper ideals is again a proper ideal (that the union is again an ideal is straightforward, to see that it is proper note that it does not contain  $B_1(e)$ ). Consequently, Zorn's lemma implies existence of a maximal element.  $\square$

Note that if  $I$  is a closed ideal, then the quotient space  $X/I$  (cf. Lemma 1.18) is again a Banach algebra if we define

$$[x][y] = [xy]. \quad (7.32)$$

Indeed  $(x + I)(y + I) = xy + I$  and hence the multiplication is well-defined and inherits the distributive and associative laws from  $X$ . Also  $[e]$  is an identity. Finally,

$$\begin{aligned} \|[xy]\| &= \inf_{a \in I} \|xy + a\| = \inf_{b, c \in I} \|(x + b)(y + c)\| \leq \inf_{b \in I} \|x + b\| \inf_{c \in I} \|y + c\| \\ &= \|[x]\| \|[y]\|. \end{aligned} \quad (7.33)$$

In particular, the quotient map  $\pi : X \rightarrow X/I$  is a Banach algebra homomorphism.

**Example 7.19.** Consider the Banach algebra  $\mathcal{L}(X)$  together with the ideal of compact operators  $\mathcal{K}(X)$ . Then the Banach algebra  $\mathcal{L}(X)/\mathcal{K}(X)$  is known as the **Calkin algebra**. Atkinson's theorem (Theorem 7.14) says that the invertible elements in the Calkin algebra are precisely the images of the Fredholm operators.  $\diamond$

**Lemma 7.18.** *Let  $X$  be a unital Banach algebra and  $m$  a character. Then  $\text{Ker}(m)$  is a maximal ideal and  $m$  is continuous with  $\|m\| = m(e) = 1$ .*

**Proof.** It is straightforward to check that  $\text{Ker}(m)$  is an ideal. Moreover, every  $x$  can be written as

$$x = m(x)e + y, \quad y \in \text{Ker}(m).$$

Let  $I$  be an ideal containing  $\text{Ker}(m)$ . If there is some  $x \in I \setminus \text{Ker}(m)$  then  $e = m(x)^{-1}(x - y) \in I$  and thus  $\text{Ker}(m)$  is maximal. Since maximal ideals are closed by the previous lemma, we conclude that  $m$  is continuous by Problem 1.40. Clearly  $\|m\| \geq m(e) = 1$ . Conversely, suppose we can find some  $x \in X$  with  $\|x\| < 1$  and  $m(x) = 1$ . Consequently  $\|x^n\| \leq \|x\|^n \rightarrow 0$  contradicting  $m(x^n) = m(x)^n = 1$ .  $\square$

In a commutative algebra the other direction is also true.

**Lemma 7.19.** *In a commutative unital Banach algebra the characters and maximal ideals are in one-to-one correspondence.*

**Proof.** We have already seen that for a character  $m$  there is a corresponding maximal ideal  $\text{Ker}(m)$ . Conversely, let  $I$  be a maximal ideal and consider the quotient map  $\pi : X \rightarrow X/I$ . We first claim that every nontrivial element in  $X/I$  is invertible. To this end suppose  $[x_0] \neq [0]$  were not invertible. Then  $J = [x_0]X/I$  is a proper ideal (if it would contain the identity,  $X/I$  would contain an inverse of  $[x_0]$ ). Moreover,  $I' = \{y \in X \mid [y] \in J\}$  is a proper ideal of  $X$  (since  $e \in I'$  would imply  $[e] \in J$ ) which contains  $I$  (since  $[x] = [0] \in J$  for  $x \in I$ ) but is strictly larger as  $x_0 \in I' \setminus I$ . This contradicts maximality and hence by the Gelfand–Mazur theorem (Theorem 5.4), every element of  $X/I$  is of the form  $\alpha[e]$ . If  $h : X/I \rightarrow \mathbb{C}$ ,  $h(\alpha[e]) \mapsto \alpha$  is the corresponding algebra isomorphism, then  $m = h \circ \pi$  is a character with  $\text{Ker}(m) = I$ .  $\square$

Now we continue with the following observation: For fixed  $x \in X$  we get a map  $X^* \rightarrow \mathbb{C}$ ,  $\ell \mapsto \ell(x)$ . Moreover, if we equip  $X^*$  with the weak-\* topology, then this map will be continuous (by the very definition of the weak-\* topology). So we have a map  $X \rightarrow C(X^*)$  and restricting this map to the set of all characters  $\mathcal{M} \subseteq X^*$  (equipped with the relative topology of the weak-\* topology) it is known as the **Gelfand transform**:

$$\Gamma : X \rightarrow C(\mathcal{M}), \quad x \mapsto \hat{x}(m) := m(x). \quad (7.34)$$

**Theorem 7.20** (Gelfand representation theorem). *Let  $X$  be a unital Banach algebra. Then the set of all characters  $\mathcal{M} \subseteq X^*$  is a compact Hausdorff space with respect to the weak-\* topology and the Gelfand transform is a continuous algebra homomorphism with  $\hat{e} = 1$ .*

Moreover,  $\hat{x}(\mathcal{M}) \subseteq \sigma(x)$  and hence  $\|\hat{x}\|_\infty \leq r(x) \leq \|x\|$  where  $r(x)$  is the spectral radius of  $x$ . If  $X$  is commutative then  $\hat{x}(\mathcal{M}) = \sigma(x)$  and hence  $\|\hat{x}\|_\infty = r(x)$ .

**Proof.** As pointed out before, for fixed  $x, y \in X$  the map  $X^* \rightarrow \mathbb{C}^3$ ,  $\ell \mapsto (\ell(x), \ell(y), \ell(xy))$  is continuous and so is the map  $X^* \rightarrow \mathbb{C}$ ,  $\ell \mapsto \ell(x)\ell(y) - \ell(xy)$  as a composition of continuous maps. Hence the kernel of this map  $M_{x,y} = \{\ell \in X^* | \ell(x)\ell(y) = \ell(xy)\}$  is weak-\* closed and so is  $\mathcal{M} = M_0 \cap \bigcap_{x,y \in X} M_{x,y}$  where  $M_0 = \{\ell \in X^* | \ell(e) = 1\}$ . So  $\mathcal{M}$  is a weak-\* closed subset of the unit ball in  $X^*$  and the first claim follows from the Banach–Alaoglu theorem (Theorem 6.10).

Next  $(x+y)^\wedge(m) = m(x+y) = m(x)+m(y) = \hat{x}(m)+\hat{y}(m)$ ,  $(xy)^\wedge(m) = m(xy) = m(x)m(y) = \hat{x}(m)\hat{y}(m)$ , and  $\hat{e}(m) = m(e) = 1$  shows that the Gelfand transform is an algebra homomorphism.

Moreover, if  $m(x) = \alpha$  then  $x - \alpha \in \text{Ker}(m)$  implying that  $x - \alpha$  is not invertible (as maximal ideals cannot contain invertible elements), that is  $\alpha \in \sigma(x)$ . Conversely, if  $X$  is commutative and  $\alpha \in \sigma(x)$ , then  $x - \alpha$  is not invertible and hence contained in some maximal ideal, which in turn is the kernel of some character  $m$ . Whence  $m(x - \alpha) = 0$ , that is  $m(x) = \alpha$  for some  $m$ .  $\square$

Of course this raises the question about injectivity or surjectivity of the Gelfand transform. Clearly

$$x \in \text{Ker}(\Gamma) \Leftrightarrow x \in \bigcap_{m \in \mathcal{M}} \text{Ker}(m) \quad (7.35)$$

and it can only be injective if  $X$  is commutative (if  $xy \neq yx$ , then  $xy - yx \in \bigcap_{m \in \mathcal{M}} \text{Ker}(m)$ ). In this case Lemma 7.19 implies

$$x \in \text{Ker}(\Gamma) \Leftrightarrow x \in \text{Rad}(X) := \bigcap_{I \text{ maximal ideal}} I, \quad (7.36)$$

where  $\text{Rad}(X)$  is known as the **Jacobson radical** of  $X$  and a Banach algebra is called **semi-simple** if the Jacobson radical is zero. So to put this result to use, one needs to understand the set of characters, or equivalently, the set of maximal ideals. Two examples where this can be done are given below. The first one is not very surprising.

**Example 7.20.** If we start with a compact Hausdorff space  $K$  and consider  $C(K)$  we get nothing new. Indeed, first of all notice that the map  $K \rightarrow \mathcal{M}$ ,  $x_0 \mapsto m_{x_0}$  which assigns each  $x_0$  the corresponding point evaluation  $m_{x_0}(f) = f(x_0)$  is injective and continuous. Hence, by compactness of  $K$ , it will be a homeomorphism once we establish surjectivity (Corollary B.17). To this end we will show that all maximal ideals are of the form  $I = \text{Ker}(m_{x_0})$  for some  $x_0 \in K$ . So let  $I$  be an ideal and suppose there is no point where all functions vanish. Then for every  $x \in K$  there is a ball  $B_{r(x)}(x)$  and a function  $f_x \in C(K)$  such that  $|f_x(y)| \geq 1$  for  $y \in B_{r(x)}(x)$ . By compactness finitely many of these balls will cover  $K$ . Now consider  $f = \sum_j f_{x_j}^* f_{x_j} \in I$ . Then  $f \geq 1$  and hence  $f$  is invertible, that is  $I = C(K)$ . Thus maximal ideals are of the form  $I_{x_0} = \{f \in C(K) | f(x_0) = 0\}$  which are precisely the kernels of the characters  $m_{x_0}(f) = f(x_0)$ . Thus  $\mathcal{M} \cong K$  as well as  $\hat{f} \cong f$ .  $\diamond$

**Example 7.21.** Consider the Wiener algebra  $\mathcal{A}$  of all periodic continuous functions which have an absolutely convergent Fourier series. As in the previous example it suffices to show that all maximal ideals are of the form  $I_{x_0} = \{f \in \mathcal{A} | f(x_0) = 0\}$ . To see this set  $e_k(x) = e^{ikx}$  and note  $\|e_k\|_{\mathcal{A}} = 1$ . Hence for every character  $m(e_k) = m(e_1)^k$  and  $|m(e_k)| \leq 1$ . Since the last claim holds for both positive and negative  $k$ , we conclude  $|m(e_k)| = 1$  and thus there is some  $x_0 \in [-\pi, \pi]$  with  $m(e_k) = e^{ikx_0}$ . Consequently  $m(f) = f(x_0)$  and point evaluations are the only characters. Equivalently, every maximal ideal is of the form  $\text{Ker}(m_{x_0}) = I_{x_0}$ .

So, as in the previous example,  $\mathcal{M} \cong [-\pi, \pi]$  (with  $-\pi$  and  $\pi$  identified) as well as  $\hat{f} \cong f$ . Moreover, the Gelfand transform is injective but not surjective since there are continuous functions whose Fourier series are not absolutely convergent. Incidentally this also shows that the Wiener algebra is *no*  $C^*$  algebra (despite the fact that we have a natural conjugation which satisfies  $\|f^*\|_{\mathcal{A}} = \|f\|_{\mathcal{A}}$  — this again underlines the special role of (5.32)) as the Gelfand–Naimark theorem below will show that the Gelfand transform is bijective for commutative  $C^*$  algebras.  $\diamond$

Since  $0 \notin \sigma(x)$  implies that  $x$  is invertible, the Gelfand representation theorem also contains a useful criterion for invertibility.

**Corollary 7.21.** *In a commutative unital Banach algebra an element  $x$  is invertible if and only if  $m(x) \neq 0$  for all characters  $m$ .*



And applying this to the last example we get the following famous theorem of Wiener:

**Theorem 7.22** (Wiener). *Suppose  $f \in C_{per}[-\pi, \pi]$  has an absolutely convergent Fourier series and does not vanish on  $[-\pi, \pi]$ . Then the function  $\frac{1}{f}$  also has an absolutely convergent Fourier series.*

If we turn to a commutative  $C^*$  algebra, the situation further simplifies. First of all note that characters respect the additional operation automatically.

**Lemma 7.23.** *If  $X$  is a unital  $C^*$  algebra, then every character satisfies  $m(x^*) = m(x)^*$ . In particular, the Gelfand transform is a continuous  $*$ -algebra homomorphism with  $\hat{e} = 1$  in this case.*

**Proof.** If  $x$  is self-adjoint then  $\sigma(x) \subseteq \mathbb{R}$  (Lemma 5.9) and hence  $m(x) \in \mathbb{R}$  by the Gelfand representation theorem. Now for general  $x$  we can write  $x = a + ib$  with  $a := \frac{x+x^*}{2}$  and  $b := \frac{x-x^*}{2i}$  self-adjoint implying

$$m(x^*) = m(a - ib) = m(a) - im(b) = (m(a) + im(b))^* = m(x)^*.$$

Consequently the Gelfand transform preserves the involution:  $(x^*)^\wedge(m) = m(x^*) = m(x)^* = \hat{x}^*(m)$ .  $\square$

**Theorem 7.24** (Gelfand–Naimark). *Suppose  $X$  is a unital commutative  $C^*$  algebra. Then the Gelfand transform is an isometric isomorphism between  $C^*$  algebras.*

**Proof.** As in a commutative  $C^*$  algebra all elements are normal, Lemma 5.8 implies that the Gelfand transform is isometric. Moreover, by the previous lemma the image of  $X$  under the Gelfand transform is a closed  $*$ -subalgebra which contains  $\hat{e} \equiv 1$  and separates points (if  $\hat{x}(m_1) = \hat{x}(m_2)$  for all  $x \in X$  we have  $m_1 = m_2$ ). Hence it must be all of  $C(\mathcal{M})$  by the Stone–Weierstraß theorem (Theorem B.42).  $\square$

The first moral from this theorem is that from an abstract point of view there is only one commutative  $C^*$  algebra, namely  $C(K)$  with  $K$  some compact Hausdorff space. Moreover, the formulation also very much reassembles the spectral theorem and in fact, we can derive the spectral theorem by applying it to  $C^*(x)$ , the  $C^*$  algebra generated by  $x$  (cf. (5.35)). This will even give us the more general version for normal elements. As a preparation we show that it makes no difference whether we compute the spectrum in  $X$  or in  $C^*(x)$ .

**Lemma 7.25** (Spectral permanence). *Let  $X$  be a  $C^*$  algebra and  $Y \subseteq X$  a closed  $*$ -subalgebra containing the identity. Then  $\sigma(y) = \sigma_Y(y)$  for every  $y \in Y$ , where  $\sigma_Y(y)$  denotes the spectrum computed in  $Y$ .*

**Proof.** Clearly we have  $\sigma(y) \subseteq \sigma_Y(y)$  and it remains to establish the reverse inclusion. If  $(y - \alpha)$  has an inverse in  $X$ , then the same is true for  $(y - \alpha)^*(y - \alpha)$ . But the last operator is self-adjoint and hence has real spectrum in  $Y$ . Thus  $((y - \alpha)^*(y - \alpha) + \frac{1}{n})^{-1} \in Y$  and letting  $n \rightarrow \infty$  shows  $((y - \alpha)^*(y - \alpha))^{-1} \in Y$  since taking the inverse is continuous and  $Y$  is closed. Whence  $(y - \alpha)^{-1} = (y - \alpha^*)((y - \alpha)^*(y - \alpha))^{-1} \in Y$ .  $\square$

Now we can show

**Theorem 7.26** (Spectral theorem). *If  $X$  is a  $C^*$  algebra and  $x$  is normal, then there is an isometric isomorphism  $\Phi : C(\sigma(x)) \rightarrow C^*(x)$  such that  $f(t) = t$  maps to  $\Phi(t) := x$  and  $f(t) := 1$  maps to  $\Phi(1) = e$ .*

Moreover, for every  $f \in C(\sigma(x))$  we have

$$\sigma(f(x)) = f(\sigma(x)), \quad (7.37)$$

where  $f(x) := \Phi(f)$ .

**Proof.** Given a normal element  $x \in X$  we want to apply the Gelfand–Naimark theorem in  $C^*(x)$ . By our lemma we have  $\sigma(x) = \sigma_{C^*(x)}(x)$ . We first show that we can identify  $\mathcal{M}$  with  $\sigma(x)$ . By the Gelfand representation theorem (applied in  $C^*(x)$ ),  $\hat{x} : \mathcal{M} \rightarrow \sigma(x)$  is continuous and surjective. Moreover, if for given  $m_1, m_2 \in \mathcal{M}$  we have  $\hat{x}(m_1) = m_1(x) = m_2(x) = \hat{x}(m_2)$ , then

$$m_1(p(x, x^*)) = p(m_1(x), m_1(x)^*) = p(m_2(x), m_2(x)^*) = m_2(p(x, x^*))$$

for any polynomial  $p : \mathbb{C}^2 \rightarrow \mathbb{C}$  and hence  $m_1(y) = m_2(y)$  for every  $y \in C^*(x)$  implying  $m_1 = m_2$ . Thus  $\hat{x}$  is injective and hence a homeomorphism as  $\mathcal{M}$  is compact. Thus we have an isometric isomorphism

$$\Psi : C(\sigma(x)) \rightarrow C(\mathcal{M}), \quad f \mapsto f \circ \hat{x},$$

and the isometric isomorphism we are looking for is  $\Phi = \Gamma^{-1} \circ \Psi$ . Finally,  $\sigma(f(x)) = \sigma_{C^*(x)}(\Phi(f)) = \sigma_{C(\sigma(x))}(f) = f(\sigma(x))$ .  $\square$

**Example 7.22.** Let  $X$  be a  $C^*$  algebra and  $x \in X$  normal. By the spectral theorem  $C^*(x)$  is isomorphic to  $C(\sigma(x))$ . Hence every  $y \in C^*(x)$  can be written as  $y = f(x)$  for some  $f \in C(\sigma(x))$  and every character is of the form  $m(y) = m(f(x)) = f(\alpha)$  for some  $\alpha \in \sigma(x)$ .  $\diamond$

**Problem 7.13.** Show that  $C^1[a, b]$  is a Banach algebra. What are the characters? Is it semi-simple?

**Problem 7.14.** Consider the subalgebra of the Wiener algebra consisting of all functions whose negative Fourier coefficients vanish. What are the characters? (Hint: Observe that these functions can be identified with holomorphic functions inside the unit disc with summable Taylor coefficients via  $f(z) = \sum_{k=0}^{\infty} \hat{f}_k z^k$  known as the **HARDY space**  $H^1$  of the disc.)



# Unbounded operators

## 8.1. Closed operators

Recall that an operator  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  between Banach spaces is called closed if its graph  $\Gamma(A) \subset X \oplus Y$  is closed. In this section we want to investigate which properties of bounded operators carry over to closed operators.

Being closed is the next option you have once an operator turns out to be unbounded. If  $A$  is closed, then  $x_n \rightarrow x$  does not guarantee you that  $Ax_n$  converges (like continuity would), but it at least guarantees that if  $Ax_n$  converges, it converges to the right thing, namely  $Ax$ :

- $A$  bounded (with  $\mathfrak{D}(A) = X$ ):  $x_n \rightarrow x$  implies  $Ax_n \rightarrow Ax$ .
- $A$  closed (with  $\mathfrak{D}(A) \subseteq X$ ):  $x_n \rightarrow x$ ,  $x_n \in \mathfrak{D}(A)$ , and  $Ax_n \rightarrow y$  implies  $x \in \mathfrak{D}(A)$  and  $y = Ax$ .

Please observe that the domain  $\mathfrak{D}(A)$  is an intrinsic part of the definition of  $A$  and that we cannot assume  $\mathfrak{D}(A) = X$  unless  $A$  is bounded (which is the content of the closed graph theorem, Theorem 4.7). Hence, if we want an unbounded operator to be closed, we have to live with domains. We will however typically assume that  $\mathfrak{D}(A)$  is dense and set

$$\mathcal{C}(X, Y) := \{A : \mathfrak{D}(A) \subseteq X \rightarrow Y \mid A \text{ is densely defined and closed}\}. \quad (8.1)$$

One writes  $B \subseteq A$  if  $\mathfrak{D}(B) \subseteq \mathfrak{D}(A)$  and  $Bx = Ax$  for  $x \in \mathfrak{D}(B)$ . In this case  $A$  is called an **extension** of  $B$ .

**Example 8.1.** Two operators having the same prescription but different domains are different. For example

$$\mathfrak{D}(A) = C^1[0, 1], \quad Af = f'$$

and

$$\mathfrak{D}(B) = \{f \in C^1[0, 1] \mid f(0) = f(1) = 0\}, \quad Bf = f'$$

are two different operators in  $X := C[0, 1]$ . Clearly  $A$  is an extension of  $B$ . Moreover, both are closed since  $f_n \rightarrow f$  and  $f'_n \rightarrow g$  implies that  $f$  is differentiable and  $f' = g$ . Note that  $A$  is densely defined while  $B$  is not.  $\diamond$

Be aware that taking sums or products of unbounded operators is tricky due to the possible different domains. Indeed, if  $A$  and  $B$  are two operators between Banach spaces  $X$  and  $Y$ , so is  $A + B$  defined on  $\mathfrak{D}(A + B) := \mathfrak{D}(A) \cap \mathfrak{D}(B)$ . The problem is that  $\mathfrak{D}(A + B)$  might contain nothing more than zero. Similarly, if  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  and  $B : \mathfrak{D}(B) \subseteq Y \rightarrow Z$ , then the composition  $BA$  is defined on  $\mathfrak{D}(BA) := \{x \in \mathfrak{D}(A) \mid Ax \in \mathfrak{D}(B)\}$ .

**Example 8.2.** Consider  $X := C[0, 1]$ . Let  $M$  be the the subspace of trigonometric polynomials and  $N$  be the subspace of piecewise linear functions. Then both  $M$  and  $N$  are dense with  $M \cap N = \{0\}$ .  $\diamond$

If an operator is not closed, you can try to take the closure of its graph, to obtain a closed operator. If  $A$  is bounded this always works (which is just the content of Theorem 1.16). However, in general, the closure of the graph might not be the graph of an operator as we might pick up points  $(x, y_1), (x, y_2) \in \overline{\Gamma(A)}$  with  $y_1 \neq y_2$ . Since  $\overline{\Gamma(A)}$  is a subspace, we also have  $(x, y_2) - (x, y_1) = (0, y_2 - y_1) \in \overline{\Gamma(A)}$  in this case and thus  $\overline{\Gamma(A)}$  is the graph of some operator if and only if

$$\overline{\Gamma(A)} \cap \{(0, y) \mid y \in Y\} = \{(0, 0)\}. \quad (8.2)$$

If this is the case,  $A$  is called **closable** and the operator  $\overline{A}$  associated with  $\overline{\Gamma(A)}$  is called the **closure** of  $A$ . Any linear subset  $\mathfrak{D} \subseteq \mathfrak{D}(A)$  with the property that  $A$  restricted to  $\mathfrak{D}$  has the same closure,  $\overline{A|_{\mathfrak{D}}} = \overline{A}$ , is called a **core** for  $A$ .

In particular,  $A$  is closable if and only if  $x_n \rightarrow 0$  and  $Ax_n \rightarrow y$  implies  $y = 0$ . In this case

$$\begin{aligned} \mathfrak{D}(\overline{A}) &= \{x \in X \mid \exists x_n \in \mathfrak{D}(A), y \in Y : x_n \rightarrow x \text{ and } Ax_n \rightarrow y\}, \\ \overline{A}x &= y. \end{aligned} \quad (8.3)$$

There is yet another way of defining the closure: Define the **graph norm** associated with  $A$  by

$$\|x\|_A := \|x\|_X + \|Ax\|_Y, \quad x \in \mathfrak{D}(A). \quad (8.4)$$

Since we have  $\|Ax\| \leq \|x\|_A$  we see that  $A : \mathfrak{D}(A) \rightarrow Y$  is bounded with norm at most one. Thus far  $(\mathfrak{D}(A), \|\cdot\|_A)$  is a normed space and it suggests itself to consider its completion  $X_A$ . Then one can check (Problem 8.5) that  $X_A$  can be regarded as a subset of  $X$  if and only if  $A$  is closable. In this case the completion can be identified with  $\mathfrak{D}(\overline{A})$  and the closure of  $A$  in  $X$

coincides with the extension from Theorem 1.16 of  $A$  in  $X_A$ . In particular,  $A$  is closed if and only if  $(\mathfrak{D}(A), \|\cdot\|_A)$  is complete.

**Example 8.3.** Consider the multiplication operator  $A$  in  $\ell^p(\mathbb{N})$  defined by  $Aa_j := ja_j$  on  $\mathfrak{D}(A) := \ell_c(\mathbb{N})$ , where  $\ell_c(\mathbb{N})$  denotes the sequences with compact support (i.e., finitely many nonzero terms).

(i).  $A$  is closable. In fact, if  $a^n \rightarrow 0$  and  $Aa^n \rightarrow b$  then we have  $a_j^n \rightarrow 0$  and thus  $ja_j^n \rightarrow 0 = b_j$  for any  $j \in \mathbb{N}$ .

(ii). The closure of  $A$  is given by

$$\mathfrak{D}(\bar{A}) = \begin{cases} \{a \in \ell^p(\mathbb{N}) | (ja_j)_{j=1}^\infty \in \ell^p(\mathbb{N})\}, & 1 \leq p < \infty, \\ \{a \in c_0(\mathbb{N}) | (ja_j)_{j=1}^\infty \in c_0(\mathbb{N})\}, & p = \infty, \end{cases}$$

and  $\bar{A}a_j = ja_j$ . In fact, if  $a^n \rightarrow a$  and  $Aa^n \rightarrow b$  then we have  $a_j^n \rightarrow a_j$  and  $ja_j^n \rightarrow b_j$  for any  $j \in \mathbb{N}$  and thus  $b_j = ja_j$  for any  $j \in \mathbb{N}$ . In particular,  $(ja_j)_{j=1}^\infty = (b_j)_{j=1}^\infty \in \ell^p(\mathbb{N})$  ( $c_0(\mathbb{N})$  if  $p = \infty$ ). Conversely, suppose  $(ja_j)_{j=1}^\infty \in \ell^p(\mathbb{N})$  ( $c_0(\mathbb{N})$  if  $p = \infty$ ) and consider

$$a_j^n := \begin{cases} a_j, & j \leq n, \\ 0, & j > n. \end{cases}$$

Then  $a^n \rightarrow a$  and  $Aa^n \rightarrow (ja_j)_{j=1}^\infty$ .

(iii). Extending the basis vectors  $\{\delta^n\}_{n \in \mathbb{N}}$  to a Hamel basis (Problem 1.7) and setting  $Aa = 0$  for every new element from this Hamel basis we obtain a (still unbounded) operator which is everywhere defined. However, this extension cannot be closed!  $\diamond$

**Example 8.4.** Here is a simple example of a nonclosable operator: Let  $X = Y := \ell^2(\mathbb{N})$  and consider  $Ba := (\sum_{j=1}^\infty a_j)\delta^1$  defined on  $\ell^1(\mathbb{N}) \subset \ell^2(\mathbb{N})$ . Let  $a_j^n := \frac{1}{n}$  for  $1 \leq j \leq n$  and  $a_j^n := 0$  for  $j > n$ . Then  $\|a^n\|_2 = \frac{1}{\sqrt{n}}$  implying  $a^n \rightarrow 0$  but  $Ba^n = \delta^1 \not\rightarrow 0$ .  $\diamond$

**Example 8.5** (Sobolev spaces). Let  $X := L^p(0, 1)$ ,  $1 \leq p < \infty$ , and consider  $Af := f'$  on  $\mathfrak{D}(A) := C^1[0, 1]$ . Then it is not hard to see that  $A$  is *not* closed (take a sequence  $g_n$  of continuous functions which converges in  $L^p$  to a non-continuous function, cf. Example 1.10, and consider its primitive  $f_n(x) = \int_0^x g_n(y)dy$ ). It is however closable. To see this suppose  $f_n \rightarrow 0$  and  $f'_n \rightarrow g$  in  $L^p$ . Then  $f_n(0) = f_n(x) - \int_0^x f'_n(y)dy \rightarrow -\int_0^x g(y)dy$ . But a sequence of constant functions can only have a constant function as a limit implying  $g \equiv 0$  as required. The domain of the closure is the **Sobolev space**  $W^{1,p}(0, 1)$  and this is one way of defining Sobolev spaces. In particular,  $W^{1,p}(0, 1)$  is a Banach space when equipped with the graph norm. In this context one chooses the  $p$ -norm for the direct sum  $X \oplus_p X$  such that the graph norm reads

$$\|f\|_{1,p} := (\|f\|_p^p + \|f'\|_p^p)^{1/p}.$$

In particular, since the case  $p = 2$  will lead to a Hilbert space usually denoted by  $H^1(0, 1) := W^{1,2}(0, 1)$ .

Note, that for  $f \in W^{1,p}(0, 1)$  there is some  $g \in L^p(0, 1)$  for which  $f(x) = f(0) + \int_0^x g(y)dy$  holds and one writes  $f' := g$  in this case. Note that  $W^{1,p}(0, 1) \subset C(0, 1)$  (Problem 8.6) and in the case  $p = 1$  these functions are also known as the absolutely continuous functions  $AC[0, 1] := W^{1,1}(0, 1)$ .  $\diamond$

**Example 8.6.** Another example are point evaluations in  $L^p(0, 1)$ ,  $1 \leq p < \infty$ : Let  $x_0 \in [0, 1]$  and consider  $\ell_{x_0} : \mathfrak{D}(\ell_{x_0}) \rightarrow \mathbb{C}$ ,  $f \mapsto f(x_0)$  defined on  $\mathfrak{D}(\ell_{x_0}) := C[0, 1] \subseteq L^p(0, 1)$ . Then  $f_n(x) := \max(0, 1 - n|x - x_0|)$  satisfies  $f_n \rightarrow 0$  but  $\ell_{x_0}(f_n) = 1$ . In fact, a linear functional is closable if and only if it is bounded (Problem 8.2).  $\diamond$

For the closure of sums and products see Problem 8.11 and Problem 8.12, respectively.

Given a subset  $\Gamma \subseteq X \oplus Y$  we can define

$$\Gamma^{-1} := \{(y, x) | (x, y) \in \Gamma\} \subseteq Y \oplus X. \quad (8.5)$$

In particular, applying this to the graph of an operator  $A$ , we will obtain the graph of its inverse (provided  $A$  is invertible). Hence we see that an invertible operator is closed if and only if its inverse is closed. Slightly more general, we have:

**Lemma 8.1.** *Suppose  $A$  is closable and  $\bar{A}$  is injective. Then  $\bar{A}^{-1} = \overline{A^{-1}}$ .*

**Proof.** This follows from  $\Gamma(A^{-1}) = \Gamma^{-1}(A)$  and

$$\overline{\Gamma(A^{-1})} = \overline{\Gamma^{-1}(A)} = \overline{\Gamma(A)}^{-1} = \Gamma^{-1}(\bar{A}) = \Gamma(\bar{A}^{-1}). \quad \square$$

Note that  $A$  injective does not imply  $\bar{A}$  injective in general.

**Example 8.7.** Let  $P_M$  be the projection in  $\ell^2(\mathbb{N})$  on  $M := \{b\}^\perp$ , where  $b := (2^{-j/2})_{j=1}^\infty$ . Explicitly we have  $P_M a = a - \langle b, a \rangle b$ . Then  $P_M$  restricted to the space of sequences with finitely many nonzero terms is injective, but its closure is not.  $\diamond$

As a consequence of the closed graph theorem we obtain:

**Corollary 8.2.** *Suppose  $A \in \mathcal{C}(X, Y)$  is injective. Then  $A^{-1}$  defined on  $\mathfrak{D}(A^{-1}) = \text{Ran}(A)$  is closed. Moreover, in this case  $\text{Ran}(A)$  is closed if and only if  $A^{-1}$  is bounded.*

**Example 8.8.** Note that in Example 8.3 the inverse  $B := A^{-1}$  is the bounded operator  $Ba_j := \frac{1}{j}a_j$  defined on  $\mathfrak{D}(B) = \text{Ran}(A) = \ell_c(\mathbb{N})$ . Hence,

the closure is  $\overline{B}a_j := \frac{1}{j}a_j$  with

$$\mathfrak{D}(\overline{B}) = \overline{\text{Ran}(A)} = \begin{cases} \ell^p(\mathbb{N}), & 1 \leq p < \infty, \\ c_0(\mathbb{N}), & p = \infty. \end{cases}$$

Hence  $\overline{B} \in \mathcal{L}(X)$  for  $1 \leq p < \infty$  but not for  $p = \infty$ .  $\diamond$

The question when  $\text{Ran}(A)$  is closed plays an important role when investigating solvability of the equation  $Ax = y$  and the last part gives us a convenient criterion. Moreover, note that  $A^{-1}$  is bounded if and only if there is some  $c > 0$  such that

$$\|Ax\| \geq c\|x\|, \quad x \in \mathfrak{D}(A). \quad (8.6)$$

Indeed, this follows upon setting  $x = A^{-1}y$  in the above inequality which also shows that  $c = \|A^{-1}\|^{-1}$  is the best possible constant. Factoring out the kernel we even get a criterion for the general case. To this end we note:

**Lemma 8.3.** *Suppose  $A \in \mathcal{C}(X, Y)$ . Then  $\text{Ker}(A)$  is closed and the quotient space  $\tilde{X} := X/\text{Ker}(A)$  is a Banach space. Moreover,  $\tilde{A} : \mathfrak{D}(\tilde{A}) \rightarrow Y$  where  $\mathfrak{D}(\tilde{A}) = \mathfrak{D}(A)/\text{Ker}(A) \subseteq \tilde{X}$ , is a closed operator with  $\text{Ker}(\tilde{A}) = \{0\}$  and  $\text{Ran}(\tilde{A}) = \text{Ran}(A)$ .*

**Proof.** It is easy to check that  $\text{Ker}(A)$  is closed (Problem 8.1). Then  $\tilde{A}$  defined via  $\tilde{A}[x] := Ax$  is well defined and injective (cf. Problem 1.46). To see that  $\tilde{A}$  is closed we use  $\tilde{\pi} : X \times Y \rightarrow \tilde{X} \times Y$ ,  $(x, y) \mapsto ([x], y)$  which is bounded, surjective and hence open. Moreover,  $\tilde{\pi}(\Gamma(A)) = \Gamma(\tilde{A})$ . In fact, we even have  $(x, y) \in \Gamma(A)$  iff  $([x], y) \in \Gamma(\tilde{A})$  and thus  $\tilde{\pi}(X \times Y \setminus \Gamma(A)) = \tilde{X} \times Y \setminus \Gamma(\tilde{A})$  implying that  $\tilde{X} \times Y \setminus \Gamma(\tilde{A})$  is open.  $\square$

Applying Corollary 8.2 to  $\tilde{A}$  gives:

**Corollary 8.4.** *Suppose  $A \in \mathcal{C}(X, Y)$ . Then  $\text{Ran}(A)$  is closed if and only if*

$$\|Ax\| \geq c \text{dist}(x, \text{Ker}(A)), \quad x \in \mathfrak{D}(A), \quad (8.7)$$

for some  $c > 0$ . The sup over all possible  $c$  is known as the (reduced) **minimum modulus** of  $A$ .

As for bounded operators we can define the **adjoint operator** using  $(A'\ell)(x) := \ell(Ax)$ . Of course the linear functional  $A'\ell$  will only be defined on  $\mathfrak{D}(A)$  and it will be unbounded in general. Hence we define the domain  $\mathfrak{D}(A')$  as the set

$$\mathfrak{D}(A') := \{\ell \in Y^* \mid \ell \circ A \text{ is bounded}\}. \quad (8.8)$$



If we assume  $\mathfrak{D}(A)$  to be dense, then Theorem 1.16 implies that  $\ell \circ A$  has a unique extension to an element of  $X^*$  and we can set

$$A'\ell := \overline{\ell \circ A}, \quad \ell \in \mathfrak{D}(A'). \quad (8.9)$$

However, note that even if  $A$  is densely defined, it can happen that  $\mathfrak{D}(A') = \{0\}$  (Problem 8.14). Clearly this extends our previous definition for bounded operators and it is not hard to see that  $A$  is bounded if and only if  $A'$  is bounded (Problem 8.13). Finally note that

$$A \subseteq B \quad \Rightarrow \quad B' \subseteq A'. \quad (8.10)$$

**Example 8.9.** Consider again  $A$  from Example 8.3 with  $1 \leq p < \infty$ . Then  $X^* \cong \ell^q(\mathbb{N})$  with  $\frac{1}{q} + \frac{1}{p} = 1$ . Now if  $b \in \mathfrak{D}(A')$  there is some  $c \in \ell^q(\mathbb{N})$  such that  $\ell_c(a) = \ell_b(Aa)$  for all  $a \in \ell_c(\mathbb{N})$ . That is,

$$\sum_{j=1}^{\infty} c_j a_j = \sum_{j=1}^{\infty} b_j (j a_j), \quad a \in \ell_c(\mathbb{N}).$$

Choosing  $a = \delta^j$  shows  $c_j = j b_j$  and hence  $A' b_j = j b_j$  with  $\mathfrak{D}(A') = \{b \in \ell^q(\mathbb{N}) \mid (j b_j)_{j=1}^{\infty} \in \ell^q(\mathbb{N})\}$ . Note that for  $1 < p < \infty$  the domain is dense, while for  $p = 1$  we have  $\overline{\mathfrak{D}(A')} = c_0(\mathbb{N}) \subset \ell^\infty(\mathbb{N})$ .  $\diamond$

**Example 8.10.** Let us compute the adjoint of  $B$  from Example 8.4. Proceeding as in the previous example we get  $c_j = b_1$  which is in  $\ell^2$  if and only if  $b_1 = 0$ . Thus  $\mathfrak{D}(B') = \{b \in \ell^2(\mathbb{N}) \mid b_1 = 0\}$  and  $B'b = 0$ . So don't expect the adjoint of a noncloseable operator to contain much information about the operator.  $\diamond$

For the closure of sums and products see Problem 8.16 and Problem 8.17, respectively.

There are two other ways of introducing the adjoint operator which are worth while mentioning. First of all an operator  $B : \mathfrak{D}(B) \subseteq Y^* \rightarrow X^*$  is adjoint to  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  if

$$(B\ell)(x) = \ell(Ax), \quad x \in \mathfrak{D}(A), \ell \in \mathfrak{D}(B). \quad (8.11)$$

Then, if  $\mathfrak{D}(A)$  is dense, there is a unique maximally adjoint operator, which is precisely the adjoint  $A'$ . The second way is using graphs. To this end recall that

$$\Gamma^\perp = \{(x', y') \in X^* \oplus Y^* \mid x'(x) + y'(y) = 0 \ \forall (x, y) \in \Gamma\} \quad (8.12)$$

and hence the last equality can be rephrased as

$$\Gamma(B) \subseteq \Gamma^{-1}(-A)^\perp = (\Gamma(-A)^\perp)^{-1} \quad (8.13)$$

with equality if  $B = A'$ . In particular, we conclude that  $A'$  is closed. Note that  $\Gamma^{-1}(-A)^\perp$  is the graph of an operator if and only if  $\mathfrak{D}(A)^\perp = \{0\}$ , that

is, if and only if  $\mathfrak{D}(A)$  is densely defined. Moreover, since  $\Gamma^\perp = \overline{\Gamma}^\perp$  we have

$$A' = \overline{A'} \quad (8.14)$$

if  $A$  is closable.

If  $A'$  is densely defined we can compute the closure by computing the doubly adjoint operator.

**Theorem 8.5.** *Suppose  $X$  and  $Y$  are Banach spaces and  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  is densely defined. Then  $A$  is closable if and only if  $\mathfrak{D}(A')_\perp = \{0\}$ . Moreover, if  $A'$  is densely defined, then  $\overline{A} = J_Y^{-1} A'' J_X$  with  $\mathfrak{D}(\overline{A}) = J_X^{-1} \mathfrak{D}(A'')$  and  $\text{Ran}(A''|_{J_X(X)}) \subseteq J_Y(Y)$ .*

**Proof.** By Lemma 4.18 we have  $\Gamma^{-1}(-A')_\perp = (\Gamma(A)^\perp)_\perp = \overline{\Gamma(A)}$ . Hence  $\overline{\Gamma(A)}$  is a graph if and only if  $\mathfrak{D}(A')_\perp = \{0\}$ . Hence  $A$  is closable if  $\mathfrak{D}(A')$  is dense and the converse holds if  $Y$  is reflexive.

Moreover, if  $\mathfrak{D}(A')$  is dense,  $A''$  exists and  $J\Gamma(\overline{A}) = J\Gamma^{-1}(-A')_\perp \subseteq \Gamma^{-1}(-A')^\perp = \Gamma(A'')$  which shows  $\overline{A} = J_Y^{-1} A'' J_X$ .  $\square$

Of course if  $Y$  is reflexive, then the last theorem says that  $A$  is closable if and only if  $A'$  is densely defined. If  $Y$  is not reflexive, this formulation still holds if dense is replaced by weak-\* dense — see Problem 6.20. Without this replacement it can fail as we have seen in Example 8.9.

Finally we want to establish the closed range theorem for closed operators. We begin by extending Lemma 4.23.

**Lemma 8.6.** *If  $A \in \mathcal{C}(X, Y)$ , then  $\text{Ran}(A)^\perp = \text{Ker}(A')$  and  $\text{Ran}(A')_\perp = \text{Ker}(A)$ .*

**Proof.** For the first claim observe:  $y' \in \text{Ker}(A')$  implies  $0 = (A'y')(x) = y'(Ax)$  for all  $x \in \mathfrak{D}(A)$  and hence  $y' \in \text{Ran}(A)^\perp$ . Conversely,  $y' \in \text{Ran}(A)^\perp$  implies  $y'(Ax) = 0$  for all  $x \in \mathfrak{D}(A)$  and hence  $y' \in \mathfrak{D}(A')$  with  $A'y' = 0$ .

For the second claim observe:  $x \in \text{Ker}(A)$  implies  $(A'y')(x) = y'(Ax) = 0$  for all  $y' \in \mathfrak{D}(A')$  and hence  $x \in \text{Ran}(A')_\perp$ . Conversely,  $x \in \text{Ran}(A')_\perp$  implies  $(A'y')(x) = 0$  for all  $y' \in \mathfrak{D}(A')$ . If we had  $(x, 0) \notin \Gamma(A)$ , we could find (Corollary 4.13) a linear functional  $(x', y') \in X^* \oplus Y^*$  which vanishes on  $\Gamma(A)$  and satisfies  $(x', y')(x, 0) = 1$ . The fact that it vanishes on  $\Gamma(A)$  implies  $x'(x) + y'(Ax) = 0$  for all  $x \in \mathfrak{D}(A)$  implying  $y' \in \mathfrak{D}(A')$  with  $A'y' = -x'$ . But this gives the contradiction  $0 = (A'y')(x) = -x'(x) = -1$  and hence  $(x, 0) \in \Gamma(A)$ , that is,  $x \in \text{Ker}(A)$ .  $\square$

Taking annihilators in these formulas we again obtain

$$\text{Ker}(A')_\perp = (\text{Ran}(A)^\perp)_\perp = \overline{\text{Ran}(A)} \quad (8.15)$$

and

$$\operatorname{Ker}(A)^\perp = (\operatorname{Ran}(A')_\perp)^\perp \supseteq \overline{\operatorname{Ran}(A')}. \quad (8.16)$$

The analog of Theorem 4.22 is

**Theorem 8.7.** *Suppose  $X, Y$  are Banach spaces. If  $A \in \mathcal{C}(X, Y)$ , then  $A^{-1}$  exists and is in  $\mathcal{L}(Y, X)$  if and only if  $(A')^{-1}$  exists and is in  $\mathcal{L}(X^*, Y^*)$ . Moreover, in this case we have*

$$(A')^{-1} = (A^{-1})'. \quad (8.17)$$

**Proof.** If  $A^{-1} \in \mathcal{L}(Y, X)$ , then Theorem 4.20 implies  $(A^{-1})' \in \mathcal{L}(X^*, Y^*)$  and  $\Gamma((A^{-1})') = \Gamma^{-1}(-A^{-1})^\perp = (\Gamma^{-1}(-A)^\perp)^{-1} = \Gamma(A')^{-1} = \Gamma((A')^{-1})$  which establishes (8.17).

Conversely, let  $(A')^{-1} \in \mathcal{L}(X^*, Y^*)$ . Then for any  $\ell \in X^*$  and  $x \in \mathfrak{D}(A)$  we have  $((A')^{-1}\ell)(Ax) = (A'(A')^{-1}\ell)(x) = \ell(x)$  and choosing  $\ell$  normalized such that  $\ell(x) = \|x\|$  (Corollary 4.11) we obtain  $\|x\| = ((A')^{-1}\ell)(Ax) \leq \|(A')^{-1}\| \|Ax\|$ . This shows that  $A$  has a bounded inverse with  $\|A^{-1}\| \leq \|(A')^{-1}\|$  and hence  $\operatorname{Ran}(A)$  is closed by Corollary 8.2. Finally  $\operatorname{Ran}(A)^\perp = \operatorname{Ker}(A')$  shows  $\operatorname{Ran}(A) = Y$ .  $\square$

Now we are ready to show:

**Theorem 8.8** (Closed range). *Suppose  $X, Y$  are Banach spaces and  $A \in \mathcal{C}(X, Y)$ . Then the following items are equivalent:*

- (i)  $\operatorname{Ran}(A)$  is closed.
- (ii)  $\operatorname{Ker}(A)^\perp = \operatorname{Ran}(A')$ .
- (iii)  $\operatorname{Ran}(A')$  is closed.
- (iv)  $\operatorname{Ker}(A')_\perp = \operatorname{Ran}(A)$ .

**Proof.** Consider  $\tilde{X} = X/\operatorname{Ker}(A)$  and  $\tilde{Y} = \overline{\operatorname{Ran}(A)}$  and the corresponding operator  $\tilde{A}$  as in Lemma 8.3. Then  $\tilde{A}$  is a closed injective operator whose range is dense. Now you can literally follow the proof of Theorem 4.24 replacing the used results for bounded operators by the corresponding one for closed operators.  $\square$

We end this section with the remark that one could try to define the concept of a weakly closed operator by replacing the norm topology in the definition of a closed operator by the weak topology. However, Theorem 6.12 implies that this gives nothing new:

**Lemma 8.9.** *For an operator  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  the following are equivalent:*

- $\Gamma(A)$  is closed.

- $x_n \in \mathfrak{D}(A)$  with  $x_n \rightarrow x$  and  $Ax_n \rightarrow y$  implies  $x \in \mathfrak{D}(A)$  and  $y = Ax$ .
- $\Gamma(A)$  is weakly closed.
- $x_n \in \mathfrak{D}(A)$  with  $x_n \rightharpoonup x$  and  $Ax_n \rightharpoonup y$  implies  $x \in \mathfrak{D}(A)$  and  $Ax = y$ .

**Problem\* 8.1.** Show that the kernel  $\text{Ker}(A)$  of a closed operator  $A$  is closed.

**Problem\* 8.2.** A linear functional defined on a dense subspace is closable if and only if it is bounded.

**Problem 8.3.** Let  $(m_j)_{j \in \mathbb{N}}$  be a sequence of complex numbers and consider the multiplication operator  $M$  in  $\ell^p(\mathbb{N})$  defined by  $Ma_j := m_j a_j$  on  $\mathfrak{D}(A) := \ell_c(\mathbb{N})$ . Under which conditions on  $m$  is  $M$  closable and what is its closure?

**Problem 8.4.** Show that the differential operator  $A = \frac{d}{dx}$  defined on  $\mathfrak{D}(A) = C^1[0, 1] \subset C[0, 1]$  (sup norm) is a closed operator. (Compare the example in Section 1.6.)

**Problem\* 8.5.** Show that the completion  $X_A$  of  $(\mathfrak{D}(A), \|\cdot\|_A)$  can be regarded as a subset of  $X$  if and only if  $A$  is closable. Show that in this case the completion can be identified with  $\mathfrak{D}(\overline{A})$  and that the closure of  $A$  in  $X$  coincides with the extension from Theorem 1.16 of  $A$  in  $X_A$ . In particular,  $A$  is closed if and only if  $(\mathfrak{D}(A), \|\cdot\|_A)$  is complete.

**Problem 8.6.** Consider the Sobolev spaces  $W^{1,p}(0, 1)$ . Show that

$$\|f\|_\infty \leq \|f\|_{1,1}$$

and conclude that functions from  $W^{1,p}(0, 1)$  are continuous. (Hint: Start with estimating  $f(a) = f(x) - \int_a^x f'(y)dy$  and integrate the result.)

**Problem 8.7.** Show that  $f \in W^{1,p}(0, 1)$ ,  $1 < p < \infty$  is Hölder continuous with

$$|f(x) - f(y)| \leq \|f'\|_p |x - y|^{1-\frac{1}{p}}.$$

Conclude that the embedding  $W^{1,p}(0, 1) \hookrightarrow L^p(0, 1)$  is compact.

**Problem 8.8.** Let  $X := \ell^2(\mathbb{N})$  and  $(Aa)_j := j a_j$  with  $\mathfrak{D}(A) := \{a \in \ell^2(\mathbb{N}) \mid (ja_j)_{j \in \mathbb{N}} \in \ell^2(\mathbb{N})\}$  and  $Ba := (\sum_{j \in \mathbb{N}} a_j) \delta^1$ . Then we have seen that  $A$  is closed while  $B$  is not closable. Show that  $A+B$ ,  $\mathfrak{D}(A+B) = \mathfrak{D}(A) \cap \mathfrak{D}(B) = \mathfrak{D}(A)$  is closed.

**Problem 8.9.** Let  $X := C[0, 1]$ . Show that the operator given by  $Ax(t) := \frac{x(t)}{t}$  with  $\mathfrak{D}(A) := \{x \in X \mid \exists \lim_{x \rightarrow 0} \frac{x(t)}{t}\}$  is closed.

**Problem 8.10.** Let  $X := C[0, 1]$ . Show that the operator given by  $Ax := x'' + x$  with  $\mathfrak{D}(A) := \{x(t) \in C^2[0, 1] \mid x(0) = x'(0) = 0\}$  is closed.

**Problem\* 8.11.** Show that if  $A$ ,  $B$ , and  $A + B$  are closable, then  $\overline{A} + \overline{B} \subseteq \overline{A + B}$  with equality if  $A$  or  $B$  is bounded. Moreover, if  $B$  is bounded, then  $A + B$  is closable if and only if  $A$  is.

Give an example where equality fails. Give an example where  $A$  and  $B$  are closable but  $A + B$  is not. (Hint: For the very last part see Problem 8.8.)

**Problem\* 8.12.** Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  and  $B : \mathfrak{D}(B) \subseteq Y \rightarrow Z$  be closable. If  $A \in \mathcal{L}(X, Y)$ , then  $BA$  is closable and  $\overline{BA} = \overline{B}\overline{A}$ . Similarly, if  $B^{-1} \in \mathcal{L}(Z, Y)$ , then  $BA$  is closable and  $\overline{BA} = \overline{B}\overline{A}$ .

**Problem\* 8.13.** Show that for  $A \in \mathcal{C}(X, Y)$  we have  $A \in \mathcal{L}(X, Y)$  if and only if  $A' \in \mathcal{L}(Y^*, X^*)$ .

**Problem 8.14.** Show that  $\mathfrak{D}(A')$  is trivial if and only if  $\Gamma(A) \subset X \times Y$  is dense. Let  $X = Y$  be a separable Hilbert space and choose an ONB  $\{u_n\}_{n \in \mathbb{N}}$  plus a dense set  $\{x_n\}_{n \in \mathbb{N}}$ . Define  $Au_n := x_n$  and extend  $A$  by linearity to  $\mathfrak{D}(A) = \text{span}\{u_n\}_{n \in \mathbb{N}}$ . Show that  $\Gamma(A)^\perp = \{(0, 0)\}$  and conclude  $\mathfrak{D}(A') = \{0\}$ .

**Problem 8.15.** Compute the adjoint of the embedding  $I : \ell^1(\mathbb{N}) \hookrightarrow \ell^2(\mathbb{N})$ .

**Problem\* 8.16.** Show that  $A' + B' \subseteq (A + B)'$  if  $A + B$  is densely defined. Show that we have equality if one of the operators is bounded. Give an example where equality fails.

**Problem\* 8.17.** Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  and  $B : \mathfrak{D}(B) \subseteq Y \rightarrow Z$  be densely defined. If  $B \in \mathcal{L}(Y, Z)$ , then  $BA$  is densely defined and  $(BA)' = A'B'$ . Similarly, if  $A^{-1} \in \mathcal{L}(Y, X)$ , then  $BA$  is densely defined and  $(BA)' = A'B'$ .

**Problem 8.18.** Let  $A$  be a closed operator. Show that for every  $\alpha \in \rho(A)$  the expression  $\|f\|_\alpha := \|(A - \alpha)x\|$  defines a norm which is equivalent to the graph norm.

**Problem 8.19.** Let  $X_j$  be Banach spaces. A sequence of operators  $A_j \in \mathcal{C}(X_j, X_{j+1})$

$$X_1 \xrightarrow{A_1} X_2 \xrightarrow{A_2} X_3 \cdots X_n \xrightarrow{A_n} X_{n+1}$$

is said to be **exact** if  $\text{Ran}(A_j) = \text{Ker}(A_{j+1})$  for  $1 \leq j \leq n$ . Show that a sequence is exact if and only if the corresponding dual sequence

$$X_1^* \xleftarrow{A_1'} X_2^* \xleftarrow{A_2'} X_3^* \cdots X_n^* \xleftarrow{A_n'} X_{n+1}^*$$

is exact.

**Problem\* 8.20.** Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow X$  be an unbounded operator and  $B \in \mathcal{L}(X)$  bounded. Then  $A$  and  $B$  are said to **commute** if

$$BA \subseteq AB.$$

Of course if  $A \in \mathcal{L}(X)$ , then we have equality and this definition reduces to the usual one. Note that the definition implies in particular, that  $B$  leaves  $\mathfrak{D}(A)$  invariant,  $B\mathfrak{D}(A) \subseteq \mathfrak{D}(A)$ .

Show that if  $A$  is invertible, then  $A$  commutes with  $B$  if and only if  $A^{-1}$  commutes with  $B$ . Conclude that if  $A$  has a nonempty resolvent set, then  $A$  commutes with  $B$  if and only if  $R_A(\alpha)$  commutes with  $B$  for one  $\alpha \in \rho(A)$ . Moreover, in this case this holds for all  $\alpha \in \rho(A)$ .

**Problem\* 8.21.** Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow X$  be closable and  $B \in \mathcal{L}(X)$  bounded. Then if  $A$  commutes with  $B$  so does  $\overline{A}$ .

## 8.2. Spectral theory for unbounded operators

As in the case of bounded operators (cf. Section 5.1) we define the **resolvent set** via

$$\rho(A) := \{\alpha \in \mathbb{C} \mid A - \alpha \text{ is bijective with a bounded inverse}\} \quad (8.18)$$

and call

$$R_A(\alpha) := (A - \alpha)^{-1}, \quad \alpha \in \rho(A) \quad (8.19)$$

the **resolvent** of  $A$ . The complement  $\sigma(A) = \mathbb{C} \setminus \rho(A)$  is called the **spectrum** of  $A$ . As in the case of Banach algebras it follows that the resolvent is analytic and that the resolvent set is open:

**Lemma 8.10.** Let  $A$  be a closed operator. Then the resolvent set is open and if  $\alpha_0 \in \rho(A)$  we have

$$R_A(\alpha) = \sum_{n=0}^{\infty} (\alpha - \alpha_0)^n R_A(\alpha_0)^{n+1}, \quad |\alpha - \alpha_0| < \|R_A(\alpha_0)\|^{-1}. \quad (8.20)$$

In particular, the resolvent is analytic and

$$\|(A - \alpha)^{-1}\| \geq \frac{1}{\text{dist}(\alpha, \sigma(A))}. \quad (8.21)$$

**Proof.** By comparing with the geometric series one sees that the series for  $R_A(\alpha)$  converges and hence defines a bounded operator. Denote this operator by  $R$ . Now fix  $x \in X$  and consider  $R_N := \sum_{n=0}^N (\alpha - \alpha_0)^n R_A(\alpha_0)^{n+1}$ . Then  $x_N := R_N x \in \mathfrak{D}(A)$  and

$$\begin{aligned} (A - \alpha)x_N &= (A - \alpha_0)x_N - (\alpha - \alpha_0)x_N \\ &= \sum_{n=0}^N (\alpha - \alpha_0)^n R_A(\alpha_0)^{n+1} x - \sum_{n=0}^N (\alpha - \alpha_0)^{n+1} R_A(\alpha_0)^{n+1} x \\ &= x - (\alpha - \alpha_0)^{N+1} R_A(\alpha_0)^{N+1} x. \end{aligned}$$

Hence  $(A - \alpha)x_N \rightarrow x$  implying  $x \in \mathfrak{D}(A)$  and  $(A - \alpha)Rx = x$ . Similarly one shows  $R(A - \alpha)x = x$  for  $x \in \mathfrak{D}(A)$ .  $\square$

As a consequence we obtain the useful

**Lemma 8.11.** *We have  $\alpha \in \sigma(A)$  if there is a sequence  $x_n \in \mathfrak{D}(A)$  such that  $\|x_n\| = 1$  and  $\|(A - \alpha)x_n\| \rightarrow 0$ . If  $\alpha$  is a boundary point of  $\rho(A)$ , then the converse is also true. Such a sequence is called a **Weyl sequence**.*

**Proof.** Let  $x_n$  be a Weyl sequence. Then  $\alpha \in \rho(A)$  is impossible by  $1 = \|x_n\| = \|R_A(\alpha)(A - \alpha)x_n\| \leq \|R_A(\alpha)\| \|(A - \alpha)x_n\| \rightarrow 0$ . Conversely, by (8.21), there is a sequence  $\alpha_n \rightarrow \alpha$  and corresponding nonzero vectors  $y_n \in X$  such that  $\|R_A(\alpha_n)y_n\| \|y_n\|^{-1} \rightarrow \infty$ . Let  $x_n := R_A(\alpha_n)y_n$  and rescale  $y_n$  such that  $\|x_n\| = 1$ . Then  $\|y_n\| \rightarrow 0$  and hence

$$\|(A - \alpha)x_n\| = \|y_n + (\alpha_n - \alpha)x_n\| \leq \|y_n\| + |\alpha - \alpha_n| \rightarrow 0$$

shows that  $x_n$  is a Weyl sequence.  $\square$

The set of all  $\alpha \in \mathbb{C}$  for which a Weyl sequence exists is called the **approximate point spectrum**

$$\sigma_{ap}(A) := \{\alpha \in \mathbb{C} \mid \exists x_n \in \mathfrak{D}(A) : \|x_n\| = 1, \|(A - \alpha)x_n\| \rightarrow 0\} \quad (8.22)$$

and the above lemma could be phrased as

$$\partial\sigma(A) \subseteq \sigma_{ap}(A) \subseteq \sigma(A). \quad (8.23)$$

Note that there are two possibilities if  $\alpha \in \sigma_{ap}(A)$ : Either  $\alpha$  is an eigenvalue or  $(A - \alpha)^{-1}$  exists but is unbounded. By the closed graph theorem the latter case is equivalent to the fact that  $\text{Ran}(A - \alpha)$  is not closed. In particular, we have

$$\sigma_p(A) \cup \sigma_c(A) \subseteq \sigma_{ap}(A) \subseteq \sigma(A). \quad (8.24)$$

**Example 8.11.** If  $\alpha$  is an eigenvalue and  $x$  a corresponding normalized eigenfunction, then  $x_n := x$  will be a Weyl sequence. Hence you should think of  $x_n$  as approximate eigenfunctions. Conversely, if a Weyl sequence converges,  $x_n \rightarrow x$ , then we also have  $Ax_n \rightarrow \alpha x$  which shows that  $x$  is an eigenfunction (assuming  $A$  is closed). However, in general a Weyl sequence does not have to converge. Indeed, consider the multiplication operator  $(Aa)_j = \frac{1}{j}a_j$  in  $\ell^p(\mathbb{N})$ . Then  $x_n := \delta^n$  is a Weyl sequence for  $\alpha = 0$ , but it clearly has no convergent subsequence.  $\diamond$

It is also straightforward to verify the **first resolvent identity**

$$\begin{aligned} R_A(\alpha_0) - R_A(\alpha_1) &= (\alpha_0 - \alpha_1)R_A(\alpha_0)R_A(\alpha_1) \\ &= (\alpha_0 - \alpha_1)R_A(\alpha_1)R_A(\alpha_0), \end{aligned} \quad (8.25)$$

for  $\alpha_0, \alpha_1 \in \rho(A)$ .

However, note that for unbounded operators the spectrum will no longer be bounded in general and both  $\sigma(A) = \emptyset$  as well as  $\sigma(A) = \mathbb{C}$  are possible.

**Example 8.12.** Consider  $X := C[0, 1]$  and  $A = \frac{d}{dx}$  with  $\mathfrak{D}(A) = C^1[0, 1]$ . We obtain the eigenvalues by solving the ordinary differential equation  $x'(t) = \alpha x(t)$  which gives  $x(t) = e^{\alpha t}$ . Hence every  $\alpha \in \mathbb{C}$  is an eigenvalue, that is,  $\sigma(A) = \mathbb{C}$ .

Now let us modify the domain and look at  $A_0 = \frac{d}{dx}$  with  $\mathfrak{D}(A_0) = \{x \in C^1[0, 1] | x(0) = 0\}$  and  $X_0 := \{x \in C[0, 1] | x(0) = 0\}$ . Then the previous eigenfunctions do not satisfy the boundary condition  $x(0) = 0$  and hence  $A_0$  has no eigenvalues. Moreover, the solution of the inhomogeneous ordinary differential equation  $x'(t) - \alpha x(t) = y(t)$  is given by  $x(t) = x(0)e^{\alpha t} + \int_0^t e^{\alpha(t-s)} y(s) ds$ . Hence  $R_{A_0}(\alpha)y(t) = \int_0^t e^{\alpha(t-s)} y(s) ds$  is the resolvent of  $A_0$ . Consequently  $\sigma(A_0) = \emptyset$ .  $\diamond$

Note that if  $A$  is closed, then bijectivity implies boundedness of the inverse (see Corollary 8.2). Moreover, by Lemma 8.1 an operator with nonempty resolvent set must be closed.

The point, continuous, and residual spectrum can be defined as in Section 7.1 and Lemma 7.1 holds for closed operators.

**Lemma 8.12.** Suppose  $A \in \mathcal{C}(X)$ . Then

$$\sigma(A) = \sigma(A') \quad (8.26)$$

and

$$R_A(\alpha)' = R_{A'}(\alpha), \quad \alpha \in \rho(A) = \rho(A'). \quad (8.27)$$

Moreover,

$$\begin{aligned} \sigma_p(A') &\subseteq \sigma_p(A) \cup \sigma_r(A), & \sigma_p(A) &\subseteq \sigma_p(A') \cup \sigma_r(A'), \\ \sigma_r(A') &\subseteq \sigma_p(A) \cup \sigma_c(A), & \sigma_r(A) &\subseteq \sigma_p(A'), \\ \sigma_c(A') &\subseteq \sigma_c(A), & \sigma_c(A) &\subseteq \sigma_r(A') \cup \sigma_c(A'). \end{aligned} \quad (8.28)$$

If in addition,  $X$  is reflexive we have  $\sigma_r(A') \subseteq \sigma_p(A)$  as well as  $\sigma_c(A') = \sigma_c(A)$ .

**Proof.** Literally follow the proof of Lemma 7.1 using the corresponding results for closed operators: Theorem 8.7, Lemma 8.6, and Theorem 8.5.  $\square$

Let us also note the following spectral mapping result.

**Lemma 8.13.** Suppose  $A \in \mathcal{C}(X)$  is injective with  $\text{Ran}(A)$  is dense. Then

$$\sigma(A^{-1}) \setminus \{0\} = (\sigma(A) \setminus \{0\})^{-1} \quad (8.29)$$

and

$$R_{A^{-1}}(\alpha^{-1}) = -\alpha A R_A(\alpha) = -\alpha - \alpha^2 R_A(\alpha), \quad \alpha \in \rho(A) \setminus \{0\}. \quad (8.30)$$

In addition, for  $\alpha \neq 0$  we have  $\text{Ker}((A - \alpha)^n) = \text{Ker}((A^{-1} - \alpha^{-1})^n)$  as well as  $\text{Ran}((A - \alpha)^n) = \text{Ran}((A^{-1} - \alpha^{-1})^n)$  for any  $n \in \mathbb{N}$ .



**Proof.** Throughout this proof  $\alpha \neq 0$ . We first show the formula for the resolvent of  $A^{-1}$ . To this end choose some  $\alpha \in \rho(A) \setminus \{0\}$  and observe that the right-hand side of (8.30) is a bounded operator from  $X \rightarrow \text{Ran}(A) = \mathfrak{D}(A^{-1})$  and

$$(A^{-1} - \alpha^{-1})(-\alpha AR_A(\alpha))x = (-\alpha + A)R_A(\alpha)x = x, \quad x \in X.$$

Conversely, if  $y \in \mathfrak{D}(A^{-1}) = \text{Ran}(A)$ , we have  $y = Ax$  and hence

$$(-\alpha AR_A(\alpha))(A^{-1} - \alpha^{-1})y = AR_A(\alpha)((A - \alpha)x) = Ax = y.$$

Thus (8.30) holds and  $\alpha^{-1} \in \rho(A^{-1})$ . Interchanging the roles of  $A$  and  $A^{-1}$  establishes the first part.

Next note that for  $x \in \mathfrak{D}(A)$  the equation  $(A - \alpha)x = y$  is equivalent to  $x = \frac{1}{\alpha}(Ax - y)$  and hence  $(A - \alpha)x \in \text{Ran}(A^m)$  implies  $x \in \text{Ran}(A^m)$  for any  $m \in \mathbb{N}$ . Consequently we get that  $x \in \text{Ker}((A - \alpha)^n)$  implies  $x \in \text{Ran}(A^m) = \mathfrak{D}(A^{-m})$  for any  $m \in \mathbb{N}$  and  $0 = A^{-n}(A - \alpha)^n \alpha = (-\alpha)^n(A^{-1} - \alpha^{-1})^n x$ . So  $\text{Ker}((A - \alpha)^n) \subseteq \text{Ker}((A^{-1} - \alpha^{-1})^n)$  and equality follows by reversing the roles of  $A$  and  $A^{-1}$ .

Similarly,  $x \in \text{Ran}((A - \alpha)^n)$  implies  $x = (A - \alpha)^n y$  for some  $y \in \mathfrak{D}(A^n)$ . Consequently  $y = A^{-n}z$  and  $x = (A - \alpha)^n y = (A - \alpha)^n A^{-n}z = (-\alpha)^n(A^{-1} - \alpha^{-1})^n z \in \text{Ran}((A^{-1} - \alpha^{-1})^n)$ . So  $\text{Ran}((A - \alpha)^n) \subseteq \text{Ran}((A^{-1} - \alpha^{-1})^n)$  and equality follows again by reversing the roles of  $A$  and  $A^{-1}$ .  $\square$

Concerning  $\alpha = 0$  note that  $0 \in \sigma(A^{-1})$  if and only if  $A$  is unbounded and vice versa.

In particular we can apply this lemma to the resolvent in case  $\alpha_0 \in \rho(A)$  which shows

$$\sigma(A) = \alpha_0 + (\sigma(R_A(\alpha_0)) \setminus \{0\})^{-1} \quad (8.31)$$

and  $\text{Ker}(R_A(\alpha_0) - \alpha)^n = \text{Ker}(A - \alpha_0 - \frac{1}{\alpha})^n$  as well as  $\text{Ran}(R_A(\alpha_0) - \alpha)^n = \text{Ran}(A - \alpha_0 - \frac{1}{\alpha})^n$  for  $\alpha \neq 0$  and  $n \in \mathbb{N}$ .

For example, this can be used to apply Theorem 7.7 to unbounded operators in case they have a compact resolvent. To this end note that if we have  $R_A(\alpha) \in \mathcal{K}(X)$  for one  $\alpha \in \rho(A)$ , then this holds in fact for all  $\alpha \in \rho(A)$  by the first resolvent identity (8.25) since compact operators form an ideal.

**Theorem 8.14.** *Suppose  $R_A(\alpha) \in \mathcal{K}(X)$  for one  $\alpha \in \rho(A)$ . Then the spectrum of  $A$  consists only of discrete eigenvalues with finite (geometric and algebraic) multiplicity and we have the splitting into closed invariant subspaces*

$$X = \text{Ker}(A - \alpha)^n \dot{+} \text{Ran}(A - \alpha)^n, \quad \alpha \in \sigma(A), \quad (8.32)$$

where  $n$  is the index of  $\alpha$ .

**Proof.** The claim follows by combining the previous lemma with the spectral theorem for compact operators (Theorem 7.7 and Lemma 7.5).  $\square$

**Example 8.13.** Consider again  $A_0$  from the previous problem. Then

$$A_0^{-1} : X_0 \rightarrow X_0, \quad x \mapsto x(t) = \int_0^t x(s) ds$$

is compact by Problem 3.6 (cf. also Problem 7.4). Hence we get again  $\sigma(A_0) = \emptyset$  without the need of computing the resolvent.  $\diamond$

**Example 8.14.** Of course another example of unbounded operators with compact resolvent are regular Sturm–Liouville operators as shown in Section 3.3.  $\diamond$

This result says in particular, that for  $\alpha \in \sigma_p(A)$  we can split  $X = X_1 \oplus X_2$  where both  $X_1 := \text{Ker}(A - \alpha)^n$  and  $X_2 := \text{Ran}(A - \alpha)^n$  are invariant subspaces for  $A$ . Consequently we can split  $A = A_1 \oplus A_2$ , where  $A_1$  is the restriction of  $A$  to the finite dimensional subspace  $X_1$  (with  $A_1 - \alpha$  a nilpotent matrix) and  $A_2$  is the restriction of  $A$  to  $X_2$ . Moreover,  $\text{Ker}(A_2 - \alpha) = \{0\}$  by construction. Now note that for  $\beta \in \rho(A)$  we must have  $\beta \in \rho(A_1) \cap \rho(A_2)$  with  $R_A(\beta) = R_{A_1}(\beta) \oplus R_{A_2}(\beta)$  (cf. Problem 8.25) which shows that  $R_{A_2}(\beta) \in \mathcal{K}(X_2)$  for  $\beta \in \rho(A)$ . Now since  $\alpha \notin \sigma_p(A_2)$  this tells us  $\alpha \in \rho(A_2)$ .

**Problem\* 8.22.** Let  $A$  be a closed operator. Show (8.25). Moreover, conclude

$$\frac{d^n}{d\alpha^n} R_A(\alpha) = n! R_A(\alpha)^{n+1}, \quad \frac{d}{d\alpha} R_A(\alpha)^n = n R_A(\alpha)^{n+1}.$$

**Problem 8.23.** Suppose  $A = \overline{A_0}$ . If  $x_n \in \mathfrak{D}(A)$  is a Weyl sequence for  $\alpha \in \sigma(A)$ , then there is also one with  $\tilde{x}_n \in \mathfrak{D}(A_0)$ .

**Problem\* 8.24.** Suppose  $A_j \mathfrak{D}(A_j) \subseteq X_j \rightarrow Y_j$ ,  $j = 1, 2$ . Then  $A_1 \oplus A_2$  is closable if and only if both  $A_1$  and  $A_2$  are closable. Moreover, in this case we have  $\overline{A_1 \oplus A_2} = \overline{A_1} \oplus \overline{A_2}$ .

**Problem\* 8.25.** Suppose  $A_j \in \mathcal{C}(X_j)$ ,  $j = 1, 2$ . Then  $A_1 \oplus A_2 \in \mathcal{C}(X_1 \oplus X_2)$  and  $\sigma(A_1 \oplus A_2) = \sigma(A_1) \cup \sigma(A_2)$  with

$$R_{A_1 \oplus A_2}(\alpha) = R_{A_1}(\alpha) \oplus R_{A_2}(\alpha), \quad \alpha \in \rho(A_1) \cap \rho(A_2).$$

### 8.3. Reducing subspaces and spectral projections

In the previous section we have seen that if  $A$  has a compact resolvent, then we can split  $X$  into two invariant subspaces such that all the spectral information pertaining to an eigenvalue is captured by restricting  $A$  to the first subspace, while the restriction to the second subspace has a bounded inverse. In this section we want to discuss such a splitting in more detail.

To begin with, recall that for a given closed subspace  $M \subseteq X$  there might not be a corresponding closed subspace  $N$  such that  $X = M \dot{+} N$ . Subspaces with this property are called complemented and  $M$  is complemented if and only if there is a projection  $P = P^2 \in \mathcal{L}(X)$  such that  $M = \text{Ran}(P)$ . In this case we have  $N = \text{Ker}(P)$  and  $1 - P$  is a projection with  $N = \text{Ran}(1 - P)$ . Recall that these problems do not arise if  $M$  is finite dimensional (or finite codimensional) or if we are in a Hilbert space (where we can choose  $N = M^\perp$ ). Moreover, in this case we have  $X \cong M \oplus N$ .

So we need to assume that we have a complemented subspace such that both  $M$  and  $N$  are invariant with respect to  $A$ . Moreover, if  $A$  is unbounded we also need to assume that the domain is compatible with this splitting, that is, if we split a vector from the domain, both pieces will be again in the domain. In this case  $M$  is said to **reduce**  $A : \mathfrak{D}(A) \subseteq X \rightarrow X$  and this is equivalent to the requirement that the corresponding projection  $P$  commutes with  $A$  (cf. Problem 8.20) in the sense that

$$PA \subseteq AP. \quad (8.33)$$

Note that the definition implies in particular, that  $P$  leaves  $\mathfrak{D}(A)$  invariant,  $P\mathfrak{D}(A) \subseteq \mathfrak{D}(A)$ . Clearly, in this case  $A$  also commutes with the complementary projection  $Q := 1 - P$ . If  $A$  has nonempty resolvent set, then  $A$  will commute with  $P$  if and only if it commutes with  $R_A(\alpha)$  for one (and hence for all)  $\alpha \in \rho(A)$ . Moreover, if  $A$  commutes with  $P$ , the same is true for  $\bar{A}$  (Problem 8.21).

**Lemma 8.15.** *Suppose  $P \in \mathcal{L}(X)$  is a projection which commutes with  $A \in \mathcal{C}(X)$ . Let  $M := \text{Ran}(P)$  and  $N := \text{Ker}(P)$ . Then we have  $A = A_1 \oplus A_2$  with  $A_1 = A|_M \in \mathcal{C}(M)$  and  $A_2 = A|_N \in \mathcal{C}(N)$ . Here  $\mathfrak{D}(A|_M) := \mathfrak{D}(A) \cap M$ .*

**Proof.** First of all note that  $A_1 : \mathfrak{D}(A_1) \subseteq M \rightarrow M$  is well defined as  $A_1x = Ax$  for  $x \in \mathfrak{D}(A_1) := M \cap \mathfrak{D}(A) = P\mathfrak{D}(A)$ . Moreover,  $A_1$  is densely defined since if  $x \in M$  we can find a sequence  $x_n \in \mathfrak{D}(A)$  with  $x_n \rightarrow x$  and hence also  $Px_n \in \mathfrak{D}(A_1) \rightarrow Px = x$ . Similarly we see that  $A_1$  is closed since  $\Gamma(A_1) = \Gamma(A) \cap (M \oplus M)$ .

Exchanging  $P$  and  $Q$  shows that the same conclusions hold for  $A_2 : \mathfrak{D}(A_2) \subseteq N \rightarrow N$  defined as  $A_2x = Ax$  for  $x \in \mathfrak{D}(A_2) := N \cap \mathfrak{D}(A) = Q\mathfrak{D}(A)$ .  $\square$

Hence in such a situation the investigation of  $A$  can be reduced to the investigation of  $A_1$  and  $A_2$  (cf. Problems 8.24 and 8.25).

In a finite dimensional case the projection onto a generalized eigenspace is given as the negative residue of the resolvent.

**Example 8.15.** Indeed consider

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

such that

$$(A - \alpha)^{-1} = \frac{1}{1 - \alpha} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} + \frac{1}{-\alpha} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} + \frac{1}{\alpha^2} \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}. \quad \diamond$$

This motivates the following definition: Suppose  $\gamma : [0, 1] \rightarrow \rho(A) \subset \mathbb{C}$  is a smooth Jordan curve. Recall that a Jordan curve is a closed curve without any self intersections and the Jordan curve theorem states that it splits the complex plane into two connected regions which both have the curve as a boundary. In particular, the interior region will be bounded. While this seems quite obvious, a proof turns out to be quite delicate. However, in the typical application of our result  $\gamma$  will just be a circle and hence you can always think of a circle.

Then we define the corresponding projector as

$$P := -\frac{1}{2\pi i} \oint_{\gamma} R_A(\alpha) d\alpha. \quad (8.34)$$

Here the integral is defined as a Riemann integral (cf. Section 9.1 for details), explicitly

$$\oint_{\gamma} R_A(\alpha) d\alpha := \lim_{n \rightarrow \infty} \sum_{j=1}^n R_A(\gamma(j/n)) \gamma'(j/n) / n. \quad (8.35)$$

Choosing some  $x \in X$  and some  $\ell \in X^*$  we get

$$\ell(Px) = -\frac{1}{2\pi i} \oint_{\gamma} \ell(R_A(\alpha)x) d\alpha, \quad (8.36)$$

where now the integral on the right is an ordinary path integral in the complex plane. In particular, this shows that all the convenient facts from complex analysis about line integrals of holomorphic functions are at our disposal. For example, we can continuously deform the curve  $\gamma$  within  $\rho(A)$  without changing the integral and the Cauchy integral theorem holds.

**Lemma 8.16.** *Let  $A \in \mathcal{C}(X)$  and  $\gamma : [0, 1] \rightarrow \rho(A)$  be a Jordan curve. Then  $P$  from (8.34) is a projection which reduces  $A$ .*

**Proof.** First of all note that  $P \in \mathcal{L}(X)$  since (see again Section 9.1 for details) the right-hand side of (8.35) converges in the operator norm. Moreover, since the integral is independent of  $\gamma$  we can take a slightly deformed

path  $\gamma'$  in the exterior of  $\gamma$  to compute:

$$\begin{aligned} P^2 &= -\frac{1}{4\pi^2} \oint_{\gamma'} \oint_{\gamma} R_A(\alpha) R_A(\beta) d\alpha d\beta \\ &= -\frac{1}{4\pi^2} \oint_{\gamma'} \left( \oint_{\gamma} R_A(\alpha) \frac{d\alpha}{\alpha - \beta} \right) d\beta + \frac{1}{4\pi^2} \oint_{\gamma'} R_A(\beta) \left( \oint_{\gamma} \frac{d\alpha}{\alpha - \beta} \right) d\beta \\ &= -\frac{1}{4\pi^2} \oint_{\gamma} R_A(\alpha) \left( \oint_{\gamma'} \frac{d\beta}{\alpha - \beta} \right) d\alpha = -\frac{1}{2\pi i} \oint_{\gamma} R_A(\alpha) d\alpha = P \end{aligned}$$

Here we have used the first resolvent identity (8.25) to obtain the second line. To obtain the third line we have used Fubini for the first double integral and the fact that in the second integral the inner integral vanishes by the Cauchy integral theorem since  $\beta$  lies in the exterior of  $\gamma$ . To obtain the last line observe that the inner integral equals  $-2\pi i$  since  $\alpha$  lies in the interior of  $\gamma'$ .

Finally, since  $P$  commutes with  $R_A(\alpha)$  it also commutes with  $A$  (Problem 8.20), that is,  $P$  reduces  $A$ .  $\square$

**Example 8.16.** Corollary 3.13 shows that the resolvent of a Sturm–Liouville operator has a simple pole at every eigenvalue whose negative residue is precisely the projector onto the corresponding eigenspace. In particular, if we choose  $\gamma(t) := E_j + \varepsilon e^{2\pi i t}$  with  $\varepsilon$  so small that the interior of  $\gamma$  contains no eigenvalues other than  $E_j$ , then  $P = \langle u_j, \cdot \rangle u_j$ .  $\diamond$

**Theorem 8.17.** Suppose  $A \in \mathcal{C}(X)$ . Let  $\gamma : [0, 1] \rightarrow \rho(A)$  be a Jordan curve and denote by  $S_1, S_2$  its interior, exterior, respectively. Then the operators in the splitting from Lemma 8.15 satisfy

$$\sigma(A_1) = \sigma(A) \cap S_1, \quad \text{and} \quad \sigma(A_2) = \sigma(A) \cap S_2. \quad (8.37)$$

Moreover,  $A_1$  is bounded.

**Proof.** Note that for  $\alpha$  in the exterior of  $\gamma$  an analogous calculation (using the first resolvent identity) as in the previous lemma shows

$$R_A(\alpha)P = \frac{1}{2\pi i} \oint_{\gamma} R_A(\alpha) \frac{d\beta}{\alpha - \beta},$$

where the right-hand side is analytic in the exterior of  $\gamma$ . Similarly, for  $\alpha$  in the interior of  $\gamma$  we have

$$R_A(\alpha)(1 - P) = -\frac{1}{2\pi i} \oint_{\gamma} R_A(\alpha) \frac{d\beta}{\alpha - \beta},$$

where the right-hand side is analytic in the interior of  $\gamma$ . Hence Problem 8.26 implies  $S_2 \subset \rho(A_1)$  and  $S_1 \subset \rho(A_1)$ . Since also have  $\rho(A_j) \subset \rho(A)$  by Problem 8.25 claim about the spectra follows.

To see that  $AP$  is bounded note that for  $x \in X$  and  $\ell \in \mathfrak{D}(A')$  we have

$$(A'\ell)(Px) = -\frac{1}{2\pi i} \oint_{\gamma} \ell(AR_A(\alpha)x) d\alpha,$$

which shows that  $AP$  is bounded since  $AR_A(\alpha) = 1 + \alpha R_A(\alpha)$  is.  $\square$

**Problem\* 8.26.** Let  $U \subseteq \mathbb{C}$  be open and connected and  $R : U \rightarrow \mathcal{L}(X)$  is weakly holomorphic (i.e.  $\alpha \rightarrow \ell(R(\alpha)x)$  is holomorphic for every  $\ell \in X^*$  and  $x \in X$ ). If  $A \in \mathcal{C}(X)$  satisfies  $U_0 \subseteq \rho(A)$ , where  $U_0$  is nonempty and open, and if  $R(\alpha) = R_A(\alpha)$  for  $\alpha \in U_0$ , then we have  $U \subseteq \rho(A)$  and  $R(\alpha) = R_A(\alpha)$  for  $\alpha \in U$ .

#### 8.4. Relatively bounded and relatively compact operators

In many applications operators have the structure  $A + B$ , where  $A$  is well-understood and  $B$  can be considered *small* with respect to  $A$ . This raises the questions about conditions on  $B$  which ensure that certain nice properties of  $A$  persist when adding  $B$ . In this section we discuss two basic conditions.

Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$ . An operator  $B : \mathfrak{D}(B) \subseteq X \rightarrow Y$  is called  **$A$ -bounded** or **relatively bounded** with respect to  $A$  if  $\mathfrak{D}(A) \subseteq \mathfrak{D}(B)$  and if there are constants  $a, b \geq 0$  such that

$$\|Bx\| \leq a\|Ax\| + b\|x\|, \quad x \in \mathfrak{D}(A). \quad (8.38)$$

The infimum of all constants  $a$  for which a corresponding  $b$  exists such that (8.38) holds is called the  **$A$ -bound** of  $B$ . The set of all  $A$ -bounded operators forms a vector space denoted by  $\mathcal{L}_A(X, Y)$  (Problem 8.27). Note that  $B$  is relatively bounded if and only if it is bounded with respect to the graph norm, that is  $\mathcal{L}_A(X, Y) \cong \mathcal{L}([\mathfrak{D}(A)], Y)$  (if we identify  $B$  with its restriction to  $\mathfrak{D}(A)$ ), where we have written  $[\mathfrak{D}(A)]$  to indicate that  $\mathfrak{D}(A)$  is equipped with the graph norm of  $A$ .

It is important to emphasize that the  $A$ -bound will in general not be attained, since  $b$  will typically increase as  $a$  approaches the  $A$ -bound. Moreover, this concept is only of interest for unbounded operators, since if  $A$  is bounded, then the relatively bounded operators are just the bounded operators.

**Example 8.17.** Let  $X := \ell^2(\mathbb{N})$  and  $(Aa)_j := j a_j$  with  $\mathfrak{D}(A) := \{a \in \ell^2(\mathbb{N}) \mid (ja_j)_{j \in \mathbb{N}} \in \ell^2(\mathbb{N})\}$  and  $Ba := (\sum_{j \in \mathbb{N}} a_j) \delta^1$ . Then we have seen that  $A$  is closed (Example 8.3) while  $B$  is not closable (Example 8.4). (Compare also Problem 8.8.)

Now

$$\sum_{j \in \mathbb{N}} |a_j| = \sum_{j \leq n} |a_j| + \sum_{j > n} \frac{1}{j} \cdot |ja_j| \leq \sqrt{n} \|a\|_2 + \sqrt{\sum_{j > n} \frac{1}{j}} \|Aa\|_2$$

shows that  $B$  is  $A$ -bounded with  $A$ -bound zero.  $\diamond$

There is also an important equivalent characterization in terms of the resolvent.

**Lemma 8.18.** *Suppose  $A \in \mathcal{C}(X)$  is closed with nonempty resolvent set. Then  $B \in \mathcal{L}_A(X)$  if and only if  $BR_A(\alpha) \in \mathcal{L}(X)$  for one (and hence for all)  $\alpha \in \rho(A)$ .*

*Moreover, the  $A$ -bound of  $B$  is no larger than  $\inf_{\alpha \in \rho(A)} \|BR_A(\alpha)\|$ .*

**Proof.** This follows since  $R_A(\alpha) : X \rightarrow [\mathfrak{D}(A)]$  is a homeomorphism by the inverse mapping theorem. In fact, this can also be seen directly since  $y = (A - \alpha)x$  implies  $\|R_A(\alpha)y\| \leq \|R_A(\alpha)\|\|Ax\| + |\alpha|\|x\|$ .

Moreover, the same idea shows that for  $x \in \mathfrak{D}(A)$  we have

$$\|Bx\| = \|BR_A(\alpha)(A - \alpha)x\| \leq a\|(A - \alpha)x\| \leq a\|Ax\| + (a|\alpha|)\|x\|,$$

where  $a := \|BR_A(\alpha)\|$ . Finally, note that if  $BR_A(\alpha)$  is bounded for one  $\alpha \in \rho(A)$ , it is bounded for all  $\alpha \in \rho(A)$  by the first resolvent formula.  $\square$

Note that if  $B$  is closable and  $A$  has nonempty resolvent set, then it is relatively bounded if and only if  $\mathfrak{D}(A) \subseteq \mathfrak{D}(B)$ . In fact, in this case Problem 8.12 implies that  $BR_A(\alpha)$  is closed and hence bounded by the closed graph theorem (Theorem 4.7).

Now we come to our first application.

**Lemma 8.19.** *Suppose  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$  and  $B \in \mathcal{L}_A(X, Y)$  with  $A$ -bound less than one. Then the graph norms of  $A$  and  $A + B$  are equivalent. In particular, in this case  $A + B$  is closable/closed if and only if  $A$  is.*

**Proof.** From  $\|Ax\| \leq \|(A + B)x\| + \|Bx\| \leq \|(A + B)x\| + a\|Ax\| + b\|x\|$  we obtain

$$\|Ax\| \leq \frac{1}{1-a}\|(A + B)x\| + \frac{b}{1-a}\|x\|$$

which shows that the graph norms of  $A$  and  $A + B$  are equivalent.  $\square$

If  $A$  is closable, then the fact that  $B$  is  $A$ -bounded always implies that there is an extension  $\bar{B}$  to  $\mathfrak{D}(\bar{A})$  given by extending  $B$  with respect to the graph norm using continuity (i.e., as in Theorem 1.16). Moreover, note that (8.38) continues to hold for  $\bar{A}$  and  $\bar{B}$  (with the same constants). This implies that we can assume  $A$  to be closed without loss of generality and that it suffices to verify (8.38) on a core of  $A$ . However, while  $\bar{A} + \bar{B}$  is bounded (and hence closed) on  $[\mathfrak{D}(\bar{A})]$ , it might not be closed on  $X$  since the graph norm of  $\bar{A} + \bar{B}$  might be weaker than the graph norm of  $\bar{A}$ . Conversely, if the graph norms are equivalent, there is no need that  $B$  has a closed extension in  $X$ . The following two examples illustrate these observations.

**Example 8.18.** If  $A$  is unbounded and  $B = -A$ , then  $A + B = 0$  and hence the condition that the  $A$ -bound is less than one cannot be dropped in general.  $\diamond$

**Example 8.19.** Consider the operators from Example 8.17. Then, since  $A$  is closed, so is  $A + B$  by our lemma (which can also be verified directly, cf. Problem 8.8). Nevertheless, note that  $B$  alone is not closable.  $\diamond$

**Example 8.20.** Let  $1 \leq p < \infty$ . Consider the differentiation operator  $A : W^{1,p}(0,1) \subset L^p(0,1) \rightarrow L^p(0,1)$  from Example 8.5 and let  $B : C[0,1] \subset L^p(0,1) \rightarrow L^p(0,1)$  be the point evaluation at 0, that is  $(Bf)(x) = f(0)$ . Then by Problem 8.6  $B$  is relatively bounded and satisfies (8.38) with  $a = b = 1$  (note that Hölder's inequality (Problem 1.27) implies  $\|f\|_1 \leq \|f\|_p$ ). This is clearly not good enough to apply our lemma and hence we need to work a bit harder in this case. If  $1 < p < \infty$  we can invoke Hölder's inequality to obtain

$$|f(0)| = \left| f(x) - \int_0^x f'(y) dy \right| \leq |f(x)| + x^{1/q} \|f'\|_p,$$

where  $q = \frac{p}{p-1}$ . Integrating from 0 to  $\varepsilon$  with respect to  $x$  further shows (using  $\|f\|_1 \leq \|f\|_p$ )

$$|f(0)| \leq \frac{q\varepsilon^{1/q}}{1+q} \|f'\|_p + \frac{1}{\varepsilon} \|f\|_p.$$

Hence the relative bound is 0 for  $1 < p < \infty$  and our lemma applies. In the case  $p = 1$  this argument breaks down. In fact, it turns out that for  $p = 1$  the  $A$ -bound is one. To see this let  $f_n(x) := \max(1 - nx, 0)$  such that  $\|f_n\|_1 = \frac{1}{2n}$  and  $\|f'_n\|_1 = |f_n(0)| = 1$ . Then letting  $n \rightarrow \infty$  in

$$1 = |f_n(0)| = \|Bf_n\|_1 \leq a\|Af_n\|_1 + b\|f_n\|_1 = a + \frac{b}{2n}$$

shows  $a \geq 1$  and hence the  $A$  bound of  $B$  is one and our lemma does not apply. We will come back to this case below.  $\diamond$

Next, there is also a convenient formula for the resolvent of  $A + B$ .

**Lemma 8.20.** Let  $A, B$  be two given operators with  $\mathfrak{D}(A) \subseteq \mathfrak{D}(B)$  such that  $A$  and  $A + B$  are closed. Then we have the **second resolvent formula**

$$R_{A+B}(\alpha) - R_A(\alpha) = -R_A(\alpha)BR_{A+B}(\alpha) = -R_{A+B}(\alpha)BR_A(\alpha) \quad (8.39)$$

for  $\alpha \in \rho(A) \cap \rho(A + B)$ .

**Proof.** We abbreviate  $C := A + B$  and compute

$$R_C(\alpha) + R_A(\alpha)BR_C(\alpha) = R_A(\alpha)(C - \alpha)R_C(\alpha) = R_A(\alpha). \quad \square$$

Moreover, we also give a criterion for a point to be in the resolvent set of  $A + B$ .



**Lemma 8.21.** *Let  $A \in \mathcal{C}(X)$  and  $B \in \mathcal{L}_A(X)$  satisfying (8.38). If  $\alpha \in \rho(A)$  and  $\|BR_A(\alpha)\| < 1$ , then  $\alpha \in \rho(A+B)$  with*

$$R_{A+B}(\alpha) = R_A(\alpha)(1 - BR_A(\alpha))^{-1}. \quad (8.40)$$

*This condition holds in particular if*

$$a\|AR_A(\alpha)\| + b\|R_A(\alpha)\| < 1. \quad (8.41)$$

**Proof.** Since we can write  $A+B+\alpha = (1+BR_A(\alpha))(A-\alpha)$  the first claim follows from Problem 8.28. The second follows from

$$\|BR_A(\alpha)x\| \leq a\|AR_A(\alpha)x\| + b\|R_A(\alpha)x\|. \quad \square$$

Now we turn to the second condition: Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$ . An operator  $B : \mathfrak{D}(B) \subseteq X \rightarrow Y$  is called **A-compact** or **relatively compact** with respect to  $A$  if  $\mathfrak{D}(A) \subseteq \mathfrak{D}(B)$  and whenever  $x_n$  and  $Ax_n$  are bounded sequences, then  $Bx_n$  has a convergent subsequence.

The set of all  $A$ -compact operators forms a vector space denoted by  $\mathcal{K}_A(X, Y)$ . As before, note that  $B$  is relatively compact if and only if it is compact with respect to the graph norm, that is  $\mathcal{K}_A(X, Y) \cong \mathcal{K}([\mathfrak{D}(A)], Y)$  (if we identify  $B$  with its restriction to  $\mathfrak{D}(A)$ ). Of course every compact operator is relatively bounded and if the embedding  $[\mathfrak{D}(A)] \hookrightarrow X$  is compact, then every bounded operator is relatively compact.

**Example 8.21.** Consider the operators from Example 8.17 or 8.20. Then,  $B$  is relatively compact since its range is one-dimensional. In fact, the same is true for any relatively bounded operator  $B$  whose range is finite dimensional (cf. Problem 8.30).  $\diamond$

**Example 8.22.** Consider the differential operator  $A : W^{1,p}(0, 1) \subseteq L^p(0, 1) \rightarrow L^p(0, 1)$ ,  $1 < p < \infty$  from Example 8.5. Then since the embedding  $W^{1,p}(0, 1) \subseteq L^p(0, 1)$  is compact (Problem 8.7), every bounded operator  $B \in \mathcal{L}(L^p(0, 1))$  is relatively compact. For example one can choose  $B$  to be a multiplication operator by a bounded function.  $\diamond$

Again there is an equivalent characterization in terms of resolvents.

**Lemma 8.22.** *Suppose  $A \in \mathcal{C}(X)$  is closed with nonempty resolvent set. Then  $B \in \mathcal{K}_A(X)$  if and only if  $BR_A(\alpha) \in \mathcal{L}(X)$  for one (and hence for all)  $\alpha \in \rho(A)$ .*

**Proof.** Again this follows since  $R_A(\alpha) : X \rightarrow [\mathfrak{D}(A)]$  is a homeomorphism. Moreover, if  $BR_A(\alpha)$  is compact for one  $\alpha \in \rho(A)$ , it is compact for all  $\alpha \in \rho(A)$  by the first resolvent formula.  $\square$

As already pointed out before, if  $A$  is closable, there is a unique extension  $\bar{B}$  of  $B$  to  $\mathfrak{D}(\bar{A})$  using the graph norm. Moreover, if  $B$  is relatively compact

with respect to  $A$ , then  $\bar{B}$  is relatively compact with respect to  $\bar{A}$ . This follows from Lemma 3.3 upon using the graph norm on  $\mathfrak{D}(A)$ . Moreover, in this case  $\bar{A} + \bar{B}$  is closed without restriction on the  $A$ -bound.

**Lemma 8.23.** *Suppose  $A \in \mathcal{C}(X, Y)$  and  $B \in \mathcal{K}_A(X, Y)$ . Then  $A + B$  is closed and  $B$  is also relatively compact with respect to  $A + B$ .*

**Proof.** Abbreviate  $C := A + B$ . We first show that  $B$  is relatively compact with respect to  $C$ . If it were not, we could find a bounded sequence  $x_n \in \mathfrak{D}(A)$  such that  $Cx_n$  is bounded but  $Ax_n$  is unbounded. After dropping some terms we can assume  $\|Ax_n\| \rightarrow \infty$ . Then  $\tilde{x}_n := \|Ax_n\|^{-1}x_n$  is a null sequence and so is  $C\tilde{x}_n$ . Moreover, since  $B$  is  $A$ -compact we can assume that  $B\tilde{x}_n \rightarrow y$  converges. Hence  $A\tilde{x}_n = (C - B)\tilde{x}_n \rightarrow -y$  implying  $y = 0$  as  $A$  is closed, contradicting  $\|A\tilde{x}_n\| = 1$ .

Now let  $x_n \in \mathfrak{D}(A)$  be a sequence with  $x_n \rightarrow x$  and  $Cx_n \rightarrow y$ . Since by the first part  $B$  is  $C$ -compact, we can pass to a subsequence such that also  $Bx_n \rightarrow z$ . Consequently  $Ax_n = (C - B)x_n \rightarrow y - z$  implying  $x \in \mathfrak{D}(A)$  and  $Ax = y - z$  since  $A$  is closed. Moreover, since  $B$  is  $A$ -bounded we also have  $Bx = z$ . Thus we have  $x \in \mathfrak{D}(C) = \mathfrak{D}(A)$  and  $Cx = y = Ax + Bx$  which shows that  $C$  is closed.  $\square$

**Example 8.23.** Consider once more the operators from Example 8.17 or 8.20. Again our lemma implies that  $A + B$  is closed and this time we don't even need the computation of the  $A$ -bound. Moreover, in the case of Example 8.20 this now covers the full range  $1 \leq p < \infty$ .  $\diamond$

Applications of these notions will be given in the next section.

**Problem 8.27.** *Suppose  $B_j$ ,  $j = 1, 2$ , are  $A$  bounded with respective  $A$ -bounds  $a_i$ ,  $i = 1, 2$ . Show that  $\alpha_1 B_1 + \alpha_2 B_2$  is also  $A$  bounded with  $A$ -bound less than  $|\alpha_1|a_1 + |\alpha_2|a_2$ .*

**Problem\* 8.28.** *Suppose  $A$  is closed and  $B$  satisfies  $\mathfrak{D}(A) \subseteq \mathfrak{D}(B)$ :*

- *Show that  $1 + B$  has a bounded inverse if  $\|B\| < 1$ .*
- *Suppose  $A$  has a bounded inverse. Then so does  $A + B$  if  $\|BA^{-1}\| < 1$ . In this case we have  $\|(A + B)^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|BA^{-1}\|}$ .*

**Problem 8.29.** *Suppose  $A, B$  are linear operators with  $\mathfrak{D}(A) \subseteq \mathfrak{D}(B)$  and  $\|Bx\| \leq a\|Ax\|^\alpha\|x\|^{1-\alpha}$  for all  $x \in \mathfrak{D}(A)$  and some  $a > 0$ ,  $\alpha \in (0, 1)$ . Then  $B$  is relatively bounded with  $A$ -bound 0. (Hint: Young's inequality (1.24).)*

**Problem 8.30.** *Show that a relatively bounded operator  $B \in \mathcal{L}_A(X, Y)$  with finite dimensional range is of the form*

$$Bx = \sum_{j=1}^n (x'_j(x) + y'_j(Ax))y_j$$

for some  $x'_j \in X^*$ ,  $y'_j \in Y^*$ , and  $y_j \in Y$ ,  $1 \leq j \leq n$ . Moreover, any such operator is relatively compact. (Hint: Start with the finite rank operator  $\hat{B} \in \mathcal{L}([\mathfrak{D}(A)], Y)$  and use  $[\mathfrak{D}(A)] \cong \Gamma(A) \subset X \oplus Y$ .)

**Problem 8.31.** Suppose  $B : \mathfrak{D}(B) \subseteq X \rightarrow Y$  is relatively compact with respect to  $A : \mathfrak{D}(A) \subseteq X \rightarrow Y$ . Then the  $A$ -bound of  $B$  is zero if either

- (i)  $B$  is closable, or
- (ii)  $A$  is closed and  $Y$  is reflexive.

(Hint: Argue by contradiction. If the claim were wrong, then for any  $a > 0$  one could find a sequence  $x_n \in \mathfrak{D}(A)$  such that  $\|x_n\|_A = 1$  and  $\|Bx_n\| \geq a\|Ax_n\| + n\|x_n\|$ . Since  $B \in \mathcal{L}([\mathfrak{D}(A)], Y)$  we must have  $x_n \rightarrow 0$ . Moreover, we can assume  $Bx_n \rightarrow y$  and we can get a contradiction if we can show  $y = 0$ . In the first case this is immediate. In the second case use that  $\Gamma(A)$  is weakly closed.)

### 8.5. Unbounded Fredholm operators

The definition of a Fredholm operator from Section 7.2 can be extended to closed operators in a straightforward manner: An operator  $A \in \mathcal{C}(X, Y)$  is called a **Fredholm operator** if both its kernel and cokernel are finite dimensional. In this case we define its **index** as

$$\text{ind}(A) := \dim \text{Ker}(A) - \dim \text{Coker}(A). \quad (8.42)$$

In fact, many results for bounded Fredholm operators carry over to the unbounded case using the following observation: Let us denote  $A$  regarded as an operator from  $\mathfrak{D}(A)$  equipped with the graph norm  $\|\cdot\|_A$  to  $Y$  by  $\hat{A}$ . Recall that  $\hat{A}$  is bounded (cf. Problem 8.5). Moreover,  $\text{Ker}(\hat{A}) = \text{Ker}(A)$  and  $\text{Ran}(\hat{A}) = \text{Ran}(A)$ . Consequently,  $A \in \mathcal{C}(X, Y)$  is Fredholm if and only if  $\hat{A}$  is.

For example, applying Lemma 7.8 to  $\hat{A}$  shows that a closed operator with finite cokernel has closed range. In particular, a Fredholm operator has closed range.

**Example 8.24.** Consider the operator  $A$  from Example 8.12. There we have seen  $\dim \text{Ker}(A - \alpha) = 1$  for every  $\alpha \in \mathbb{C}$ . Moreover, the solution of the inhomogeneous differential equation  $f' - \alpha f = g$  is given by

$$f(x) = f(0)e^{\alpha x} + \int_0^x e^{\alpha(x-y)} g(y) dy,$$

which shows  $\text{Ran}(A - \alpha) = X$ . Thus  $A - \alpha$  is Fredholm with  $\text{ind}(A - \alpha) = 1$ .  $\diamond$

Literally following the proof of Lemma 7.9 (using Lemma 8.6 and Theorem 8.8) gives

**Lemma 8.24.** *Suppose  $A \in \mathcal{C}(X, Y)$  with  $\text{Ran}(A)$  closed. Then*

$$\text{Ker}(A') \cong \text{Coker}(A)^*, \quad \text{Coker}(A') \cong \text{Ker}(A)^*. \quad (8.43)$$

*In particular, we have*

$$\dim \text{Ker}(A') = \dim \text{Coker}(A), \quad \dim \text{Ker}(A) = \dim \text{Coker}(A'). \quad (8.44)$$

And hence Riesz' theorem continues to hold:

**Theorem 8.25.** *A closed operator  $A$  is Fredholm if and only if  $A'$  is and in this case*

$$\text{ind}(A') = -\text{ind}(A). \quad (8.45)$$

Applying Dieudonné's and Yood's theorem to  $\hat{A}$  gives:

**Theorem 8.26.** *Suppose  $A \in \mathcal{C}(X, Y)$  is Fredholm. Then there is an  $\varepsilon > 0$  such that for every  $B \in \mathcal{L}_A(X, Y)$  with norm smaller than  $\varepsilon$ , we have that  $A + B$  is Fredholm with index  $\text{ind}(A + B) = \text{ind}(A)$ . Similarly, for every  $B \in \mathcal{K}_A(X, Y)$  the operator  $A + B$  is Fredholm with index  $\text{ind}(A + B) = \text{ind}(A)$ .*

The structure of Fredholm operators is essentially the same as in the bounded case. As discussed in Section 7.2 there are closed subspaces  $X_0 \subseteq X$  and  $Y_0 \subseteq Y$  such that  $X = \text{Ker}(A) \dot{+} X_0$  and  $Y = Y_0 \dot{+} \text{Ran}(A)$ , respectively. Moreover, we have corresponding projections  $P \in \mathcal{L}(X)$  with  $\text{Ran}(P) = \text{Ker}(A)$  and  $Q \in \mathcal{L}(Y)$  with  $\text{Ran}(Q) = Y_0$ . With respect to the decomposition  $\text{Ker}(A) \oplus X_0 \rightarrow Y_0 \oplus \text{Ran}(A)$  our Fredholm operator is given by

$$A = \begin{pmatrix} 0 & 0 \\ 0 & A_0 \end{pmatrix}, \quad (8.46)$$

where  $A_0$  is the restriction of  $A$  to  $\mathfrak{D}(A) \cap X_0 \rightarrow \text{Ran}(A)$ . By construction  $A_0$  is closed and surjective implying that it has a bounded inverse (Corollary 8.2). Defining

$$B := \begin{pmatrix} 0 & 0 \\ 0 & A_0^{-1} \end{pmatrix} \quad (8.47)$$

we get

$$AB = \mathbb{I} - Q, \quad BA = (\mathbb{I} - P)|_{\mathfrak{D}(A)} \quad (8.48)$$

and hence  $A$  is invertible up to finite rank operators. The converse can be shown as in the bounded case which gives the analog of Theorem 7.14:

**Theorem 8.27.** *An operator  $A \in \mathcal{C}(X, Y)$  is Fredholm if and only if there exist  $B_1, B_2 \in \mathcal{L}(Y, X)$  such that  $\overline{B_1 A - \mathbb{I}} \in \mathcal{K}(X)$  and  $AB_2 - \mathbb{I} \in \mathcal{K}(Y)$ .*

Moreover, the essential and Fredholm spectrum can be defined as in the bounded case and the same argument shows that

$$\sigma_{\text{ess}}(A) = \bigcap_{K \in \mathcal{K}(X)} \sigma(A + K). \quad (8.49)$$

In particular, Theorem 7.16 continues to hold.

**Theorem 8.28** (Weyl). *Let  $A \in \mathcal{C}(X)$ , then*

$$\sigma_\Phi(A + K) = \sigma_\Phi(A), \quad \sigma_{ess}(A + K) = \sigma_{ess}(A), \quad K \in \mathcal{K}_A(X). \quad (8.50)$$

**Example 8.25.** Consider again the operator  $A$  from Example 8.12. Since  $A - \alpha$  is Fredholm for every  $\alpha$  we have  $\sigma_{ess}(A) = \emptyset$ . Moreover, since the embedding  $C^1[0, 1] \hookrightarrow C[0, 1]$  is compact, every bounded operator  $B \in \mathcal{L}(C[0, 1])$  is relatively compact and we have  $\sigma_{ess}(A + B) = \sigma_{ess}(A) = \emptyset$  in this case.  $\diamond$

We end this section with a useful criterion extending Lemma 8.11. To this end we call a Weyl sequence  $x_n$  a **singular Weyl sequence** if in addition  $x_n$  has no convergent subsequence.

**Lemma 8.29** (Weyl's criterion). *We have  $\alpha \in \sigma_\Phi(A)$  if there is a singular Weyl sequence  $x_n \in \mathfrak{D}(A)$  such that  $\|x_n\| = 1$ ,  $\|(A - \alpha)x_n\| \rightarrow 0$ , and  $x_n$  has no convergent subsequence.*

**Proof.** By Lemma 8.11 we have  $\alpha \in \sigma(A)$  and hence we need to show that  $A - \alpha$  is not Fredholm. If  $\dim(\text{Ker}(A - \alpha)) = \infty$  there is nothing to do. Hence we can assume  $\dim(\text{Ker}(A - \alpha)) < \infty$  and there is a projector  $P = P^2 \in \mathcal{L}(X)$  with  $\text{Ker}(P) = \text{Ker}(A - \alpha)$ . Consider  $y_n := (1 - P)x_n$ . Then  $(A - \alpha)y_n = (A - \alpha)x_n \rightarrow 0$  and  $y_n$  has no convergent subsequence. In indeed, if it had a convergent subsequence, then by compactness of  $P$ , there is a further subsequence which works for both  $y_n$ ,  $Px_n$ , and hence also for  $x_n$ , a contradiction. Moreover,  $\|y_n\| = \|x_n - Px_n\| \rightarrow 1$  and hence we can assume  $\|y_n\| = 1$  after rescaling. Moreover, since  $X \cong \text{Ker}(A - \alpha) \oplus X_0$ , where  $X_0 := \text{Ran}(P)$  and  $X/\text{Ker}(A - \alpha) \cong X_0$  we see that  $\text{dist}(y_n, \text{Ker}(A - \alpha)) \geq \varepsilon > 0$  and hence the  $\text{Ran}(A - \alpha)$  is not closed by Corollary 8.4.  $\square$

Note that one way of ensuring that  $x_n$  has no convergent subsequence is to require  $x_n \rightharpoonup 0$ . For example, in a Hilbert space one could try to choose the elements orthogonal.

**Example 8.26.** Let us take  $X := L^p(\mathbb{R})$ ,  $1 \leq p < \infty$ , and start with  $A_0 = \frac{1}{i} \frac{d}{dx}$  on  $\mathfrak{D}(A_0) = C_c^\infty(\mathbb{R})$ . Then the closure  $A := \overline{A_0}$  is defined on  $\mathfrak{D}(A) = W^{1,p}(\mathbb{R})$ , where the Sobolev space  $W^{1,p}(\mathbb{R})$  is defined as the closure of  $C_c^\infty(\mathbb{R})$  with respect to the norm

$$\|f\|_{1,p}^p = \int_{\mathbb{R}} (|f'(x)|^p + |f(x)|^p) dx.$$

(Similarly  $L^p(\mathbb{R})$  can be understood as the closure of  $C_c^\infty(\mathbb{R})$  with respect to the norm  $\|f\|_p^p = \int_{\mathbb{R}} |f(x)|^p dx$ .) Then any eigenfunction must be of the form

$$f(x) = f(0)e^{-i\alpha x}$$

and such a function will never be square integrable unless  $f(0) = 0$ . However, for  $\alpha \in \mathbb{R}$  this function is at least bounded and we could try to get approximating eigenfunctions by restricting  $e^{-i\alpha x}$  to compact intervals. More precisely, choose  $\varphi_n \in C_c^\infty(\mathbb{R})$  such that  $\varphi_n$  is symmetric with  $\varphi_n(x) = 1$  for  $0 \leq x \leq n$ ,  $\varphi_n(x) = 1$  for  $x \geq n+1$  and such that the piece on  $[n, n+1]$  is independent of  $n$ . Set  $u_n(x) := \varphi_n(x)e^{-i\alpha x}$ . Then  $\|u_n\|_p = (2(C_0 + n))^{1/p}$  while  $\|(A - \alpha)u_n\|_p = (2C_1)^{1/p}$  and thus  $u_n/\|u_n\|_p$  is a Weyl sequence. Since the same is true for  $u_n(x + r_n)$  we get a singular Weyl sequence by choosing  $r_n$  such that the supports of  $\varphi_n$  and  $\varphi_m$  are disjoint for  $n \neq m$ . Hence we conclude  $\mathbb{R} \subset \sigma_\Phi(A)$ .

If  $\alpha \in \mathbb{C} \setminus \mathbb{R}$  we can write down the resolvent of  $A$ . To this end note that the solution of the inhomogeneous equation  $-if' - \alpha f = g$  is given by

$$f(x) = f(a)e^{i\alpha(x-a)} + i \int_a^x e^{i\alpha(x-y)} g(y) dy.$$

If  $g$  has compact support, we can shift  $a$  beyond the support of  $g$  and hence we expect

$$R_A(\alpha)g(x) := i \int_{\pm\infty}^x e^{i\alpha(x-y)} g(y) dy, \quad \mp \operatorname{Im}(\alpha) > 0.$$

Note that here we have shifted  $a$  to the left/right of the support depending on the sign of  $\operatorname{Im}(\alpha)$  such that the above expression at least formally makes sense without the compact support assumption. Then Young's inequality for convolutions implies

$$\|R_A(\alpha)g\|_p \leq \frac{1}{|\operatorname{Im}(\alpha)|} \|g\|_p,$$

which shows  $\alpha \in \rho(A)$  with  $\|R_A(\alpha)\| \leq |\operatorname{Im}(\alpha)|^{-1}$ .

In summary, we have  $\sigma(A) = \sigma_\Phi(A) = \sigma_{ess}(A) = \mathbb{R}$ .

Moreover, since the embedding  $W^{1,p}(a, b) \hookrightarrow L^p(a, b)$  is compact (Problem 8.7) and the embedding  $L^p(a, b) \hookrightarrow L^p(\mathbb{R})$  is bounded for every bounded interval  $(a, b)$  every multiplication operator  $Q$  with a function  $q \in L_c^\infty(\mathbb{R})$  (bounded with compact support) is relatively compact. Moreover, every bounded function  $q \in L_0^\infty(\mathbb{R})$  which vanishes as  $|x| \rightarrow \infty$  can be approximated by bounded functions with compact support  $q_n$  in the sup norm. Hence we also have  $Q_n \rightarrow Q$  in the operator norm, and for any  $q \in L_0^\infty(\mathbb{R})$  we have  $\sigma_\Phi(A + Q) = \sigma_{ess}(A + Q) = \mathbb{R}$ .

However, note that while the spectrum of  $A + Q$  could be different from the spectrum of  $A$ , this is not the case here. Indeed one can verify that the resolvent is given by

$$R_{A+Q}(\alpha)g(x) := i \int_{\pm\infty}^x e^{i\alpha(x-y) - i \int_y^x q(t) dt} g(y) dy, \quad \mp \operatorname{Im}(\alpha) > 0. \quad \diamond$$

**Problem 8.32.** Consider  $X := C[0, 1]$  and  $A = \frac{d}{dx}$  with  $\mathfrak{D}(A) = \{f \in C^1[0, 1] | f(0) = f(1) = 0\}$ . Investigate when  $A - \alpha$  is Fredholm and compute the essential spectrum of  $A$ .

---

*Part 3*

# Nonlinear Functional Analysis





# Analysis in Banach spaces

## 9.1. Single variable calculus in Banach spaces

As a warmup we will look at mappings from an interval to a Banach space. This case is somewhat simpler than the case of mappings between Banach spaces but nevertheless is sufficient for many applications.

Let  $X$  be a Banach space. Let  $I \subseteq \mathbb{R}$  be some interval and denote by  $C(I, X)$  the set of continuous functions from  $I$  to  $X$ . Given  $t \in I$  we call  $f : I \rightarrow X$  differentiable at  $t$  if the limit

$$\dot{f}(t) := \lim_{\varepsilon \rightarrow 0} \frac{f(t + \varepsilon) - f(t)}{\varepsilon} \quad (9.1)$$

exists. If  $t$  is a boundary point, the limit/derivative is understood as the corresponding one-sided limit/derivative.

The set of functions  $f : I \rightarrow X$  which are differentiable at all  $t \in I$  and for which  $\dot{f} \in C(I, X)$  is denoted by  $C^1(I, X)$ . Clearly  $C^1(I, X) \subset C(I, X)$ . As usual we set  $C^{k+1}(I, X) := \{f \in C^1(I, X) \mid \dot{f} \in C^k(I, X)\}$ . Note that if  $A \in \mathcal{L}(X, Y)$  and  $f \in C^k(I, X)$ , then  $Af \in C^k(I, Y)$  and  $\frac{d}{dt} Af = A\dot{f}$ .

The following version of the mean value theorem will be crucial.

**Theorem 9.1** (Mean value theorem). *Suppose  $f \in C^1(I, X)$ . Then*

$$\|f(t) - f(s)\| \leq M|t - s|, \quad M := \sup_{\tau \in [s, t]} \|\dot{f}(\tau)\|, \quad (9.2)$$

for  $s \leq t \in I$ .

**Proof.** Fix  $\tilde{M} > M$  and consider  $d(\tau) := \|f(\tau) - f(s)\| - \tilde{M}(\tau - s)$  for  $\tau \in [s, t]$ . Suppose  $\tau_0$  is the largest  $\tau$  for which  $d(\tau) \leq 0$  holds. Then there must be a sequence  $\varepsilon_n \downarrow 0$  such that

$$\begin{aligned} 0 < d(\tau_0 + \varepsilon_n) &\leq \|f(\tau_0 + \varepsilon_n) - f(\tau_0)\| - \tilde{M}\varepsilon_n + d(\tau_0) \\ &= \|\dot{f}(\tau_0)\varepsilon_n + o(\varepsilon_n)\| - \tilde{M}\varepsilon_n \leq (M - \tilde{M} + o(1))\varepsilon_n < 0. \end{aligned}$$

Taking  $n \rightarrow \infty$  contradicts our assumption.  $\square$

In particular,

**Corollary 9.2.** *For  $f \in C^1(I, X)$  we have  $\dot{f} = 0$  if and only if  $f$  is constant.*

Next we turn to integration. Let  $I := [a, b]$  be compact. A function  $f : I \rightarrow X$  is called a **step function** provided there are numbers

$$t_0 = a < t_1 < t_2 < \cdots < t_{n-1} < t_n = b \quad (9.3)$$

such that  $f(t)$  is constant on each of the open intervals  $(t_{j-1}, t_j)$ . The set of all step functions  $S(I, X)$  forms a linear space and can be equipped with the sup norm. The corresponding Banach space obtained after completion is called the set of **regulated functions**  $R(I, X)$ . In other words, a regulated function is the uniform limit of a step function.

Observe that  $C(I, X) \subset R(I, X)$ . In fact, consider the functions  $f_n := \sum_{j=0}^{n-1} f(t_j)\chi_{[t_j, t_{j+1})} \in S(I, X)$ , where  $t_j = a + j\frac{b-a}{n}$  and  $\chi$  is the characteristic function. Since  $f \in C(I, X)$  is uniformly continuous, we infer that  $f_n$  converges uniformly to  $f$ . Slightly more general, note that piecewise continuous functions are regulated since every piecewise continuous function is the sum of a continuous function and a step function.

For a step function  $f \in S(I, X)$  we can define a linear map  $\int : S(I, X) \rightarrow X$  by

$$\int_a^b f(t)dt := \sum_{j=1}^n x_j(t_j - t_{j-1}), \quad (9.4)$$

where  $x_i$  is the value of  $f$  on  $(t_{j-1}, t_j)$ . This map satisfies

$$\left\| \int_a^b f(t)dt \right\| \leq \|f\|_\infty(b-a). \quad (9.5)$$

and hence it can be extended uniquely to a linear map  $\int : R(I, X) \rightarrow X$  with the same norm  $(b-a)$ . We even have

$$\left\| \int_a^b f(t)dt \right\| \leq \int_a^b \|f(t)\|dt \quad (9.6)$$

since this holds for simple functions by the triangle inequality and hence for all functions by approximation.

We remark that it is possible to extend the integral to a larger class of functions in various ways. The first generalization is to replace step functions by simple functions (and at the same time one could also replace the Lebesgue measure on  $I$  by an arbitrary finite measure). Then the same approach defines the integral for uniform limits of simple functions. However, things only get interesting when you also replace the sup norm by an  $L^1$  type seminorm:  $\|f\|_1 := \int \|f(x)\| d\mu(x)$ . As before the integral can be extended to all functions which can be approximated by simple functions with respect to this seminorm. This is known as the Bochner integral and we refer to Section 5.5 from [48] for details.

In addition, if  $A \in \mathcal{L}(X, Y)$ , then  $f \in R(I, X)$  implies  $Af \in R(I, Y)$  and

$$A \int_a^b f(t) dt = \int_a^b A f(t) dt. \quad (9.7)$$

Again this holds for step functions and thus extends to all regulated functions by continuity. In particular, if  $\ell \in X^*$  is a continuous linear functional, then

$$\ell\left(\int_a^b f(t) dt\right) = \int_a^b \ell(f(t)) dt, \quad f \in R(I, X). \quad (9.8)$$

Moreover, we will use the usual conventions  $\int_{t_1}^{t_2} f(s) ds := \int_I \chi_{(t_1, t_2)}(s) f(s) ds$  and  $\int_{t_2}^{t_1} f(s) ds := -\int_{t_1}^{t_2} f(s) ds$ . Note that we could replace  $(t_1, t_2)$  by a closed or half-open interval with the same endpoints (why?) and hence  $\int_{t_1}^{t_3} f(s) ds = \int_{t_1}^{t_2} f(s) ds + \int_{t_2}^{t_3} f(s) ds$ .

**Theorem 9.3** (Fundamental theorem of calculus). *Suppose  $F \in C^1(I, X)$ , then*

$$F(t) = F(a) + \int_a^t \dot{F}(s) ds. \quad (9.9)$$

*Conversely, if  $f \in C(I, X)$ , then  $F(t) = \int_a^t f(s) ds \in C^1(I, X)$  and  $\dot{F}(t) = f(t)$ .*

**Proof.** Let  $f \in C(I, X)$  and set  $G(t) := \int_a^t f(s) ds$ . Then  $G \in C^1(I, X)$  with  $\dot{G}(t) = f(t)$  as can be seen from

$$\begin{aligned} \left\| \int_a^{t+\varepsilon} f(s) ds - \int_a^t f(s) ds - f(t)\varepsilon \right\| &= \left\| \int_t^{t+\varepsilon} (f(s) - f(t)) ds \right\| \\ &\leq |\varepsilon| \sup_{s \in [t, t+\varepsilon]} \|f(s) - f(t)\|. \end{aligned}$$

Hence if  $F \in C^1(I, X)$  then  $G(t) := \int_a^t (\dot{F}(s)) ds$  satisfies  $\dot{G} = \dot{F}$  and hence  $F(t) = C + G(t)$  by Corollary 9.2. Choosing  $t = a$  finally shows  $F(a) = C$ .  $\square$

**Problem\* 9.1** (Product rule). Let  $X$  be a Banach algebra. Show that if  $f, g \in C^1(I, X)$  then  $fg \in C^1(I, X)$  and  $\frac{d}{dt}fg = \dot{f}g + f\dot{g}$ .

**Problem 9.2.** Let  $X$  be a Banach algebra and  $\mathcal{G}(X)$  the group of invertible elements. Show that if  $f \in C^1(I, \mathcal{G}(X))$ , then  $f^{-1} \in C^1(I, X)$  with

$$\frac{d}{dt}f^{-1}(t) = -f^{-1}(t)\dot{f}(t)f^{-1}(t).$$

(Hint: Corollary 5.2)

**Problem\* 9.3.** Let  $f \in R(I, X)$  and  $\tilde{I} := I + t_0$ . then  $f(t - t_0) \in R(\tilde{I}, X)$  and

$$\int_I f(t)dt = \int_{\tilde{I}} f(t - t_0)dt.$$

**Problem\* 9.4.** Let  $A : \mathfrak{D}(A) \subseteq X \rightarrow X$  be a closed operator. Show that (9.7) holds for  $f \in C(I, X)$  with  $\text{Ran}(f) \subseteq \mathfrak{D}(A)$  and  $Af \in C(I, X)$ .

**Problem 9.5.** Let  $I = [a, b]$  and  $J = [c, d]$  be two compact intervals. Suppose  $f(s, t) : I \times J \rightarrow X$  is regulated in the sense that it is a uniform limit of step functions being constant on disjoint open rectangles  $(s_{j-1}, s_j) \times (t_{k-1}, t_k)$  whose closure cover  $I \times J$ . Show that

$$\int_J \left( \int_I f(s, t)ds \right) dt = \int_I \left( \int_J f(s, t)dt \right) ds.$$

(Hint: One way is to use linear functionals and reduce it to the classical Fubini theorem.)

## 9.2. Multivariable calculus in Banach spaces

We now turn to calculus in Banach spaces. Most facts will be similar to the situation of multivariable calculus for functions from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ . To emphasize this we will use  $|\cdot|$  for the norm in this section.

Let  $X$  and  $Y$  be two Banach spaces and let  $U$  be an open subset of  $X$ . Denote by  $C(U, Y)$  the set of continuous functions from  $U \subseteq X$  to  $Y$  and by  $\mathcal{L}(X, Y) \subset C(X, Y)$  the Banach space of bounded linear functions equipped with the operator norm

$$\|L\| := \sup_{|u|=1} |Lu|. \quad (9.10)$$

Then a function  $F : U \rightarrow Y$  is called differentiable at  $x \in U$  if there exists a linear function  $dF(x) \in \mathcal{L}(X, Y)$  such that

$$F(x + u) = F(x) + dF(x)u + o(u), \quad (9.11)$$

where  $o, O$  are the **Landau symbols**. Explicitly

$$\lim_{u \rightarrow 0} \frac{|F(x + u) - F(x) - dF(x)u|}{|u|} = 0. \quad (9.12)$$

The linear map  $dF(x)$  is called the **Fréchet derivative** of  $F$  at  $x$ . It is uniquely defined since if  $dG(x)$  were another derivative we had  $(dF(x) - dG(x))u = o(u)$  implying that for every  $\varepsilon > 0$  we can find a  $\delta > 0$  such that  $|(dF(x) - dG(x))u| \leq \varepsilon|u|$  whenever  $|u| \leq \delta$ . By homogeneity of the norm we conclude  $\|dF(x) - dG(x)\| \leq \varepsilon$  and since  $\varepsilon > 0$  is arbitrary  $dF(x) = dG(x)$ . Note that for this argument to work it is crucial that we can approach  $x$  from arbitrary directions  $u$ , which explains our requirement that  $U$  should be open.

If  $I \subseteq \mathbb{R}$ , we have an isomorphism  $\mathcal{L}(I, X) \equiv X$  and if  $F : I \rightarrow X$  we will write  $\dot{F}(t)$  instead of  $dF(t)$  if we regard  $dF(t)$  as an element of  $X$ . Clearly this is consistent with the definition (9.1) from the previous section.

**Example 9.1.** Let  $X$  be a Hilbert space and consider  $F : X \rightarrow \mathbb{R}$  given by  $F(x) := |x|^2$ . Then

$$F(x+u) = \langle x+u, x+u \rangle = |x|^2 + 2\operatorname{Re}\langle x, u \rangle + |u|^2 = F(x) + 2\operatorname{Re}\langle x, u \rangle + o(u).$$

Hence if  $X$  is a real Hilbert space, then  $F$  is differentiable with  $dF(x)u = 2\langle x, u \rangle$ . However, if  $X$  is a complex Hilbert space, then  $F$  is not differentiable.  $\diamond$

The previous example emphasizes that for  $F : U \subseteq X \rightarrow Y$  it makes a big difference whether  $X$  is a real or a complex Banach space. In fact, in case of a complex Banach space  $X$ , we obtain a version of complex differentiability which of course is much stronger than real differentiability. Note that in this respect it makes no difference whether  $Y$  is real or complex.

**Example 9.2.** Suppose  $f \in C^1(\mathbb{R})$  with  $f(0) = 0$ . Let  $X := \ell_{\mathbb{R}}^p(\mathbb{N})$ , then

$$F : X \rightarrow X, \quad (x_n)_{n \in \mathbb{N}} \mapsto (f(x_n))_{n \in \mathbb{N}}$$

is differentiable for every  $x \in X$  with derivative given by the multiplication operator

$$(dF(x)u)_n = f'(x_n)u_n.$$

First of all note that the mean value theorem implies  $|f(t)| \leq M_R|t|$  for  $|t| \leq R$  with  $M_R := \sup_{|t| \leq R} |f'(t)|$ . Hence, since  $\|x\|_{\infty} \leq \|x\|_p$ , we have  $\|F(x)\|_p \leq M_{\|x\|_{\infty}}\|x\|_p$  and  $F$  is well defined. This also shows that multiplication by  $f'(x_n)$  is a bounded linear map. To establish differentiability we use

$$f(t+s) - f(t) - f'(t)s = s \int_0^1 (f'(t+s\tau) - f'(t))d\tau$$

and since  $f'$  is uniformly continuous on every compact interval, we can find a  $\delta > 0$  for every given  $R > 0$  and  $\varepsilon > 0$  such that

$$|f'(t+s) - f'(t)| < \varepsilon \quad \text{if} \quad |s| < \delta, \quad |t| < R.$$

Now for  $x, u \in X$  with  $\|x\|_\infty < R$  and  $\|u\|_\infty < \delta$  we have  $|f(x_n + u_n) - f(x_n) - f'(x_n)u_n| < \varepsilon|u_n|$  and hence

$$\|F(x + u) - F(x) - dF(x)u\|_p < \varepsilon\|u\|_p$$

which establishes differentiability. Moreover, using uniform continuity of  $f$  on compact sets a similar argument shows that  $dF : X \rightarrow \mathcal{L}(X, X)$  is continuous (observe that the operator norm of a multiplication operator by a sequence is the sup norm of the sequence) and hence one writes  $F \in C^1(X, X)$  as usual.  $\diamond$

Differentiability implies existence of directional derivatives

$$\delta F(x, u) := \lim_{\varepsilon \rightarrow 0} \frac{F(x + \varepsilon u) - F(x)}{\varepsilon}, \quad \varepsilon \in \mathbb{R} \setminus \{0\}, \quad (9.13)$$

which are also known as **Gâteaux derivative** or **variational derivative**. Indeed, if  $F$  is differentiable at  $x$ , then (9.11) implies

$$\delta F(x, u) = dF(x)u. \quad (9.14)$$

In particular, we call  $F$  Gâteaux differentiable at  $x \in U$  if the limit on the right-hand side in (9.13) exists for all  $u \in X$ . However, note that Gâteaux differentiability does not imply differentiability. In fact, the Gâteaux derivative might be unbounded or it might even fail to be linear in  $u$ . Some authors require the Gâteaux derivative to be a bounded linear operator and in this case we will write  $\delta F(x, u) = \delta F(x)u$  but even this additional requirement does not imply differentiability in general. Note that in any case the Gâteaux derivative is homogenous, that is, if  $\delta F(x, u)$  exists, then  $\delta F(x, \lambda u)$  exists for every  $\lambda \in \mathbb{R}$  and

$$\delta F(x, \lambda u) = \lambda \delta F(x, u), \quad \lambda \in \mathbb{R}. \quad (9.15)$$

**Example 9.3.** The function  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $F(x, y) := \frac{x^3}{x^2 + y^2}$  for  $(x, y) \neq 0$  and  $F(0, 0) = 0$  is Gâteaux differentiable at 0 with Gâteaux derivative

$$\delta F(0, (u, v)) = \lim_{\varepsilon \rightarrow 0} \frac{F(\varepsilon u, \varepsilon v)}{\varepsilon} = F(u, v),$$

which is clearly nonlinear.

The function  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $F(x, y) = x$  for  $y = x^2$  and  $F(x, y) := 0$  else is Gâteaux differentiable at 0 with Gâteaux derivative  $\delta F(0) = 0$ , which is clearly linear. However,  $F$  is not differentiable.

If you take a linear function  $L : X \rightarrow Y$  which is unbounded, then  $L$  is everywhere Gâteaux differentiable with derivative equal to  $Lu$ , which is linear but, by construction, not bounded.  $\diamond$

**Example 9.4.** Let  $X := L^2_{\mathbb{R}}(0, 1)$  and consider

$$F : X \rightarrow X, \quad x \mapsto \sin(x).$$

First of all note that by  $|\sin(t)| \leq |t|$  our map is indeed from  $X$  to  $X$  and since sine is Lipschitz continuous we get the same for  $F$ :  $\|F(x) - F(y)\|_2 \leq \|x - y\|_2$ . Moreover,  $F$  is Gâteaux differentiable at  $x = 0$  with derivative given by

$$\delta F(0) = \mathbb{I}$$

but it is not differentiable at  $x = 0$ .

To see that the Gâteaux derivative is the identity note that

$$\lim_{\varepsilon \rightarrow 0} \frac{\sin(\varepsilon u(t))}{\varepsilon} = u(t)$$

pointwise and hence

$$\lim_{\varepsilon \rightarrow 0} \left\| \frac{\sin(\varepsilon u(\cdot))}{\varepsilon} - u(\cdot) \right\|_2 = 0$$

by dominated convergence since  $|\frac{\sin(\varepsilon u(t))}{\varepsilon}| \leq |u(t)|$ .

To see that  $F$  is not differentiable let

$$u_n = \pi \chi_{[0, 1/n]}, \quad \|u_n\|_2 = \frac{\pi}{\sqrt{n}}$$

and observe that  $F(u_n) = 0$ , implying that

$$\frac{\|F(u_n) - u_n\|_2}{\|u_n\|_2} = 1$$

does not converge to 0. Note that this problem does not occur in  $X := C[0, 1]$  (Problem 9.7).  $\diamond$

**Example 9.5.** Consider  $L^p(X, d\mu)$ ,  $1 \leq p < \infty$  and let  $G : \mathbb{C} \rightarrow \mathbb{R}$  be (real) differentiable with

$$|G(z)| \leq C|z|^p, \quad \sqrt{|\partial_x G(z)|^2 + |\partial_y G(z)|^2} \leq C|z|^{p-1}, \quad z = x + iy,$$

or, if  $\mu$  is finite,

$$|G(z)| \leq C(1 + |z|^p), \quad \sqrt{|\partial_x G(z)|^2 + |\partial_y G(z)|^2} \leq C(1 + |z|^{p-1}).$$

Note that the first condition comes for free from the second in the finite case and also in the general case if  $G(0) = 0$ . We only write down the estimates in the first case and leave the easy adaptations for the second case as an exercise.

Then

$$N(f) := \int_X G(f) d\mu$$

is Gâteaux differentiable and we have

$$\delta N(f)g = \int_X ((\partial_x G)(f)\operatorname{Re}(g) + (\partial_y G)(f)\operatorname{Im}(g)) d\mu.$$



In fact, by the chain rule  $h(\varepsilon) := G(f + \varepsilon g)$  is differentiable with  $h'(0) = (\partial_x G)(f)\operatorname{Re}(g) + (\partial_y G)(f)\operatorname{Im}(g)$ . Moreover, by the mean value theorem

$$\begin{aligned} \left| \frac{h(\varepsilon) - h(0)}{\varepsilon} \right| &\leq \sup_{0 \leq \tau \leq \varepsilon} \sqrt{(\partial_x G)(f + \tau g)^2 + (\partial_y G)(f + \tau g)^2} |g| \\ &\leq C 2^{p-1} (|f|^{p-1} + |g|^{p-1}) |g| \end{aligned}$$

and hence we can invoke dominated convergence to interchange differentiation and integration. Note that using Hölder's inequality this last estimate also shows local Lipschitz continuity on bounded domains:

$$|N(f) - N(g)| \leq C(\|f\|_p + \|g\|_p)^{p-1} \|f - g\|_p.$$

In particular, for  $1 < p < \infty$  the norm

$$N(f) := \int_X |f|^p d\mu$$

is Gâteaux differentiable with

$$\delta N(f)g = p \int_X |f|^{p-2} \operatorname{Re}(fg^*) d\mu. \quad \diamond$$

We will mainly consider Fréchet derivatives in the remainder of this chapter as it will allow a theory quite close to the usual one for multivariable functions. First of all we of course have linearity (which is easy to check):

**Lemma 9.4.** *Suppose  $F, G : U \rightarrow Y$  are differentiable at  $x \in U$  and  $\alpha, \beta \in \mathbb{C}$ . Then  $\alpha F + \beta G$  is differentiable at  $x$  with  $d(\alpha F + \beta G)(x) = \alpha dF(x) + \beta dG(x)$ . Similarly, if the Gâteaux derivatives  $\delta F(x, u)$  and  $\delta G(x, u)$  exist, then so does  $\delta(F + G)(x, u) = \delta F(x, u) + \delta G(x, u)$ .*

Next, Fréchet differentiability implies continuity:

**Lemma 9.5.** *Suppose  $F : U \rightarrow Y$  is differentiable at  $x \in U$ . Then  $F$  is continuous at  $x$ . Moreover, we can find constants  $M, \delta > 0$  such that*

$$|F(x + u) - F(x)| \leq M|u|, \quad |u| \leq \delta. \quad (9.16)$$

**Proof.** For every  $\varepsilon > 0$  we can find a  $\delta > 0$  such that  $|F(x + u) - F(x) - dF(x)u| \leq \varepsilon|u|$  for  $|u| \leq \delta$ . Now choose  $M = \|dF(x)\| + \varepsilon$ .  $\square$

**Example 9.6.** Note that this lemma fails for the Gâteaux derivative as the example of an unbounded linear function shows. In fact, it already fails in  $\mathbb{R}^2$  as the function  $F : \mathbb{R}^2 \rightarrow \mathbb{R}$  given by  $F(x, y) = 1$  for  $y = x^2 \neq 0$  and  $F(x, y) = 0$  else shows: It is Gâteaux differentiable at 0 with  $\delta F(0) = 0$  but it is not continuous since  $\lim_{\varepsilon \rightarrow 0} F(\varepsilon, \varepsilon^2) = 1 \neq 0 = F(0, 0)$ .  $\diamond$

If  $F$  is differentiable for all  $x \in U$  we call  $F$  **differentiable**. In this case we get a map

$$\begin{aligned} dF : U &\rightarrow \mathcal{L}(X, Y) \\ x &\mapsto dF(x) \end{aligned} \quad (9.17)$$

If  $dF : U \rightarrow \mathcal{L}(X, Y)$  is continuous, we call  $F$  continuously differentiable and write  $F \in C^1(U, Y)$ .

If  $X$  or  $Y$  has a (finite) product structure, then the computation of the derivatives can be reduced as usual. The following facts are simple and can be shown as in the case of  $X = \mathbb{R}^n$  and  $Y = \mathbb{R}^m$ .

Let  $Y := \times_{j=1}^m Y_j$  and let  $F : X \rightarrow Y$  be given by  $F = (F_1, \dots, F_m)$  with  $F_j : X \rightarrow Y_j$ . Then  $F \in C^1(X, Y)$  if and only if  $F_j \in C^1(X, Y_j)$ ,  $1 \leq j \leq m$ , and in this case  $dF = (dF_1, \dots, dF_m)$ . Similarly, if  $X = \times_{i=1}^n X_i$ , then one can define the **partial derivative**  $\partial_i F \in \mathcal{L}(X_i, Y)$ , which is the derivative of  $F$  considered as a function of the  $i$ -th variable alone (the other variables being fixed). We have  $dF u = \sum_{i=1}^n \partial_i F u_i$ ,  $u = (u_1, \dots, u_n) \in X$ , and  $F \in C^1(X, Y)$  if and only if all partial derivatives exist and are continuous.

**Example 9.7.** In the case of  $X = \mathbb{R}^n$  and  $Y = \mathbb{R}^m$ , the matrix representation of  $dF$  with respect to the canonical basis in  $\mathbb{R}^n$  and  $\mathbb{R}^m$  is given by the partial derivatives  $\partial_i F_j(x)$  and is called **Jacobi matrix** of  $F$  at  $x$ .  $\diamond$

Given  $F \in C^1(U, Y)$  we have  $dF \in C(U, \mathcal{L}(X, Y))$  and we can define the second derivative (provided it exists) via

$$dF(x + v) = dF(x) + d^2F(x)v + o(v). \quad (9.18)$$

In this case  $d^2F : U \rightarrow \mathcal{L}(X, \mathcal{L}(X, Y))$  which maps  $x$  to the linear map  $v \mapsto d^2F(x)v$  which for fixed  $v$  is a linear map  $u \mapsto (d^2F(x)v)u$ . Equivalently, we could regard  $d^2F(x)$  as a map  $d^2F(x) : X^2 \rightarrow Y$ ,  $(u, v) \mapsto (d^2F(x)v)u$  which is linear in both arguments. That is,  $d^2F(x)$  is a bilinear map  $X^2 \rightarrow Y$ . The corresponding norm on  $\mathcal{L}(X, \mathcal{L}(X, Y))$  explicitly spelled out reads

$$\|d^2F(x)\| = \sup_{|v|=1} \|d^2F(x)v\| = \sup_{|u|=|v|=1} \|(d^2F(x)v)u\|. \quad (9.19)$$

**Example 9.8.** Note that if  $F \in \mathcal{L}(X, Y)$ , then  $dF(x) = F$  (independent of  $x$ ) and  $d^2F(x) = 0$ .  $\diamond$

**Example 9.9.** Let  $X$  be a real Hilbert space and  $F(x) = |x|^2$ . Then we have already seen  $dF(x)u = 2\langle x, u \rangle$  and hence

$$dF(x + v)u = 2\langle x + v, u \rangle = 2\langle x, u \rangle + 2\langle v, u \rangle = dF(x)u + 2\langle v, u \rangle$$

which shows  $(d^2F(x)v)u = 2\langle v, u \rangle$ .  $\diamond$

**Example 9.10.** Suppose  $f \in C^2(\mathbb{R})$  with  $f(0) = 0$  and continue Example 9.2. Then we have  $F \in C^2(X, X)$  with  $d^2F(x)v$  the multiplication operator by the sequence  $f''(x_n)v_n$ , that is,

$$((d^2F(x)v)u)_n = f''(x_n)v_nu_n.$$

Indeed, arguing in a similar fashion we can find a  $\delta_1$  such that  $|f'(x_n + v_n) - f'(x_n) - f''(x_n)v_n| \leq \varepsilon|v_n|$  whenever  $\|x\|_\infty < R$  and  $\|v\|_\infty < \delta_1$ . Hence

$$\|dF(x+v) - dF(x) - d^2F(x)v\| < \varepsilon\|v\|_p$$

which shows differentiability. Moreover, since  $\|d^2F(x)\| = \|f''(x)\|_\infty$  one also easily verifies that  $F \in C^2(X, X)$  using uniform continuity of  $f''$  on compact sets.  $\diamond$

We can iterate the procedure of differentiation and write  $F \in C^r(U, Y)$ ,  $r \geq 1$ , if the  $r$ -th derivative of  $F$ ,  $d^rF$  (i.e., the derivative of the  $(r-1)$ -th derivative of  $F$ ), exists and is continuous. Note that  $d^rF(x)$  will be a multilinear map in  $r$  arguments as we will show below. Finally, we set  $C^\infty(U, Y) = \bigcap_{r \in \mathbb{N}} C^r(U, Y)$  and, for notational convenience,  $C^0(U, Y) = C(U, Y)$  and  $d^0F = F$ .

**Example 9.11.** Let  $X$  be a Banach algebra. Consider the multiplication  $M : X \times X \rightarrow X$ . Then

$$\partial_1 M(x, y)u = uy, \quad \partial_2 M(x, y)u = xu$$

and hence

$$dM(x, y)(u_1, u_2) = u_1y + xu_2.$$

Consequently  $dM$  is linear in  $(x, y)$  and hence

$$(d^2M(x, y)(v_1, v_2))(u_1, u_2) = u_1v_2 + v_1u_2$$

Consequently all differentials of order higher than two will vanish and in particular  $M \in C^\infty(X \times X, X)$ .  $\diamond$

If  $F$  is bijective and  $F, F^{-1}$  are both of class  $C^r$ ,  $r \geq 1$ , then  $F$  is called a **diffeomorphism** of class  $C^r$ .

For the composition of mappings we have the usual **chain rule**.

**Lemma 9.6** (Chain rule). *Let  $U \subseteq X$ ,  $V \subseteq Y$  and  $F \in C^r(U, V)$  and  $G \in C^r(V, Z)$ ,  $r \geq 1$ . Then  $G \circ F \in C^r(U, Z)$  and*

$$d(G \circ F)(x) = dG(F(x)) \circ dF(x), \quad x \in X. \quad (9.20)$$

**Proof.** Fix  $x \in U$ ,  $y = F(x) \in V$  and let  $u \in X$  such that  $v = dF(x)u$  with  $x+u \in U$  and  $y+v \in V$  for  $|u|$  sufficiently small. Then  $F(x+u) = y+v+o(u)$  and, with  $\tilde{v} = v + o(u)$ ,

$$G(F(x+u)) = G(y + \tilde{v}) = G(y) + dG(y)\tilde{v} + o(\tilde{v}).$$

Using  $|\tilde{v}| \leq \|dF(x)\| |u| + |o(u)|$  we see that  $o(\tilde{v}) = o(u)$  and hence

$$G(F(x+u)) = G(y) + dG(y)v + o(u) = G(F(x)) + dG(F(x)) \circ dF(x)u + o(u)$$

as required. This establishes the case  $r = 1$ . The general case follows from induction.  $\square$

In particular, if  $\ell \in Y^*$  is a bounded linear functional, then  $d(\ell \circ F) = d\ell \circ dF = \ell \circ dF$ . As an application of this result we obtain

**Theorem 9.7** (Schwarz). *Suppose  $F \in C^2(\mathbb{R}^n, Y)$ . Then*

$$\partial_i \partial_j F = \partial_j \partial_i F$$

for any  $1 \leq i, j \leq n$ .

**Proof.** First of all note that  $\partial_j F(x) \in \mathcal{L}(\mathbb{R}, Y)$  and thus it can be regarded as an element of  $Y$ . Clearly the same applies to  $\partial_i \partial_j F(x)$ . Let  $\ell \in Y^*$  be a bounded linear functional, then  $\ell \circ F \in C^2(\mathbb{R}^2, \mathbb{R})$  and hence  $\partial_i \partial_j (\ell \circ F) = \partial_j \partial_i (\ell \circ F)$  by the classical theorem of Schwarz. Moreover, by our remark preceding this lemma  $\partial_i \partial_j (\ell \circ F) = \partial_i \ell(\partial_j F) = \ell(\partial_i \partial_j F)$  and hence  $\ell(\partial_i \partial_j F) = \ell(\partial_j \partial_i F)$  for every  $\ell \in Y^*$  implying the claim.  $\square$

Now we let  $F \in C^2(X, Y)$  and look at the function  $G : \mathbb{R}^2 \rightarrow Y$ ,  $(t, s) \mapsto G(t, s) = F(x + tu + sv)$ . Then one computes

$$\partial_t G(t, s) \Big|_{t=0} = dF(x + sv)u$$

and hence

$$\partial_s \partial_t G(t, s) \Big|_{(s,t)=0} = \partial_s dF(x + sv)u \Big|_{s=0} = (d^2 F(x)u)v.$$

Since by the previous theorem the order of the derivatives is irrelevant, we obtain

$$d^2 F(u, v) = d^2 F(v, u), \quad (9.21)$$

that is,  $d^2 F$  is a symmetric bilinear form. This result easily generalizes to higher derivatives. To this end we introduce some notation first.

A function  $L : \times_{j=1}^n X_j \rightarrow Y$  is called **multilinear** if it is linear with respect to each argument. It is not hard to see that  $L$  is continuous if and only if

$$\|L\| = \sup_{x: |x_1|=\dots=|x_n|=1} |L(x_1, \dots, x_n)| < \infty. \quad (9.22)$$

If we take  $n$  copies of the same space, the set of multilinear functions  $L : X^n \rightarrow Y$  will be denoted by  $\mathcal{L}^n(X, Y)$ . A multilinear function is called **symmetric** provided its value remains unchanged if any two arguments are switched. With the norm from above it is a Banach space and in fact there is a canonical isometric isomorphism between  $\mathcal{L}^n(X, Y)$

and  $\mathcal{L}(X, \mathcal{L}^{n-1}(X, Y))$  given by  $L : (x_1, \dots, x_n) \mapsto L(x_1, \dots, x_n)$  maps to  $x_1 \mapsto L(x_1, \cdot)$ .

**Lemma 9.8.** *Suppose  $F \in C^r(X, Y)$ . Then for every  $x \in X$  we have that*

$$d^r F(x)(u_1, \dots, u_r) = \partial_{t_1} \cdots \partial_{t_r} F(x + \sum_{i=1}^r t_i u_i) \Big|_{t_1=\dots=t_r=0}. \quad (9.23)$$

Moreover,  $d^r F(x) \in \mathcal{L}^r(X, Y)$  is a bounded symmetric multilinear form.

**Proof.** The representation (9.23) follows using induction as before. Symmetry follows since the order of the partial derivatives can be interchanged by Lemma 9.7.  $\square$

Finally, note that to each  $L \in \mathcal{L}^n(X, Y)$  we can assign its polar form  $L \in C(X, Y)$  using  $L(x) = L(x, \dots, x)$ ,  $x \in X$ . If  $L$  is symmetric it can be reconstructed using polarization (Problem 9.9):

$$L(u_1, \dots, u_n) = \frac{1}{n!} \partial_{t_1} \cdots \partial_{t_n} L\left(\sum_{i=1}^n t_i u_i\right). \quad (9.24)$$

We also have the following version of the **product rule**: Suppose  $L \in \mathcal{L}^2(X, Y)$ , then  $L \in C^1(X^2, Y)$  with

$$dL(x)u = L(u_1, x_2) + L(x_1, u_2) \quad (9.25)$$

since

$$\begin{aligned} L(x_1 + u_1, x_2 + u_2) - L(x_1, x_2) &= L(u_1, x_2) + L(x_1, u_2) + L(u_1, u_2) \\ &= L(u_1, x_2) + L(x_1, u_2) + O(|u|^2) \end{aligned} \quad (9.26)$$

as  $|L(u_1, u_2)| \leq \|L\| |u_1| |u_2| = O(|u|^2)$ . If  $X$  is a Banach algebra and  $L(x_1, x_2) = x_1 x_2$  we obtain the usual form of the product rule.

Next we have the following mean value theorem.

**Theorem 9.9** (Mean value). *Suppose  $U \subseteq X$  and  $F : U \rightarrow Y$  is Gâteaux differentiable at every  $x \in U$ . If  $U$  is convex, then*

$$|F(x) - F(y)| \leq M |x - y|, \quad M := \sup_{0 \leq t \leq 1} \left| \delta F((1-t)x + ty, \frac{x-y}{|x-y|}) \right|. \quad (9.27)$$

Conversely, (for any open  $U$ ) if

$$|F(x) - F(y)| \leq M |x - y|, \quad x, y \in U, \quad (9.28)$$

then

$$\sup_{x \in U, |e|=1} |\delta F(x, e)| \leq M. \quad (9.29)$$

**Proof.** Abbreviate  $f(t) = F((1-t)x + ty)$ ,  $0 \leq t \leq 1$ , and hence  $df(t) = \delta F((1-t)x + ty, y - x)$  implying  $|df(t)| \leq \tilde{M} := M|x - y|$  by (9.15). Hence the first part follows from Theorem 9.1.

To prove the second claim suppose we can find an  $e \in X$ ,  $|e| = 1$  such that  $|\delta F(x_0, e)| = M + \delta$  for some  $\delta > 0$  and hence

$$\begin{aligned} M\varepsilon &\geq |F(x_0 + \varepsilon e) - F(x_0)| = |\delta F(x_0, e)\varepsilon + o(\varepsilon)| \\ &\geq (M + \delta)\varepsilon - |o(\varepsilon)| > M\varepsilon \end{aligned}$$

since we can assume  $|o(\varepsilon)| < \varepsilon\delta$  for  $\varepsilon > 0$  small enough, a contradiction.  $\square$

**Corollary 9.10.** *Suppose  $U \subseteq X$  and  $F \in C^1(U, Y)$ . Then  $F$  is locally Lipschitz continuous.*

**Proof.** Fix  $x_0 \in U$  and note that by continuity there is a neighborhood  $U_0$  of  $x_0$  such that  $\|dF(x)\| \leq M$  for  $x \in U_0$ . Hence the claim follows from the mean value theorem.  $\square$

Note, however, that a  $C^1$  function is in general not Lipschitz on arbitrary bounded sets since in the infinite dimensional case continuity of  $dF$  does not suffice to conclude boundedness on bounded closed sets.

**Example 9.12.** Let  $X$  be an infinite Hilbert space and  $\{u_n\}_{n \in \mathbb{N}}$  some orthonormal set. Then the functions  $F_n(x) := \max(0, 1 - 2\|x - u_n\|)$  are continuous with disjoint supports. Hence  $F(x) := \sum_{n \in \mathbb{N}} nF_n(x)$  is also continuous (show this). But  $F$  is not bounded on the unit ball since  $F(u_n) = n$ .  $\diamond$

As an immediate consequence we obtain

**Corollary 9.11.** *Suppose  $U$  is a connected subset of a Banach space  $X$ . A Gâteaux differentiable mapping  $F : U \rightarrow Y$  is constant if and only if  $\delta F = 0$ . In addition, if  $F_{1,2} : U \rightarrow Y$  and  $\delta F_1 = \delta F_2$ , then  $F_1$  and  $F_2$  differ only by a constant.*

As an application of the fundamental theorem of calculus (Theorem 9.3) we obtain a generalization of the well-known fact that continuity of the directional derivatives implies continuous differentiability.

**Lemma 9.12.** *Suppose  $F : U \subseteq X \rightarrow Y$  is Gâteaux differentiable such that the Gâteaux derivative is linear and continuous,  $\delta F \in C(U, \mathcal{L}(X, Y))$ . Then  $F \in C^1(U, Y)$  and  $dF = \delta F$ .*

**Proof.** By assumption  $f(t) := F(x + tu)$  is in  $C^1([0, 1], Y)$  for  $u$  with sufficiently small norm. Moreover, by definition we have  $\dot{f} = \delta F(x + tu)u$  and

using the fundamental theorem of calculus we obtain

$$\begin{aligned} F(x+u) - F(x) &= f(1) - f(0) = \int_0^1 \dot{f}(t) dt = \int_0^1 \delta F(x+tu) u dt \\ &= \left( \int_0^1 \delta F(x+tu) dt \right) u, \end{aligned}$$

where the last equality follows from continuity of the integral since it clearly holds for simple functions. Consequently

$$\begin{aligned} |F(x+u) - F(x) - \delta F(x)u| &= \left| \left( \int_0^1 (\delta F(x+tu) - \delta F(x)) dt \right) u \right| \\ &\leq \left( \int_0^1 \|\delta F(x+tu) - \delta F(x)\| dt \right) |u| \\ &\leq \max_{t \in [0,1]} \|\delta F(x+tu) - \delta F(x)\| |u|. \end{aligned}$$

By the continuity assumption on  $\delta F$ , the right-hand side is  $o(u)$  as required.  $\square$

As another consequence we obtain **Taylor's theorem**.

**Theorem 9.13** (Taylor). *Suppose  $U \subseteq X$  and  $F \in C^{r+1}(U, Y)$ . Then*

$$\begin{aligned} F(x+u) &= F(x) + dF(x)u + \frac{1}{2}d^2F(x)u^2 + \cdots + \frac{1}{r!}d^rF(x)u^r \\ &\quad + \left( \frac{1}{r!} \int_0^1 (1-t)^r d^{r+1}F(x+tu) dt \right) u^{r+1}, \end{aligned} \quad (9.30)$$

where  $u^k := (u, \dots, u) \in X^k$ .

**Proof.** As in the proof of the previous lemma, the case  $r = 0$  is just the fundamental theorem of calculus applied to  $f(t) := F(x+tu)$ . For the induction step we use integration by parts. To this end let  $f_j \in C^1([0, 1], X_j)$ ,  $L \in \mathcal{L}^2(X_1 \times X_2, Y)$  bilinear. Then the product rule (9.25) and the fundamental theorem of calculus imply

$$\int_0^1 L(\dot{f}_1(t), f_2(t)) dt = L(f_1(1), f_2(1)) - L(f_1(0), f_2(0)) - \int_0^1 L(f_1(t), \dot{f}_2(t)) dt.$$

Hence applying integration by parts with  $L(y, t) = ty$ ,  $f_1(t) = d^{r+1}F(x+ut)$ , and  $f_2(t) = \frac{(1-t)^{r+1}}{(r+1)!}$  establishes the induction step.  $\square$

Of course this also gives the Peano form for the remainder:

**Corollary 9.14.** *Suppose  $U \subseteq X$  and  $F \in C^r(U, Y)$ . Then*

$$F(x+u) = F(x) + dF(x)u + \frac{1}{2}d^2F(x)u^2 + \cdots + \frac{1}{r!}d^rF(x)u^r + o(|u|^r). \quad (9.31)$$

**Proof.** Just estimate

$$\begin{aligned} & \left| \left( \frac{1}{(r-1)!} \int_0^1 (1-t)^{r-1} d^r F(x+tu) dt - \frac{1}{r!} d^r F(x) \right) u^r \right| \\ & \leq \frac{|u|^r}{(r-1)!} \int_0^1 (1-t)^{r-1} \|d^r F(x+tu) - d^r F(x)\| dt \\ & \leq \frac{|u|^r}{r!} \sup_{0 \leq t \leq 1} \|d^r F(x+tu) - d^r F(x)\|. \quad \square \end{aligned}$$

Finally we remark that it is often necessary to equip  $C^r(U, Y)$  with a norm. A suitable choice is

$$\|F\| = \sum_{0 \leq j \leq r} \sup_{x \in U} \|d^j F(x)\|. \quad (9.32)$$

The set of all  $r$  times continuously differentiable functions for which this norm is finite forms a Banach space which is denoted by  $C_b^r(U, Y)$ .

In the definition of differentiability we have required  $U$  to be open. Of course there is no stringent reason for this and (9.12) could simply be required for all sequences from  $U \setminus \{x\}$  converging to  $x$ . However, note that the derivative might not be unique in case you miss some directions (the ultimate problem occurring at an isolated point). Our requirement avoids all these issues. Moreover, there is usually another way of defining differentiability at a boundary point: By  $C^r(\overline{U}, Y)$  we denote the set of all functions in  $C^r(U, Y)$  all whose derivatives of order up to  $r$  have a continuous extension to  $\overline{U}$ . Note that if you can approach a boundary point along a half-line then the fundamental theorem of calculus shows that the extension coincides with the Gâteaux derivative.

**Problem 9.6.** Let  $X$  be a real Hilbert space,  $A \in \mathcal{L}(X)$  and  $F(x) := \langle x, Ax \rangle$ . Compute  $d^n F$ .

**Problem\* 9.7.** Let  $X := C([0, 1], \mathbb{R})$  and suppose  $f \in C^1(\mathbb{R})$ . Show that

$$F : X \rightarrow X, \quad x \mapsto f \circ x$$

is differentiable for every  $x \in X$  with derivative given by

$$(dF(x)y)(t) = f'(x(t))y(t).$$

**Problem 9.8.** Let  $X := \ell^2(\mathbb{N})$ ,  $Y := \ell^1(\mathbb{N})$  and  $F : X \rightarrow Y$  given by  $F(x)_j := x_j^2$ . Show  $F \in C^\infty(X, Y)$  and compute all derivatives.

**Problem\* 9.9.** Show (9.24).

**Problem 9.10.** Let  $X$  be a Banach algebra,  $I \subseteq \mathbb{R}$  an open interval, and  $x, y \in C^1(I, X)$ . Show that  $xy \in C^1(I, X)$  with  $(xy)' = \dot{x}y + x\dot{y}$ .



**Problem 9.11.** Let  $X$  be a Banach algebra and  $\mathcal{G}(X)$  the group of invertible elements. Show that  $I : \mathcal{G}(X) \rightarrow \mathcal{G}(X)$ ,  $x \mapsto x^{-1}$  is differentiable with

$$dI(x)y = -x^{-1}yx.$$

(Hint: Corollary 5.2.)

### 9.3. Minimizing nonlinear functionals via calculus

Many problems in applications lead to finding the minimum (or maximum) of a given (nonlinear) functional  $F : X \rightarrow \mathbb{R}$ . For example, many physical problems can be described by an energy functional and one seeks a solution which minimizes this energy. Since the minima of  $-F$  are the maxima of  $F$  and vice versa, we will restrict our attention to minima only. Of course if  $X = \mathbb{R}$  (or  $\mathbb{R}^n$ ) we can find the local extrema by searching for the zeros of the derivative and then checking the second derivative to determine if it is a minimum or maximum. In fact, by virtue of our version of Taylor's theorem (9.31) we see that  $F$  will take values above and below  $F(x)$  in a vicinity of  $x$  if we can find some  $u$  such that  $dF(x)u \neq 0$ . Hence  $dF(x) = 0$  is clearly a necessary condition for a local extremum. Moreover, if  $dF(x) = 0$  we can go one step further and conclude that all values in a vicinity of  $x$  will lie above  $F(x)$  provided the second derivative  $d^2F(x)$  is positive in the sense that there is some  $c > 0$  such that  $d^2F(x)u^2 > c$  for all directions  $u \in \partial B_1(0)$ . While this gives a viable solution to the problem of finding local extrema, we can easily do a bit better. To this end we look at the variations of  $f$  along lines through  $x$ , that is, we look at the behavior of the function

$$f(t) := F(x + tu) \tag{9.33}$$

for a fixed direction  $u \in B_1(0)$ . Hence this approach is also known as **calculus of variations**. Then, if  $F$  has a local extremum at  $x$  the same will be true for  $f$  and hence a necessary condition for an extremum is that the Gâteaux derivative vanishes in every direction:  $\delta F(x, u) = 0$  for all unit vectors  $u$ . Similarly, a necessary condition for a local minimum at  $x$  is that  $f$  has a local minimum at 0 for all unit vectors  $u$ . For example  $\delta^2 F(x, u) > 0$  for all unit vectors  $u$ . Here the higher order Gâteaux derivatives are defined as

$$\delta^n F(x, u) := \left( \frac{d}{dt} \right)^n F(x + tu) \Big|_{t=0} \tag{9.34}$$

with the derivative defined as a limit as in (9.13). That is we have the recursive definition  $\delta^n F(x, u) = \lim_{\varepsilon \rightarrow 0} \varepsilon^{-1} (\delta^{n-1} F(x + \varepsilon u, u) - \delta^{n-1} F(x, u))$ . Note that if  $\delta^n F(x, u)$  exists, then  $\delta^n F(x, \lambda u)$  exists for every  $\lambda \in \mathbb{R}$  and

$$\delta^n F(x, \lambda u) = \lambda^n \delta^n F(x, u), \quad \lambda \in \mathbb{R}. \tag{9.35}$$

However, the condition  $\delta^2 F(x, u) > 0$  for all unit vectors  $u$  is not sufficient as there are certain features you might miss when you only look at the function along rays through a fixed point. This is demonstrated by the following example:

**Example 9.13.** Let  $X = \mathbb{R}^2$  and consider the points  $(x_n, y_n) := (\frac{1}{n}, \frac{1}{n^2})$ . For each point choose a radius  $r_n$  such that the balls  $B_n := B_{r_n}(x_n, y_n)$  are disjoint and lie between two parabolas  $B_n \subset \{(x, y) | x \geq 0, \frac{x^2}{2} \leq y \leq 2x^2\}$ . Moreover, choose a smooth nonnegative bump function  $\phi(r^2)$  with support in  $[-1, 1]$  and maximum 1 at 0. Now consider  $F(x, y) = x^2 + y^2 - 2 \sum_{n \in \mathbb{N}} \rho_n \phi(\frac{(x-x_n)^2 + (y-y_n)^2}{r_n^2})$ , where  $\rho_n = x_n^2 + y_n^2$ . By construction  $F$  is smooth away from zero. Moreover, at zero  $F$  is continuous and Gâteaux differentiable of arbitrary order with  $F(0, 0) = 0$ ,  $\delta F((0, 0), (u, v)) = 0$ ,  $\delta^2 F((0, 0), (u, v)) = 2(u^2 + v^2)$ , and  $\delta^k F((0, 0), (u, v)) = 0$  for  $k \geq 3$ .

In particular,  $F(ut, vt)$  has a strict local minimum at  $t = 0$  for every  $(u, v) \in \mathbb{R}^2 \setminus \{0\}$ , but  $F$  has no local minimum at  $(0, 0)$  since  $F(x_n, y_n) = -\rho_n$ . Clearly  $F$  is not differentiable at 0. In fact, note that the Gâteaux derivatives are not continuous at 0 (the derivatives in  $B_n$  grow like  $r_n^{-2}$ ).  $\diamond$

**Lemma 9.15.** Suppose  $F : U \rightarrow \mathbb{R}$  has Gâteaux derivatives up to the order of two. A necessary condition for  $x \in U$  to be a local minimum is that  $\delta F(x, u) = 0$  and  $\delta^2 F(x, u) \geq 0$  for all  $u \in X$ . A sufficient condition for a strict local minimum is if in addition  $\delta^2 F(x, u) \geq c > 0$  for all  $u \in \partial B_1(0)$  and  $\delta^2 F$  is continuous at  $x$  uniformly with respect to  $u \in \partial B_1(0)$ .

**Proof.** The necessary conditions have already been established. To see the sufficient conditions note that the assumptions on  $\delta^2 F$  imply that there is some  $\varepsilon > 0$  such that  $\delta^2 F(y, u) \geq \frac{c}{2}$  for all  $y \in B_\varepsilon(x)$  and all  $u \in \partial B_1(0)$ . Equivalently,  $\delta^2 F(y, u) \geq \frac{c}{2}|u|^2$  for all  $y \in B_\varepsilon(x)$  and all  $u \in X$ . Hence applying Taylor's theorem to  $f(t)$  using  $\ddot{f}(t) = \delta^2 F(x + tu, u)$  gives

$$F(x + u) = f(1) = f(0) + \int_0^1 (1-s)\ddot{f}(s)ds \geq F(x) + \frac{c}{4}|u|^2$$

for  $u \in B_\varepsilon(0)$ .  $\square$

Note that if  $F \in C^2(U, \mathbb{R})$  then  $\delta^2 F(x, u) = d^2 F(x)u^2$  and we obtain

**Corollary 9.16.** Suppose  $F \in C^2(U, \mathbb{R})$ . A sufficient condition for  $x \in U$  to be a strict local minimum is  $dF(x) = 0$  and  $d^2 F(x)u^2 \geq c|u|^2$  for all  $u \in X$ .

**Proof.** Observe that by  $|\delta^2 F(x, u) - \delta^2 F(y, u)| \leq \|d^2 F(x) - d^2 F(y)\||u|^2$  the continuity requirement from the previous lemma is satisfied.  $\square$

**Example 9.14.** If  $X$  is a real Hilbert space, then the symmetric bilinear form  $d^2 F$  has a corresponding self-adjoint operator  $A \in \mathcal{L}(X)$  such that

$d^2F(u, v) = \langle u, Av \rangle$  and the condition  $d^2F(x)u^2 \geq c|u|^2$  is equivalent to the spectral condition  $\sigma(A) \subset [c, \infty)$ . In the finite dimensional case  $A$  is of course the Jacobi matrix and the spectral conditions says that all eigenvalues must be positive.  $\diamond$

**Example 9.15.** Let  $X := \ell^2(\mathbb{N})$  and consider

$$F(x) := \sum_{n \in \mathbb{N}} \left( \frac{x_n^2}{2n^2} - x_n^4 \right).$$

Then  $F \in C^2(X, \mathbb{R})$  with  $dF(x)u = \sum_{n \in \mathbb{N}} (\frac{x_n}{n^2} - 4x_n^3)u_n$  and  $d^2F(x)(u, v) = \sum_{n \in \mathbb{N}} (\frac{1}{n^2} - 12x_n^2)v_n u_n$ . In particular,  $F(0) = 0$ ,  $dF(0) = 0$  and  $d^2F(0)u^2 = \sum_n n^{-2}u_n^2 > 0$  for  $u \neq 0$ . However,  $F(\delta^m/m) < 0$  shows that 0 is no local minimum. So the condition  $d^2F(x)u^2 > 0$  is not sufficient in infinite dimensions. It is however, sufficient in finite dimensions since compactness of the unit ball leads to the stronger condition  $d^2F(x, u) \geq c > 0$  for all  $u \in \partial B_1(0)$ .  $\diamond$

**Example 9.16.** Consider a classical particle whose location at time  $t$  is given by  $q(t)$ . Then the **least action principle** states that, if the particle moves from  $q(a)$  to  $q(b)$ , the path of the particle will make the action functional

$$S(q) := \int_a^b L(t, q(t), \dot{q}(t)) dt$$

stationary, that is

$$\delta S(q) = 0.$$

Here  $L : \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  is the **Lagrangian** of the system. The name suggests that the action should attain a minimum, but this is not always the case and hence it is also referred to as **stationary action principle**.

More precisely, let  $L \in C^2(\mathbb{R}^{2n+1}, \mathbb{R})$  and in order to incorporate the requirement that the initial and end points are fixed, we take  $X = \{x \in C^2([a, b], \mathbb{R}^n) | x(a) = x(b) = 0\}$  and consider

$$q(t) := q(a) + \frac{t-a}{b-a}(q(b) - q(a)) + x(t), \quad x \in X.$$

Hence we want to compute the Gâteaux derivative of  $F(x) := S(q)$ , where  $x$  and  $q$  are related as above with  $q(a)$ ,  $q(b)$  fixed. Then

$$\begin{aligned} \delta F(x, u) &= \frac{d}{dt} \Big|_{t=0} \int_a^b L(s, q(s) + t u(s), \dot{q}(s) + t \dot{u}(s)) ds \\ &= \int_a^b (L_q(s, q(s), \dot{q}(s))u(s) + L_{\dot{q}}(s, q(s), \dot{q}(s))\dot{u}(s)) ds \\ &= \int_a^b (L_q(s, q(s), \dot{q}(s))u(s) - \frac{d}{ds} L_{\dot{q}}(s, q(s), \dot{q}(s))) u(s) ds, \end{aligned}$$

where we have used integration by parts (including the boundary conditions) to obtain the last equality. Here  $L_q, L_{\dot{q}}$  are the gradients with respect to  $q, \dot{q}$ , respectively, and products are understood as scalar products in  $\mathbb{R}^n$ .

If we want this to vanish for all  $u \in X$  we obtain the corresponding **Euler–Lagrange equation**

$$\frac{d}{ds} L_{\dot{q}}(s, q(s), \dot{q}(s)) = L_q(s, q(s), \dot{q}(s)).$$

For example, for a classical particle of mass  $m > 0$  moving in a conservative force field described by a potential  $V \in C^1(\mathbb{R}^n, \mathbb{R})$  the Lagrangian is given by the difference between kinetic and potential energy

$$L(t, q, \dot{q}) := \frac{m}{2} \dot{q}^2 - V(q)$$

and the Euler–Lagrange equations read

$$m\ddot{q} = -V_q(q),$$

which are just Newton’s equations of motion.  $\diamond$

Finally we note that the situation simplifies a lot when  $F$  is convex. Our first observation is that a local minimum is automatically a global one.

**Lemma 9.17.** *Suppose  $C \subseteq X$  is convex and  $F : C \rightarrow \mathbb{R}$  is convex. Every local minimum is a global minimum. Moreover, if  $F$  is strictly convex then the minimum is unique.*

**Proof.** Suppose  $x$  is a local minimum and  $F(y) < F(x)$ . Then  $F(\lambda y + (1 - \lambda)x) \leq \lambda F(y) + (1 - \lambda)F(x) < F(x)$  for  $\lambda \in (0, 1)$  contradicts the fact that  $x$  is a local minimum. If  $x, y$  are two global minima, then  $F(\lambda y + (1 - \lambda)x) < F(y) = F(x)$  yielding a contradiction unless  $x = y$ .  $\square$

Moreover, to find the global minimum it suffices to find a point where the Gâteaux derivative vanishes.

**Lemma 9.18.** *Suppose  $C \subseteq X$  is convex and  $F : C \rightarrow \mathbb{R}$  is convex. If the Gâteaux derivative exists at an interior point  $x \in C$  and satisfies  $\delta F(x, u) = 0$  for all  $u \in X$ , then  $x$  is a global minimum.*

**Proof.** By assumption  $f(t) := F(x + tu)$  is a convex function defined on an interval containing 0 with  $f'(0) = 0$ . If  $y$  is another point we can choose  $u = y - x$  and Lemma 3.2 from [48] (iii) implies  $F(y) = f(1) \geq f(0) = F(x)$ .  $\square$

As in the one-dimensional case, convexity can be read off from the second derivative.

**Lemma 9.19.** *Suppose  $C \subseteq X$  is open and convex and  $F : C \rightarrow \mathbb{R}$  has Gâteaux derivatives up to order two. Then  $F$  is convex if and only if  $\delta^2 F(x, u) \geq 0$  for all  $x \in C$  and  $u \in X$ . Moreover,  $F$  is strictly convex if  $\delta^2 F(x, u) > 0$  for all  $x \in C$  and  $u \in X \setminus \{0\}$ .*

**Proof.** We consider  $f(t) := F(x + tu)$  as before such that  $f'(t) = \delta F(x + tu, u)$ ,  $f''(t) = \delta^2 F(x + tu, u)$ . Moreover, note that  $f$  is (strictly) convex for all  $x \in C$  and  $u \in X \setminus \{0\}$  if and only if  $F$  is (strictly) convex. Indeed, if  $F$  is (strictly) convex so is  $f$  as is easy to check. To see the converse note

$$F(\lambda y + (1 - \lambda)x) = f(\lambda) \leq \lambda f(1) - (1 - \lambda)f(0) = \lambda F(y) - (1 - \lambda)F(x)$$

with strict inequality if  $f$  is strictly convex. The rest follows from Problem 3.6 from [48].  $\square$

There is also a version using only first derivatives plus the concept of a monotone operator. A map  $F : U \subseteq X \rightarrow X^*$  is **monotone** if

$$(F(x) - F(y))(x - y) \geq 0, \quad x, y \in U.$$

It is called **strictly monotone** if we have strict inequality for  $x \neq y$ . Monotone operators will be the topic of Chapter 14.

**Lemma 9.20.** *Suppose  $C \subseteq X$  is open and convex and  $F : C \rightarrow \mathbb{R}$  has Gâteaux derivatives  $\delta F(x) \in X^*$  for every  $x \in C$ . Then  $F$  is (strictly) convex if and only if  $\delta F$  is (strictly) monotone.*

**Proof.** Note that by assumption  $\delta F : C \rightarrow X^*$  and the claim follows as in the previous lemma from Problem 3.6 from [48] since  $f'(t) = \delta F(x + tu)u$  which shows that  $\delta F$  is (strictly) monotone if and only if  $f'$  is (strictly) increasing.  $\square$

**Example 9.17.** The length of a curve  $q : [a, b] \rightarrow \mathbb{R}^n$  is given by

$$\int_a^b |q'(s)| ds.$$

Of course we know that the shortest curve between two given points  $q_0$  and  $q_1$  is a straight line. Notwithstanding that this is evident, defining the length as the total variation, let us show this by seeking the minimum of the following functional

$$F(x) := \int_a^b |q'(s)| ds, \quad q(t) = x(t) + q_0 + \frac{t - a}{b - a}(q_1 - q_0)$$

for  $x \in X := \{x \in C^1([a, b], \mathbb{R}^n) | x(a) = x(b) = 0\}$ . Unfortunately our integrand will not be differentiable unless  $|q'| \geq c$ . However, since the absolute

value is convex, so is  $F$  and it will suffice to search for a local minimum within the convex open set  $C := \{x \in X \mid |x'| < \frac{|q_1 - q_0|}{2(b-a)}\}$ . We compute

$$\delta F(x, u) = \int_a^b \frac{q'(s)u'(s)}{|q'(s)|} ds$$

which shows (Lemma 3.24 from [48]) that  $q'/|q'|$  must be constant. Hence the local minimum in  $C$  is indeed a straight line and this must also be a global minimum in  $X$ . However, since the length of a curve is independent of its parametrization, this minimum is not unique!  $\diamond$

**Example 9.18.** Let us try to find a curve  $y(x)$  from  $y(0) = 0$  to  $y(x_1) = y_1$  which minimizes

$$F(y) := \int_0^{x_1} \sqrt{\frac{1 + y'(x)^2}{x}} dx.$$

Note that since the function  $t \mapsto \sqrt{1 + t^2}$  is convex, we obtain that  $F$  is convex. Hence it suffices to find a zero of

$$\delta F(y, u) = \int_0^{x_1} \frac{y'(x)u'(x)}{\sqrt{x(1 + y'(x)^2)}} dx,$$

which shows (Lemma 3.24 from [48]) that  $\frac{y'}{\sqrt{x(1 + y'^2)}} = C^{-1/2}$  is constant or equivalently

$$y'(x) = \sqrt{\frac{x}{C - x}}$$

and hence

$$y(x) = C \arctan \left( \sqrt{\frac{x}{C - x}} \right) - \sqrt{x(C - x)}.$$

The constant  $C$  has to be chosen such that  $y(x_1)$  matches the given value  $y_1$ . Note that  $C \mapsto y(x_1)$  decreases from  $\frac{\pi x_1}{2}$  to 0 and hence there will be a unique  $C > x_1$  for  $0 < y_1 < \frac{\pi x_1}{2}$ .  $\diamond$

**Problem 9.12.** Consider the least action principle for a classical one-dimensional particle. Show that

$$\delta^2 F(x, u) = \int_a^b (m \dot{u}(s)^2 - V''(q(s))u(s)^2) ds.$$

Moreover, show that we have indeed a minimum if  $V'' \leq 0$ .

## 9.4. Minimizing nonlinear functionals via compactness

Another approach for minimizing a nonlinear functional  $F : M \subseteq X \rightarrow \mathbb{R}$  is based on compactness. If  $M$  is compact and  $F$  is continuous, then we can proceed as in the finite-dimensional case to show that there is a minimizer: Start with a sequence  $x_n$  such that  $F(x_n) \rightarrow \inf_M F$ . By compactness we can assume that  $x_n \rightarrow x_0$  after passing to a subsequence and by continuity

$F(x_n) \rightarrow F(x_0) = \inf_M F$ . Now in the infinite dimensional case we will use weak convergence to get compactness and hence we will also need weak (sequential) continuity of  $F$ . However, since there are more weakly than strongly convergent subsequences, weak (sequential) continuity is in fact a stronger property than just continuity!

**Example 9.19.** By Lemma 4.25 (ii) the norm is weakly sequentially lower semicontinuous but it is in general not weakly sequentially continuous as any infinite orthonormal set in a Hilbert space converges weakly to 0. However, note that this problem does not occur for linear maps. This is an immediate consequence of the very definition of weak convergence (Problem 4.29).  $\diamond$

Hence weak continuity might be too much to hope for in concrete applications. In this respect note that, for our argument to work lower semicontinuity (cf. Problem B.19) will already be sufficient. This is frequently referred to as the **direct method in the calculus of variations** due to Zaremba and Hilbert:

**Theorem 9.21** (Variational principle). *Let  $X$  be a reflexive Banach space and let  $F : M \subseteq X \rightarrow (-\infty, \infty]$ . Suppose  $M$  is nonempty, weakly sequentially closed and that either  $F$  is **weakly coercive**, that is  $F(x) \rightarrow \infty$  whenever  $\|x\| \rightarrow \infty$ , or that  $M$  is bounded. Then, if  $F$  is weakly sequentially lower semicontinuous, there exists some  $x_0 \in M$  with  $F(x_0) = \inf_M F$ .*

*If  $F$  is Gâteaux differentiable, then*

$$\delta F(x_0) = 0.$$

**Proof.** Without loss of generality we can assume  $F(x) < \infty$  for some  $x \in M$ . As above we start with a sequence  $x_n \in M$  such that  $F(x_n) \rightarrow \inf_M F < \infty$ . If  $M$  is unbounded, then the fact that  $F$  is coercive implies that  $x_n$  is bounded. Otherwise, if  $M$  is bounded, it is obviously bounded. Hence by Theorem 4.28 we can pass to a subsequence such that  $x_n \rightharpoonup x_0$  with  $x_0 \in M$  since  $M$  is assumed sequentially closed. Now since  $F$  is weakly sequentially lower semicontinuous we finally get  $\inf_M F = \lim_{n \rightarrow \infty} F(x_n) = \liminf_{n \rightarrow \infty} F(x_n) \geq F(x_0)$ .  $\square$

Of course in a metric space the definition of closedness in terms of sequences agrees with the corresponding topological definition. In the present situation sequentially weakly closed implies (sequentially) closed and the converse holds at least for convex sets.

**Lemma 9.22.** *Suppose  $M \subseteq X$  is convex. Then  $M$  is closed if and only if it is sequentially weakly closed.*

**Proof.** Suppose  $M$  is closed and let  $x$  be in the weak sequential closure of  $M$ , that is, there is a sequence  $x_n \rightharpoonup x$ . If  $x \notin M$ , then by Corollary 6.4 we can

find a linear functional  $\ell$  which separates  $\{x\}$  and  $M$ . But this contradicts  $\ell(x) = d < c \leq \ell(x_n) \rightarrow \ell(x)$ .  $\square$

Similarly, the same is true with lower semicontinuity. In fact, a slightly weaker assumption suffices. Let  $X$  be a vector space and  $M \subseteq X$  a convex subset. A function  $F : M \rightarrow \overline{\mathbb{R}}$  is called **quasiconvex** if

$$F(\lambda x + (1 - \lambda)y) \leq \max\{F(x), F(y)\}, \quad \lambda \in (0, 1), \quad x, y \in M. \quad (9.36)$$

It is called **strictly quasiconvex** if the inequality is strict for  $x \neq y$ . By  $\lambda F(x) + (1 - \lambda)F(y) \leq \max\{F(x), F(y)\}$  every (strictly) convex function is (strictly) quasiconvex. The converse is not true as the following example shows.

**Example 9.20.** Every (strictly) monotone function on  $\mathbb{R}$  is (strictly) quasiconvex. Moreover, the same is true for symmetric functions which are (strictly) monotone on  $[0, \infty)$ . Hence the function  $F(x) = \sqrt{|x|}$  is strictly quasiconvex. But it is clearly not convex on  $M = \mathbb{R}$ .  $\diamond$

Note that we can extend a (quasi-)convex function  $F : M \rightarrow \overline{\mathbb{R}}$  to all of  $X$  by setting  $F(x) = \infty$  for  $x \in X \setminus M$  and the resulting function will still be (quasi-)convex and will have the same infimum.

Now we are ready for the next

**Lemma 9.23.** *Suppose  $M \subseteq X$  is a closed convex set and suppose  $F : M \rightarrow \overline{\mathbb{R}}$  is quasiconvex. Then  $F$  is weakly sequentially lower semicontinuous if and only if it is (sequentially) lower semicontinuous.*

**Proof.** Suppose  $F$  is lower semicontinuous. If it were not weakly sequentially lower semicontinuous we could find a sequence  $x_n \rightharpoonup x_0$  with  $F(x_n) \rightarrow a < F(x_0)$ . But then  $x_n \in F^{-1}((-\infty, a])$  for  $n$  sufficiently large implying  $x_0 \in F^{-1}((-\infty, a])$  as this set is convex (Problem 9.15) and closed (Problem B.19). But this gives the contradiction  $a < F(x_0) \leq a$ .  $\square$

**Example 9.21.** Let  $U \subseteq \mathbb{R}^n$  and  $K : U \times \mathbb{C} \rightarrow [0, \infty)$ . Suppose  $u \mapsto K(x, u)$  is quasi-convex and continuous for fixed  $x \in U$ . Then

$$F(u) := \int_U K(x, u(x)) d^n x$$

is weakly sequentially lower semicontinuous on  $L^p(U)$  for  $1 \leq p \leq \infty$ . Since  $F$  is obviously quasiconvex, it suffices to show lower semicontinuity. Assume the contrary, then we can find some  $u \in L^p$  and a sequence  $u_n \rightarrow u$  such that  $F(u) > \liminf F(u_n)$ . After passing to a subsequence we can assume that  $u_n(x) \rightarrow u(x)$  a.e. and hence  $K(x, u_n(x)) \rightarrow K(x, u(x))$  a.e. Finally applying Fatou's lemma (Theorem 2.4 from [48]) gives the contradiction  $F(u) \leq \liminf F(u_n)$ . Note that this result generalizes to  $\mathbb{C}^n$ -valued functions in a straightforward manner.  $\diamond$



Moreover, in this case our variational principle reads as follows:

**Corollary 9.24.** *Let  $X$  be a reflexive Banach space and let  $M$  be a nonempty closed convex subset. If  $F : M \subseteq X \rightarrow \overline{\mathbb{R}}$  is quasiconvex, lower semicontinuous, and, if  $M$  is unbounded, weakly coercive, then there exists some  $x_0 \in M$  with  $F(x_0) = \inf_M F$ . If  $F$  is strictly quasiconvex then  $x_0$  is unique.*

**Proof.** It remains to show uniqueness. Let  $x_0$  and  $x_1$  be two different minima. Then  $F(\lambda x_0 + (1 - \lambda)x_1) < \max\{F(x_0), F(x_1)\} = \inf_M F$ , a contradiction.  $\square$

**Example 9.22.** Let  $X$  be a reflexive Banach space. Suppose  $M \subseteq X$  is a nonempty closed convex set. Then for every  $x \in X$  there is a point  $x_0 \in M$  with minimal distance,  $\|x - x_0\| = \text{dist}(x, M)$ . Indeed,  $F(z) = \text{dist}(x, z)$  is convex, continuous and, if  $M$  is unbounded, weakly coercive. Hence the claim follows from Corollary 9.24. Note that the assumption that  $X$  is reflexive is crucial (Problem 9.13). Moreover, we also get that  $x_0$  is unique if  $X$  is strictly convex (see Problem 1.13).  $\diamond$

**Example 9.23.** Let  $\mathfrak{H}$  be a Hilbert space and  $\ell \in \mathfrak{H}^*$  a linear functional. We will give a variational proof of the Riesz lemma (Theorem 2.10). Since we already need to know that Hilbert spaces are reflexive, it should not be taken too serious. To this end consider

$$F(x) = \frac{1}{2}\|x\|^2 - \text{Re}(\ell(x)), \quad x \in \mathfrak{H}.$$

Then  $F$  is convex, continuous, and weakly coercive. Hence there is some  $x_0 \in \mathfrak{H}$  with  $F(x_0) = \inf_{x \in \mathfrak{H}} F(x)$ . Moreover, for fixed  $x \in \mathfrak{H}$ ,

$$\mathbb{R} \rightarrow \mathbb{R}, \quad \varepsilon \mapsto F(x_0 + \varepsilon x) = F(x_0) + \varepsilon \text{Re}(\langle x_0, x \rangle - \ell(x)) + \frac{\varepsilon^2}{2}\|x\|^2$$

is a smooth map which has a minimum at  $\varepsilon = 0$ . Hence its derivative at  $\varepsilon = 0$  must vanish:  $\text{Re}(\langle x_0, x \rangle - \ell(x)) = 0$  for all  $x \in \mathfrak{H}$ . Replacing  $x \rightarrow -ix$  we also get  $\text{Im}(\langle x_0, x \rangle - \ell(x)) = 0$  and hence  $\ell(x) = \langle x_0, x \rangle$ .  $\diamond$

**Example 9.24.** Let  $\mathfrak{H}$  be a Hilbert space and let us consider the problem of finding the lowest eigenvalue of a positive operator  $A \geq 0$ . Of course this is bound to fail since the eigenvalues could accumulate at 0 without 0 being an eigenvalue (e.g. the multiplication operator with the sequence  $\frac{1}{n}$  in  $\ell^2(\mathbb{N})$ ). Nevertheless it is instructive to see how things can go wrong (and it underlines the importance of our various assumptions).

To this end consider its quadratic form  $q_A(f) = \langle f, Af \rangle$ . Then, since  $q_A^{1/2}$  is a seminorm (Problem 1.23) and taking squares is convex,  $q_A$  is convex. If we consider it on  $M = \bar{B}_1(0)$  we get existence of a minimum from Theorem 9.21. However this minimum is just  $q_A(0) = 0$  which is not very interesting. In order to obtain a minimal eigenvalue we would need to take

$M = S_1 = \{f \mid \|f\| = 1\}$ , however, this set is not weakly closed (its weak closure is  $\bar{B}_1(0)$  as we will see in the next section). In fact, as pointed out before, the minimum is in general not attained on  $M$  in this case.

Note that our problem with the trivial minimum at 0 would also disappear if we would search for a maximum instead. However, our lemma above only guarantees us weak sequential lower semicontinuity but not weak sequential upper semicontinuity. In fact, note that not even the norm (the quadratic form of the identity) is weakly sequentially upper continuous (cf. Lemma 4.25 (ii) versus Lemma 4.26). If we make the additional assumption that  $A$  is compact, then  $q_A$  is weakly sequentially continuous as can be seen from Theorem 4.29. Hence for compact operators the maximum is attained at some vector  $f_0$ . Of course we will have  $\|f_0\| = 1$  but is it an eigenvalue? To see this we resort to a small ruse: Consider the real function

$$\phi(t) = \frac{q_A(f_0 + tf)}{\|f_0 + tf\|^2} = \frac{\alpha_0 + 2t\operatorname{Re}\langle f, Af_0 \rangle + t^2 q_A(f)}{1 + 2t\operatorname{Re}\langle f, f_0 \rangle + t^2 \|f\|^2}, \quad \alpha_0 = q_A(f_0),$$

which has a maximum at  $t = 0$  for any  $f \in \mathfrak{H}$ . Hence we must have  $\phi'(0) = 2\operatorname{Re}\langle f, (A - \alpha_0)f_0 \rangle = 0$  for all  $f \in \mathfrak{H}$ . Replacing  $f \rightarrow if$  we get  $2\operatorname{Im}\langle f, (A - \alpha_0)f_0 \rangle = 0$  and hence  $\langle f, (A - \alpha_0)f_0 \rangle = 0$  for all  $f$ , that is  $Af_0 = \alpha_0 f_0$ . So we have recovered Theorem 3.6.  $\diamond$

**Example 9.25.** The classical **Poisson problem** asks for a solution of

$$-\Delta u = f$$

in a bounded domain  $U \subset \mathbb{R}^n$  attaining given boundary values  $g$  on  $\partial U$ . We are going to look for weak solutions, that is, solutions  $u \in H_{\mathbb{R}}^1(\mathbb{R}^n)$  satisfying

$$\int_{\mathbb{R}^n} (\partial u \cdot \partial \phi - f\phi) d^n x = 0, \quad \phi \in C_c^\infty(\mathbb{R}^n).$$

We start by introducing the functional

$$F(u) := \int_{\mathbb{R}^n} \left( \frac{1}{2} |\partial u|^2 - uf \right) d^n x$$

on  $H_{\mathbb{R}}^1(\mathbb{R}^n)$ . To incorporate the boundary values we introduce

$$M := \{v \in H_{\mathbb{R}}^1(\mathbb{R}^n) \mid v|_{\partial U} = g\}.$$

Here the equality  $v|_{\partial U} = g$  has to be understood in the sense of traces and hence we need to require  $U$  to have a  $C^1$  boundary such that the trace operator is well-defined. Moreover, we assume  $f \in L_{\mathbb{R}}^2(\mathbb{R}^n)$  and  $g$  in the range of the trace operator, such that  $M$  is nonempty. In particular, there is some  $\bar{g} \in H_{\mathbb{R}}^1(\mathbb{R}^n)$  with  $\bar{g}|_{\partial U} = g$ . For example, we can assume that  $g$  is Lipschitz continuous in which case it has a Lipschitz continuous extension  $\bar{g}$  to  $\bar{U}$  (Lemma B.31). In this case  $\bar{g} \in W_{\mathbb{R}}^{1,\infty}(U) \subset H_{\mathbb{R}}^1(U)$ . Moreover, by continuity of the trace operator,  $M$  is closed and convexity is obvious.

To see that  $F$  is weakly coercive, let  $u = \bar{g} + v$ , where  $v \in H_0^1(U)$  vanishes on the boundary, then

$$F(u) \geq \frac{1}{2} \|\partial v\|_2^2 - \|\partial \bar{g}\|_2 \|\partial v\|_2 - \|f\|_2 \|v\|_2 - C,$$

with  $C$  depending on  $f$  and  $g$  only. Now the Poincaré inequality (Theorem 7.38 from [48])  $\|v\|_2 \leq C_0 \|\partial v\|_2$  implies that  $F(u) \rightarrow \infty$  if  $\|v\|_{1,2} \rightarrow \infty$ . Finally, since  $F$  is convex and continuous, Corollary 9.24 implies existence of a unique minimizer. This minimizer solves the weak formulation of our boundary value problem.  $\diamond$

**Example 9.26.** Let us consider the following nonlinear elliptic problem

$$-\Delta u + u|u| + u = f$$

in  $L_{\mathbb{R}}^2(\mathbb{R}^n)$  for a given function  $f \in L_{\mathbb{R}}^2(\mathbb{R}^n)$ . We are going to look for weak solutions, that is, solutions  $u \in H_{\mathbb{R}}^1(\mathbb{R}^n)$  satisfying

$$\int_{\mathbb{R}^n} (\partial u \cdot \partial \phi + (|u|u + u - f)\phi) d^n x = 0, \quad \phi \in C_c^\infty(\mathbb{R}^n).$$

We start by introducing the functional

$$F(u) := \int_{\mathbb{R}^n} \left( \frac{1}{2} |\partial u|^2 + \frac{1}{3} |u|^3 + \frac{1}{2} u^2 - uf \right) d^n x$$

on  $L_{\mathbb{R}}^2(\mathbb{R}^n)$  and set  $F(u) = \infty$  if  $u \notin H_{\mathbb{R}}^1(\mathbb{R}^n) \cap L_{\mathbb{R}}^3(\mathbb{R}^n)$ . One checks that for  $u \in H_{\mathbb{R}}^1(\mathbb{R}^n) \cap L_{\mathbb{R}}^3(\mathbb{R}^n)$  and  $\phi \in C_c^\infty(\mathbb{R}^n)$  this functional has a variational derivative

$$\delta F(u, \phi) = \int_{\mathbb{R}^n} (\partial u \cdot \partial \phi + (|u|u + u - f)\phi) d^n x = 0$$

which coincides with the weak formulation of our problem. Hence a minimizer (which is necessarily in  $H_{\mathbb{R}}^1(\mathbb{R}^n) \cap L_{\mathbb{R}}^3(\mathbb{R}^n)$ ) is a weak solution of our nonlinear elliptic problem and it remains to show existence of a minimizer.

First of all note that

$$F(u) \geq \frac{1}{2} \|u\|_2^2 - \|u\|_2 \|f\|_2 \geq \frac{1}{4} \|u\|_2^2 - \|f\|_2^2$$

and hence  $F$  is coercive. To see that it is weakly sequentially lower continuous, observe that for the second terms this follows from Example 9.21 and the last two are easy. For the first term let  $u_n \rightharpoonup u$  in  $L^2$  and observe

$$\begin{aligned} \|\partial u\|_2 &= \sup_{\|\phi\|_2=1, \phi \in C_c^\infty} |\langle u, \partial \phi \rangle| = \sup_{\|\phi\|_2=1, \phi \in C_c^\infty} \lim_n |\langle u_n, \partial \phi \rangle| \\ &\leq \liminf_n \sup_{\|\phi\|_2=1, \phi \in C_c^\infty} |\langle u_n, \partial \phi \rangle| \\ &= \liminf_n \|\partial u_n\|_2. \end{aligned}$$

Hence the claim follows.  $\diamond$

If we look at the previous problem in the case  $f = 0$ , our approach will only give us the trivial solution. In fact, for a linear problem one expects nontrivial solutions for the homogenous problem only at an eigenvalue. Since the Laplace operator has no eigenvalues on  $\mathbb{R}^n$  (as is not hard to see using the Fourier transform), we look at a bounded domain  $U$  instead. To avoid the trivial solution we will add a constraint. As a preparation we note

**Lemma 9.25.** *Let  $X, Y$  be Banach spaces such that  $X$  is compactly embedded into  $Y$  and let  $N : Y \rightarrow \mathbb{R}$  be continuous. Then  $M := \{x \in X | N(x) = N_0\}$  is weakly sequentially closed for any  $N_0 \in \mathbb{R}$ . The same holds for  $M := \{x \in X | N(x) \leq N_0\}$ .*

**Proof.** This follows from Theorem 4.29 since every weakly convergent sequence in  $X$  is convergent in  $Y$ .  $\square$

**Theorem 9.26** (Variational principle with constraints). *Let  $X$  be a reflexive Banach space and let  $F : X \rightarrow \mathbb{R}$  be weakly sequentially lower semicontinuous and weakly coercive. Let  $Y$  be another Banach space such that  $X$  is compactly embedded into  $Y$  and let  $N : Y \rightarrow \mathbb{R}$  be continuous. Fix  $N_0 \in \mathbb{R}$  and suppose that  $M := \{x \in X | N(x) = N_0\}$  is nonempty. Then there exists some  $x_0 \in M$  with  $F(x_0) = \inf_M F$ .*

*If in addition  $F$  is differentiable,  $N$  is Gâteaux differentiable and  $\delta N$  does not vanish on  $M$ , then there is a constant  $\lambda \in \mathbb{R}$  (the **Lagrange multiplier**) such that*

$$dF(x_0) = \lambda \delta N(x_0). \quad (9.37)$$

**Proof.** Existence follows from Theorem 9.21 which is applicable thanks to our previous lemma. Now choose some  $x_1 \in X$  such that  $\delta N(x_0)x_1 \neq 0$  and  $x \in X$  arbitrary. Then the function

$$f(t, s) := N(x_0 + tx + x_1s)$$

is  $C^2(\mathbb{R}^2)$  and satisfies

$$\partial_t f(t, s) = \delta N(x_0 + tx + x_1s)x, \quad \partial_s f(t, s) = \delta N(x_0 + tx + x_1s)x_1$$

and since  $\partial_s f(0, 0) \neq 0$  the implicit function theorem implies existence of a function  $\sigma \in C^1(-\varepsilon, \varepsilon)$  such that  $\sigma(0) = 0$  and  $f(t, \sigma(t)) = f(0, 0)$ , that is,  $x(t) := x_0 + tx + \sigma(t)x_1 \in M$  for  $|t| < \varepsilon$ . Moreover,

$$\sigma'(0) = -\frac{\partial_t f(0, 0)}{\partial_s f(0, 0)} = -\frac{\delta N(x_0)x}{\delta N(x_0)x_1}.$$

Hence by the chain rule

$$\frac{d}{dt} F(x_0 + tx + \sigma(t)x_1)|_{t=0} = dF(x_0)(x + \sigma'(0)x_1) = dF(x_0)x - \lambda \delta N(x_0)x = 0,$$

where

$$\lambda := \frac{dF(x_0)x_1}{\delta N(x_0)x_1}. \quad \square$$

**Example 9.27.** Let  $U \subset \mathbb{R}^n$  be a bounded domain and consider

$$F(u) := \frac{1}{2} \int_U |\partial u|^2 d^n x, \quad u \in H_0^1(U, \mathbb{R})$$

subject to the constraint

$$N(u) := \int_U G(u) d^n x = N_0,$$

where  $G : \mathbb{R} \rightarrow \mathbb{R}$  is differentiable and satisfies

$$|G'(x)| \leq C(1 + |x|).$$

This condition implies

$$|G(x)| \leq \tilde{C}(1 + |x|^2)$$

and ensures that  $N(u)$  is well-defined for all  $u \in L^2(U)$ .

In order to apply the theorem we set  $X = H_0^1(U)$  and  $Y = L^2(U)$ . That  $X$  is compactly embedded into  $Y$  is the Rellich–Kondrachov theorem (Theorem 7.35 from [48]). Moreover, by the Poincaré inequality (Theorem 7.38 from [48]) we can choose  $\|x\| := F(x)$  as an equivalent norm on  $X$ . In particular,  $F$  satisfies the requirements of our theorem and so does  $N$  by Example 9.5. Consequently, if  $N_0$  is such that

$$M := \{f \in H_0^1(U, \mathbb{R}) \mid \int_U G(u) d^n x = N_0\}$$

is nonempty, there is a minimizer  $u_0$ . By Example 9.5

$$\delta N(u_0)u = \int_U G'(u_0)u d^n x$$

and if we can find some  $u \in H_0^1(U)$  such that this derivative is nonzero, then  $u_0$  satisfies

$$\int_U (\partial u_0 \cdot \partial u - \lambda G'(u_0)u) d^n x = 0, \quad u \in H_0^1(\mathbb{R}),$$

and hence is a weak solution of the nonlinear eigenvalue problem

$$-\Delta u_0 = \lambda G'(u_0).$$

Note that this last condition is for example satisfied if  $G(0) = 0$ ,  $G'(x)x > 0$  for  $x \neq 0$ , and  $N_0 > 0$ . Indeed, in this case  $\delta N(u_0)u_0 = \int_U G'(u_0)u_0 d^n x > 0$  since otherwise we would have  $u_0 = 0$  contradicting  $0 < N_0 = N(u_0) = N(0) = 0$ .

Of course in the case  $G(x) = \frac{1}{2}|x|^2$  and  $N_0 = 1$  this gives us the lowest eigenvalue of the Laplacian on  $U$  with Dirichlet boundary conditions.

Note that using continuous embeddings  $L^2 \hookrightarrow L^p$  with  $2 \leq p \leq \infty$  for  $n = 1$ ,  $2 \leq p < \infty$  for  $n = 2$ , and  $2 \leq p \leq \frac{2n}{n-2}$  for  $n \geq 3$  one can improve this result to the case

$$|G'(x)| \leq C(1 + |x|^{p-1}). \quad \diamond$$

**Problem 9.13.** Consider  $X = C[0, 1]$  and  $M = \{f \mid \int_0^1 f(x)dx = 1, f(0) = 0\}$ . Show that  $M$  is closed and convex. Show that  $d(0, M) = 1$  but there is no minimizer. If we replace the boundary condition by  $f(0) = 1$  there is a unique minimizer and for  $f(0) = 2$  there are infinitely many minimizers.

**Problem 9.14.** Show that  $F : M \rightarrow \overline{\mathbb{R}}$  is convex if and only if its **epigraph**  $\text{epi } F := \{(x, a) \in M \times \mathbb{R} \mid F(x) \leq a\} \subset X \times \mathbb{R}$  is convex.

**Problem\* 9.15.** Show that  $F : M \rightarrow \overline{\mathbb{R}}$  is quasiconvex if and only if the sublevel sets  $F^{-1}((-\infty, a])$  are convex for every  $a \in \mathbb{R}$ .

## 9.5. Contraction principles

Let  $X$  be a Banach space. A fixed point of a mapping  $F : C \subseteq X \rightarrow C$  is an element  $x \in C$  such that  $F(x) = x$ . Moreover,  $F$  is called a contraction if there is a contraction constant  $\theta \in [0, 1)$  such that

$$|F(x) - F(\tilde{x})| \leq \theta|x - \tilde{x}|, \quad x, \tilde{x} \in C. \quad (9.38)$$

Note that a contraction is continuous. We also recall the notation  $F^n(x) = F(F^{n-1}(x))$ ,  $F^0(x) = x$ .

**Theorem 9.27** (Contraction principle). *Let  $C$  be a nonempty closed subset of a Banach space  $X$  and let  $F : C \rightarrow C$  be a contraction, then  $F$  has a unique fixed point  $\bar{x} \in C$  such that*

$$|F^n(x) - \bar{x}| \leq \frac{\theta^n}{1 - \theta}|F(x) - x|, \quad x \in C. \quad (9.39)$$

**Proof.** If  $x = F(x)$  and  $\tilde{x} = F(\tilde{x})$ , then  $|x - \tilde{x}| = |F(x) - F(\tilde{x})| \leq \theta|x - \tilde{x}|$  shows that there can be at most one fixed point.

Concerning existence, fix  $x_0 \in C$  and consider the sequence  $x_n = F^n(x_0)$ . We have

$$|x_{n+1} - x_n| \leq \theta|x_n - x_{n-1}| \leq \cdots \leq \theta^n|x_1 - x_0|$$

and hence by the triangle inequality (for  $n > m$ )

$$\begin{aligned} |x_n - x_m| &\leq \sum_{j=m+1}^n |x_j - x_{j-1}| \leq \theta^m \sum_{j=0}^{n-m-1} \theta^j |x_1 - x_0| \\ &\leq \frac{\theta^m}{1 - \theta} |x_1 - x_0|. \end{aligned} \quad (9.40)$$

Thus  $x_n$  is Cauchy and tends to a limit  $\bar{x}$ . Moreover,

$$|F(\bar{x}) - \bar{x}| = \lim_{n \rightarrow \infty} |x_{n+1} - x_n| = 0$$

shows that  $\bar{x}$  is a fixed point and the estimate (9.39) follows after taking the limit  $m \rightarrow \infty$  in (9.40).  $\square$

**Example 9.28.** Consider  $X := C[0, 1]$  and let  $C := \{x \in X | x(0) = 0, x(1) = 1, 0 \leq x \leq 1\}$ . Then the map  $F(x)(t) := tx(t)$  maps  $C$  to  $C$  and satisfies  $\|F(x) - F(\tilde{x})\|_\infty < \|x - \tilde{x}\|_\infty$  for  $x \neq \tilde{x}$  and  $x, \tilde{x} \in C$ . Nevertheless it has no fixed point in  $C$ .  $\diamond$

Note that we can replace  $\theta^n$  by any other summable sequence  $\theta_n$  (Problem 9.17):

**Theorem 9.28** (Weissinger). *Let  $C$  be a nonempty closed subset of a Banach space  $X$ . Suppose  $F : C \rightarrow C$  satisfies*

$$|F^n(x) - F^n(y)| \leq \theta_n |x - y|, \quad x, y \in C, \quad (9.41)$$

*with  $\sum_{n=1}^\infty \theta_n < \infty$ . Then  $F$  has a unique fixed point  $\bar{x}$  such that*

$$|F^n(x) - \bar{x}| \leq \left( \sum_{j=n}^\infty \theta_j \right) |F(x) - x|, \quad x \in C. \quad (9.42)$$

Next, we want to investigate how fixed points of contractions vary with respect to a parameter. Let  $X, Y$  be Banach spaces,  $U \subseteq X$ ,  $V \subseteq Y$  be open and consider  $F : \bar{U} \times V \rightarrow U$ . The mapping  $F$  is called a uniform contraction if there is a  $\theta \in [0, 1)$  such that

$$|F(x, y) - F(\tilde{x}, y)| \leq \theta |x - \tilde{x}|, \quad x, \tilde{x} \in \bar{U}, y \in V, \quad (9.43)$$

that is, the contraction constant  $\theta$  is independent of  $y$ .

**Theorem 9.29** (Uniform contraction principle). *Let  $U, V$  be nonempty open subsets of Banach spaces  $X, Y$ , respectively. Let  $F : \bar{U} \times V \rightarrow U$  be a uniform contraction and denote by  $\bar{x}(y) \in U$  the unique fixed point of  $F(\cdot, y)$ . If  $F \in C^r(U \times V, U)$ ,  $r \geq 0$ , then  $\bar{x}(\cdot) \in C^r(V, U)$ . If  $F$  is Lipschitz with respect to the parameter, so is the fixed point.*

**Proof.** Let us first show that  $\bar{x}(y)$  is continuous. From

$$\begin{aligned} |\bar{x}(y+v) - \bar{x}(y)| &= |F(\bar{x}(y+v), y+v) - F(\bar{x}(y), y+v) \\ &\quad + F(\bar{x}(y), y+v) - F(\bar{x}(y), y)| \\ &\leq \theta |\bar{x}(y+v) - \bar{x}(y)| + |F(\bar{x}(y), y+v) - F(\bar{x}(y), y)| \end{aligned}$$

we infer

$$|\bar{x}(y+v) - \bar{x}(y)| \leq \frac{1}{1-\theta} |F(\bar{x}(y), y+v) - F(\bar{x}(y), y)| \quad (9.44)$$

and hence  $\bar{x}(y) \in C(V, U)$ . Now let  $r := 1$  and let us formally differentiate  $\bar{x}(y) = F(\bar{x}(y), y)$  with respect to  $y$ ,

$$d\bar{x}(y) = \partial_x F(\bar{x}(y), y)d\bar{x}(y) + \partial_y F(\bar{x}(y), y). \quad (9.45)$$

Considering this as a fixed point equation  $T(x', y) = x'$ , where  $T(\cdot, y) : \mathcal{L}(Y, X) \rightarrow \mathcal{L}(Y, X)$ ,  $x' \mapsto \partial_x F(\bar{x}(y), y)x' + \partial_y F(\bar{x}(y), y)$  is a uniform contraction since we have  $\|\partial_x F(\bar{x}(y), y)\| \leq \theta$  by Theorem 9.9. Hence we get a unique continuous solution  $\bar{x}'(y)$ . It remains to show

$$\bar{x}(y + v) - \bar{x}(y) - \bar{x}'(y)v = o(v).$$

Let us abbreviate  $u := \bar{x}(y + v) - \bar{x}(y)$ , then using (9.45) and the fixed point property of  $\bar{x}(y)$  we see

$$\begin{aligned} (1 - \partial_x F(\bar{x}(y), y))(u - \bar{x}'(y)v) &= \\ &= F(\bar{x}(y) + u, y + v) - F(\bar{x}(y), y) - \partial_x F(\bar{x}(y), y)u - \partial_y F(\bar{x}(y), y)v \\ &= o(u) + o(v) \end{aligned}$$

since  $F \in C^1(U \times V, U)$  by assumption. Moreover,  $\|(1 - \partial_x F(\bar{x}(y), y))^{-1}\| \leq (1 - \theta)^{-1}$  and  $u = O(v)$  (by (9.44)) implying  $u - \bar{x}'(y)v = o(v)$  as desired.

Finally, suppose that the result holds for some  $r - 1 \geq 1$ . Thus, if  $F$  is  $C^r$ , then  $\bar{x}(y)$  is at least  $C^{r-1}$  and the fact that  $d\bar{x}(y)$  satisfies (9.45) shows  $d\bar{x}(y) \in C^{r-1}(V, U)$  and hence  $\bar{x}(y) \in C^r(V, U)$ .  $\square$

As an important consequence we obtain the implicit function theorem.

**Theorem 9.30** (Implicit function). *Let  $X, Y$ , and  $Z$  be Banach spaces and let  $U, V$  be open subsets of  $X, Y$ , respectively. Let  $F \in C^r(U \times V, Z)$ ,  $r \geq 0$ , and fix  $(x_0, y_0) \in U \times V$ . Suppose  $\partial_x F \in C(U \times V, \mathcal{L}(X, Z))$  exists (if  $r = 0$ ) and  $\partial_x F(x_0, y_0) \in \mathcal{L}(X, Z)$  is an isomorphism. Then there exists an open neighborhood  $U_1 \times V_1 \subseteq U \times V$  of  $(x_0, y_0)$  such that for each  $y \in V_1$  there exists a unique point  $(\xi(y), y) \in U_1 \times V_1$  satisfying  $F(\xi(y), y) = F(x_0, y_0)$ . Moreover,  $\xi$  is in  $C^r(V_1, Z)$  and fulfills (for  $r \geq 1$ )*

$$d\xi(y) = -(\partial_x F(\xi(y), y))^{-1} \circ \partial_y F(\xi(y), y). \quad (9.46)$$

**Proof.** Using the shift  $F \rightarrow F - F(x_0, y_0)$  we can assume  $F(x_0, y_0) = 0$ . Next, the fixed points of  $G(x, y) = x - (\partial_x F(x_0, y_0))^{-1}F(x, y)$  are the solutions of  $F(x, y) = 0$ . The function  $G$  has the same smoothness properties as  $F$  and since  $\partial_x G(x_0, y_0) = 0$ , we can find balls  $U_1$  and  $V_1$  around  $x_0$  and  $y_0$  such that  $\|\partial_x G(x, y)\| \leq \theta < 1$  for  $(x, y) \in U_1 \times V_1$ . Thus by the mean value theorem (Theorem 9.9)  $G(\cdot, y)$  is a uniform contraction on  $U_1$  for  $y \in V_1$ . Moreover, choosing the radius of  $V_1$  sufficiently small such that  $|G(x_0, y) - G(x_0, y_0)| < (1 - \theta)r$  for  $y \in V_1$ , where  $r$  is the radius of  $U_1$ , we get

$$|G(x, y) - x_0| = |G(x, y) - G(x_0, y_0)| \leq \theta|x - x_0| + (1 - \theta)r < r$$



for  $(x, y) \in \overline{U_1} \times V_1$ , that is,  $G : \overline{U_1} \times V_1 \rightarrow U_1$ . The rest follows from the uniform contraction principle. Formula (9.46) follows from differentiating  $F(\xi(y), y) = 0$  using the chain rule.  $\square$

Note that our proof is constructive, since it shows that the solution  $\xi(y)$  can be obtained by iterating  $x - (\partial_x F(x_0, y_0))^{-1}(F(x, y) - F(x_0, y_0))$ .

Moreover, as a corollary of the implicit function theorem we also obtain the inverse function theorem.

**Theorem 9.31** (Inverse function). *Suppose  $F \in C^r(U, Y)$ ,  $r \geq 1$ ,  $U \subseteq X$ , and let  $dF(x_0)$  be an isomorphism for some  $x_0 \in U$ . Then there are neighborhoods  $U_1, V_1$  of  $x_0, F(x_0)$ , respectively, such that  $F \in C^r(U_1, V_1)$  is a diffeomorphism.*

**Proof.** Apply the implicit function theorem to  $G(x, y) = y - F(x)$ .  $\square$

**Example 9.29.** It is important to emphasize that invertibility of  $dF$  on all of  $U$  does not imply injectivity on  $U$  as the following example in  $X := \mathbb{R}^2$  shows:  $F(x, y) = (e^{2x} - y^2 + 3, 4e^{2x}y - y^3)$ . Note that  $\det \frac{\partial F}{\partial(x,y)} = 8e^{4x} + 10e^{2x}y^2$  and  $F(0, 2) = (0, 0) = F(0, -2)$ .  $\diamond$

**Example 9.30.** Let  $X$  be a Banach algebra and  $\mathcal{G}(X)$  the group of invertible elements. We have seen that multiplication is  $C^\infty(X \times X, X)$  and hence taking the inverse is also  $C^\infty(\mathcal{G}(X), \mathcal{G}(X))$ . Consequently,  $\mathcal{G}(X)$  is an (in general infinite-dimensional) **Lie group**.  $\diamond$

Further applications will be given in the next section.

**Problem 9.16.** *Derive Newton's method for finding the zeros of a twice continuously differentiable function  $f(x)$ ,*

$$x_{n+1} = F(x_n), \quad F(x) = x - \frac{f(x)}{f'(x)},$$

*from the contraction principle by showing that if  $\bar{x}$  is a zero with  $f'(\bar{x}) \neq 0$ , then there is a corresponding closed interval  $C$  around  $\bar{x}$  such that the assumptions of Theorem 9.27 are satisfied.*

**Problem\* 9.17.** *Prove Theorem 9.28. Moreover, suppose  $F : C \rightarrow C$  and that  $F^n$  is a contraction. Show that the fixed point of  $F^n$  is also one of  $F$ . Hence Theorem 9.28 (except for the estimate) can also be considered as a special case of Theorem 9.27 since the assumption implies that  $F^n$  is a contraction for  $n$  sufficiently large.*

## 9.6. Ordinary differential equations

As a first application of the implicit function theorem, we prove (local) existence and uniqueness for solutions of ordinary differential equations in Banach spaces. Let  $X$  be a Banach space,  $U \subseteq X$  a (nonempty) open subset, and  $I \subseteq \mathbb{R}$  a compact interval. Denote by  $C(I, U)$  the Banach space of bounded continuous functions equipped with the sup norm.

The following lemma, known as **omega lemma**, will be needed in the proof of the next theorem.

**Lemma 9.32.** *Suppose  $I \subseteq \mathbb{R}$  is a compact interval and  $f \in C^r(U, Y)$ . Then  $f_* \in C^r(C(I, U), C(I, Y))$ , where*

$$(f_*x)(t) := f(x(t)). \quad (9.47)$$

**Proof.** Fix  $x_0 \in C(I, U)$  and  $\varepsilon > 0$ . For each  $t \in I$  we have a  $\delta(t) > 0$  such that  $\overline{B_{2\delta(t)}(x_0(t))} \subset U$  and  $|f(x) - f(x_0(t))| \leq \varepsilon/2$  for all  $x$  with  $|x - x_0(t)| \leq 2\delta(t)$ . The balls  $B_{\delta(t)}(x_0(t))$ ,  $t \in I$ , cover the set  $\{x_0(t)\}_{t \in I}$  and since  $I$  is compact, there is a finite subcover  $B_{\delta(t_j)}(x_0(t_j))$ ,  $1 \leq j \leq n$ . Let  $\|x - x_0\| \leq \delta := \min_{1 \leq j \leq n} \delta(t_j)$ . Then for each  $t \in I$  there is a  $t_j$  such that  $|x_0(t) - x_0(t_j)| \leq \delta(t_j)$  and hence  $|f(x(t)) - f(x_0(t))| \leq |f(x(t)) - f(x_0(t_j))| + |f(x_0(t_j)) - f(x_0(t))| \leq \varepsilon$  since  $|x(t) - x_0(t_j)| \leq |x(t) - x_0(t)| + |x_0(t) - x_0(t_j)| \leq 2\delta(t_j)$ . This settles the case  $r = 0$ .

Next let us turn to  $r = 1$ . We claim that  $df_*$  is given by  $(df_*(x_0)x)(t) := df(x_0(t))x(t)$ . To show this we use Taylor's theorem (cf. the proof of Corollary 9.14) to conclude that

$$|f(x_0(t) + x) - f(x_0(t)) - df(x_0(t))x| \leq |x| \sup_{0 \leq s \leq 1} \|df(x_0(t) + sx) - df(x_0(t))\|.$$

By the first part  $(df)_*$  is continuous and hence for a given  $\varepsilon$  we can find a corresponding  $\delta$  such that  $|x(t) - y(t)| \leq \delta$  implies  $\|df(x(t)) - df(y(t))\| \leq \varepsilon$  and hence  $\|df(x_0(t) + sx) - df(x_0(t))\| \leq \varepsilon$  for  $|x_0(t) + sx - x_0(t)| \leq |x| \leq \delta$ . But this shows differentiability of  $f_*$  as required and it remains to show that  $df_*$  is continuous. To see this we use the linear map

$$\begin{array}{ccc} \lambda : C(I, \mathcal{L}(X, Y)) & \rightarrow & \mathcal{L}(C(I, X), C(I, Y)) \\ T & \mapsto & T_* \end{array}$$

where  $(T_*x)(t) := T(t)x(t)$ . Since we have

$$\|T_*x\| = \sup_{t \in I} |T(t)x(t)| \leq \sup_{t \in I} \|T(t)\| |x(t)| \leq \|T\| \|x\|,$$

we infer  $|\lambda| \leq 1$  and hence  $\lambda$  is continuous. Now observe  $df_* = \lambda \circ (df)_*$ .

The general case  $r > 1$  follows from induction.  $\square$

Now we come to our existence and uniqueness result for the initial value problem in Banach spaces.

**Theorem 9.33.** *Let  $I$  be an open interval,  $U$  an open subset of a Banach space  $X$  and  $\Lambda$  an open subset of another Banach space. Suppose  $F \in C^r(I \times U \times \Lambda, X)$ ,  $r \geq 1$ , then the initial value problem*

$$\dot{x} = F(t, x, \lambda), \quad x(t_0) = x_0, \quad (t_0, x_0, \lambda) \in I \times U \times \Lambda, \quad (9.48)$$

*has a unique solution  $x(t, t_0, x_0, \lambda) \in C^r(I_1 \times I_2 \times U_1 \times \Lambda_1, U)$ , where  $I_{1,2}$ ,  $U_1$ , and  $\Lambda_1$  are open subsets of  $I$ ,  $U$ , and  $\Lambda$ , respectively. The sets  $I_2$ ,  $U_1$ , and  $\Lambda_1$  can be chosen to contain any point  $t_0 \in I$ ,  $x_0 \in U$ , and  $\lambda_0 \in \Lambda$ , respectively.*

**Proof.** Adding  $t$  and  $\lambda$  to the dependent variables  $x$ , that is considering  $(\tau, x, \lambda) \in \mathbb{R} \times X \times \Lambda$  and augmenting the differential equation according to  $(\dot{\tau}, \dot{x}, \dot{\lambda}) = (1, F(\tau, x, \lambda), 0)$ , we can assume that  $F$  is independent of  $t$  and  $\lambda$ . Moreover, by a translation we can even assume  $t_0 = 0$ .

Our goal is to invoke the implicit function theorem. In order to do this we introduce an additional parameter  $\varepsilon \in \mathbb{R}$  and consider

$$\dot{x} = \varepsilon F(x_0 + x), \quad x \in D^1 := \{x \in C^1([-1, 1], B_\delta(0)) | x(0) = 0\}, \quad (9.49)$$

such that we know the solution for  $\varepsilon = 0$ . The implicit function theorem will show that solutions still exist as long as  $\varepsilon$  remains small. At first sight this doesn't seem to be good enough for us since our original problem corresponds to  $\varepsilon = 1$ . But since  $\varepsilon$  corresponds to a scaling  $t \rightarrow \varepsilon t$ , the solution for one  $\varepsilon > 0$  suffices. Now let us turn to the details.

Our problem (9.49) is equivalent to looking for zeros of the function

$$\begin{aligned} G: D^1 \times U_0 \times \mathbb{R} &\rightarrow C([-1, 1], X), \\ (x, x_0, \varepsilon) &\mapsto \dot{x} - \varepsilon F(x_0 + x), \end{aligned} \quad (9.50)$$

where  $U_0$  is a neighborhood of  $x_0$  and  $\delta$  sufficiently small such that  $U_0 + B_\delta(0) \subseteq U$ . Lemma 9.32 ensures that this function is  $C^1$ . Now fix  $x_0$ , then  $G(0, x_0, 0) = 0$  and  $\partial_x G(0, x_0, 0) = T$ , where  $Tx := \dot{x}$ . Since  $(T^{-1}x)(t) = \int_0^t x(s)ds$  we can apply the implicit function theorem to conclude that there is a unique solution  $x(x_0, \varepsilon) \in C^1(U_1 \times (-\varepsilon_0, \varepsilon_0), D^1) \hookrightarrow C^1([-1, 1] \times U_1 \times (-\varepsilon_0, \varepsilon_0), X)$ . In particular, the map  $(t, x_0) \mapsto x_0 + x(x_0, \varepsilon)(t/\varepsilon)$  is in  $C^1((-\varepsilon, \varepsilon) \times U_1, X)$ . Hence it is the desired solution of our original problem. This settles the case  $r = 1$ .

For  $r > 1$  we use induction. Suppose  $F \in C^{r+1}$  and let  $x(t, x_0)$  be the solution which is at least  $C^r$ . Moreover,  $y(t, x_0) := \partial_{x_0} x(t, x_0)$  satisfies

$$\dot{y} = \partial_x F(x(t, x_0))y, \quad y(0) = \mathbb{I},$$

and hence  $y(t, x_0) \in C^r$ . Moreover, the differential equation shows  $\partial_t x(t, x_0) = F(x(t, x_0)) \in C^r$  which shows  $x(t, x_0) \in C^{r+1}$ .  $\square$

**Example 9.31.** The simplest example is a linear equation

$$\dot{x} = Ax, \quad x(0) = x_0,$$

where  $A \in \mathcal{L}(X)$ . Then it is easy to verify that the solution is given by

$$x(t) = \exp(tA)x_0,$$

where

$$\exp(tA) := \sum_{k=0}^{\infty} \frac{t^k}{k!} A^k.$$

It is easy to check that the last series converges absolutely (cf. also Problem 1.39) and solves the differential equation (Problem 9.18).  $\diamond$

**Example 9.32.** The classical example  $\dot{x} = x^2$ ,  $x(0) = x_0$ , in  $X := \mathbb{R}$  with solution

$$x(t) = \frac{x_0}{1 - x_0 t}, \quad t \in \begin{cases} (-\infty, \frac{1}{x_0}), & x_0 > 0, \\ \mathbb{R}, & x_0 = 0, \\ (\frac{1}{x_0}, \infty), & x_0 < 0. \end{cases}$$

shows that solutions might not exist for all  $t \in \mathbb{R}$  even though the differential equation is defined for all  $t \in \mathbb{R}$ .  $\diamond$

This raises the question about the maximal interval on which a solution of the initial value problem (9.48) can be defined.

Suppose that solutions of the initial value problem (9.48) exist locally and are unique (as guaranteed by Theorem 9.33). Let  $\phi_1, \phi_2$  be two solutions of (9.48) defined on the open intervals  $I_1, I_2$ , respectively. Let  $I := I_1 \cap I_2 = (T_-, T_+)$  and let  $(t_-, t_+)$  be the maximal open interval on which both solutions coincide. I claim that  $(t_-, t_+) = (T_-, T_+)$ . In fact, if  $t_+ < T_+$ , both solutions would also coincide at  $t_+$  by continuity. Next, considering the initial value problem with initial condition  $x(t_+) = \phi_1(t_+) = \phi_2(t_+)$  shows that both solutions coincide in a neighborhood of  $t_+$  by local uniqueness. This contradicts maximality of  $t_+$  and hence  $t_+ = T_+$ . Similarly,  $t_- = T_-$ .

Moreover, we get a solution

$$\phi(t) := \begin{cases} \phi_1(t), & t \in I_1, \\ \phi_2(t), & t \in I_2, \end{cases} \quad (9.51)$$

defined on  $I_1 \cup I_2$ . In fact, this even extends to an arbitrary number of solutions and in this way we get a (unique) solution defined on some maximal interval.

**Theorem 9.34.** *Suppose the initial value problem (9.48) has a unique local solution (e.g. the conditions of Theorem 9.33 are satisfied). Then there exists a unique maximal solution defined on some maximal interval  $I_{(t_0, x_0)} = (T_-(t_0, x_0), T_+(t_0, x_0))$ .*

**Proof.** Let  $\mathcal{S}$  be the set of all solutions  $\phi$  of (9.48) which are defined on an open interval  $I_\phi$ . Let  $\mathcal{I} := \bigcup_{\phi \in \mathcal{S}} I_\phi$ , which is again open. Moreover, if  $t_1 > t_0 \in \mathcal{I}$ , then  $t_1 \in I_\phi$  for some  $\phi$  and thus  $[t_0, t_1] \subseteq I_\phi \subseteq \mathcal{I}$ . Similarly for  $t_1 < t_0$  and thus  $\mathcal{I}$  is an open interval containing  $t_0$ . In particular, it is of the form  $\mathcal{I} = (T_-, T_+)$ . Now define  $\phi_{\max}(t)$  on  $\mathcal{I}$  by  $\phi_{\max}(t) := \phi(t)$  for some  $\phi \in \mathcal{S}$  with  $t \in I_\phi$ . By our above considerations any two  $\phi$  will give the same value, and thus  $\phi_{\max}(t)$  is well-defined. Moreover, for every  $t_1 > t_0$  there is some  $\phi \in \mathcal{S}$  such that  $t_1 \in I_\phi$  and  $\phi_{\max}(t) = \phi(t)$  for  $t \in (t_0 - \varepsilon, t_1 + \varepsilon)$  which shows that  $\phi_{\max}$  is a solution. By construction there cannot be a solution defined on a larger interval.  $\square$

The solution found in the previous theorem is called the **maximal solution**. A solution defined for all  $t \in \mathbb{R}$  is called a **global solution**. Clearly every global solution is maximal.

The next result gives a simple criterion for a solution to be global.

**Lemma 9.35.** *Suppose  $F \in C^1(\mathbb{R} \times X, X)$  and let  $x(t)$  be a maximal solution of the initial value problem (9.48). Suppose  $|F(t, x(t))|$  is bounded on finite  $t$ -intervals. Then  $x(t)$  is a global solution.*

**Proof.** Let  $(T_-, T_+)$  be the domain of  $x(t)$  and suppose  $T_+ < \infty$ . Then  $|F(t, x(t))| \leq C$  for  $t \in (t_0, T_+)$  and for  $t_0 < s < t < T_+$  we have

$$|x(t) - x(s)| \leq \int_s^t |\dot{x}(\tau)| d\tau = \int_s^t |F(\tau, x(\tau))| d\tau \leq C|t - s|.$$

Thus  $x(t_n)$  is Cauchy whenever  $t_n$  is and hence  $\lim_{t \rightarrow T_+} x(t) = x_+$  exists. Now let  $y(t)$  be the solution satisfying the initial condition  $y(T_+) = x_+$ . Then

$$\tilde{x}(t) = \begin{cases} x(t), & t < T_+, \\ y(t), & t \geq T_+, \end{cases}$$

is a larger solution contradicting maximality of  $T_+$ .  $\square$

**Example 9.33.** Finally, we want to apply this to a famous example, the so-called **FPU lattices** (after Enrico Fermi, John Pasta, and Stanislaw Ulam who investigated such systems numerically). This is a simple model of a linear chain of particles coupled via nearest neighbor interactions. Let us assume for simplicity that all particles are identical and that the interaction

is described by a potential  $V \in C^2(\mathbb{R})$ . Then the equation of motions are given by

$$\ddot{q}_n(t) = V'(q_{n+1} - q_n) - V'(q_n - q_{n-1}), \quad n \in \mathbb{Z},$$

where  $q_n(t) \in \mathbb{R}$  denotes the position of the  $n$ 'th particle at time  $t \in \mathbb{R}$  and the particle index  $n$  runs through all integers. If the potential is quadratic,  $V(r) = \frac{k}{2}r^2$ , then we get the discrete linear wave equation

$$\ddot{q}_n(t) = k(q_{n+1}(t) - 2q_n(t) + q_{n-1}(t)).$$

If we use the fact that the Jacobi operator  $Aq_n = -k(q_{n+1} - 2q_n + q_{n-1})$  is a bounded operator in  $X = \ell_{\mathbb{R}}^p(\mathbb{Z})$  we can easily solve this system as in the case of ordinary differential equations. In fact, if  $q^0 = q(0)$  and  $p^0 = \dot{q}(0)$  are the initial conditions then one can easily check (cf. Problem 9.18) that the solution is given by

$$q(t) = \cos(tA^{1/2})q^0 + \frac{\sin(tA^{1/2})}{A^{1/2}}p^0.$$

In the Hilbert space case  $p = 2$  these functions of our operator  $A$  could be defined via the spectral theorem but here we just use the more direct definition

$$\cos(tA^{1/2}) := \sum_{k=0}^{\infty} \frac{t^{2k}}{(2k)!} A^k, \quad \frac{\sin(tA^{1/2})}{A^{1/2}} := \sum_{k=0}^{\infty} \frac{t^{2k+1}}{(2k+1)!} A^k.$$

In the general case an explicit solution is no longer possible but we are still able to show global existence under appropriate conditions. To this end we will assume that  $V$  has a global minimum at 0 and hence looks like  $V(r) = V(0) + \frac{k}{2}r^2 + o(r^2)$ . As  $V(0)$  does not enter our differential equation we will assume  $V(0) = 0$  without loss of generality. Moreover, we will also introduce  $p_n := \dot{q}_n$  to have a first order system

$$\dot{q}_n = p_n, \quad \dot{p}_n = V'(q_{n+1} - q_n) - V'(q_n - q_{n-1}).$$

Since  $V' \in C^1(\mathbb{R})$  with  $V'(0) = 0$  it gives rise to a  $C^1$  map on  $\ell_{\mathbb{R}}^p(\mathbb{N})$  (see Example 9.2). Since the same is true for shifts, the chain rule implies that the right-hand side of our system is a  $C^1$  map and hence Theorem 9.33 gives us existence of a local solution. To get global solutions we will need a bound on solutions. This will follow from the fact that the energy of the system

$$H(p, q) := \sum_{n \in \mathbb{Z}} \left( \frac{p_n^2}{2} + V(q_{n+1} - q_n) \right)$$

is conserved. To ensure that the above sum is finite we will choose  $X := \ell_{\mathbb{R}}^2(\mathbb{Z}) \oplus \ell_{\mathbb{R}}^2(\mathbb{Z})$  as our underlying Banach (in this case even Hilbert) space. Recall that since we assume  $V$  to have a minimum at 0 we have  $|V(r)| \leq$

$C_R r^2$  for  $|r| < R$  and hence  $H(p, q) < \infty$  for  $(p, q) \in X$ . Under these assumptions it is easy to check that  $H \in C^1(X, \mathbb{R})$  and that

$$\begin{aligned} \frac{d}{dt} H(p(t), q(t)) &= \sum_{n \in \mathbb{Z}} \left( \dot{p}_n(t) p_n(t) + V'(q_{n+1}(t) - q_n(t)) (\dot{q}_{n+1}(t) - \dot{q}_n(t)) \right) \\ &= \sum_{n \in \mathbb{Z}} \left( (V'(q_{n+1} - q_n) - V'(q_n - q_{n-1})) p_n(t) \right. \\ &\quad \left. + V'(q_{n+1}(t) - q_n(t)) (p_{n+1}(t) - p_n(t)) \right) \\ &= \sum_{n \in \mathbb{Z}} \left( -V'(q_n - q_{n-1}) p_n(t) + V'(q_{n+1}(t) - q_n(t)) p_{n+1}(t) \right) \\ &= 0 \end{aligned}$$

provided  $(p(t), q(t))$  solves our equation. Consequently, since  $V \geq 0$ ,

$$\|p(t)\|_2^2 \leq 2H(p(t), q(t)) = 2H(p(0), q(0)).$$

Moreover,  $q_n(t) = q_n(0) + \int_0^t p_n(s) ds$  (note that since the  $\ell^2$  norm is stronger than the  $\ell^\infty$  norm,  $q_n(t)$  is differentiable for fixed  $n$ ) implies

$$\|q(t)\|_2 \leq \|q(0)\|_2 + \int_0^t \|p_n(s)\|_2 ds \leq \|q(0)\|_2 + \sqrt{2H(p(0), q(0))} t.$$

So Lemma 9.35 ensures that solutions are global in  $X$ . Of course every solution from  $X$  is also a solution from  $Y = \ell_{\mathbb{R}}^p(\mathbb{Z}) \oplus \ell_{\mathbb{R}}^p(\mathbb{Z})$  for all  $p \geq 2$  (since the  $\|\cdot\|_2$  norm is stronger than the  $\|\cdot\|_p$  norm for  $p \geq 2$ ).

Examples include the original FPU  $\beta$ -model  $V_\beta(r) := \frac{1}{2}r^2 + \frac{\beta}{4}r^4$ ,  $\beta > 0$ , and the famous Toda lattice  $V(r) := e^{-r} + r - 1$ .  $\diamond$

**Example 9.34.** Consider the **discrete nonlinear Schrödinger equation (dNLS)**

$$i\dot{u}(t) = Hu(t) \pm |u(t)|^{2p}u(t), \quad t \in \mathbb{R},$$

where  $Hu_n = u_{n+1} + u_{n-1} + q_n u_n$  is the Jacobi operator, in  $X = \ell^2(\mathbb{Z}) \cong \ell_{\mathbb{R}}^2(\mathbb{Z}) \oplus \ell_{\mathbb{R}}^2(\mathbb{Z})$ . Here  $q \in \ell^\infty(\mathbb{Z})$  is a real-valued sequence corresponding to an external potential and  $q = 0$  (or  $q = -2$ , depending on your preferences) is the free discrete Schrödinger operator. Clearly the right-hand side is  $C^1$  for  $p \geq 0$  and hence there is a unique local solution by Theorem 9.33. Please note that even though  $X$  is a complex Banach space we consider it as a real Banach space

Moreover, for a solution we have

$$\frac{d}{dt} \|u(t)\|_2^2 = 2\operatorname{Re}\langle \dot{u}(t), u(t) \rangle = 2\operatorname{Im}(\langle Hu, u \rangle \pm \langle |u(t)|^{2p}u(t), u(t) \rangle) = 0$$

and hence the dNLS has a unique global norm preserving solution  $u \in C^1(\mathbb{R}, \ell^2(\mathbb{Z}))$ . Note that this in fact works for any self-adjoint  $H \in \mathcal{L}(X)$ .  $\diamond$

It should be mentioned that the above theory does not suffice to cover partial differential equations. In fact, if we replace the difference operator by a differential operator we run into the problem that differentiation is not a continuous process!

**Problem\* 9.18.** *Let*

$$f(z) := \sum_{j=0}^{\infty} f_j z^j, \quad |z| < R,$$

*be a convergent power series with convergence radius  $R > 0$ . Suppose  $X$  is a Banach space and  $A \in \mathcal{L}(X)$  is a bounded operator with  $\|A\| < R$ . Show that*

$$f(tA) := \sum_{j=0}^{\infty} f_j t^j A^j$$

*is in  $C^\infty(I, \mathcal{L}(X))$ ,  $I = (-R\|A\|^{-1}, R\|A\|^{-1})$  and*

$$\frac{d^n}{dt^n} f(tA) = A^n f^{(n)}(tA), \quad n \in \mathbb{N}_0.$$

*(Compare also Problem 1.39.)*

**Problem 9.19.** *Consider the FPU  $\alpha$ -model  $V_\alpha(r) := \frac{1}{2}r^2 + \frac{\alpha}{3}r^3$ . Show that solutions satisfying  $\|q_{n+1}(0) - q_n(0)\|_\infty < \frac{1}{|\alpha|}$  and  $H(p(0), q(0)) < \frac{1}{6\alpha^2}$  are global in  $X := \ell^2(\mathbb{Z}) \oplus \ell^2(\mathbb{Z})$ . (Hint: Of course local solutions follow from our considerations above. Moreover, note that  $V_\alpha(r)$  has a maximum at  $r = -\frac{1}{\alpha}$ . Now use conservation of energy to conclude that the solution cannot escape the region  $|r| < \frac{1}{|\alpha|}$ .)*

## 9.7. Bifurcation theory

One of the most basic tasks is finding the zeros of a given function  $F \in C^k(U, Y)$ , where  $U \subseteq X$  and  $X, Y$  are Banach spaces. Frequently the equation will depend on a parameter  $\mu \in \mathbb{R}$  (of course we could also consider the case where the parameter is again in some Banach space, but we will only consider this basic case here). That is, we are looking for solutions  $x(\mu)$  of the equation

$$F(\mu, x) = 0, \tag{9.52}$$

where  $F \in C^k(I \times U, Y)$  for some suitable  $k \in \mathbb{N}$  and some open interval  $I \subseteq \mathbb{R}$ . Moreover, we are interested in the case of values  $\mu_0 \in I$ , where there is a change in the number of solutions (i.e. where new solutions appear or old solutions disappear as  $\mu$  increases). Such points  $\mu_0 \in I$  will be called **bifurcation points**. Clearly this cannot happen at a point where the implicit function theorem is applicable and hence a necessary condition for a



bifurcation point is that

$$\partial_x F(\mu_0, x_0) \quad (9.53)$$

is not invertible at some point  $x_0$  satisfying the equation  $F(\mu_0, x_0) = 0$ .

**Example 9.35.** Consider  $f(\mu, x) = x^2 - \mu$  in  $C^\infty(\mathbb{R} \times \mathbb{R}, \mathbb{R})$ . Then  $f(\mu, x) = x^2 - \mu = 0$ ,  $\partial_x f(\mu, x) = 2x = 0$  shows that  $\mu_0 = 0$  is the only possible bifurcation point. Since there are no solutions for  $\mu < 0$  and two solutions  $x_0(\mu) = \pm\sqrt{\mu}$  for  $\mu > 0$ , we see that  $\mu_0 = 0$  is a bifurcation point. Consider  $f(\mu, x) = x^3 - \mu$  in  $C^\infty(\mathbb{R} \times \mathbb{R}, \mathbb{R})$ . Then  $f(\mu, x) = x^3 - \mu = 0$ ,  $\partial_x f(\mu, x) = 3x^2 = 0$  shows that again  $\mu_0 = 0$  is the only possible bifurcation point. However, this time there is only one solution  $x(\mu) = \text{sign}(\mu)|\mu|^{1/3}$  for all  $\mu \in \mathbb{R}$  and hence there is no bifurcation occurring at  $\mu_0 = 0$ .  $\diamond$

So the derivative  $\partial_x F$  tells us where to look for bifurcation points while further derivatives can be used to determine what kind of bifurcation (if any) occurs. Here we want to show how this can be done in the infinite dimensional case.

Suppose we have an abstract problem  $F(\mu, x) = 0$  with  $\mu \in \mathbb{R}$  and  $x \in X$  some Banach space. We assume that  $F \in C^1(\mathbb{R} \times X, X)$  and that there is a trivial solution  $x = 0$ , that is,  $F(\mu, 0) = 0$ .

The first step is to split off the trivial solution and reduce it effectively to a finite-dimensional problem. To this end we assume that we have found a point  $\mu_0 \in \mathbb{R}$  such that the derivative  $A := \partial_x F(\mu_0, 0)$  is not invertible. Moreover, we will assume that  $A$  is a Fredholm operator such that there exists (cf. Section 7.2) continuous linear projections

$$P : X = \text{Ker}(A) \dot{+} X_0 \rightarrow \text{Ker}(A), \quad Q : X = X_1 \dot{+} \text{Ran}(A) \rightarrow X_1. \quad (9.54)$$

Now split our equation into a system of two equations according to the above splitting of the underlying Banach space:

$$F(\mu, x) = 0 \quad \Leftrightarrow \quad F_1(\mu, u, v) = 0, \quad F_2(\mu, u, v) = 0, \quad (9.55)$$

where  $x = u + v$  with  $u = Px \in \text{Ker}(A)$ ,  $v = (1 - P)x \in X_0$  and  $F_1(\mu, u, v) = QF(\mu, u + v)$ ,  $F_2(\mu, u, v) = (1 - Q)F(\mu, u + v)$ .

Since  $P, Q$  are bounded, this system is still  $C^1$  and the derivatives are given by (recall the block structure of  $A$  from (7.23))

$$\begin{aligned} \partial_u F_1(\mu_0, 0, 0) &= 0, & \partial_v F_1(\mu_0, 0, 0) &= 0, \\ \partial_u F_2(\mu_0, 0, 0) &= 0, & \partial_v F_2(\mu_0, 0, 0) &= A_0. \end{aligned} \quad (9.56)$$

Moreover, since  $A_0$  is an isomorphism, the implicit function theorem tells us that we can (locally) solve  $F_2$  for  $v$ . That is, there exists a neighborhood  $U$  of  $(\mu_0, 0) \in \mathbb{R} \times \text{Ker}(A)$  and a unique function  $\psi \in C^1(U, X_0)$  such that

$$F_2(\mu, u, \psi(\mu, u)) = 0, \quad (\mu, u) \in U. \quad (9.57)$$

In particular, by the uniqueness part we have  $\psi(\mu, 0) = 0$ . Moreover,  $\partial_u \psi(\mu_0, 0) = -A_0^{-1} \partial_u F_2(\mu_0, 0, 0) = 0$ .

Plugging this into the first equation reduces to the original system to the finite dimensional system

$$\tilde{F}_1(\mu, u) = F_1(\mu, u, \psi(\mu, u)) = 0. \quad (9.58)$$

Of course the chain rule tells us that  $\tilde{F} \in C^1$ . Moreover, we still have  $\tilde{F}_1(\mu, 0) = F_1(\mu, 0, \psi(\mu, 0)) = QF(\mu, 0) = 0$  as well as

$$\partial_u \tilde{F}_1(\mu_0, 0) = \partial_u F_1(\mu_0, 0, 0) + \partial_v F_1(\mu_0, 0, 0) \partial_u \psi(\mu_0, 0) = 0. \quad (9.59)$$

This is known as **Lyapunov–Schmidt reduction**.

Now that we have reduced the problem to a finite-dimensional system, it remains to find conditions such that the finite dimensional system has a nontrivial solution. For simplicity we make the requirement

$$\dim \text{Ker}(A) = \dim \text{Coker}(A) = 1 \quad (9.60)$$

such that we actually have a problem in  $\mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ .

Explicitly, let  $u_0$  span  $\text{Ker}(A)$  and let  $u_1$  span  $X_1$ . Then we can write

$$\tilde{F}_1(\mu, \lambda u_0) = f(\mu, \lambda) u_1, \quad (9.61)$$

where  $f \in C^1(V, \mathbb{R})$  with  $V = \{(\mu, \lambda) | (\mu, \lambda u_0) \in U\} \subseteq \mathbb{R}^2$  a neighborhood of  $(\mu_0, 0)$ . Of course we still have  $f(\mu, 0) = 0$  for  $(\mu, 0) \in V$  as well as

$$\partial_\lambda f(\mu_0, 0) u_1 = \partial_u \tilde{F}_1(\mu_0, 0) u_0 = 0. \quad (9.62)$$

It remains to investigate  $f$ . To split off the trivial solution it suggests itself to write

$$f(\mu, \lambda) = \lambda g(\mu, \lambda) \quad (9.63)$$

We already have

$$g(\mu_0, 0) = \partial_\lambda f(\mu_0, 0) = 0 \quad (9.64)$$

and hence if

$$0 \neq \partial_\mu g(\mu_0, 0) = \partial_\mu \partial_\lambda f(\mu_0, 0) \neq 0 \quad (9.65)$$

the implicit function theorem implies existence of a function  $\mu(\lambda)$  with  $\mu(0) = \mu_0$  and  $g(\mu(\lambda), \lambda) = 0$ . Moreover,  $\mu'(0) = -\frac{\partial_\lambda g(\mu_0, 0)}{\partial_\mu g(\mu_0, 0)} = -\frac{\partial_\lambda^2 f(\mu_0, 0)}{2\partial_\mu \partial_\lambda f(\mu_0, 0)}$ .

Of course this last condition is a bit problematic since up to this point we only have  $f \in C^1$  and hence  $g \in C^0$ . However, if we change our original assumption to  $F \in C^2$  we get  $f \in C^2$  and thus  $g \in C^1$ .

So all we need to do is to trace back our definitions and compute

$$\begin{aligned} \partial_\lambda^2 f(\mu_0, 0) u_1 &= \partial_\lambda^2 \tilde{F}_1(\mu_0, \lambda u_0) \Big|_{\lambda=0} = \partial_\lambda^2 F_1(\mu_0, \lambda u_0, \psi(\mu_0, \lambda u_0)) \Big|_{\lambda=0} \\ &= \partial_u^2 F_1(\mu_0, 0, 0)(u_0, u_0) = Q \partial_x^2 F(\mu_0, 0)(u_0, u_0) \end{aligned}$$

(recall  $\partial_u \psi(\mu_0, 0) = 0$ ) and

$$\begin{aligned} \partial_\mu \partial_\lambda f(\mu_0, 0)u_1 &= \partial_\mu \partial_\lambda \tilde{F}_1(\mu, \lambda u_0) \Big|_{\lambda=0, \mu=\mu_0} \\ &= \partial_\mu \partial_\lambda F_1(\mu, \lambda u_0, \psi(\mu, \lambda u_0)) \Big|_{\lambda=0, \mu=\mu_0} \\ &= Q \partial_\mu \partial_x F(\mu_0, 0)u_0. \end{aligned}$$

**Theorem 9.36** (Crandall–Rabinowitz). *Assume  $F \in C^2(\mathbb{R} \times X, X)$  with  $F(\mu, 0) = 0$  for all  $\mu \in \mathbb{R}$ . Suppose that for some  $\mu_0 \in \mathbb{R}$  we have that  $\partial_x F(\mu_0, 0)$  is a Fredholm operator of index zero with a one-dimensional kernel spanned by  $u_0 \in X$ . Then, if*

$$\partial_\mu \partial_x F(\mu_0, 0)u_0 \notin \text{Ran}(\partial_x F(\mu_0, 0)) \quad (9.66)$$

*there are open neighborhoods  $I \subseteq \mathbb{R}$  of 0,  $J \subseteq \mathbb{R}$  of  $\mu_0$ , and  $U \subseteq \text{span}\{u_0\}$  of 0 plus corresponding functions  $\mu \in C^1(I, J)$  and  $\psi \in C^2(J \times U, X_0)$  such that every nontrivial solution of  $F(\mu, x) = 0$  in a neighborhood of  $(\mu_0, 0)$  is given by*

$$x(\lambda) = \lambda u_0 + \psi(\mu(\lambda), \lambda u_0). \quad (9.67)$$

Moreover,

$$\mu(\lambda) = \mu_0 - \frac{\ell_1(\partial_x^2 F(\mu_0, 0)(u_0, u_0))}{2\ell_1(\partial_\mu \partial_x F(\mu_0, 0)u_0)} \lambda + o(\lambda), \quad x(\lambda) = \lambda u_0 + o(\lambda). \quad (9.68)$$

where  $\ell_1$  is any nontrivial functional which vanishes on  $\text{Ran}(\partial_x F(\mu_0, 0))$ .

Note that if  $Q \partial_x^2 F(\mu_0, 0)(u_0, u_0) \neq 0$  we could have also solved for  $\lambda$  obtaining a function  $\lambda(\mu)$  with  $\lambda(\mu_0) = 0$ . However, in this case it is not obvious that  $\lambda(\mu) \neq 0$  for  $\mu \neq \mu_0$ , and hence that we get a nontrivial solution, unless we also require  $Q \partial_\mu \partial_x F(\mu_0, 0)u_0 \neq 0$  which brings us back to our previous condition. If both conditions are met, then  $\mu'(0) \neq 0$  and there is a unique nontrivial solution  $x(\mu)$  which crosses the trivial solution non transversally at  $\mu_0$ . This is known as **transcritical bifurcation**. If  $\mu'(0) = 0$  but  $\mu''(0) \neq 0$  (assuming this derivative exists), then two solutions will branch off (either for  $\mu > \mu_0$  or for  $\mu < \mu_0$  depending on the sign of the second derivative). This is known as a **pitchfork bifurcation** and is typical in case of an odd function  $F(\mu, -x) = -F(\mu, x)$ .

**Example 9.36.** Now we can establish existence of a stationary solution of the dNLS of the form

$$u_n(t) = e^{-it\omega} \phi_n(\omega)$$

Plugging this ansatz into the dNLS we get the stationary dNLS

$$H\phi - \omega\phi \pm |\phi|^{2p}\phi = 0.$$

Of course we always have the trivial solution  $\phi = 0$ .

Applying our analysis to

$$F(\omega, \phi) = (H - \omega)\phi \pm |\phi|^{2p}\phi, \quad p > \frac{1}{2},$$

we have (with respect to  $\phi = \phi_r + i\phi_i \cong (\phi_r, \phi_i)$ )

$$\partial_\phi F(\omega, \phi)u = (H - \omega)u \pm 2p|\phi|^{2(p-1)} \begin{pmatrix} \phi_r^2 & \phi_r\phi_i \\ \phi_r\phi_i & \phi_i^2 \end{pmatrix} u \pm |\phi|^{2p}u$$

and in particular  $\partial_\phi F(\omega, 0) = H - \omega$  and hence  $\omega$  must be an eigenvalue of  $H$ . In fact, if  $\omega_0$  is a discrete eigenvalue, then self-adjointness implies that  $H - \omega_0$  is Fredholm of index zero. Moreover, if there are two eigenfunction  $u$  and  $v$ , then one checks that the Wronskian  $W(u, v) = u(n)v(n+1) - u(n+1)v(n)$  is constant. But square summability implies that the Wronskian must vanish and hence  $u$  and  $v$  must be linearly dependent (note that a solution of  $Hu = \omega_0 u$  vanishing at two consecutive points must vanish everywhere). Hence eigenvalues are always simple for our Jacobi operator  $H$ . Finally, if  $u_0$  is the eigenfunction corresponding to  $\omega_0$  we have

$$\partial_\omega \partial_\phi F(\omega_0, 0)u_0 = -u_0 \notin \text{Ran}(H - \omega_0) = \text{Ker}(H - \omega_0)^\perp$$

and the Crandall–Rabinowitz theorem ensures existence of a stationary solution  $\phi$  for  $\omega$  in a neighborhood of  $\omega_0$ . Note that

$$\partial_\phi^2 F(\omega, \phi)(u, v) = \pm 2p(2p+1)|\phi|^{2p-1} \text{sign}(\phi)uv$$

and hence  $\partial_\phi^2 F(\omega, 0) = 0$ . This is of course not surprising and related to the symmetry  $F(\omega, -\phi) = -F(\omega, \phi)$  which implies that zeros branch off in symmetric pairs.

Of course this leaves the problem of finding a discrete eigenvalue open. One can show that for the free operator  $H_0$  (with  $q = 0$ ) the spectrum is  $\sigma(H_0) = [-2, 2]$  and that there are no eigenvalues (in fact, the discrete Fourier transform will map  $H_0$  to a multiplication operator in  $L^2[-\pi, \pi]$ ). If  $q \in c_0(\mathbb{Z})$ , then the corresponding multiplication operator is compact and  $\sigma_{\text{ess}}(H) = \sigma(H_0)$  by Weyl's theorem (Theorem 7.16). Hence every point in  $\sigma(H) \setminus [-2, 2]$  will be an isolated eigenvalue.  $\diamond$

**Problem 9.20.** Show that if  $F(\mu, -x) = -F(\mu, x)$ , then  $\psi(\mu, -u) = -\psi(\mu, u)$  and  $\mu(-\lambda) = \mu(\lambda)$ .



# Operator semigroups

In this chapter we want to look at (semi)linear ordinary linear differential equations in Banach spaces. We will need a few relevant facts about differentiation and integration for Banach space valued functions. Section 9.1 will be sufficient.

## 10.1. Uniformly continuous operator groups

Our aim is to investigate the abstract Cauchy problem

$$\dot{u} = Au, \quad u(0) = u_0 \quad (10.1)$$

in some Banach space  $X$ . Here  $A$  is some linear operator and we will assume that  $A \in \mathcal{L}(X)$  to begin with. Note that in the simplest case  $X = \mathbb{R}^n$  this is just a linear first order system with constant coefficient matrix  $A$ . In this case the solution is given by

$$u(t) = T(t)u_0, \quad (10.2)$$

where

$$T(t) := \exp(tA) := \sum_{j=0}^{\infty} \frac{t^j}{j!} A^j \quad (10.3)$$

is the exponential of  $tA$ . It is not difficult to see that this also gives the solution in our Banach space setting.

**Theorem 10.1.** *Let  $A \in \mathcal{L}(X)$ . Then the series in (10.3) converges and defines a **uniformly continuous operator group**:*

- (i) *The map  $t \mapsto T(t)$  is continuous,  $T \in C(\mathbb{R}, \mathcal{L}(X))$ .*
- (ii)  *$T(0) = \mathbb{I}$  and  $T(t+s) = T(t)T(s)$  for all  $t, s \in \mathbb{R}$ .*

Moreover,  $T \in C^\infty(\mathbb{R}, \mathcal{L}(X))$  is the unique solution of  $\dot{T}(t) = AT(t)$  with  $T(0) = \mathbb{I}$  and it commutes with  $A$ ,  $AT(t) = T(t)A$ .

**Proof.** Set

$$T_n(t) := \sum_{j=0}^n \frac{t^j}{j!} A^j.$$

Then (for  $m \leq n$ )

$$\|T_n(t) - T_m(t)\| = \left\| \sum_{j=m+1}^n \frac{t^j}{j!} A^j \right\| \leq \sum_{j=m+1}^n \frac{|t|^j}{j!} \|A\|^j \leq \frac{|t|^{m+1}}{(m+1)!} \|A\|^{m+1} e^{|t|\|A\|}.$$

In particular,

$$\|T(t)\| \leq e^{|t|\|A\|}$$

and  $AT(t) = \lim_{n \rightarrow \infty} AT_n(t) = \lim_{n \rightarrow \infty} T_n(t)A = T(t)A$ . Furthermore we have  $\dot{T}_{n+1} = AT_n$  and thus

$$T_{n+1}(t) = \mathbb{I} + \int_0^t AT_n(s) ds.$$

Taking limits shows

$$T(t) = \mathbb{I} + \int_0^t AT(s) ds$$

or equivalently  $T \in C^1(\mathbb{R}, \mathcal{L}(X))$  and  $\dot{T}(t) = AT(t)$ ,  $T(0) = \mathbb{I}$ . Differentiating this last equation shows  $(\frac{d}{dt})^{k+1}T(t) = A(\frac{d}{dt})^kT(t)$  and establishes  $T \in C^\infty(\mathbb{R}, \mathcal{L}(X))$ .

Suppose  $S(t)$  is another solution,  $\dot{S} = AS$ ,  $S(0) = \mathbb{I}$ . Then, by the product rule (Problem 9.1),  $\frac{d}{dt}T(-t)S(t) = T(-t)AS(t) - AT(-t)S(t) = 0$  implying  $T(-t)S(t) = T(0)S(0) = \mathbb{I}$ . In the special case  $T = S$  this shows  $T(-t) = T^{-1}(t)$  and in the general case it hence proves uniqueness  $S = T$ . Finally,  $T(t+s)$  and  $T(t)T(s)$  both satisfy our differential equation and coincide at  $t = 0$ . Hence they coincide for all  $t$  by uniqueness.  $\square$

Note that choosing  $s = -t$  in (ii) shows

$$T(t)^{-1} = T(-t). \quad (10.4)$$

Clearly  $A$  is uniquely determined by  $T(t)$  via  $A = \dot{T}(0)$ . Moreover, from this we also easily get uniqueness for our original Cauchy problem. We will in fact be slightly more general and consider the inhomogeneous problem

$$\dot{u} = Au + g, \quad u(0) = u_0, \quad (10.5)$$

where  $g \in C(I, X)$ . A solution necessarily satisfies

$$\frac{d}{dt}T(-t)u(t) = -AT(-t)u(t) + T(-t)\dot{u}(t) = T(-t)g(t)$$

and integrating this equation (fundamental theorem of calculus) shows the **Duhamel formula**

$$u(t) = T(t) \left( u_0 + \int_0^t T(-s)g(s)ds \right) = T(t)u_0 + \int_0^t T(t-s)g(s)ds. \quad (10.6)$$

**Lemma 10.2.** *Let  $A \in \mathcal{L}(X)$  and  $g \in C^k(I, X)$  for some  $k \in \mathbb{N}_0$ . Then (10.5) has a unique solution  $u \in C^{k+1}(I, X)$  given by (10.6).*

**Proof.** Using Problem 10.2 it is straightforward to verify that this is indeed a solution for any given  $g \in C(I, X)$ . This also shows that  $u \in C^{k+1}(I, X)$  since  $T \in C^\infty(\mathbb{R}, \mathcal{L}(X))$ .  $\square$

**Example 10.1.** For example look at the discrete linear wave equation

$$\ddot{q}_n(t) = k(q_{n+1}(t) - 2q_n(t) + q_{n-1}(t)), \quad n \in \mathbb{Z}.$$

Factorizing this equation according to

$$\dot{q}_n(t) = p_n(t), \quad \dot{p}_n(t) = k(q_{n+1}(t) - 2q_n(t) + q_{n-1}(t)),$$

we can write this as a first order system

$$\frac{d}{dt} \begin{pmatrix} q_n \\ p_n \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ kA_0 & 0 \end{pmatrix} \begin{pmatrix} q_n \\ p_n \end{pmatrix}$$

with the Jacobi operator  $(A_0q)_n = q_{n+1} - 2q_n + q_{n-1}$ . Since  $A_0$  is a bounded operator on  $X = \ell^p(\mathbb{Z})$ , we obtain a well-defined uniformly continuous operator group in  $\ell^p(\mathbb{Z}) \oplus \ell^p(\mathbb{Z})$ .  $\diamond$

**Problem 10.1.** *Show that if  $A, B \in \mathcal{L}(X)$  commute,  $[A, B] := AB - BA = 0$ , then so do their associated groups and we have*

$$\exp(sA + tB) = \exp(sA)\exp(tB) = \exp(tB)\exp(sA), \quad [A, B] = 0.$$

**Problem\* 10.2** (Product rule). *Suppose  $f \in C^1(I, X)$  and  $T \in C^1(I, \mathcal{L}(X, Y))$ . Show that  $Tf \in C^1(I, Y)$  and  $\frac{d}{dt}Tf = \dot{T}f + T\dot{f}$ .*

**Problem 10.3.** *Let  $X$  be a Hilbert space and  $A \in \mathcal{L}(X)$ . Show that  $T(t)^*$  is a uniformly continuous operator group whose generator is  $A^*$ . Conclude that if  $A$  is skew adjoint, that is,  $A^* = -A$ , then  $T$  is unitary.*

**Problem 10.4.** *Discuss the discrete Schrödinger equation*

$$i\ddot{u} = Hu, \quad (Hu)_n := u_{n+1} + u_{n-1} + q_n u_n,$$

in  $\ell^2(\mathbb{Z})$ , where  $q \in \ell^\infty(\mathbb{Z}, \mathbb{R})$ . In particular, show  $\|u(t)\| = \|u(0)\|$  and  $\langle u(t), Hu(t) \rangle = \langle u(0), Hu(0) \rangle$ .

**Problem 10.5.** *Let  $X := C_0(\mathbb{R})$  and consider  $(A_\alpha f)(x) := \frac{1}{\alpha}(f(x + \alpha) - f(x))$  for  $\alpha > 0$ . Show that  $A_\alpha$  is bounded and compute its norm. Compute the corresponding group  $T$  as well as its norm.*



## 10.2. Strongly continuous semigroups

In the previous section we have found a quite complete solution of the abstract Cauchy problem (10.5) in the case when  $A$  is bounded. However, since differential operators are typically unbounded, this assumption is too strong for applications to partial differential equations. Since it is unclear what the conditions on  $A$  should be, we will go the other way and impose conditions on  $T$ . First of all, even rather simple equations like the heat equation are only solvable for positive times and hence we will only assume that the solutions give rise to a semigroup. Moreover, continuity in the operator topology is too much to ask for (in fact, it is equivalent to boundedness of  $A$  — Problem 10.6) and hence we go for the next best option, namely strong continuity. In this sense, our problem is still well-posed.

A **strongly continuous operator semigroup** (also  $C_0$ -semigroup) is a family of bounded operators  $T(t) \in \mathcal{L}(X)$ ,  $t \geq 0$ , such that

- (i)  $T(t)g \in C([0, \infty), X)$  for every  $g \in X$  (strong continuity) and
- (ii)  $T(0) = \mathbb{I}$ ,  $T(t+s) = T(t)T(s)$  for all  $t, s \geq 0$  (semigroup property).

If  $T(t) \in \mathcal{L}(X)$  is defined for  $t \in \mathbb{R}$  and item (ii) holds for all  $t, s \in \mathbb{R}$  it is called a **strongly continuous operator group**.

We first note that  $\|T(t)\|$  is uniformly bounded on compact time intervals.

**Lemma 10.3.** *Let  $T(t)$  be a  $C_0$ -semigroup. Then there are constants  $M \geq 1$ ,  $\omega \geq 0$  such that*

$$\|T(t)\| \leq Me^{\omega t}, \quad t \geq 0. \quad (10.7)$$

*In case of a  $C_0$ -group we have  $\|T(t)\| \leq Me^{\omega|t|}$ ,  $t \in \mathbb{R}$ .*

**Proof.** Since  $\|T(\cdot)g\| \in C[0, 1]$  for every  $g \in X$  we have  $\sup_{t \in [0, 1]} \|T(t)g\| \leq M_g$ . Hence by the uniform boundedness principle  $\sup_{t \in [0, 1]} \|T(t)\| \leq M$  for some  $M \geq 1$ . Setting  $\omega = \log(M)$  the claim follows by induction using the semigroup property. For the group case apply the semigroup case to both  $T(t)$  and  $S(t) := T(-t)$ .  $\square$

The infimum over all possible  $\omega$  for which a corresponding  $M$  exists such that (10.7) holds is known as the **growth bound**  $\omega_0(T)$  of the semigroup. It is given by

$$\omega_0(T) := \limsup_{t \rightarrow \infty} \frac{\log(\|T(t)\|)}{t} \quad (10.8)$$

and it can be shown that the limsup is actually a limit (Problem 10.7). However, there will not be a corresponding constant  $M$

Inspired by the previous section we define the **generator**  $A$  of a strongly continuous semigroup as the linear operator

$$Af := \lim_{t \downarrow 0} \frac{1}{t} (T(t)f - f), \quad (10.9)$$

where the domain  $\mathfrak{D}(A)$  is precisely the set of all  $f \in X$  for which the above limit exists. By linearity of limits  $\mathfrak{D}(A)$  is a linear subspace of  $X$  (and  $A$  is a linear operator) but at this point it is unclear whether it contains any nontrivial elements. We will however postpone this issue and begin with the observation that a  $C_0$ -semigroup is the solution of the abstract Cauchy problem associated with its generator  $A$ :

**Lemma 10.4.** *Let  $T(t)$  be a  $C_0$ -semigroup with generator  $A$ . If  $f \in \mathfrak{D}(A)$  then  $T(t)f \in \mathfrak{D}(A)$  and  $AT(t)f = T(t)Af$ . Moreover, suppose  $g \in X$  with  $u(t) := T(t)g \in \mathfrak{D}(A)$  for  $t > 0$ . Then  $u(t) \in C([0, \infty), X) \cap C^1((0, \infty), X)$  and  $u(t)$  is the unique solution of the abstract Cauchy problem*

$$\dot{u}(t) = Au(t), \quad u(0) = g. \quad (10.10)$$

*This is, for example, the case if  $g \in \mathfrak{D}(A)$  in which case we even have  $u(t) \in C^1([0, \infty), X)$ .*

*Similarly, if  $T(t)$  is a  $C_0$ -group and  $g \in \mathfrak{D}(A)$ , then  $u(t) := T(t)g \in C^1(\mathbb{R}, X)$  is the unique solution of (10.10) for all  $t \in \mathbb{R}$ .*

**Proof.** Let  $f \in \mathfrak{D}(A)$  and  $t > 0$  (respectively  $t \in \mathbb{R}$  for a group), then

$$\lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} (u(t + \varepsilon) - u(t)) = \lim_{\varepsilon \downarrow 0} T(t) \frac{1}{\varepsilon} (T(\varepsilon)f - f) = T(t)Af.$$

This shows the first part. To show that  $u(t)$  is differentiable it remains to compute

$$\begin{aligned} \lim_{\varepsilon \downarrow 0} \frac{1}{-\varepsilon} (u(t - \varepsilon) - u(t)) &= \lim_{\varepsilon \downarrow 0} T(t - \varepsilon) \frac{1}{\varepsilon} (T(\varepsilon)f - f) \\ &= \lim_{\varepsilon \downarrow 0} T(t - \varepsilon) (Af + o(1)) = T(t)Af \end{aligned}$$

since  $\|T(t)\|$  is bounded on compact  $t$  intervals. Hence  $u(t) \in C^1([0, \infty), X)$  (respectively  $u(t) \in C^1(\mathbb{R}, X)$  for a group) solves (10.10). In the general case  $f = T(t_0)g \in \mathfrak{D}(A)$  and  $u(t) = T(t)g = T(t - t_0)f$  solves our differential equation for every  $t > t_0$ . Since  $t_0 > 0$  is arbitrary it follows that  $u(t)$  solves (10.10) by the first part. To see that it is the only solution, let  $v(t)$  be a

solution corresponding to the initial condition  $v(0) = 0$ . For  $s < t$  we have

$$\begin{aligned} \frac{d}{ds}T(t-s)v(s) &= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (T(t-s-\varepsilon)v(s+\varepsilon) - T(t-s)v(s)) \\ &= \lim_{\varepsilon \rightarrow 0} T(t-s-\varepsilon) \frac{1}{\varepsilon} (v(s+\varepsilon) - v(s)) \\ &\quad - \lim_{\varepsilon \rightarrow 0} T(t-s-\varepsilon) \frac{1}{\varepsilon} (T(\varepsilon)v(s) - v(s)) \\ &= T(t-s)Av(s) - T(t-s)Av(s) = 0. \end{aligned}$$

Whence,  $v(t) = T(t-t)v(t) = T(t-s)v(s) = T(t)v(0) = 0$ .  $\square$

Note that our proof in fact even shows a bit more: If  $g \in \mathfrak{D}(A)$  we have  $u \in C^1([0, \infty), X)$  and hence not only  $u \in C([0, \infty), X)$  but also  $Au = \dot{u} \in C([0, \infty), X)$ . Hence, if we regard  $\mathfrak{D}(A)$  as a normed space equipped with the graph norm  $\|f\|_A := \|f\| + \|Af\|$ , in which case we will write  $[\mathfrak{D}(A)]$ , then  $g \in \mathfrak{D}(A)$  implies  $u \in C([0, \infty), [\mathfrak{D}(A)])$ . In particular,  $T$  restricted to  $[\mathfrak{D}(A)]$  is again a  $C_0$ -semigroup and it is straightforward to check that its generator is  $A$  restricted to  $\mathfrak{D}(A^2) = \{f \in \mathfrak{D}(A) | Af \in \mathfrak{D}(A)\}$ .

Similarly,  $u(t) = T(t)g \in \mathfrak{D}(A)$  for  $t > 0$  implies  $u \in C((0, \infty), [\mathfrak{D}(A)])$ . Moreover, recall that  $[\mathfrak{D}(A)]$  will be a Banach space if and only if  $A$  is a closed operator (cf. Section ??) and the latter fact will be established in Corollary 10.7 below.

Also observe that if one assumes  $g \in \mathfrak{D}(A^k)$ , one can apply Lemma 10.4 recursively to obtain:

**Corollary 10.5.** *Let  $T(t)$  be a  $C_0$ -semigroup with generator  $A$ . If  $g \in \mathfrak{D}(A^k)$  then  $T(t)g \in \mathfrak{D}(A^k)$  and  $T(t)g \in C^k([0, \infty), X)$  with*

$$\left(\frac{d}{dt}\right)^k T(t)g = A^j T(t)A^{k-j}g, \quad t \geq 0, \quad (10.11)$$

for any  $j = 0, \dots, k$ . In case of a  $C_0$ -group we have  $T(t)g \in C^k(\mathbb{R}, X)$  and the above formula holds for all  $t \in \mathbb{R}$ .

If we have  $T(t)g \in \mathfrak{D}(A)$  for all  $t > 0$ , then  $T(t)g \in \mathfrak{D}(A^k)$  for all  $k \in \mathbb{N}$ ,  $t > 0$  and  $T(t)g \in C^\infty((0, \infty), X)$ .

**Proof.** The case  $k = 1$  is established in Lemma 10.4. Suppose the claim holds for  $k \geq 1$  and  $g \in \mathfrak{D}(A^{k+1})$ . Then, applying  $A$  to  $T(t)A^k g = A^k T(t)g \in \mathfrak{D}(A)$  shows  $T(t)g \in \mathfrak{D}(A^{k+1})$  and  $AT(t)A^k g = A^{k+1}T(t)g$ . Finally,  $(\frac{d}{dt})^{k+1}T(t)g = \frac{d}{dt}T(t)A^k g = AT(t)A^k g$  finishes the induction step.

To see the second claim we use again induction to show  $T(t)g \in \mathfrak{D}(A^k)$ . The case  $k = 1$  holds by assumption. Now suppose the claim holds for some  $k \in \mathbb{N}$ , then  $AT(t)g = T(t/2)\tilde{g}$ , where  $\tilde{g} := AT(t/2)g \in \mathfrak{D}(A^{k-1})$  by

the induction hypothesis. Hence  $AT(t)g \in \mathfrak{D}(A^k)$  by assumption implying  $T(t)g \in \mathfrak{D}(A^{k+1})$ .  $\square$

Extending our remark from above we can use the differential equation show that for  $g \in \mathfrak{D}(A^k)$  we have  $T(t)g \in C^{k-j}([0, \infty), [\mathfrak{D}(A^j)])$ ,  $0 \leq j \leq k$ , where we equip  $\mathfrak{D}(A^j)$  with the norm  $\|f\|_{A^j} := \sum_{i=0}^j \|A^i f\|$ . Then  $T$  restricted to  $[\mathfrak{D}(A^j)]$  is a  $C_0$ -semigroup whose generator is  $A$  restricted to  $\mathfrak{D}(A^{j+1})$ .

A  $C_0$ -semigroup for which we have  $T(t)g \in \mathfrak{D}(A)$  for all  $g \in X$  and all  $t > 0$  is called **differentiable**.

Before turning to some examples, we establish a useful criterion for a semigroup to be strongly continuous.

**Lemma 10.6.** *A (semi)group of bounded operators is strongly continuous if and only if  $\limsup_{\varepsilon \downarrow 0} \|T(\varepsilon)g\| < \infty$  for every  $g \in X$  and  $\lim_{\varepsilon \downarrow 0} T(\varepsilon)f = f$  for  $f$  in a dense subset.*

**Proof.** We first show that  $\limsup_{\varepsilon \downarrow 0} \|T(\varepsilon)g\| < \infty$  for every  $g \in X$  implies that  $T(t)$  is bounded in a small interval  $[0, \delta]$ . Otherwise there would exist a sequence  $\varepsilon_n \downarrow 0$  with  $\|T(\varepsilon_n)\| \rightarrow \infty$ . Hence  $\|T(\varepsilon_n)g\| \rightarrow \infty$  for some  $g$  by the uniform boundedness principle, a contradiction. Thus there exists some  $M$  such that  $\sup_{t \in [0, \delta]} \|T(t)\| \leq M$ . Setting  $\omega = \frac{\log(M)}{\delta}$  we even obtain (10.7). Moreover, boundedness of  $T(t)$  shows that  $\lim_{\varepsilon \downarrow 0} T(\varepsilon)f = f$  for all  $f \in X$  by a simple approximation argument (Lemma 4.30 (iv)).

In case of a group this also shows  $\|T(-t)\| \leq \|T(\delta - t)\| \|T(-\delta)\| \leq M \|T(-\delta)\|$  for  $0 \leq t \leq \delta$ . Choosing  $\tilde{M} = \max(M, M \|T(-\delta)\|)$  we conclude  $\|T(t)\| \leq \tilde{M} \exp(\tilde{\omega}|t|)$ .

Finally, right continuity is implied by the semigroup property:  $\lim_{\varepsilon \downarrow 0} T(t + \varepsilon)g = \lim_{\varepsilon \downarrow 0} T(\varepsilon)T(t)g = T(t)g$ . Left continuity follows from  $\|T(t - \varepsilon)g - T(t)g\| = \|T(t - \varepsilon)(T(\varepsilon)g - g)\| \leq \|T(t - \varepsilon)\| \|T(\varepsilon)g - g\|$ .  $\square$

**Example 10.2.** Let  $X := C_0(\mathbb{R})$  be the continuous functions vanishing as  $|x| \rightarrow \infty$ . Then it is straightforward to check that

$$(T(t)f)(x) := f(x + t)$$

defines a group of continuous operators on  $X$ . Since shifting a function does not alter its supremum we have  $\|T(t)f\|_\infty = \|f\|_\infty$  and hence  $\|T(t)\| = 1$ . Moreover, strong continuity is immediate for uniformly continuous functions. Since every function with compact support is uniformly continuous and since such functions are dense, we get that  $T$  is strongly continuous. Moreover, for  $f \in \mathfrak{D}(A)$  we have

$$\lim_{\varepsilon \rightarrow 0} \frac{f(t + \varepsilon) - f(t)}{\varepsilon} = (Af)(t)$$

uniformly. In particular,  $f \in C^1(\mathbb{R})$  with  $f, f' \in C_0(\mathbb{R})$ . Conversely, for  $f \in C^1(\mathbb{R})$  with  $f, f' \in C_0(\mathbb{R})$  we have

$$\frac{f(t+\varepsilon) - f(t) - \varepsilon f'(t)}{\varepsilon} = \frac{1}{\varepsilon} \int_0^\varepsilon (f'(t+s) - f'(t)) ds \leq \sup_{0 \leq s \leq \varepsilon} \|T(s)f' - f'\|_\infty$$

which converges to zero as  $\varepsilon \downarrow 0$  by strong continuity of  $T$ . Whence

$$A = \frac{d}{dx}, \quad \mathfrak{D}(A) = \{f \in C^1(\mathbb{R}) \cap C_0(\mathbb{R}) \mid f' \in C_0(\mathbb{R})\}.$$

It is not hard to see that  $T$  is not uniformly continuous or, equivalently, that  $A$  is not bounded (cf. Problem 10.6).

Note that this group is not strongly continuous when considered on  $X := C_b(\mathbb{R})$ . Indeed for  $f(x) = \cos(x^2)$  we can choose  $x_n = \sqrt{2\pi n}$  and  $t_n = \sqrt{2\pi}(\sqrt{n + \frac{1}{4}} - \sqrt{n}) = \frac{1}{4}\sqrt{\frac{\pi}{2n}} + O(n^{-3/2})$  such that  $\|T(t_n)f - f\|_\infty \geq |f(x_n + t_n) - f(x_n)| = 1$ .  $\diamond$

Next consider

$$u(t) = T(t)g, \quad v(t) := \int_0^t u(s)ds, \quad g \in X. \quad (10.12)$$

Then  $v \in C^1([0, \infty), X)$  with  $\dot{v}(t) = u(t)$  and (Problem 9.3)

$$T(\varepsilon)v(t) = \int_0^t u(\varepsilon + s)ds = \int_\varepsilon^{t+\varepsilon} u(s)ds = v(t + \varepsilon) - v(t) \quad (10.13)$$

implying

$$\lim_{\varepsilon \downarrow 0} \frac{1}{\varepsilon} (T(\varepsilon)v(t) - v(t)) = \lim_{\varepsilon \downarrow 0} \left( -\frac{1}{\varepsilon}v(\varepsilon) + \frac{1}{\varepsilon}(v(t + \varepsilon) - v(t)) \right) = -g + u(t). \quad (10.14)$$

Consequently  $v(t) \in \mathfrak{D}(A)$  and  $Av(t) = -g + u(t)$  implying that  $u(t)$  solves the following integral version of our abstract Cauchy problem

$$u(t) = g + A \int_0^t u(s)ds. \quad (10.15)$$

Note that while in the case of a bounded generator both versions are equivalent, this will not be the case in general. So while  $u(t) = T(t)g$  always solves the integral version, it will only solve the differential version if  $u(t) \in \mathfrak{D}(A)$  for  $t > 0$  (which is clearly also necessary for the differential version to make sense). In the latter case  $u(t)$  is sometimes called a **strong solution** (also **classical solution**), while otherwise it is called a **mild solution**.

Two further consequences of these considerations are also worth while noticing:

**Corollary 10.7.** *Let  $T(t)$  be a  $C_0$ -semigroup with generator  $A$ . Then  $A$  is a densely defined and closed operator.*

**Proof.** Since  $v(t) \in \mathfrak{D}(A)$  and  $\lim_{t \downarrow 0} \frac{1}{t}v(t) = g$  for arbitrary  $g$ , we see that  $\mathfrak{D}(A)$  is dense. Moreover, if  $f_n \in \mathfrak{D}(A)$  and  $f_n \rightarrow f$ ,  $Af_n \rightarrow g$  then

$$T(t)f_n - f_n = \int_0^t T(s)Af_n ds.$$

Taking  $n \rightarrow \infty$  and dividing by  $t$  we obtain

$$\frac{1}{t}(T(t)f - f) = \frac{1}{t} \int_0^t T(s)g ds.$$

Taking  $t \downarrow 0$  finally shows  $f \in \mathfrak{D}(A)$  and  $Af = g$ .  $\square$

Note that by the closed graph theorem we have  $\mathfrak{D}(A) = X$  if and only if  $A$  is bounded. Moreover, since a  $C_0$ -semigroup provides the unique solution of the abstract Cauchy problem for  $A$ , we obtain

**Corollary 10.8.** *A  $C_0$ -semigroup is uniquely determined by its generator.*

**Proof.** Suppose  $T$  and  $S$  have the same generator  $A$ . Then by uniqueness for (10.10) we have  $T(t)f = S(t)f$  for all  $f \in \mathfrak{D}(A)$ . Since  $\mathfrak{D}(A)$  is dense this implies  $T(t) = S(t)$  as both operators are continuous.  $\square$

Finally, as in the uniformly continuous case, the inhomogeneous problem can be solved by Duhamel's formula. However, now it is not so clear when this will actually be a solution.

**Lemma 10.9.** *Let  $A$  be the generator of a  $C_0$ -semigroup and  $f \in C([0, \infty), X)$ . If the inhomogeneous problem*

$$\dot{u} = Au + f, \quad u(0) = g, \quad (10.16)$$

*has a solution it is necessarily given by **Duhamel's formula***

$$u(t) = T(t)g + \int_0^t T(t-s)f(s)ds. \quad (10.17)$$

*Conversely, for  $g = 0$ , the function  $u$  given by (10.17) satisfies  $u \in C^1([0, \infty), X)$  if and only if  $u \in C([0, \infty), [\mathfrak{D}(A)])$ . Moreover, in this case it will be a solution.*

*Specifically, (10.17) gives a solution if either one of the following conditions is satisfied:*

- $g \in \mathfrak{D}(A)$  and  $f \in C([0, \infty), [\mathfrak{D}(A)])$ .
- $g \in \mathfrak{D}(A)$  and  $f \in C^1([0, \infty), X)$ .

*If  $A$  is the generator of a strongly continuous group, we can replace  $[0, \infty)$  by  $\mathbb{R}$ .*

**Proof.** Let  $u(t)$  be a solution of (10.16) and set  $v(s) := T(t-s)u(s)$ ,  $0 \leq s \leq t$ , then one shows as in the proof of Lemma 10.4 that

$$\begin{aligned}\dot{v}(s) &= -AT(t-s)u(s) + T(t-s)\dot{u}(s) \\ &= -AT(t-s)u(s) + T(t-s)(Au(s) + f(s)) \\ &= T(t-s)f(s), \quad 0 < s < t.\end{aligned}$$

Hence the fundamental theorem of calculus (taking limits towards the boundary points) gives (10.17).

For the converse observe that  $T(t)g$  is a solution of the homogenous equation if  $g \in \mathfrak{D}(A)$ . Hence it remains to investigate the integral, which we will denote by  $u(t)$ . We first note that

$$\frac{1}{\varepsilon}(u(t+\varepsilon) - u(t)) = \frac{1}{\varepsilon} \int_0^\varepsilon T(\varepsilon-s)f(t+s)ds + \frac{1}{\varepsilon}(T(\varepsilon) - \mathbb{I})u(t),$$

where the integral term on the right converges to  $f(t)$  thanks to our assumption  $f \in C([0, \infty), X)$ . Hence, if one of the remaining two expressions has a limit, so has the other. In particular, if  $u(t)$  is differentiable, we see that the limit on the right exists implying  $u(t) \in \mathfrak{D}(A)$  and  $\dot{u}(t) = f(t) + Au(t)$ . Similarly if  $u(t) \in \mathfrak{D}(A)$ , then the limit on the right exists and we see that  $u(t)$  is differentiable.

From this the first case is immediate since  $u \in C([0, \infty), [\mathfrak{D}(A)])$  provided  $f \in C([0, \infty), [\mathfrak{D}(A)])$  by Problem 9.4.

In case of the second condition we note that

$$u(t) = \int_0^t T(s)f(t-s)ds$$

by a change of variables (Problem 9.3) and hence

$$\begin{aligned}\frac{1}{\varepsilon}(u(t+\varepsilon) - u(t)) &= \frac{1}{\varepsilon} \int_0^t T(s)(f(t+\varepsilon-s) - f(t-s))ds \\ &\quad + \frac{1}{\varepsilon} \int_0^\varepsilon T(t+s)f(\varepsilon-s)ds \\ &\xrightarrow{\varepsilon \rightarrow 0} \int_0^t T(s)\dot{f}(t-s)ds + T(t)f(0)\end{aligned}$$

since  $f \in C^1$ . □

The function  $u(t)$  defined by (10.17) is called the **mild solution** of the inhomogeneous problem. In general a mild solution is not a solution:

**Example 10.3.** Let  $T(t)$  be a strongly continuous group with an unbounded generator  $A$  (e.g. the one from Example 10.2). Choose  $f_0 \in X \setminus \mathfrak{D}(A)$  and

set  $g := 0$ ,  $f(t) := T(t)f_0$ . Then  $f \in C(\mathbb{R}, X)$  and the mild solution is given by

$$u(t) = T(t) \int_0^t T(-s)f(s)ds = T(t) \int_0^t f_0 ds = t T(t)f_0.$$

Since  $T(t)$  leaves  $\mathfrak{D}(A)$  invariant, we have  $u(t) \notin \mathfrak{D}(A)$  for all  $t \in \mathbb{R}$  and hence  $u(t)$  is not a solution.  $\diamond$

**Problem\* 10.6.** Show that a uniformly continuous semigroup has a bounded generator. (Hint: Write  $T(t) = V(t_0)^{-1}V(t_0)T(t) = \dots$  with  $V(t) := \int_0^t T(s)ds$  and conclude that it is  $C^1$ .)

**Problem 10.7.** Let  $f : [0, \infty) \rightarrow \mathbb{R}$  be bounded from above on every compact interval and subadditive, that is,  $f(t_1 + t_2) \leq f(t_1) + f(t_2)$ . Then

$$\lim_{t \rightarrow \infty} \frac{f(t)}{t} = \inf_{t \geq 0} \frac{f(t)}{t}.$$

**Problem 10.8.** Show that the growth bound of a semigroup is given by

$$\omega_0(T) = \inf_{t \geq 0} \frac{\log(\|T(t)\|)}{t} = \lim_{t \rightarrow \infty} \frac{\log(\|T(t)\|)}{t}.$$

Moreover, show that the spectral radius of  $T(t)$  is given by

$$r(T(t)) = e^{\omega_0(T)t}.$$

(Hint: The spectral radius of an operator  $T \in \mathcal{L}(X)$  is defined as  $r(T) := \sup_{z \in \sigma(T)} |z| = \lim_{n \rightarrow \infty} \|T^n\|^{1/n} \leq \|T(t)\|.$ )

**Problem 10.9.** Let  $T(t)$  be a  $C_0$ -semigroup. Show that if  $T(t_0)$  has a bounded inverse for one  $t_0 > 0$  then this holds for all  $t > 0$  and it extends to a strongly continuous group via  $T(t) := T(-t)^{-1}$  for  $t < 0$ .

**Problem 10.10.** Consider the translation group  $T(t) := T_t$  on  $L^p(\mathbb{R})$ ,  $1 \leq p < \infty$ . Show that this is a strongly continuous group and compute its generator. Show that it is not strongly continuous for  $p = \infty$ . (Hint: Problem ??.)

**Problem 10.11.** Consider the translation semigroup  $T(t) := T_t$  on  $L^p(0, 1)$ ,  $1 \leq p < \infty$ . Show that this is a strongly continuous group and compute its generator. Show that it is nilpotent:  $T_t g = 0$  for  $t \geq 1$ . (Hint: Problem ??.)

**Problem 10.12.** Let  $U \subseteq \mathbb{R}^n$  and let  $m : U \rightarrow \mathbb{C}$  be a measurable function with  $\sup_{x \in U} \operatorname{Re}(m(x)) < \infty$ . Consider the multiplication semigroup  $T(t)g(x) := e^{tm(x)}g(x)$  on  $L^p(U)$ ,  $1 \leq p < \infty$ . Show that this is a strongly continuous group and compute its generator.

**Problem 10.13.** Define a semigroup on  $L^1(-1, 1)$  via

$$(T(t)f)(s) = \begin{cases} 2f(s-t), & 0 < s \leq t, \\ f(s-t), & \text{else,} \end{cases}$$



where we set  $f(s) = 0$  for  $s < 0$ . Show that the estimate from Lemma 10.3 does not hold with  $M < 2$ .

**Problem 10.14.** Let  $A$  be the generator of a  $C_0$ -semigroup  $T(t)$ . Show

$$T(t)f = f + tAf + \int_0^t (t-s)T(s)A^2f \, ds, \quad f \in \mathfrak{D}(A^2).$$

**Problem 10.15.** Let  $A$  be the generator of a  $C_0$ -semigroup  $T(t)$ . Show that  $\bigcap_{k \in \mathbb{N}} \mathfrak{D}(A^k)$  is dense. (Hint: Set  $g_m := m \int_0^1 \phi(ms)T(s)g \, ds$ , where  $\phi \in C_c^\infty(0,1)$  with  $\int_0^1 \phi(s)ds = 1$ .)

**Problem 10.16.** Let  $T(t)$  be a differentiable  $C_0$ -semigroup with generator  $A$ . Show  $T \in C^\infty((0, \infty), \mathcal{L}(X))$  with  $\frac{d^k}{dt} T(t) = A^k T(t)$ ,  $t > 0$ . (Hint: Show that  $A^k T(t)$  is bounded and use Problem 10.14.)

**Problem 10.17** (Landau inequality). Let  $A$  be the generator of a  $C_0$ -semigroup  $T(t)$  satisfying  $\|T(t)\| \leq M$ . Derive the **abstract Landau inequality**

$$\|Af\| \leq 2M\|A^2f\|^{1/2}\|f\|^{1/2}.$$

(Hint: Problem 10.14.)

**Problem 10.18.** Let  $A$  be the generator of a  $C_0$ -semigroup. Consider the integral version of our inhomogeneous problem (10.16):

$$u(t) = g + A \int_0^t u(s)ds + \int_0^t f(s)ds$$

for given  $g \in X$ ,  $f \in C([0,1), X)$ . Show that this problem has a unique solution  $u \in C([0,1), X)$  such that  $\int_0^t u(s)ds \in \mathfrak{D}(A)$  for  $t \geq 0$  which is given by Duhamel's formula (10.17). (Hint: Problem 9.5.)

**Problem 10.19.** A bounded operator  $P \in \mathcal{L}(X)$  is said to commute with a closed operator  $A$  if

$$PA \subseteq AP.$$

That is, if  $\mathfrak{D}(PA) = \mathfrak{D}(A) \subseteq \mathfrak{D}(AP) = \{x \in X | Px \in \mathfrak{D}(A)\}$  and both operators agree on the smaller set  $\mathfrak{D}(A)$ . Note that this in particular requires that  $P$  leaves the domain invariant,  $P\mathfrak{D}(A) \subseteq \mathfrak{D}(A)$ .

Show that if  $P \in \mathcal{L}(X)$  commutes with the generator of a (semi)group  $A$ , then it also **commutes** with the (semi)group,  $PT(t) = T(t)P$ . (Hint: Uniqueness of solutions.)

### 10.3. Generator theorems

Of course in practice the abstract Cauchy problem, that is the operator  $A$ , is given and the question is if  $A$  generates a corresponding  $C_0$ -semigroup.

Corollary 10.7 already gives us some necessary conditions but this alone is not enough.

It turns out that it is crucial to understand the resolvent of  $A$  (see Section 8.2). Using an operator-valued version of the elementary integral  $\int_0^\infty e^{t(a-z)} dt = -(a-z)^{-1}$  (for  $\operatorname{Re}(a-z) < 0$ ) we can make the connection between the resolvent and the semigroup.

**Lemma 10.10.** *Let  $T$  be a  $C_0$ -semigroup with generator  $A$  satisfying (10.7). Then  $\{z | \operatorname{Re}(z) > \omega\} \subseteq \rho(A)$  and*

$$R_A(z) = - \int_0^\infty e^{-zt} T(t) dt, \quad \operatorname{Re}(z) > \omega, \quad (10.18)$$

where the right-hand side is defined as

$$\left( \int_0^\infty e^{-zt} T(t) dt \right) g := \lim_{s \rightarrow \infty} \int_{1/s}^s e^{-zt} T(t) g dt. \quad (10.19)$$

Moreover,

$$\|R_A(z)\| \leq \frac{M}{\operatorname{Re}(z) - \omega}, \quad \operatorname{Re}(z) > \omega. \quad (10.20)$$

**Proof.** Let us abbreviate  $R_s(z)f := - \int_0^s e^{-zt} T(t) f dt$ . Then, by virtue of (10.7),  $\|e^{-zt} T(t) f\| \leq M e^{(\omega - \operatorname{Re}(z))t} \|f\|$  shows that  $R_s(z)$  is a bounded operator satisfying  $\|R_s(z)\| \leq M(\operatorname{Re}(z) - \omega)^{-1}$ . Moreover, this estimate on the integrand also shows that the limit  $R(z) := \lim_{s \rightarrow \infty} R_s(z)$  exists (and still satisfies  $\|R(z)\| \leq M(\operatorname{Re}(z) - \omega)^{-1}$ ). Next note that  $S(t) := e^{-zt} T(t)$  is a semigroup with generator  $A - z$  (Problem 10.20) and hence for  $f \in \mathfrak{D}(A)$  we have

$$R_s(z)(A - z)f = - \int_0^s S(t)(A - z)f dt = - \int_0^s \dot{S}(t)f dt = f - S(s)f.$$

In particular, taking the limit  $s \rightarrow \infty$ , we obtain  $R(z)(A - z)f = f$  for  $f \in \mathfrak{D}(A)$ . Similarly, still for  $f \in \mathfrak{D}(A)$ , by Problem 9.4

$$(A - z)R_s(z)f = - \int_0^s (A - z)S(t)f dt = - \int_0^s \dot{S}(t)f dt = f - S(s)f$$

and taking limits, using closedness of  $A$ , implies  $(A - z)R(z)f = f$  for  $f \in \mathfrak{D}(A)$ . Finally, if  $g \in X$  choose  $f_n \in \mathfrak{D}(A)$  with  $f_n \rightarrow g$ . Then  $R(z)f_n \rightarrow R(z)g$  and  $(A - z)R(z)f_n = f_n \rightarrow g$  proving  $R(z)g \in \mathfrak{D}(A)$  and  $(A - z)R(z)g = g$  for  $g \in X$ .  $\square$

The number

$$s(A) := \sup\{\operatorname{Re}(z) | z \in \sigma(A)\} \quad (10.21)$$

is known as the **spectral bound** of  $A$ . The above lemma shows  $s(A) \leq \omega_0(T)$ . Moreover, for matrices it is not hard to see that we have equality.

However, this is not true in general. In particular, knowledge of  $\sigma(A)$  alone is in general not sufficient to estimate the growth of the associated semigroup.

**Corollary 10.11.** *Let  $T$  be a  $C_0$ -semigroup with generator  $A$  satisfying (10.7). Then*

$$R_A(z)^{n+1} = \frac{(-1)^{n+1}}{n!} \int_0^\infty t^n e^{-zt} T(t) dt, \quad \operatorname{Re}(z) > \omega, \quad (10.22)$$

and

$$\|R_A(z)^n\| \leq \frac{M}{(\operatorname{Re}(z) - \omega)^n}, \quad \operatorname{Re}(z) > \omega, \quad n \in \mathbb{N}. \quad (10.23)$$

**Proof.** Abbreviate  $R_n(z) := \int_0^\infty t^n e^{-zt} T(t) dt$  and note that

$$\frac{R_n(z + \varepsilon) - R_n(z)}{\varepsilon} = -R_{n+1}(z) + \varepsilon \int_0^\infty t^{n+2} \phi(\varepsilon t) e^{-zt} T(t) dt$$

where  $|\phi(\varepsilon)| \leq \sum_{j=0}^\infty \frac{|\varepsilon|^j}{(j+2)!} \leq \frac{1}{2} e^{|\varepsilon|}$  from which we see  $\frac{d}{dz} R_n(z) = -R_{n+1}(z)$  and hence  $\frac{d^n}{dz^n} R_A(z) = -\frac{d^n}{dz^n} R_0(z) = (-1)^{n+1} R_n(z)$ . Now the first claim follows using  $R_A(z)^{n+1} = \frac{1}{n!} \frac{d^n}{dz^n} R_A(z)$  (Problem 8.22). Estimating the integral using (10.7) establishes the second claim.  $\square$

Given these preparations we can now try to answer the question when  $A$  generates a semigroup. In fact, we will be constructive and obtain the corresponding semigroup by approximation. To this end we introduce the **Yosida approximation**

$$A_n := -nAR_A(\omega + n) = -n - n(\omega + n)R_A(\omega + n) \in \mathcal{L}(X). \quad (10.24)$$

Of course this is motivated by the fact that this is a valid approximation for numbers  $\lim_{n \rightarrow \infty} \frac{-n}{a - \omega - n} = 1$ . That we also get a valid approximation for operators is the content of the next lemma.

**Lemma 10.12.** *Suppose  $A$  is a densely defined closed operator with  $(\omega, \infty) \subset \rho(A)$  satisfying*

$$\|R_A(\omega + n)\| \leq \frac{M}{n}. \quad (10.25)$$

Then

$$\lim_{n \rightarrow \infty} -nR_A(\omega + n)f = f, \quad f \in X, \quad \lim_{n \rightarrow \infty} A_n f = Af, \quad f \in \mathfrak{D}(A). \quad (10.26)$$

**Proof.** If  $f \in \mathfrak{D}(A)$  we have  $-nR_A(\omega + n)f = f - R_A(\omega + n)(A - \omega)f$  which shows  $-nR_A(\omega + n)f \rightarrow f$  if  $f \in \mathfrak{D}(A)$ . Since  $\mathfrak{D}(A)$  is dense and  $\|nR_A(\omega + n)\| \leq M$  this even holds for all  $f \in X$ . Moreover, for  $f \in \mathfrak{D}(A)$  we have  $A_n f = -nAR_A(\omega + n)f = -nR_A(\omega + n)(Af) \rightarrow Af$  by the first part.  $\square$

Moreover,  $A_n$  can also be used to approximate the corresponding semigroup under suitable assumptions.

**Theorem 10.13** (Feller–Miyadera–Phillips). *A linear operator  $A$  is the generator of a  $C_0$ -semigroup  $T$  satisfying (10.7) if and only if it is densely defined, closed,  $(\omega, \infty) \subseteq \rho(A)$ , and*

$$\|R_A(\lambda)^n\| \leq \frac{M}{(\lambda - \omega)^n}, \quad \lambda > \omega, n \in \mathbb{N}. \quad (10.27)$$

Moreover, if  $A_n$  is the Yosida approximation (10.24) and

$$T_n(t) := \exp(tA_n) = e^{-tn} \exp(-tn(\omega + n)R_A(\omega + n)) \quad (10.28)$$

are the corresponding groups, we have

$$T(t)g = \lim_{n \rightarrow \infty} T_n(t)g, \quad g \in X. \quad (10.29)$$

**Proof.** Necessity has already been established in Corollaries 10.7 and 10.11.

For the converse we note

$$\|T_n(t)\| \leq e^{-tn} \sum_{j=0}^{\infty} \frac{(tn(\omega + n))^j}{j!} \|R_A(\omega + n)^j\| \leq Me^{-tn} e^{t(\omega + n)} = Me^{\omega t}.$$

Moreover, since  $R_A(\omega + m)$  and  $R_A(\omega + n)$  commute by the first resolvent identity (Problem 8.22), we conclude that the same is true for  $A_m$ ,  $A_n$  as well as for  $T_m(t)$ ,  $T_n(t)$  (by the very definition as a power series). Consequently

$$\begin{aligned} \|T_n(t)f - T_m(t)f\| &= \left\| \int_0^1 \frac{d}{ds} T_n(st)T_m((1-s)t)f ds \right\| \\ &\leq t \int_0^1 \|T_n(st)T_m((1-s)t)(A_n - A_m)f\| ds \\ &\leq tM^2 e^{\omega t} \|(A_n - A_m)f\|. \end{aligned}$$

Thus, for  $f \in \mathfrak{D}(A)$  we have a Cauchy sequence and can define a linear operator by  $T(t)f := \lim_{n \rightarrow \infty} T_n(t)f$ . Since  $\|T(t)f\| = \lim_{n \rightarrow \infty} \|T_n(t)f\| \leq Me^{\omega t}\|f\|$ , we see that  $T(t)$  is bounded and has a unique extension to all of  $X$ . Moreover,  $T(0) = \mathbb{I}$  and

$$\begin{aligned} \|T_n(t)T_n(s)f - T(t)T(s)f\| &\leq \\ &Me^{\omega t}\|T_n(s)f - T(s)f\| + \|(T_n(t) - T(t))T(s)f\| \end{aligned}$$

implies  $T(t+s)f = \lim_{n \rightarrow \infty} T_n(t+s)f = \lim_{n \rightarrow \infty} T_n(t)T_n(s)f = T(t)T(s)f$ , that is, the semigroup property holds. Finally, by

$$\begin{aligned} \|T(\varepsilon)f - f\| &\leq \|T(\varepsilon)f - T_n(\varepsilon)f\| + \|T_n(\varepsilon)f - f\| \\ &\leq \varepsilon M^2 e^{\omega \varepsilon} \|(A - A_n)f\| + \|T_n(\varepsilon)f - f\| \end{aligned}$$

we see  $\lim_{\varepsilon \downarrow 0} T(\varepsilon)f = f$  for  $f \in \mathfrak{D}(A)$  and Lemma 10.6 shows that  $T$  is a  $C_0$ -semigroup. It remains to show that  $A$  is its generator. To this end let  $f \in \mathfrak{D}(A)$ , then

$$\begin{aligned} T(t)f - f &= \lim_{n \rightarrow \infty} T_n(t)f - f = \lim_{n \rightarrow \infty} \int_0^t T_n(s)A_n f \, ds \\ &= \lim_{n \rightarrow \infty} \left( \int_0^t T_n(s)A f \, ds + \int_0^t T_n(s)(A_n - A)f \, ds \right) \\ &= \int_0^t T(s)A f \, ds \end{aligned}$$

which shows  $\lim_{t \downarrow 0} \frac{1}{t}(T(t)f - f) = Af$  for  $f \in \mathfrak{D}(A)$ . Finally, note that the domain of the generator cannot be larger, since  $A - \omega - 1$  is bijective and adding a vector to its domain would destroy injectivity. But then  $\omega + 1$  would not be in the resolvent set contradicting Lemma 10.10.  $\square$

Note that in combination with the following lemma this also answers the question when  $A$  generates a  $C_0$ -group.

**Lemma 10.14.** *An operator  $A$  generates a  $C_0$ -group if and only if both  $A$  and  $-A$  generate  $C_0$ -semigroups.*

**Proof.** Clearly, if  $A$  generates a  $C_0$ -group  $T(t)$ , then  $S(t) := T(-t)$  is a  $C_0$ -group with generator  $-A$ . Conversely, let  $T(t)$ ,  $S(t)$  be the  $C_0$ -semigroups generated by  $A$ ,  $-A$ , respectively. Then a short calculation shows

$$\frac{d}{dt}T(t)S(t)g = -T(t)AS(t)g + T(t)AS(t)g = 0, \quad t \geq 0.$$

Consequently,  $T(t)S(t) = T(0)S(0) = \mathbb{I}$  and similarly  $S(t)T(t) = \mathbb{I}$ , that is,  $S(t) = T(t)^{-1}$ . Hence it is straightforward to check that  $T$  extends to a group via  $T(-t) := S(t)$ ,  $t \geq 0$ .  $\square$

The following examples show that the spectral conditions are indeed crucial. Moreover, they also show that an operator might give rise to a Cauchy problem which is uniquely solvable for a dense set of initial conditions, without generating a strongly continuous semigroup.

**Example 10.4.** Let

$$A = \begin{pmatrix} 0 & A_0 \\ 0 & 0 \end{pmatrix}, \quad \mathfrak{D}(A) = X \times \mathfrak{D}(A_0).$$

Then  $u(t) = \begin{pmatrix} 1 & tA_0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} f_0 \\ f_1 \end{pmatrix} = \begin{pmatrix} f_0 + tA_0f_1 \\ f_1 \end{pmatrix}$  is the unique solution of the corresponding abstract Cauchy problem for given  $f \in \mathfrak{D}(A)$ . Nevertheless, if  $A_0$  is unbounded, the corresponding semigroup is not strongly continuous.

Note that in this case we have  $\sigma(A) = \{0\}$  if  $A_0$  is bounded and  $\sigma(A) = \mathbb{C}$  else. In fact, since  $A$  is not injective we must have  $\{0\} \subseteq \sigma(A)$ . For  $z \neq 0$  the inverse of  $A - z$  is given by

$$(A - z)^{-1} = -\frac{1}{z} \begin{pmatrix} 1 & \frac{1}{z}A_0 \\ 0 & 1 \end{pmatrix}, \quad \mathfrak{D}((A - z)^{-1}) = \text{Ran}(A - z) = X \times \mathfrak{D}(A_0),$$

which is bounded if and only if  $A$  is bounded.  $\diamond$

**Example 10.5.** Let  $X_0 = C_0(\mathbb{R})$  and  $m(x) = ix$ . Then we can regard  $m$  as a multiplication operator on  $X_0$  when defined maximally, that is,  $f \mapsto mf$  with  $\mathfrak{D}(m) = \{f \in X_0 \mid mf \in X_0\}$ . Note that since  $C_c(\mathbb{R}) \subseteq \mathfrak{D}(m)$  we see that  $m$  is densely defined. Moreover, it is easy to check that  $m$  is closed.

Now consider  $X = X_0 \oplus X_0$  with  $\|f\| = \max(\|f_0\|, \|f_1\|)$  and note that

$$A = \begin{pmatrix} m & m \\ 0 & m \end{pmatrix}, \quad \mathfrak{D}(A) = \mathfrak{D}(m) \oplus \mathfrak{D}(m),$$

is closed. Moreover, for  $z \notin i\mathbb{R}$  the resolvent is given by the multiplication operator

$$R_A(z) = \frac{1}{m - z} \begin{pmatrix} 1 & -\frac{m}{m-z} \\ 0 & 1 \end{pmatrix}.$$

For  $\lambda > 0$  we compute

$$\|R_A(\lambda)f\| \leq \left( \sup_{x \in \mathbb{R}} \frac{1}{|ix - \lambda|} + \sup_{x \in \mathbb{R}} \frac{|x|}{|ix - \lambda|^2} \right) \|f\| = \frac{3}{2\lambda} \|f\|$$

and hence  $A$  satisfies (10.27) with  $M = \frac{3}{2}$ ,  $\omega = 0$  and  $n = 1$ . However, by

$$\|R_A(\lambda + in)\| \geq \|R_A(\lambda + in)(0, f_n)\| \geq \left| \frac{inf_n(n)}{(\lambda - in + in)^2} \right| = \frac{n}{\lambda^2},$$

where  $f_n$  is chosen such that  $f_n(n) = 1$  and  $\|f_n\|_\infty = 1$ , it does not satisfy (10.23). Hence  $A$  does not generate a  $C_0$ -semigroup. Indeed, the solution of the corresponding Cauchy problem is

$$T(t) = e^{tm} \begin{pmatrix} 1 & tm \\ 0 & 1 \end{pmatrix}, \quad \mathfrak{D}(T) = X_0 \oplus \mathfrak{D}(m),$$

which is unbounded.  $\diamond$

When it comes to applying this theorem, the main difficulty will be establishing the resolvent estimate (10.27). Moreover, while it might be already difficult to estimate the resolvent, it will in general be even more challenging to get estimates on its powers. In this connection note that the trivial estimate  $\|R_A(z)^n\| \leq \|R_A(z)\|^n$  will do the job if and only if  $M = 1$ . Hence we finally look at the special case of **contraction semigroups** satisfying

$$\|T(t)\| \leq 1. \tag{10.30}$$

By a simple transform the case  $M = 1$  in Lemma 10.3 can always be reduced to this case (Problem 10.20). Moreover, as already anticipated, in the case  $M = 1$  the estimate (10.20) immediately implies the general estimate (10.23) and it suffices to establish (10.27) for  $n = 1$ :

**Corollary 10.15** (Hille–Yosida). *A linear operator  $A$  is the generator of a contraction semigroup if and only if it is densely defined, closed,  $(0, \infty) \subseteq \rho(A)$ , and*

$$\|R_A(\lambda)\| \leq \frac{1}{\lambda}, \quad \lambda > 0. \quad (10.31)$$

**Example 10.6.** If  $A$  is the generator of a contraction, then clearly all eigenvalues  $z$  must satisfy  $\operatorname{Re}(z) \leq 0$ . Moreover, for

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

we have

$$R_A(z) = -\frac{1}{z} \begin{pmatrix} 1 & 1/z \\ 0 & 1 \end{pmatrix}, \quad T(t) = \begin{pmatrix} 1 & t \\ 0 & 1 \end{pmatrix},$$

which shows that the bound on the resolvent is crucial.  $\diamond$

**Example 10.7.** If  $X$  is a Hilbert space and  $A$  is a self-adjoint operator, then we have the required estimate if and only if  $A$  is bounded from above,

$$E := \sup \sigma(A) = \sup_{f \in \mathfrak{D}(A), \|f\|=1} \langle f, Af \rangle < \infty.$$

Indeed in this case we have (cf. [47, Theorem 2.19])

$$\|R_A(\lambda)\| \leq \frac{1}{\lambda - E}, \quad \lambda > E,$$

such that  $A - E$  generates a contraction semigroup. In particular, note that in this case the growth bound equals the spectral bound.  $\diamond$

However, for a given operator even the simple estimate (10.31) might be difficult to establish directly. Hence we outline another criterion.

**Example 10.8.** Let  $X$  be a Hilbert space and observe that for a contraction semigroup the expression  $\|T(t)f\|$  must be nonincreasing. Consequently, for  $f \in \mathfrak{D}(A)$  we must have

$$\left. \frac{d}{dt} \|T(t)f\|^2 \right|_{t=0} = 2\operatorname{Re}(\langle f, Af \rangle) \leq 0.$$

Operators satisfying  $\operatorname{Re}(\langle f, Af \rangle) \leq 0$  are called dissipative and this clearly suggests to replace the resolvent estimate by dissipativity.  $\diamond$

To formulate this condition for Banach spaces, we first introduce the **duality set**

$$\mathcal{J}(x) := \{x' \in X^* \mid x'(x) = \|x\|^2 = \|x'\|^2\} \quad (10.32)$$

of a given vector  $x \in X$ . In other words, the elements from  $\mathcal{J}(x)$  are those linear functionals which attain their norm at  $x$  and are normalized to have the same norm as  $x$ . As a consequence of the Hahn–Banach theorem (Corollary 4.11) note that  $\mathcal{J}(x)$  is nonempty. Moreover, it is also easy to see that  $\mathcal{J}(x)$  is convex and weak-\* closed.

**Example 10.9.** Let  $X$  be a Hilbert space and identify  $X$  with  $X^*$  via  $x \mapsto \langle x, \cdot \rangle$  as usual. Then  $\mathcal{J}(x) = \{x\}$ . Indeed since we have equality  $\langle x', x \rangle = \|x'\| \|x\|$  in the Cauchy–Schwarz inequality, we must have  $x' = \alpha x$  for some  $\alpha \in \mathbb{C}$  with  $|\alpha| = 1$  and  $\alpha^* \|x\|^2 = \langle x', x \rangle = \|x\|^2$  shows  $\alpha = 1$ .  $\diamond$

**Example 10.10.** Recall (cf. Problem 1.13) that a Banach space  $X$  is called strictly convex, if equality in the triangle inequality,  $\|x + y\| = \|x\| + \|y\|$ , can only occur if the vectors are parallel,  $y = \alpha x$  for some  $\alpha \geq 0$  or  $x = 0$ .

If  $X^*$  is strictly convex, then the duality set contains only one point. In fact, suppose  $x', y' \in \mathcal{J}(x)$ , then  $z' = \frac{1}{2}(x' + y') \in \mathcal{J}(x)$  and  $\frac{\|x\|}{2} \|x' + y'\| = z'(x) = \frac{\|x\|}{2} (\|x'\| + \|y'\|)$  implying  $x' = y'$  by strict convexity. Note that the converse is also true: If  $x', y' \in \mathcal{J}(x)$  for some  $x \in B_1^X(0)$ , then  $x'(x) + y'(x) = 2$  implies  $\|x' + y'\| = 2$  contradicting strict convexity.

This applies for example to  $X := \ell^p(\mathbb{N})$  if  $1 < p < \infty$  (cf. Problem 1.13) in which case  $X^* \cong \ell^q(\mathbb{N})$  with  $q = \frac{p}{p-1}$ . In fact, for  $a \in X$  we have  $\mathcal{J}(a) = \{a'\}$  with  $a'_j = \|a\|_p^{2-p} \text{sign}(a_j^*) |a_j|^{p-1}$ .  $\diamond$

**Example 10.11.** Let  $X$  be a measurable space with a  $\sigma$ -finite measure  $\mu$ . The previous example can be generalized to  $L^p(X, d\mu)$  if  $1 < p < \infty$  (which are strictly, in fact even uniformly, convex by Theorem 3.11 from [48]). In this case we have  $L^p(X, d\mu)^* \cong L^q(X, d\mu)$  and for  $f \in L^p(X, d\mu)$  we have  $\mathcal{J}(f) = \{g\}$  with  $g = \|f\|_p^{2-p} \text{sign}(f^*) |f|^{p-1}$ .  $\diamond$

**Example 10.12.** Let  $X := C[0, 1]$  and choose  $x \in X$ . If  $t_0$  is chosen such that  $|x(t_0)| = \|x\|_\infty$ , then the functional  $y \mapsto x'(y) := x(t_0)^* y(t_0)$  satisfies  $x' \in \mathcal{J}(x)$ . Clearly  $\mathcal{J}(x)$  will contain more than one element in general.

Note that for  $X = C_b(\mathbb{R})$  the situation is more complicated since the supremum might not be attained. However, we can choose a sequence  $t_n \in \mathbb{R}$  such that  $x(t_n) \rightarrow x_0$  with  $|x_0| = \|x\|_\infty$  and set  $x'(y) = x_0^* L(y(t_n))$ , where  $L$  is the Banach limit from Problem 4.19.  $\diamond$

Now a given operator  $\mathfrak{D}(A) \subseteq X \rightarrow X$  is called **dissipative** if

$$\text{Re}(x'(Ax)) \leq 0 \quad \text{for one } x' \in \mathcal{J}(x) \text{ and all } x \in \mathfrak{D}(A). \quad (10.33)$$

**Lemma 10.16.** Let  $x, y \in X$ . Then  $\|x\| \leq \|x - \alpha y\|$  for all  $\alpha > 0$  if and only if there is an  $x' \in \mathcal{J}(x)$  such that  $\text{Re}(x'(y)) \leq 0$ .



**Proof.** Without loss of generality we can assume  $x \neq 0$ . If  $\operatorname{Re}(x'(y)) \leq 0$  for some  $x' \in \mathcal{J}(x)$ , then for  $\alpha > 0$  we have

$$\|x'\|\|x\| = x'(x) \leq \operatorname{Re}(x'(x - \alpha y)) \leq \|x'\|\|x - \alpha y\|$$

implying  $\|x\| \leq \|x - \alpha y\|$ .

Conversely, if  $\|x\| \leq \|x - \alpha y\|$  for all  $\alpha > 0$ , let  $x'_\alpha \in \mathcal{J}(x - \alpha y)$  and set  $y'_\alpha = \|x'_\alpha\|^{-1}x'_\alpha$ . Then

$$\begin{aligned} \|x\| &\leq \|x - \alpha y\| = y'_\alpha(x - \alpha y) = \operatorname{Re}(y'_\alpha(x)) - \alpha \operatorname{Re}(y'_\alpha(y)) \\ &\leq \|x\| - \alpha \operatorname{Re}(y'_\alpha(y)). \end{aligned}$$

Now choose a sequence  $\alpha_j \rightarrow 0$  such that both  $y'_{\alpha_j}(x)$  and  $y'_{\alpha_j}(y)$  converge. This defines a linear functional on the two dimensional subspace spanned by  $x$  and  $y$  which can be extended to a functional  $y'_0 \in X^*$  using Hahn–Banach. Taking the limit in the above inequality yields  $y'_0(x) = \|x\|$ . Moreover, the above inequality also shows  $\operatorname{Re}(y'_\alpha(y)) \leq 0$  and hence  $\operatorname{Re}(y'_0(y)) \leq 0$ . Whence  $x'_0 = \|x\|y'_0 \in \mathcal{J}(x)$  and  $\operatorname{Re}(x'_0(y)) \leq 0$ .  $\square$

As a straightforward consequence we obtain:

**Corollary 10.17.** *A linear operator is dissipative if and only if*

$$\|(A - \lambda)x\| \geq \lambda\|x\|, \quad \lambda > 0, x \in \mathfrak{D}(A). \quad (10.34)$$

In particular, for a dissipative operator  $A - \lambda$  is injective for  $\lambda > 0$  and  $(A - \lambda)^{-1}$  is bounded with  $\|(A - \lambda)^{-1}\| \leq \lambda^{-1}$ . However, this does not imply that  $\lambda$  is in the resolvent set of  $A$  since  $\mathfrak{D}((A - \lambda)^{-1}) = \operatorname{Ran}(A - \lambda)$  might not be all of  $X$ .

Now we are ready to show

**Theorem 10.18** (Lumer–Phillips). *A linear operator  $A$  is the generator of a contraction semigroup if and only if it is densely defined, dissipative, and  $A - \lambda_0$  is surjective for one  $\lambda_0 > 0$ . Moreover, in this case (10.33) holds for all  $x' \in \mathcal{J}(x)$ .*

**Proof.** Let  $A$  generate a contraction semigroup  $T(t)$  and let  $x \in \mathfrak{D}(A)$ ,  $x' \in \mathcal{J}(x)$ . Then

$$\operatorname{Re}(x'(T(t)x - x)) \leq |x'(T(t)x)| - \|x\|^2 \leq \|x'\|\|x\| - \|x\|^2 = 0$$

and dividing by  $t$  and letting  $t \downarrow 0$  shows  $\operatorname{Re}(x'(Ax)) \leq 0$ . Hence  $A$  is dissipative and by Corollary 10.15  $(0, \infty) \subseteq \rho(A)$ , that is,  $A - \lambda$  is bijective for  $\lambda > 0$ .

Conversely, by Corollary 10.17  $A - \lambda$  has a bounded inverse satisfying  $\|(A - \lambda)^{-1}\| \leq \lambda^{-1}$  for all  $\lambda > 0$ . In particular, for  $\lambda_0$  the inverse is defined on all of  $X$  and hence closed. Thus  $A$  is also closed and  $\lambda_0 \in \rho(A)$ . Moreover,

from  $\|R_A(\lambda_0)\| \leq \lambda_0^{-1}$  (cf. Lemma 8.10) we even get  $(0, 2\lambda_0) \subseteq \rho(A)$  and iterating this argument shows  $(0, \infty) \subseteq \rho(A)$  as well as  $\|R_A(\lambda)\| \leq \lambda^{-1}$ ,  $\lambda > 0$ . Hence the requirements from Corollary 10.15 are satisfied.  $\square$

Note that generators of contraction semigroups are maximal dissipative in the sense that they do not have any dissipative extensions. In fact, if we extend  $A$  to a larger domain we must destroy injectivity of  $A - \lambda$  and thus the extension cannot be dissipative.

**Example 10.13.** Let  $X$  be a Hilbert space. An equation of the form

$$iu = Hu$$

with  $H$  a self-adjoint operator, is known as abstract **Schrödinger equation**. If  $H$  is bounded, it is easy to see that  $\exp(-itH)$  is a uniformly continuous unitary group (cf. Problem 10.3). The Lumer–Phillips theorem also allows us to handle the unbounded case.

In this context a densely defined operator  $H$  is called symmetric if

$$\langle f, Hg \rangle = \langle Hf, g \rangle, \quad f, g \in \mathfrak{D}(H).$$

In this case  $\langle f, Hf \rangle$  is real-valued or equivalently, both  $A = -iH$  and  $-A = iH$  are dissipative. Hence if we assume  $\text{Ran}(H + i) = \text{Ran}(H - i) = X$ , then both  $A$  and  $-A$  will generate contraction semigroups from which it is not hard to see that  $T(t)$  is a strongly continuous group which preserves the norm (cf. Problem 10.22). But an operator preserving the norm is unitary,  $T(t)^{-1} = T(t)^*$ . Since for a symmetric operator  $\text{Ran}(H + i) = \text{Ran}(H - i) = X$  is equivalent to self-adjointness ([47, Lemma 2.3]), we see that a self-adjoint operators give rise to a strongly continuous unitary group. This is known as **Stone’s theorem**.

In fact, the converse is also true. To see this observe that  $H$  is symmetric if and only if  $\langle f, Hf \rangle$  is real-valued which in turn is equivalent to both  $-iH$  and  $iH$  being dissipative.  $\diamond$

**Example 10.14.** Let  $X := C_0[0, 1] = \{f \in C[0, 1] \mid f(0) = f(1) = 0\}$  and consider the one-dimensional **heat equation**

$$\frac{\partial}{\partial t} u(t, x) = \frac{\partial^2}{\partial x^2} u(t, x)$$

on a finite interval  $x \in [0, 1]$  with Dirichlet boundary conditions  $u(0) = u(1) = 0$  and the initial condition  $u(0, x) = g(x)$ . The corresponding operator is

$$Af = f'', \quad \mathfrak{D}(A) = \{f \in C_0[0, 1] \mid f \in C^2[0, 1]\} \subset X.$$

Note that  $\mathfrak{D}(A)$  is dense. For  $\ell \in \mathcal{J}(f)$  we can choose  $\ell(g) = f(x_0)^* g(x_0)$ , where  $x_0$  is chosen such that  $|f(x_0)| = \|f\|_\infty$ . Then  $\text{Re}(f(x_0)^* f(x))$  has a global maximum at  $x = x_0$  and if  $f \in \mathfrak{D}(A)$  we must have  $\text{Re}(f(x_0)^* f''(x_0)) \leq$

0 provided this maximum is in the interior of  $(0, 1)$ . If  $x_0$  is at the boundary this holds trivially and consequently  $A$  is dissipative. That  $A - \lambda$  is surjective follows using the Green's function as in Section 3.3: For  $g \in X$  the function

$$f(x) := (R_A(\lambda)g)(x) = \int_0^1 G(\lambda, x, y)g(y)dy,$$

where

$$G(\lambda, x, y) := \frac{-1}{\sqrt{\lambda} \sinh(\lambda)} \begin{cases} \sinh(\sqrt{\lambda}(1-x)) \sinh(\sqrt{\lambda}y), & y \leq x, \\ \sinh(\sqrt{\lambda}(1-y)) \sinh(\sqrt{\lambda}x), & x \leq y, \end{cases}$$

is in  $\mathfrak{D}(A)$  and satisfies  $(A - \lambda)f = g$ . Note that alternatively one could compute the norm of the resolvent

$$\|R_A(\lambda)\| = \frac{1}{\lambda} \left( 1 - \frac{1}{\cosh(\sqrt{\lambda}/2)} \right)$$

(equality is attained for constant functions; while these are not in  $X$ , you can approximate them by choosing functions which are constant on  $[\varepsilon, 1 - \varepsilon]$ ).  $\diamond$

**Example 10.15.** Let us consider the heat equation on  $x := C_{buc}(\mathbb{R})$  the bounded uniformly continuous functions. Since the uniform limit of uniformly continuous functions is again uniformly continuous, this is a closed subspace of  $C_b(\mathbb{R})$  and hence a Banach space (it will become clear why we do not choose  $C_b(\mathbb{R})$  in a moment). In this case we choose  $\mathfrak{D}(A) := C_{buc}^2(\mathbb{R}) := \{f \in C^2(\mathbb{R}) | f, f', f'' \in X\}$ . Now dissipativity does not follow as in the previous example since the maximum might not be attained. Hence we go directly for the resolvent whose kernel is given by

$$G(\lambda, x, y) := \frac{-1}{2\sqrt{\lambda}} e^{-\sqrt{\lambda}|x-y|}.$$

One checks that  $R_A(\lambda)$  is a bounded map on  $X$  whose norm is given by  $\|R_A(\lambda)\| = \frac{1}{\lambda}$  with equality for constant functions. Moreover, for given  $g \in X$  we have  $g := R_A(\lambda)f \in \mathfrak{D}(A)$  with  $(A - \lambda)f = g$ . Conversely, if  $f \in \mathfrak{D}(A)$  one checks that  $R_A(\lambda)(A - \lambda)f = f$  and hence  $R_A(\lambda)$  is indeed the resolvent of  $A$ . Up to this point everything would work on  $C_b(\mathbb{R})$ . Moreover, note that the mollification of a function  $g \in X$  will be in  $\mathfrak{D}(A)$  and converge uniformly to  $g$ . Hence  $\mathfrak{D}(A)$  is dense in  $X$  (but not in  $C_b(\mathbb{R})$ ).

Note that since the heat group is given by mollification with the heat kernel, the same argument also shows directly that the heat group is not strongly continuous on  $C_b(\mathbb{R})$ .  $\diamond$

**Example 10.16.** Another neat example is the following linear **delay differential equation**

$$\dot{u}(t) = \int_{t-1}^t u(s) d\nu(s), \quad t > 0, \quad u(s) = g(s), \quad -1 \leq s \leq 0,$$

where  $\nu$  is a complex measure. To this end we introduce the following operator

$$Af := f', \quad \mathfrak{D}(A) := \{f \in C^1[-1, 0] \mid f'(0) = \int_{-1}^0 f(s) d\nu(s)\} \subset C[0, 1].$$

Suppose that we can show that it generates a semigroup  $T$  on  $X = C[0, 1]$  and set  $u(t) := (T(t)f)(0)$  for  $f \in \mathfrak{D}(A)$ . Then, since  $T$  leaves  $\mathfrak{D}(A)$  invariant, the function  $r \mapsto (T(t+r)f)(s-r)$  is differentiable with

$$\frac{d}{dr}(T(t+r)f)(s-r) = (T(t+r)Af)(s-r) - (T(t+r)f')(s-r) = 0$$

and we conclude  $(T(t+r)f)(s-r) = (T(t)f)(s)$  for  $-1+r \leq s \leq 0$ . In particular, for  $r = s$  we obtain  $u(t+s) = (T(t)f)(s)$ . Hence we obtain

$$\begin{aligned} \dot{u}(t) &= \frac{d}{dt}(T(t)f)(0) = (AT(t)f)(0) = \int_{-1}^0 (T(t)f)(s) d\nu(s) \\ &= \int_{-1}^0 u(t+s) d\nu(s) \end{aligned}$$

and  $u$  solves our delay differential equation. Now if  $g \in C[0, 1]$  is given we can approximate it by a sequence  $f_n \in \mathfrak{D}(A)$ . Then  $u_n(t) := (T_n(t)f_n)(0)$  will converge uniformly on compact sets to  $u(t) := (T(t)g)(0)$  and taking the limit in the differential equation shows that  $u$  is differentiable and satisfies the differential equation.

Hence it remains to show that  $A$  generates a semigroup. First of all we claim that  $\tilde{A} := A - \|\nu\|$  is dissipative, where  $\|\nu\|$  is the total variation of  $\nu$ . As in the previous example, for  $\ell \in \mathcal{J}(f)$  we can choose  $\ell(g) = f(x_0)^*g(x_0)$  where  $x_0$  is chosen such that  $|f(x_0)| = \|f\|_\infty$ . Then  $\operatorname{Re}(f(x_0)^*f(x))$  has a global maximum at  $x = x_0$  and if  $f \in \mathfrak{D}(A)$  we must have  $\operatorname{Re}(f(x_0)^*f'(x_0)) = 0$  provided  $x_0$  is in the interior. If  $x_0 = -1$  we still must have  $\operatorname{Re}(f(x_0)^*f'(x_0)) \leq 0$ . In both cases  $\operatorname{Re}(\ell(\tilde{A}f)) \leq -\|\nu\|\|f(x_0)\|^2 \leq 0$ . If  $x_0 = 0$  we compute

$$\operatorname{Re}(\ell(\tilde{A}f)) = \operatorname{Re}\left(f^*(0) \int_{-1}^0 f(s) d\nu(s)\right) - \|\nu\|\|f(0)\|^2 \leq 0$$

since  $|f(s)| \leq |f(0)|$ . Thus  $\tilde{A}$  is dissipative. Moreover, it is straightforward to verify that the differential equation  $(\tilde{A} - \lambda)f = g$  has a unique solution  $f \in \mathfrak{D}(A)$  for  $\lambda > 0$  since  $|\int_{-1}^0 e^{(\lambda + \|\nu\|)s} d\nu(s)| \leq \|\nu\|$ .  $\diamond$

Finally, we note that the condition that  $A - \lambda_0$  is surjective can be weakened to the condition that  $\operatorname{Ran}(A - \lambda_0)$  is dense. To this end we need:

**Lemma 10.19.** *Suppose  $A$  is a densely defined dissipative operator. Then  $A$  is closable and the closure  $\bar{A}$  is again dissipative.*

**Proof.** Recall that  $A$  is closable if and only if for every  $x_n \in \mathfrak{D}(A)$  with  $x_n \rightarrow 0$  and  $Ax_n \rightarrow y$  we have  $y = 0$ . So let  $x_n$  be such a sequence and chose another sequence  $y_n \in \mathfrak{D}(A)$  such that  $y_n \rightarrow y$  (which is possible since  $\mathfrak{D}(A)$  is assumed dense). Then by dissipativity (specifically Corollary 10.17)

$$\|(A - \lambda)(\lambda x_n + y_m)\| \geq \lambda \|\lambda x_n + y_m\|, \quad \lambda > 0$$

and letting  $n \rightarrow \infty$  and dividing by  $\lambda$  shows

$$\|y + (\lambda^{-1}A - 1)y_m\| \geq \|y_m\|.$$

Finally  $\lambda \rightarrow \infty$  implies  $\|y - y_m\| \geq \|y_m\|$  and  $m \rightarrow \infty$  yields  $0 \geq \|y\|$ , that is,  $y = 0$  and  $A$  is closable. To see that  $\bar{A}$  is dissipative choose  $x \in \mathfrak{D}(\bar{A})$  and  $x_n \in \mathfrak{D}(A)$  with  $x_n \rightarrow x$  and  $Ax_n \rightarrow \bar{A}x$ . Then (again using Corollary 10.17) taking the limit in  $\|(A - \lambda)x_n\| \geq \lambda\|x_n\|$  shows  $\|(\bar{A} - \lambda)x\| \geq \lambda\|x\|$  as required.  $\square$

Consequently:

**Corollary 10.20.** *Suppose the linear operator  $A$  is densely defined, dissipative, and  $\text{Ran}(A - \lambda_0)$  is dense for one  $\lambda_0 > 0$ . Then  $A$  is closable and  $\bar{A}$  is the generator of a contraction semigroup.*

**Proof.** By the previous lemma  $A$  is closable with  $\bar{A}$  again dissipative. In particular,  $\bar{A}$  is injective and by Lemma 8.1 we have  $(\bar{A} - \lambda_0)^{-1} = (\overline{A - \lambda_0})^{-1}$ . Since  $(A - \lambda_0)^{-1}$  is bounded its closure is defined on the closure of its domain, that is,  $\text{Ran}(\bar{A} - \lambda_0) = \overline{\text{Ran}(A - \lambda_0)} = X$ . The rest follows from the Lumer–Phillips theorem.  $\square$

**Problem\* 10.20.** *Let  $T(t)$  be a  $C_0$ -semigroup and  $\alpha > 0$ ,  $\lambda \in \mathbb{C}$ . Show that  $S(t) := e^{\lambda t}T(\alpha t)$  is a  $C_0$ -semigroup with generator  $B = \alpha A + \lambda$ ,  $\mathfrak{D}(B) = \mathfrak{D}(A)$ .*

**Problem 10.21.** *Let  $T$  be a bounded  $C_0$ -semigroup satisfying  $\|T(t)\| \leq M$ . Then*

$$\|g\|_T := \sup_{t \geq 0} \|T(t)g\|$$

*defines an equivalent norm on  $X$  satisfying*

$$\|g\| \leq \|g\|_T \leq M\|g\|.$$

**Problem 10.22.** *Show that  $A$  generates a  $C_0$  group of isometries, that is,  $\|T(t)g\| = \|g\|$  for all  $g \in X$  if and only if both  $A$  and  $-A$  generate contraction semigroups. That is, both  $A$  and  $-A$  satisfy the hypothesis of either the Hill–Yosida or the Lumer–Phillips theorem.*

**Problem 10.23.** *Let  $T(t)$  be a  $C_0$ -semigroup satisfying  $\|T(t)\| \leq Me^{\omega t}$  with generator  $A$  and  $B \in \mathcal{L}(X)$ . Show that  $A + B$  generates a  $C_0$ -semigroup*

$S(t)$  satisfying  $\|S(t)\| \leq Me^{(\omega+M\|B\|)t}$ . (Hint: First use Problem 10.21 to reduce it to the case of a contraction. Then use Problem 8.28.)

**Problem 10.24.** Let  $X = \ell^2(\mathbb{N})$  and  $(Aa)_n := in^2a_n$ ,  $(Ba)_n := na_n$  both defined maximally. Show that  $A$  generates a  $C_0$ -semigroup but  $A + \varepsilon B$  does not for any  $\varepsilon > 0$ .

**Problem 10.25.** Consider the heat equation (Example 10.14) on  $[0, 1]$  with Neumann boundary conditions  $u'(0) = u'(1) = 0$ .

**Problem 10.26.** Consider the heat equation (Example 10.15) on  $C_0(\mathbb{R})$ .

#### section Semilinear equations

Linear problems are often only a first approximation and adding a non-linear perturbation leads to the following semilinear problem

$$\dot{u} = Au + F(u), \quad u(0) = g, \quad (10.35)$$

where  $A$  is supposed to generate a semigroup  $T(t)$  and  $F \in C(X, X)$  such that we can recast this problem as

$$u(t) = T(t)g + \int_0^t T(t-s)F(u(s))ds. \quad (10.36)$$

In fact, if we have a solution  $u \in C([0, t_+), [\mathfrak{D}(A)]) \cap C^1([0, t_+), X)$  of (10.35), then Duhamel's formula shows that (10.36) holds. In the other direction you need a stronger assumption on  $F$ . However, it will be more convenient to work with (10.36) and we will call a solution a mild solution of (10.35). In fact, (10.36) is of fixed point type and hence begs us to apply the contraction principle. As always with nonlinear equations, we expect the solution to be only defined on a finite time interval  $[0, t_+)$  in general.

**Theorem 10.21.** Suppose  $F \in C(X, X)$  is Lipschitz continuous on bounded sets. Then for every  $g \in X$  there is a  $t_0 = t_0(\|g\|) > 0$ , such that there is a unique mild solution  $u \in C([0, t_0], X)$ . Moreover, the solution map  $g \mapsto u(t)$  will be Lipschitz continuous from every ball  $\|g\| \leq \rho$  to  $C([0, t_0(\rho)], X)$ .

**Proof.** We will consider  $0 \leq t \leq 1$  and set  $M := \sup_{0 \leq t \leq 1} \|T(t)\|$ . Let  $r := 1 + M\|g\|$  and consider the closed ball  $\bar{B}_r(0) \subset X$ . Let  $L = L(r)$  be the Lipschitz constant of  $F$  on  $\bar{B}_r(0)$ . Set

$$K(u)(t) := T(t)g + \int_0^t T(t-s)F(u(s))ds$$

and note that

$$\begin{aligned} \|K(u)(t)\| &\leq M\|g\| + M \int_0^t (\|F(0)\| + L\|u(s)\|)ds \\ &\leq M\|g\| + M\|F(0)\|t + MLt \sup_{0 \leq s \leq t} \|u(s)\| \end{aligned}$$

and

$$\|K(u)(t) - K(v)(t)\| \leq M \int_0^t L(\|u(s) - v(s)\|) ds \leq MLt \sup_{0 \leq s \leq t} \|u(s) - v(s)\|$$

Hence if we choose  $t_0 \leq 1$  such that

$$M(\|F(0)\| + Lr)t_0 < 1$$

then  $\theta := MLt_0 < 1$  and  $K$  will be a contraction on  $\bar{B}_r(0) \subset C([0, t_0], X)$ . In particular, for two solutions  $u_j$  corresponding to  $g_j$  with  $\|g_j\| \leq \|g\|$  we will have  $\|u_1 - u_2\|_\infty \leq \frac{1}{1-\theta} \|g_1 - g_2\|$ .

This establishes the theorem except for the fact that it only shows uniqueness for solutions which stay within  $\bar{B}_r(0)$ . However, since  $K$  maps from  $\bar{B}_r(0)$  to its interior  $B_r(0)$ , a potential different solution starting at  $g \in B_r(0)$  would need to branch off at the boundary, which is impossible since our solution does not reach the boundary.  $\square$

**Corollary 10.22.** *Suppose that  $F \in C([\mathfrak{D}(A)], [\mathfrak{D}(A)])$  is Lipschitz continuous on bounded sets. Then for every  $g \in \mathfrak{D}(A)$  there is a  $t_1 = t_1(\|g\|_A) > 0$ , such that there is a unique strong solution  $u \in C^1([0, t_1], X) \cap C([0, t_1], [\mathfrak{D}(A)])$ .*

**Proof.** Since  $T$  restricted to  $[\mathfrak{D}(A)]$  generates a  $C_0$ -semigroup (see the discussion after Lemma 10.4), we can apply the previous result to this semigroup giving a solution  $u \in C([0, t_1], [\mathfrak{D}(X)])$ . This solution is in  $C^1([0, t_1], X)$  by Lemma 10.9.  $\square$

**Corollary 10.23.** *If  $F$  is globally Lipschitz, then solutions are global.*

**Proof.** In this case we can consider  $K$  on all of  $C([0, t_0], X)$  and set  $M := \sup_{0 \leq t \leq t_0} \|T(t)\|$ . By induction we get for the iterates

$$\|K^n(u)(t) - K^n(v)(t)\| \leq \frac{(MLt)^n}{n!} \sup_{0 \leq s \leq t} \|u(s) - v(s)\|$$

and Weissinger's fixed point theorem (Theorem 9.28) gives a solution on  $C([0, t_0], X)$ . Since  $t_0 > 0$  is arbitrary, the claim follows.  $\square$

If solutions are not global, there is still a unique maximal solution: Fix  $g \in X$  and let  $u_j$  be two solutions on  $[0, t_j]$  with  $0 < t_1 < t_2$ . By the uniqueness part of our theorem, we will have  $u_1(t) = u_2(t)$  for  $0 \leq t < \tau$  for some  $\tau > 0$ . Suppose  $\tau < t_1$  and  $\tau$  is chosen maximal. Let  $r := \max_{0 \leq t \leq \tau} \|u_1(t)\|$  and  $0 < \varepsilon < \min(\tau, t_0(r)/2)$  with  $t_0(r)$  from our theorem. Then there is a solution  $v$  starting with initial condition  $u_1(\tau - \varepsilon)$  which is defined on  $[0, 2\varepsilon]$ . Moreover, again by the uniqueness part of our theorem  $u_1(t) = v(t - (\tau - \varepsilon)) = u_2(t)$  for  $\tau - \varepsilon \leq t \leq \tau + \varepsilon$  contradiction our assumption that  $\tau$  is maximal. Hence taking the union (with respect to their domain) over all mild solutions starting at  $g$ , we get a unique solution

defined on a maximal domain  $[0, t_+(g))$ . Note that if  $t_+(g) < \infty$ , then  $\|u(t)\|$  must blow up as  $t \rightarrow t_+(g)$ :

**Lemma 10.24.** *Let  $t_+(g)$  be the maximal time of existence for the mild solution starting at  $g$ . If  $t_+(g) < \infty$ , then  $\liminf_{t \rightarrow t_+(g)} \|u(t)\| = \infty$ .*

**Proof.** Assume that  $\rho := \sup_{0 \leq t < t_+(g)} \|u(t)\| < \infty$ . As above, choose  $0 < \varepsilon < \min(t_+(g), t_0(\rho)/2)$  with  $t_0(\rho)$  from our theorem. Then the solution  $v$  starting with initial condition  $u(t_+(g) - \varepsilon)$  extends  $u$  to the interval  $[0, t_+(g) + \varepsilon)$ , contradicting maximality.  $\square$

In many applications it will happen that the local Lipschitz constant depends only on a weaker norm. In such a situation also the weaker norm will have to blow up.

**Lemma 10.25.** *Let  $\|\cdot\|_0$  be a norm, which is weaker than the standard norm on  $X$ , that is,  $\|x\|_0 \leq C_0\|x\|$  for all  $x \in X$ . Suppose that there is a nondecreasing function  $L : [0, \infty) \rightarrow [0, \infty)$  and a constant  $C$  such that*

$$\|F(x)\| \leq C + L(\|x\|_0)\|x\|. \quad (10.37)$$

*If  $t_+(g) < \infty$ , then  $\liminf_{t \rightarrow t_+(g)} \|u(t)\|_0 = \infty$ .*

**Proof.** Starting from (10.36) we obtain

$$\begin{aligned} \|u(t)\| &\leq Me^{\omega t}\|g\| + M \int_0^t e^{\omega(t-s)} \|F(u(s))\| ds \\ &\leq Me^{\omega t}\|g\| + MC \frac{e^{\omega t} - 1}{\omega} + M \int_0^t e^{\omega(t-s)} L(\|u(s)\|_0) \|u(s)\| ds \end{aligned}$$

and hence Gronwall's inequality ([?, Lemma 2.7]) implies

$$\|u(t)\| \leq M(\|g\| + Ct) \exp \left( \omega t + M \int_0^t L(\|u(s)\|_0) ds \right).$$

This shows that the  $\|\cdot\|$  norm cannot blow up before the  $\|\cdot\|_0$  norm.  $\square$

So the key to proving global existence of solutions is an a priori bound on the norm of the solution. Typically such a bound will come from a conservation law.

**Example 10.17.** Consider the **discrete nonlinear Schrödinger equation (dNLS)**

$$i\dot{u}(t) = Hu(t) + f(|u(t)|)u(t), \quad t \in \mathbb{R},$$

in  $X = \ell^2(\mathbb{Z})$ . Here  $H$  could be any bounded self-adjoint operator any locally Lipschitz continuous function. In applications  $Hu_n = u_{n+1} + u_{n-1} + q_n u_n$  is the Jacobi operator, with  $q \in \ell^\infty(\mathbb{Z})$  a real-valued sequence corresponding to an external potential and  $q = 0$  (or  $q = -2$ , depending on your preferences)



is the free discrete Schrödinger operator. The function  $f$  is typically an even polynomial.

Clearly we have

$$|f(|x|)x - f(|y|)y| \leq |f(|x|) - f(|y|)||x| + |f(|y|)||x - y| \leq L(\max(|x|, |y|))|x - y|$$

for  $x, y \in \mathbb{C}$ , where

$$L(r) := r \max_{[0, r]} |f'| + \max_{[0, r]} |f|.$$

Consequently

$$\|f(|u|)u - f(|v|)v\|_2 \leq L(\max(\|u\|_\infty, \|v\|_\infty))\|u - v\|_2$$

is the required Lipschitz estimate (recall  $\|u\|_\infty \leq \|u\|_2$ ) to apply Theorem 10.21 to conclude existence of local solutions. Note that since our generator is bounded, there is no difference between mild and strong solutions. Moreover, Lemma 10.25 implies that if solutions are not global, then  $\|u(t)\|_\infty$  must blow up. In this respect note that while  $\ell^2(\mathbb{Z})$  is the most natural choice from a quantum mechanical point of view, our analysis still applies if we replace  $\ell^2(\mathbb{Z})$  by  $\ell^p(\mathbb{Z})$  for any  $1 \leq p \leq \infty$ . Then by  $\ell^{p_1}(\mathbb{Z}) \subset \ell^{p_2}(\mathbb{Z})$  for  $p_2 \leq p_1$  and for a solution starting in  $\ell^{p_1}(\mathbb{Z}) \subset \ell^{p_2}(\mathbb{Z})$  the existence interval in  $\ell^{p_2}(\mathbb{Z})$  could be larger than in  $\ell^{p_1}(\mathbb{Z})$ . However, by Lemma 10.25 this is not the case and the solutions does not just loose decay but will always blow up pointwise (if it blows up at all).

Finally, if we assume that  $f$  is real-valued then solutions satisfy

$$\frac{d}{dt} \|u(t)\|_2^2 = 2\operatorname{Re}\langle \dot{u}(t), u(t) \rangle = 2\operatorname{Im}(\langle Hu, u \rangle + \langle f(|u(t)|)u(t), u(t) \rangle) = 0$$

and hence the dNLS equation has a unique global norm preserving solution  $u \in C^1(\mathbb{R}, \ell^2(\mathbb{Z}))$ .  $\diamond$

Let me close with a few remarks: First of all, it is straightforward to extend these results to the situation where  $F$  depends on  $t$  or to the case where  $T$  is a group. Details are left to the reader. Moreover, if  $A$  is bounded, then it is Lipschitz continuous and could be absorbed in  $F$ . In fact, in this case our theorem just gives the Picard–Lindelöf theorem for ordinary differential equations in Banach spaces (in particular, in this case the differential equation (10.35) and the integral equation (10.36) are equivalent).

**Problem 10.27.** *Show that solutions are global if  $\|F(x)\| \leq C(1 + \|x\|)$  for some constant  $C$ . (Hint: Use Gronwall's inequality to bound  $\|u(t)\|$ .)*

# The nonlinear Schrödinger equation

The purpose of this chapter is to investigate a prototypical example, the initial value problem for the **nonlinear Schrödinger equation** (NLS)

$$iu_t + \Delta u = \pm |u|^{\alpha-1}u, \quad u(0) = g. \quad (11.1)$$

The two cases  $-$  and  $+$  are known as **focusing** and **defocusing**, respectively. Of particular importance in applications are the cubic ( $\alpha = 3$ ) and quintic ( $\alpha = 5$ ) case. Note that if  $u$  is a solution, then so will be  $v(t, x) = u(-t, x)^*$  and hence it suffices to look at positive times only.

## 11.1. Local well-posedness in $H^r$ for $r > \frac{n}{2}$

Equation (11.1) is a semilinear equation of the type considered in Section 10.3 and hence we need to look at the linear **Schrödinger equation**

$$iu_t + \Delta u = 0, \quad u(0) = g \quad (11.2)$$

first. We recall that the solution for  $g \in H^2(\mathbb{R}^n)$  can be obtained using the Fourier transform and is given by

$$u(t) = T_S(t)g, \quad T_S(t) = \mathcal{F}^{-1}e^{-i|p|^2t}\mathcal{F}. \quad (11.3)$$

Note that  $T_S(t) : L^2(\mathbb{R}^n) \rightarrow L^2(\mathbb{R}^n)$  is a unitary operator (since  $|e^{-i|p|^2t}| = 1$ ):

$$\|u(t)\|_2 = \|g\|_2. \quad (11.4)$$

In fact, we even have that  $T_S(t) : H^r(\mathbb{R}^n) \rightarrow H^r(\mathbb{R}^n)$  is unitary.

**Theorem 11.1.** *The family  $T_S(t)$  is a  $C_0$ -group in  $H^r(\mathbb{R}^n)$  whose generator is  $i\Delta$ ,  $\mathfrak{D}(i\Delta) = H^{r+2}(\mathbb{R}^n)$ .*

Note that we have

$$(T_S(t)g)^* = T_S(-t)g^*, \quad g \in L^2(\mathbb{R}^n) \quad (11.5)$$

and

$$\partial_j T_S(t)g = T_S(t)\partial_j g, \quad 1 \leq j \leq n, g \in H^1(\mathbb{R}^n). \quad (11.6)$$

Next we turn to the nonlinear Schrödinger equation. If we assume that  $u, |u|^{\alpha-1}u \in C([0, T], L^2(\mathbb{R}^n))$  we can use Duhamel's formula to rewrite the nonlinear Schrödinger equation as

$$u(t) = T_S(t)g \mp i \int_0^t T_S(t-s)|u(s)|^{\alpha-1}u(s)ds \quad (11.7)$$

just as we did in Section 10.3. In order to apply our theory, we need that the nonlinearity  $F(u) = \mp i|u|^{\alpha-1}u$  is Lipschitz on  $X$ . Clearly for  $X = L^2(\mathbb{R}^n)$  this will not be the case, as the image of a square integrable function will not be square integrable. However, the key observation is that for  $r > \frac{n}{2}$  the space  $H^r(\mathbb{R}^n)$  is a Banach algebra (Lemma 8.33 from [48]) and hence, if we assume our nonlinearity to be of the form  $F(u) = \mp i|u|^{\alpha-1}u$  with  $\alpha-1 = 2k$  where  $k \in \mathbb{N}$ , then  $F : H^r(\mathbb{R}^n) \rightarrow H^r(\mathbb{R}^n)$  is Lipschitz on bounded sets since

$$F(u) - F(v) = u^{k+1}Q_{k-1}(u^*, v^*)(u-v)^* + (v^*)^k Q_k(u, v)(u-v), \quad (11.8)$$

where  $Q_k(x, y) = \mp i \sum_{j=0}^k x^{k-j} y^j$ . Another algebra which is natural in this context is the Wiener algebra.

$$\mathcal{A}(\mathbb{R}^n) := \{f | \hat{f} \in L^1(\mathbb{R}^n)\}, \quad \|f\|_{\mathcal{A}} := \|\hat{f}\|_1. \quad (11.9)$$

Just as with  $H^r(\mathbb{R}^n)$ , the Schrödinger group  $T_S$  leaves  $\mathcal{A}(\mathbb{R}^n)$  invariant and preserves its norm. Note that we have  $H^r(\mathbb{R}^n) \subset \mathcal{A}(\mathbb{R}^n)$  for  $r > \frac{n}{2}$  since  $(1 + |p|^2)^{-r} \in L^2(\mathbb{R}^n)$  for such  $r$ . The embedding being continuous,  $\|f\|_{\mathcal{A}} \leq \|(1 + |\cdot|^2)^{-r}\|_2 \|f\|_{H^r}$ .

Hence Theorem 10.21 applies and we get:

**Theorem 11.2.** *Let  $\alpha = 2k+1$  be an odd integer and  $X = H^r(\mathbb{R}^n)$  for  $r > \frac{n}{2}$  or  $X = \mathcal{A}(\mathbb{R}^n)$ . Then for every  $g \in X$  there is a  $t_0 = t_0(\|g\|) > 0$ , such that there is a unique solution  $u \in C([-t_0, t_0], X)$  of (11.7). Moreover, the solution map  $g \mapsto u(t)$  will be Lipschitz continuous from every ball  $\|g\| \leq \rho$  to  $C([-t_0(\rho), t_0(\rho)], X)$ .*

Note that the mild solution will be a strong solution for  $g \in H^{r+2}$  since  $F : H^{r+2} \rightarrow H^{r+2}$  is Lipschitz continuous on bounded sets. Moreover, for each initial condition there is a maximal solution and Lemma 10.24 implies:

**Lemma 11.3.** *This solution exists on a maximal time interval  $(t_-(g), t_+(g))$  and if  $|t_{\pm}(g)| < \infty$  we must have  $\liminf_{t \rightarrow t_{\pm}(g)} \|u(t)\| = \infty$ .*

An interesting observation is that the maximal existence time does not depend on  $r$ . This is known as persistence of regularity:

**Lemma 11.4.** *Let  $g \in H^r(\mathbb{R}^n)$  with  $r > \frac{n}{2}$  or  $g \in \mathcal{A}(\mathbb{R}^n)$ . Let  $t_{+,r}(g)$ ,  $t_{+,\mathcal{A}}(g)$  be the maximal existence time of the solution with initial condition  $g$  with respect to these cases. Then  $t_{+,r}(g) = t_{+,\mathcal{A}}(g)$ .*

**Proof.** Using (Lemma 8.33 from [48] )

$$\|fg\|_{H^r} \leq C_{n,r}(\|f\|_{H^r}\|g\|_{\mathcal{A}} + \|f\|_{\mathcal{A}}\|g\|_{H^r})$$

recursively we obtain  $\| |u|^{\alpha-1}u \|_{H^r} \leq C\|u\|_{\mathcal{A}}^{\alpha-1}\|u\|_{H^r}$ . Now the claim follows from Lemma 10.25.  $\square$

Of course up to this point we can replace the nonlinearity by an arbitrary polynomial in  $u$  and  $u^*$ . In fact, it is even possible to replace the nonlinearity by a (sufficiently smooth) function, but in this case the required Lipschitz estimate is more tedious to derive since we cannot just simply rely on the algebra structure.

In order to get global solutions the following conservation laws will be crucial: **Momentum**

$$M(t) := \frac{1}{2}\|u(t)\|_2^2$$

and **energy**

$$E(t) := \frac{1}{2}\|\nabla u(t)\|_2^2 \pm \frac{1}{\alpha+1}\|u(t)\|_{\alpha+1}^{\alpha+1}.$$

**Lemma 11.5.** *Let  $r > \frac{n}{2}$  and  $g \in H^r(\mathbb{R}^n)$ . Then  $M(t) = M(0)$  for all  $t \in (t_-(g), t_+(g))$ . If in addition,  $r \geq 1$  then also  $E(t) = E(0)$  for all  $t \in (t_-(g), t_+(g))$ .*

**Proof.** If  $u$  is a sufficiently smooth solution this can be verified directly (Problem 11.1). For the general case approximate by smooth solutions (using local Lipschitz continuity of the solution map) and conclude that  $M(t)$  is locally constant and hence constant on its interval of existence. Similarly for  $E(t)$ .  $\square$

So in the focusing case we get existence of global solutions in  $H^1$  if  $n = 1$  such that our local results holds for  $r = 1$ . In the defocusing case the energy is not positive and we cannot immediately control the  $H^1$  norm using  $E$  and  $M$ .

**Problem 11.1.** *Let  $u \in C([-t_0, t_0], H^{r+2}(\mathbb{R}^n)) \cap C^1([-t_0, t_0], H^r(\mathbb{R}^n))$  be a strong solution of the NLS equation (with  $r > \frac{n}{2}$ ). Show that momentum and energy are independent of  $t \in [-t_0, t_0]$ .*

### 11.2. Strichartz estimates

In order to improve the results from the previous section we need a better understanding of the linear Schrödinger equation. Unlike for example the heat equation, the Schrödinger equation does only preserve but not improve the regularity of the initial condition. For example, choosing  $f \in L^2 \setminus L^p$  (for some  $p \neq 2$ ) and considering  $g = T_S(-t_0)f$  shows that there are initial conditions in  $L^2$  which are not in  $L^p$  at a given later time  $t_0$ . However, our aim in this section is to show that we still have  $T(t)g \in L^p$  most of the time.

To this end we first need an explicit expression for the solution. As in the case of the heat equation, we would like to express our solution as a convolution with the initial condition. However, here we run into the problem that  $e^{-i|p|^2 t}$  is not integrable. To overcome this problem we consider

$$f_\varepsilon(p) = e^{-(it+\varepsilon)p^2}, \quad \varepsilon > 0. \quad (11.10)$$

Then using the well-known formula (Lemma 8.6 from [48])

$$\mathcal{F}(e^{-z|x|^2/2})(p) = \frac{1}{z^{n/2}} e^{-|p|^2/(2z)}, \quad \operatorname{Re}(z) > 0, \quad (11.11)$$

where  $z^{n/2}$  is the standard branch with branch cut along the negative real axis, together with the fact that the Fourier transform maps convolutions into products (Corollary 8.15 from [48]) we obtain

$$(f_\varepsilon \hat{g})^\vee(x) = \frac{1}{(4\pi(it+\varepsilon))^{n/2}} \int_{\mathbb{R}^n} e^{-\frac{|x-y|^2}{4(it+\varepsilon)}} g(y) d^n y. \quad (11.12)$$

Taking the limit  $\varepsilon \downarrow 0$  we finally arrive at

$$T_S(t)g(x) = \frac{1}{(4\pi it)^{n/2}} \int_{\mathbb{R}^n} e^{i\frac{|x-y|^2}{4t}} g(y) d^n y \quad (11.13)$$

for  $t \neq 0$  and  $g \in L^2(\mathbb{R}^n) \cap L^1(\mathbb{R}^n)$ . In fact, the left-hand side converges to  $T_S(t)g$  in  $L^2$  and the limit of the right-hand side exists pointwise by dominated convergence and its pointwise limit must thus be equal to its  $L^2$  limit.

Using this explicit form, we can again draw some further consequences. For example, if  $g \in L^2(\mathbb{R}^n) \cap L^1(\mathbb{R}^n)$ , then  $u(t) := T_S(t)g \in C_0(\mathbb{R}^n)$  for  $t \neq 0$  (Problem 11.3) and satisfies

$$\|u(t)\|_\infty \leq \frac{1}{|4\pi t|^{n/2}} \|g\|_1. \quad (11.14)$$

Moreover, we even have  $u \in C(\mathbb{R} \setminus \{0\}, C_0(\mathbb{R}^n))$  (Problem 11.3).

Thus we have spreading of wave functions in this case. In fact, invoking the Riesz–Thorin interpolation theorem (Theorem 9.2 from [48]) we even

get

$$\|u(t)\|_p \leq \frac{1}{|4\pi t|^{n/2-n/p}} \|g\|_{p'}. \quad (11.15)$$

for any  $p \in [2, \infty]$  with  $\frac{1}{p} + \frac{1}{p'} = 1$ . This also gives  $u \in C(\mathbb{R} \setminus \{0\}, L^p(\mathbb{R}^n))$ .

Next we look at average decay in an  $L^p$  sense instead of pointwise estimates with respect to  $t$ . To this end we will consider functions  $f \in L^r(\mathbb{R}, L^p(\mathbb{R}^n))$  and we will denote the corresponding norm by

$$\|f\|_{L^r(L^p)} := \begin{cases} \left( \int_{\mathbb{R}} \|f(t)\|_p^r dt \right)^{1/r}, & r < \infty, \\ \sup_{t \in \mathbb{R}} \|f(t)\|_p, & r = \infty. \end{cases} \quad (11.16)$$

Please recall that  $L^r(\mathbb{R}, L^p(\mathbb{R}^n))$  is a Banach space defined with the help of the Bochner integral (cf. Theorem 5.32 from [48]). It consists of (equivalence classes with respect to equality a.e. of) strongly measurable functions  $f(t)$  for which  $\|f(t)\|_p$  is in  $L^r$ . Here strongly measurable means, that  $f(t)$  is a limit of simple functions  $s_n(t)$ . It turns out that a function is strongly measurable if and only if it is measurable and its range is separable. In our situation this latter condition will come for free in the case  $p < \infty$  and similarly in the case  $p = \infty$  if the range is contained in  $C_0(\mathbb{R}^n)$ . We will also need the following variational characterization of our space-time norms (Problem 5.25 from [48]) for a given strongly measurable function  $f$ :

$$\|f\|_{L^r(L^p)} = \sup_{\|g\|_{L^{r'}(L^{p'})}=1} \left| \int_{\mathbb{R}} \int_{\mathbb{R}^n} f(x, t) g(x, t) d^n x dt \right|. \quad (11.17)$$

Moreover, it suffices to take the sup over functions which have support in a compact rectangle.

We call a pair  $(p, r)$  **admissible** if

$$\begin{cases} 2 \leq p \leq \infty, & n = 1 \\ 2 \leq p < \frac{2n}{n-2}, & n \geq 2 \end{cases}, \quad \frac{2}{r} = \frac{n}{2} - \frac{n}{p}. \quad (11.18)$$

Note  $r \in [4, \infty]$  for  $n = 1$  and  $r \in (\frac{2}{n-1}, \infty]$  for  $n \geq 2$ .

**Lemma 11.6.** *Let  $T_S$  be the Schrödinger group and let  $(p, r)$  be admissible with  $p > 2$ . Then we have*

$$\left( \int_{\mathbb{R}} \left( \int_{\mathbb{R}} \|T_S(t-s)g(s)\|_p ds \right)^r dt \right)^{1/r} \leq C \|g\|_{L^{r'}(L^{p'})}, \quad (11.19)$$

where a prime denotes the corresponding dual index. Moreover,  $s \mapsto T(t-s)g(s) \in L^p(\mathbb{R}^n)$  is integrable for a.e.  $t \in \mathbb{R}$ .

**Proof.** Of course  $T_S(t-s)g(s)$  is measurable. Applying our interpolation estimate we obtain

$$\int \|T_S(t-s)g(s)\|_p ds \leq C \int \frac{1}{|t-s|^{1-\alpha}} \|g(s)\|_{p'} ds,$$

where  $\alpha = 1 - n(1/2 - 1/p) \in (0, 1)$  by our restriction on  $p$ .

Furthermore, our choice for  $r$  implies  $\alpha = 1 - \frac{2}{r} = \frac{1}{r'} - \frac{1}{r}$  with  $r' = \frac{2}{1+\alpha} \in (1, \alpha^{-1})$ . So taking the  $\|\cdot\|_{L^r}$  norm on both sides and using the Hardy–Littlewood–Sobolev inequality (Theorem 9.10 from [48]) gives the estimate.

Hence the claim about integrability follows from Minkowski’s integral inequality (Theorem 5.30 from [48]).  $\square$

Note that the case  $p = 2$  (and  $r = \infty$ ) the above lemma holds by unitarity and does not provide much new insight.

**Theorem 11.7** (Strichartz estimates). *Let  $T_S$  be the Schrödinger group and let  $(p, r)$  be admissible. Suppose  $g \in L^{r'}(\mathbb{R}, L^{p'}(\mathbb{R}^n))$  and  $f \in L^2(\mathbb{R}^n)$ . Then we have the following estimates:*

$$\|T_S(t)f\|_{L^r(L^p)} \leq C\|f\|_2, \quad (11.20)$$

$$\left\| \int_{\mathbb{R}} T_S(s)g(s)ds \right\|_2 \leq C\|g\|_{L^{r'}(L^{p'})}, \quad (11.21)$$

$$\left\| \int_{\mathbb{R}} T_S(t-s)g(s)ds \right\|_{L^r(L^p)} \leq C\|g\|_{L^{r'}(L^{p'})}, \quad (11.22)$$

where a prime denotes the corresponding dual index.

Here  $s \mapsto T_S(t-s)g(s) \in L^p(\mathbb{R}^n)$  is integrable for a.e.  $t \in \mathbb{R}$  and the integral in (11.21) has to be understood as a limit in  $L^2$  when taking an approximating sequence of functions  $g$  with support in compact rectangles.

**Proof.** Since the case  $p = 2$  follows from unitarity, we can assume  $p > 2$ . The claims about integrability and the last estimate follow from the lemma.

Using unitarity of  $T_S$  and Fubini we get

$$\int_{\mathbb{R}} \int_{\mathbb{R}^n} (T_S(t)f)(x)g(t, x)d^n x dt = \int_{\mathbb{R}^n} f(x) \int_{\mathbb{R}} (T_S(t)g(t))(x)dt d^n x,$$

for  $g \in L^{r'}(\mathbb{R}, L^{p'}(\mathbb{R}^n))$  with support in a compact rectangle. Note that in this case we have  $g(t) \in L^2(\mathbb{R}^n)$  since  $p' \leq 2$ . This shows that the first and second estimate are equivalent upon using the above characterization (11.17) as well as the analogous characterization for the  $L^2$  norm.

Similarly, using again unitarity of  $T_S$  and Fubini

$$\begin{aligned} & \left\| \int T_S(t)g(t)dt \right\|_2^2 \\ &= \int_{\mathbb{R}^n} \int_{\mathbb{R}} (T_S(t)g(t))(x)dt \int_{\mathbb{R}} (T_S(s)g(s))(x)^* ds d^n x \\ &= \int_{\mathbb{R}^n} \int_{\mathbb{R}} g(t, x) \int_{\mathbb{R}} T_S(t-s)g(s, x)^* ds dt d^n x, \end{aligned}$$

which shows that the second and the third estimate are equivalent with a similar argument as before.  $\square$

Note that using the scaling  $f(x) \rightarrow f(\lambda x)$  for  $\lambda > 0$  shows that the left-hand side of (11.20) scales like  $\lambda^{-n/p-2/r}$  while the right-hand side scales like  $\lambda^{-n/2}$ . So (11.20) can only hold if  $\frac{n}{p} + \frac{2}{r} = \frac{n}{2}$ .

In connection with the Duhamel formula the following easy consequence is also worth while noticing:

**Corollary 11.8.** *We also have*

$$\left\| \int_0^t T_S(t-s)g(s)ds \right\|_2 \leq C \|g\|_{L^{r'}(L^p)}, \quad (11.23)$$

$$\left\| \int_0^t T_S(t-s)g(s)ds \right\|_{L^r(L^p)} \leq C \|g\|_{L^{r'}(L^p)}. \quad (11.24)$$

**Proof.** The second estimate is immediate from the lemma and the first estimate follows from (11.21) upon restricting to functions  $g$  supported in  $[0, t]$  and using a simple change of variables  $\int_0^t T(t-s)g(s)ds = \int_0^t T(s)g(t-s)ds$ .  $\square$

Note that, apart from unitarity of  $T_S$ , only (11.14) was used to derive these estimates. Moreover, since  $T_S$  commutes with derivatives, we can also get analogous estimates for derivatives:

**Corollary 11.9.** *We have the following estimates for  $k \in \mathbb{N}_0$ :*

$$\|T_S(t)f\|_{L^r(W^{k,p})} \leq C \|f\|_{H^k}, \quad (11.25)$$

$$\left\| \int_{\mathbb{R}} T_S(s)g(s)ds \right\|_{H^k} \leq C \|g\|_{L^{r'}(W^{k,p'})}, \quad (11.26)$$

$$\left\| \int_{\mathbb{R}} T_S(t-s)g(s)ds \right\|_{L^r(W^{k,p})} \leq C \|g\|_{L^{r'}(W^{k,p'})}, \quad (11.27)$$



as well as

$$\left\| \int_0^t T_S(t-s)g(s)ds \right\|_{H^k} \leq C \|g\|_{L^{r'}(W^{k,p'})}, \quad (11.28)$$

$$\left\| \int_0^t T_S(t-s)g(s)ds \right\|_{L^r(W^{k,p})} \leq C \|g\|_{L^{r'}(W^{k,p'})}. \quad (11.29)$$

**Proof.** Consider dense sets  $f \in \mathcal{S}(\mathbb{R}^n)$  and  $g \in C_c(\mathbb{R}, \mathcal{S}(\mathbb{R}^n))$ . Then we have for example

$$\|\partial_j T_S(t)f\|_{L^r(L^p)} = \|T_S(t)\partial_j f\|_{L^r(L^p)} \leq C \|\partial_j f\|_2$$

by applying (11.20) to  $\partial_j f$ . Combining the estimates for  $f$  and its derivatives gives (11.25). Similarly for the other estimates.  $\square$

**Problem 11.2.** Does the translation group  $T(t)g(x) := g(x-t)$  satisfy (11.14)?

**Problem 11.3.** Let  $u(t) := T_S(t)g$  for some  $g \in L^1(\mathbb{R}^n)$ . Show that  $u \in C(\mathbb{R} \setminus \{0\}, C_0(\mathbb{R}^n))$ . (Hint: Lemma 4.30 (iv).)

**Problem 11.4.** Prove that there is no triple  $p, q, t$  with  $1 \leq q < p < \infty$ ,  $t \in \mathbb{R}$  such that

$$\|T_S(t)g\|_q \leq C \|g\|_p.$$

(Hint: The translation operator  $T_a f(x) := f(x-a)$  commutes with  $T_S(t)$ . Moreover, we have

$$\lim_{|a| \rightarrow \infty} \|f + T_a f\|_p = 2^{1/p} \|f\|_p, \quad 1 \leq p < \infty.$$

Now apply this to the claimed estimate.)

### 11.3. Well-posedness in $L^2$ and $H^1$

The main obstacle to proving a local existence result in  $L^2$  is the fact that our nonlinearity does not map  $L^2$  to  $L^2$  (and this was precisely the reason for choosing  $H^r$  in the previous section). On the other hand, the time evolution conserves the  $L^2$  norm and hence we expect global solutions in this case.

So let us make two observations: First of all our nonlinearity  $F(u) = |u|^{\alpha-1}u$  maps  $L^p$  to  $L^{p/\alpha}$ , so the only chance is that the linear time evolution improves this behavior. Now we know, since our evolution is unitary, there is no hope to get this for fixed  $t$ , but this is true in some averaged sense by the Strichartz estimate (11.20). Hence, if we add such a space-time norm to the  $L^2$  norm, we might be able to control our singularity. In fact, the estimates (11.23) and (11.24) allow us to control the Duhamel part in (11.7) both in the  $L^2$  and the space-time norm, respectively (the linear part being taken care of by unitarity and (11.20)). Since the spatial parts of the space-time norms must match up, we need  $p'\alpha = p$ , that is,  $p = 1 + \alpha$ . For the

time part an inequality  $r'\alpha \leq r$  is sufficient since in this case  $L^{r'\alpha} \subseteq L^r$  by Hölder's inequality. This imposes the restriction  $\alpha \leq 1 + \frac{4}{n}$ . In fact, we will impose a strict inequality since we will use the contribution from Hölder's inequality to get a contraction. Moreover, note that the dependence on the initial condition  $g$  is controlled by the  $L^2$  norm alone and this will imply that our contraction is uniform (in fact Lipschitz on bounded domains) with respect to the initial condition in  $L^2$ , and so will be the solution.

**Theorem 11.10.** *Suppose  $1 < \alpha < 1 + \frac{4}{n}$  and consider the Banach space*

$$X := C([-t_0, t_0], L^2(\mathbb{R}^n)) \cap L^r([-t_0, t_0], L^{\alpha+1}(\mathbb{R}^n)), \quad r = \frac{4(\alpha+1)}{n(\alpha-1)}, \quad (11.30)$$

with norm

$$\|f\| := \sup_{t \in [-t_0, t_0]} \|f(t)\|_2 + \left( \int_{-t_0}^{t_0} \|f(t)\|_{\alpha+1}^r dt \right)^{1/r}. \quad (11.31)$$

Then for every  $g \in L^2(\mathbb{R}^n)$  there is a  $t_0 = t_0(\|g\|_2) > 0$ , such that there is a unique solution  $u \in X$  of (11.7). Moreover, the solution map  $g \mapsto u(t)$  will be Lipschitz continuous from every ball  $\|g\|_2 \leq \rho$  to  $X$  defined with  $t_0(\rho)$ .

**Proof.** We take  $[0, t_0]$  as an interval for notational simplicity. We will show that (11.7) gives rise to a contraction on the closed ball  $\bar{B}_a(0) \subset X$  provided  $a$  and  $t_0$  are chosen accordingly. Denote the right-hand side of (11.7) by  $K(u) \equiv K_g(u)$ . We will first show that  $K : \bar{B}_a(0) \rightarrow \bar{B}_a(0)$  for a suitable  $a$  depending on  $\|g\|_2$ . To this end we first invoke (11.20), (11.23), and (11.24) with  $p = \alpha + 1$  ( $p' = \frac{\alpha+1}{\alpha}$ ) to obtain

$$\begin{aligned} \|K(u)\| &\leq (1+C)\|g\|_2 + 2C \left( \int_0^{t_0} \| |u|^\alpha(t) \|_{(\alpha+1)/\alpha}^{r'} dt \right)^{1/r'} \\ &\leq (1+C)\|g\|_2 + 2C \left( \int_0^{t_0} \|u(t)\|_{\alpha+1}^{\alpha r'} dt \right)^{1/r'}. \end{aligned}$$

Next, since  $\frac{1}{r'} = \theta + \frac{\alpha-1}{r} + \frac{1}{r}$ , where  $\theta = 1 - \frac{\alpha+1}{r} = 1 - \frac{n(\alpha-1)}{4} > 0$  we can use the generalized Hölder inequality in the form

$$\|1 \cdot f_1^{\alpha-1} f_2\|_{r'} \leq \|1\|_{1/\theta} \|f_1^{\alpha-1}\|_{r/(\alpha-1)} \|f_2\|_r = t_0^\theta \|f_1\|_r^{\alpha-1} \|f_2\|_r$$

(with  $f_1(t) = f_2(t) = \|u(t)\|_{\alpha+1}$ ) to obtain

$$\begin{aligned} \|K(u)\| &\leq (1+C)\|g\|_2 + 2C t_0^\theta \left( \int_0^{t_0} \|u(t)\|_{\alpha+1}^r dt \right)^{\alpha/r} \\ &\leq (1+C)\|g\|_2 + 2C t_0^\theta a^\alpha \end{aligned}$$

for  $u \in \bar{B}_a(0)$ . Now we choose  $a = (2 + C)\|g\|_2$  and  $2C(2 + C)t_0^\theta a^{\alpha-1} < 1$  such that

$$\|K(u)\| \leq (1 + C)\|g\|_2 + 2Ct_0^\theta(2 + C)^\alpha\|g\|_2^\alpha < (2 + C)\|g\|_2 = a.$$

Similarly we can show that  $K$  is a contraction. Invoking (11.23) and (11.24) we have

$$\|K(u) - K(v)\| \leq 2C \left( \int_0^{t_0} \| |u(t)|^{\alpha-1}u(t) - |v(t)|^{\alpha-1}v(t) \|_{(\alpha+1)/\alpha}^{r'} dt \right)^{1/r'}$$

Now using (Problem 11.5)

$$||u|^{\alpha-1}u - |v|^{\alpha-1}v| \leq \alpha(|u|^{\alpha-1} + |v|^{\alpha-1})|u - v|, \quad u, v \in \mathbb{C},$$

and invoking the generalized Hölder inequality in the form

$$\||u|^{\alpha-1}|u - v|\|_{(\alpha+1)/\alpha} \leq \| |u|^{\alpha-1} \|_{(\alpha+1)/(\alpha-1)} \|u - v\|_{\alpha+1} = \|u\|_{\alpha+1}^{\alpha-1} \|u - v\|_{\alpha+1}$$

and then in the previous form with  $f_1 = \|u\|_{\alpha+1}$ ,  $f_2 = \|u - v\|_{\alpha+1}$ , we obtain

$$\begin{aligned} \|K(u) - K(v)\| &\leq 2\alpha C \left( \int_0^{t_0} ((\|u\|_{\alpha+1}^{\alpha-1} + \|v\|_{\alpha+1}^{\alpha-1}) \|u - v\|_{\alpha+1})^{r'} dt \right)^{1/r'} \\ &\leq 2\alpha C t_0^\theta 2a^{\alpha-1} \left( \int_0^{t_0} \|u - v\|_{\alpha+1}^r dt \right)^{1/r} \\ &\leq 4\alpha C t_0^\theta a^{\alpha-1} \|u - v\|. \end{aligned}$$

Hence, decreasing  $t_0$  further (if necessary), such that we also have  $4\alpha C t_0^\theta a^{\alpha-1} < 1$ , we get a contraction. Moreover, since  $\|K_g(u) - K_f(u)\| = \|K_{g-f}(0)\| \leq (1 + C)\|g - f\|_2$ , the uniform contraction principle establishes the theorem.  $\square$

By interpolation (Problem 11.7) we also have:

**Corollary 11.11.** *The solution  $u$  is also in*

$$L^{r/\theta}([-t_0, t_0], L^{2(\alpha+1)/(\alpha+1-\theta(\alpha-1))}(\mathbb{R}^n))$$

for any  $\theta \in (0, 1)$ .

Moreover, as in the previous section we obtain:

**Corollary 11.12.** *The maximal solution  $u$  is global in  $C(\mathbb{R}, L^2(\mathbb{R}^n))$  and preserves the  $L^2$  norm:  $\|u(t)\|_2 = \|g\|_2$ . In addition, it has the properties stated in the theorem for any  $t_0 > 0$ .*

Let me remark that it is possible to cover the case  $\alpha = 1 + \frac{4}{n}$ . The main difference is that the Hölder-type estimate in terms of  $t^\theta$  for the integral in (11.7) is useless since  $\theta = 0$ . However, the integral still tends to zero as  $t \rightarrow 0$ . This will be true locally in a sufficiently small neighborhood, but we cannot control this neighborhood in terms of  $\|g\|_2$ .

However, we will turn to the case of initial conditions in  $H^1$  instead.

**Theorem 11.13.** *Suppose  $n \geq 3$  and  $2 \leq \alpha < \frac{n+2}{n-2}$ . Consider the Banach space*

$$X := C([-t_0, t_0], H^1(\mathbb{R}^n)) \cap L^r([-t_0, t_0], W^{1,p}(\mathbb{R}^n)), \quad (11.32)$$

where

$$p = \frac{n(\alpha+1)}{n+\alpha-1}, \quad r = \frac{4(\alpha+1)}{(n-2)(\alpha-1)}, \quad (11.33)$$

with norm

$$\|f\| := \sup_{t \in [-t_0, t_0]} \|f(t)\|_{1,2} + \left( \int_{-t_0}^{t_0} \|f(t)\|_{1,p}^r dt \right)^{1/r}. \quad (11.34)$$

Then for every  $g \in H^1(\mathbb{R}^n)$  there is a  $t_0 = t_0(\|g\|_{1,2}) > 0$ , such that there is a unique solution  $u \in X$  of (11.7). Moreover, the solution map  $g \mapsto u(t)$  will be Lipschitz continuous from every ball  $\|g\|_{1,2} \leq \rho$  to  $X$  defined with  $t_0(\rho)$ .

**Proof.** We begin with estimating the nonlinearity. For  $u, v \in W^{1,p}$  and  $w \in L^p$  we obtain

$$\| |u|^{\alpha-2} v w \|_{p'} \leq \|u\|_q^{\alpha-2} \|v\|_q \|w\|_p \leq C \|\partial u\|_p^{\alpha-2} \|\partial v\|_p \|w\|_p,$$

where we have applied the generalized Hölder inequality with  $\frac{1}{p'} = \frac{\alpha-2}{q} + \frac{1}{q} + \frac{1}{p}$  in the first step (requiring  $\alpha \geq 2$ ) and the Gagliardo–Nirenberg–Sobolev inequality (Theorem 7.26 from [48] – since we need  $p < n$ , we need to require  $n > 2$ ) with  $\frac{1}{q} = \frac{1}{p} - \frac{1}{n}$  in the second step. In particular, this imposes

$$1 - \frac{2}{p} = \frac{\alpha-1}{q} = \frac{\alpha-1}{p} - \frac{\alpha-1}{n}$$

and explains our choice for  $p$ . The choice of  $r$  is of course dictated by (11.18) such that we can apply our Strichartz estimates. At this point a weaker upper bound (namely  $\alpha < \frac{n}{n-4}$  for  $n \geq 4$ ) is still sufficient.

Now using this estimate we see (cf. Problem 11.5)

$$\| |u|^{\alpha-1} u \|_{p'} \leq C \|\partial u\|_p^{\alpha-1} \|u\|_p, \quad \|\partial |u|^{\alpha-1} u\|_{p'} \leq \alpha \| |u|^{\alpha-1} \partial u \|_{p'} \leq \alpha C \|\partial u\|_p^\alpha$$

and hence

$$\| |u|^{\alpha-1} u \|_{1,p'} \leq \tilde{C} \|u\|_{1,p}^\alpha.$$

Similarly we obtain

$$\begin{aligned} \| |u|^{\alpha-1} u - |v|^{\alpha-1} v \|_{p'} &\leq \alpha \left( \| |u|^{\alpha-1} \| + \| |v|^{\alpha-1} \| \right) \|u - v\|_{p'} \\ &\leq \alpha C \left( \|\partial u\|_p^{\alpha-1} + \|\partial v\|_p^{\alpha-1} \right) \|u - v\|_p \end{aligned}$$

and

$$\begin{aligned}
& \|\partial|u|^{\alpha-1}u - \partial|v|^{\alpha-1}v\|_{p'} \\
& \leq (\alpha-1)(\alpha+2)\|(|u|^{\alpha-2} + |v|^{\alpha-2})|u-v|\partial u\|_{p'} \\
& \quad + \alpha\||v|^{\alpha-1}\partial u - \partial v\|_{p'} \\
& \leq (\alpha-1)(\alpha+2)C(\|\partial u\|_p^{\alpha-2} + \|\partial v\|_p^{\alpha-2})\|\partial(u-v)\|_p\|\partial u\|_p \\
& \quad + \alpha C\|\partial v\|_p^{\alpha-1}\|\partial u - \partial v\|_p.
\end{aligned}$$

In summary,

$$\||u|^{\alpha-1}u - |v|^{\alpha-1}v\|_{1,p'} \leq \bar{C}(\|u\|_{1,p}^{\alpha-1} + \|v\|_{1,p}^{\alpha-1})\|u-v\|_{1,p}.$$

Now the rest follows as in the proof of Theorem 11.10. Note that in this case  $\theta = 1 - \frac{\alpha+1}{r} = \frac{2+n+(2-n)\alpha}{4}$  explaining our upper limit for  $\alpha$ .  $\square$

Note that since we have  $H^1(\mathbb{R}^n) \subseteq L^{\alpha+1}(\mathbb{R}^n)$  for  $n \geq 3$  and  $\alpha < \frac{n+2}{n-2}$  by the Gagliardo–Nirenberg–Sobolev inequality (Theorem 7.26 from [48]), both the momentum and the energy are finite and preserved by our solutions. Moreover, in the defocusing case the momentum and the energy control the  $H^1$  norm and hence we obtain:

**Corollary 11.14.** *In the defocusing case the maximal solution  $u$  is global in  $C(\mathbb{R}, H^1(\mathbb{R}^n))$  and preserves both momentum and energy. In addition, it has the properties stated in the theorem for any  $t_0 > 0$ .*

In the focusing case we need to control the  $L^{\alpha+1}$  norm in terms of the  $H^1$  norm using the Gagliardo–Nirenberg–Sobolev inequality.

**Corollary 11.15.** *In the focusing case the maximal solution  $u$  is global in  $C(\mathbb{R}, H^1(\mathbb{R}^n))$  and preserves both momentum and energy if one of the following conditions hold:*

- (i)  $\alpha < 1 + \frac{4}{n}$ .
- (ii)  $\alpha = 1 + \frac{4}{n}$  and  $\|g\|_2 < (\frac{2(n-1)}{(n+2)(n-2)})^{n/4}$ .
- (iii)  $\alpha > 1 + \frac{4}{n}$  and  $\|g\|_{1,2}$  is sufficiently small such that  $\|\partial g\|_2 < 1$  and  $2E(0) + \frac{4(n-1)}{(n(n-2)(\alpha+1))}\|g\|_2^{\alpha+1-n(\alpha-1)/2} < 1$ .

**Proof.** Using the Gagliardo–Nirenberg–Sobolev inequality and the Lyapunov inequality (Problem 3.12 from [48]) with  $\frac{1}{1+\alpha} = \theta(\frac{1}{2} - \frac{1}{n}) + \frac{1-\theta}{2}$  (i.e.  $\theta = \frac{n(\alpha-1)}{2(\alpha+1)}$ ) we obtain

$$\|u(t)\|_{\alpha+1}^{\alpha+1} \leq \frac{2(n-1)}{n(n-2)}\|u(t)\|_2^{\alpha+1-n(\alpha-1)/2}\|\partial u(t)\|_2^{n(\alpha-1)/2}. \quad (11.35)$$

Thus

$$\begin{aligned}\|\partial u(t)\|_2^2 &= 2E(0) + \frac{2}{\alpha+1} \|u(t)\|_{\alpha+1}^{\alpha+1} \\ &\leq 2E(0) + 2C\|g\|_2^{\alpha+1-n(\alpha-1)/2} \|\partial u(t)\|_2^{n(\alpha-1)/2},\end{aligned}\quad (11.36)$$

where we have set  $C := \frac{2(n-1)}{n(n-2)(\alpha+1)}$ .

(i). Now if  $\alpha < 1 + \frac{4}{n}$ , then  $\frac{n(\alpha-1)}{2} < 2$  and  $\|\partial u(t)\|_2$  remains bounded.

(ii). In the case  $\alpha = 1 + \frac{4}{n}$  this remains still true if  $2C\|g\|_2^{4/n} < 1$ .

(iii). If  $\alpha > 1 + \frac{4}{n}$  we can choose  $\|g\|_{1,2}$  so small such that the given conditions hold. Note that this is possible since our above calculation shows

$$E(0) \leq \frac{1}{2} \|\partial g\|_2^2 + C\|g\|_2^{\alpha+1-n(\alpha-1)/2} \|\partial g\|_2^{n(\alpha-1)/2}.$$

Now if we start with  $\|\partial u(0)\|_2^2 \leq 1$  and assume  $\|\partial u(t)\|_2^2 = 1$  we get the contradiction  $1 = \|\partial u(t)\|_2^2 \leq 2E(0) + 2C\|g\|_2^{\alpha+1-n(\alpha-1)/2} < 1$ . Hence  $\|\partial u(t)\|_2^2 < 1$  as desired.  $\square$

**Problem 11.5.** Show that the real derivative (with respect to the identification  $\mathbb{C} \cong \mathbb{R}^2$ ) of  $F(u) = |u|^{\alpha-1}u$  is given by

$$F'(u)v = |u|^{\alpha-1}v + (\alpha-1)|u|^{\alpha-3}u\operatorname{Re}(u^*v).$$

Conclude in particular,

$$|F'(u)v| \leq \alpha|u|^{\alpha-1}|v|, \quad |F(u) - F(v)| \leq \alpha(|u|^{\alpha-1} + |v|^{\alpha-1})|u - v|.$$

Moreover, the second derivative is given by

$$vF''(u)w = (\alpha-1)|u|^{\alpha-5}u((\alpha+1)\operatorname{Re}(u^*v)\operatorname{Re}(u^*w) - u^2v^*w^*).$$

and hence

$$|vF''(u)w| \leq (\alpha-1)(\alpha+2)|u|^{\alpha-2}|v||w|.$$

**Problem 11.6.** Show that (11.30) is a Banach space. (Hint: Work with test functions from  $C_c^\infty$ .)

**Problem 11.7.** Suppose  $f \in L^{p_0}(I, L^{q_0}(U)) \cap L^{p_1}(I, L^{q_1}(U))$ . Show that  $f \in L^{p_\theta}(I, L^{q_\theta}(U))$  for  $\theta \in [0, 1]$ , where

$$\frac{1}{p_\theta} = \frac{1-\theta}{p_0} + \frac{\theta}{p_1}, \quad \frac{1}{q_\theta} = \frac{1-\theta}{q_0} + \frac{\theta}{q_1}.$$

(Hint: Lyapunov and generalized Hölder inequality — Problem 3.12 from [48] and Problem 3.9 from [48].)

### 11.4. Blowup in $H^1$

In this section we will show that solutions are not always global in the focusing case. For simplicity we will only consider the one-dimensional case. We first complement Theorem 11.13 with a result for the one-dimensional case.

**Theorem 11.16.** *Let  $n = 1$  and  $\alpha \geq 2$ . For every  $g \in H^1(\mathbb{R})$  there is a  $t_0 = t_0(\|g\|_{1,2}) > 0$ , such that there is a unique solution  $u \in C([-t_0, t_0], H^1(\mathbb{R}))$  of (11.7). Moreover, the solution map  $g \mapsto u(t)$  will be Lipschitz continuous from every ball  $\|g\|_{1,2} \leq \rho$  to  $C([-t_0(\rho), t_0(\rho)], H^1(\mathbb{R}))$ .*

**Proof.** It suffices to verify that  $F : H^1(\mathbb{R}) \rightarrow H^1(\mathbb{R})$  is locally Lipschitz on bounded sets. But this follows using Problem 11.5 since

$$\|F(u) - F(v)\|_2 \leq \alpha(\|u\|_\infty^{\alpha-1} + \|v\|_\infty^{\alpha-1})\|u - v\|_2$$

and

$$\begin{aligned} \|\partial(F(u) - F(v))\|_2 &\leq (\alpha - 1)(\alpha + 2)(\|u\|_\infty^{\alpha-2} + \|v\|_\infty^{\alpha-2})\|\partial u\|_2\|u - v\|_2 \\ &\quad + \alpha\|v\|_\infty^{\alpha-1}\|\partial(u - v)\|_2 \end{aligned}$$

together with  $\|f\|_\infty \leq \|f\|_{1,2}$  (Problem 11.8).  $\square$

In the one-dimensional case we did not need a space-time norm to establish the above theorem. However, it is worthwhile mentioning that they still hold.

**Corollary 11.17.** *Let  $u$  be the solution for given  $g \in H^1$ . Then*

$$u \in L^r([-t_0, t_0], L^p(\mathbb{R}^n)) \quad (11.37)$$

*whenever  $(p, r)$  is an admissible pair satisfying (11.18).*

**Proof.** Just observe that the solution satisfies

$$\|F(u)\|_{L^{r'}(L^{p'})} \leq (2t_0)^{1/r'} \sup_{|t| \leq t_0} \|u\|_\infty^{\alpha - \frac{2}{p'}} \|u\|_2^{\frac{2}{p'}} \leq (2t_0)^{1/r'} \sup_{|t| \leq t_0} \|u\|_{1,2}^\alpha$$

and apply the Strichartz estimates (11.20) and (11.24).  $\square$

As in the previous section one gets global existence in the defocusing case and, with a little more effort, also in the focusing case if  $\alpha < 5$  (Problem 11.9).

**Lemma 11.18.** *Let  $n = 1$  and  $\alpha \geq 2$ . Suppose  $g \in H^1(\mathbb{R})$  satisfies  $\|xg(x)\|_2 < \infty$  and let  $u \in C((t_-, t_+), H^1(\mathbb{R}))$  be the maximal solution of (11.7). Then*

$$M_1(t) := \int_{\mathbb{R}} x^2 |u(t, x)|^2 dx \quad (11.38)$$

remains finite as long as  $u$  exists and satisfies

$$\dot{M}_1(t) = 4\text{Im} \int_{\mathbb{R}} x u(t, x)^* \partial u(t, x) dx, \quad (11.39)$$

$$\ddot{M}_1(t) = 16E(t) \pm \frac{4(\alpha - 5)}{\alpha + 1} \int_{\mathbb{R}} |u(t, x)|^{\alpha+1} dx, \quad (11.40)$$

known as **virial** and **Morawetz identity**, respectively.

**Proof.** Consider  $H^{1,1}(\mathbb{R}) := H^1(\mathbb{R}) \cap L^2(\mathbb{R}, x^2 dx)$  together with the norm  $\|f\|^2 = \|f\|_2^2 + \|f'\|_2^2 + \|xf(x)\|_2^2$ . Then  $T_S(t)$  is a  $C^0$  group satisfying  $\|T_S(t)f\| \leq (1 + 2|t|)\|f\|$ . Moreover, as in the previous theorem one verifies that  $F : H^{1,1}(\mathbb{R}) \rightarrow H^{1,1}(\mathbb{R})$  is locally Lipschitz on bounded sets. In fact, note that by

$$\|x(F(u)(x) - F(v)(x))\|_2 \leq \alpha(\|u\|_{\infty}^{\alpha-1} + \|v\|_{\infty}^{\alpha-1})\|x(u(x) - v(x))\|_2$$

the Lipschitz constant depends only on the  $H^1$  norm. Hence we get existence of local solutions. Moreover, using (11.7) we obtain

$$\|u(t)\| \leq (1 + 2t)\|g\| + \alpha \int_0^t (1 + 2(t-s))\|u(s)\|_{1,2}^{\alpha-1}\|u(s)\| ds,$$

and Gronwall's inequality

$$\|u(t)\| \leq (1 + 2t)\|g\| \exp \left( \int_0^t (1 + 2(t-s))\|u(s)\|_{1,2}^{\alpha-1} ds \right)$$

shows that our norm cannot blow up before the  $H^1$  norm.

The formulas for  $\dot{M}_1$  and  $\ddot{M}_1$  are straightforward computations.  $\square$

Now we are ready to establish blowup for the focusing NLS equation.

**Theorem 11.19.** *Consider the one-dimensional focusing NLS equation with  $\alpha \geq 5$ . Let  $g \in H^1(\mathbb{R}) \cap L^2(\mathbb{R}, x^2 dx)$  with negative energy  $E < 0$ . Then the corresponding maximal mild solution  $u$  satisfies  $t_+(g) < \infty$ .*

**Proof.** Due to our assumption  $\alpha \geq 5$  we obtain  $\ddot{M}_1(t) \leq 16E$  implying  $M_1(t) \leq 8Et^2 + \dot{M}_1(0)t + M_1(0)$ . Hence

$$t_+(g) < \frac{-1}{16E} \left( \dot{M}_1(0) + \sqrt{\dot{M}_1(0)^2 - 32EM_1(0)} \right)$$

since  $M_1(t)$  must remain positive. Note that this also shows  $\dot{M}_1(0) > 0$  since otherwise  $M_1(t)$  would be decreasing and hence would remain bounded.  $\square$

Notice that there are initial conditions with negative energy, since the two contributions to the energy scale differently. In particular, the energy will become negative if we scale  $g$  with a sufficiently large factor.



**Problem 11.8.** Let  $f \in H^1(\mathbb{R})$ . Show  $\|f\|_\infty^2 \leq 2\|f\|_2\|f'\|_2$  and hence  $\|f\|_\infty \leq \|f\|_{1,2}$ .

**Problem 11.9.** Show that the one-dimensional focusing NLS equation has global solutions in  $H^1(\mathbb{R})$  if either  $\alpha < 5$  or  $\alpha = 5$  and  $\|g\|_2 \leq (\frac{3}{4})^{1/4}$  or  $\alpha > 5$  and  $\|g\|_{1,2}$  sufficiently small. (Hint: Use the estimate from Problem 11.8.)

### 11.5. Standing waves

A solution of the form

$$u(x, t) = \varphi_\omega(x)e^{i\omega t}, \quad \omega > 0, \quad (11.41)$$

of the focusing NLS equation is called a **standing wave**. Inserting this ansatz into the equation shows that  $\varphi_\omega$  must be a solution of the following nonlinear elliptic problem

$$-\Delta\varphi_\omega + \omega\varphi_\omega = |\varphi_\omega|^{\alpha-1}\varphi_\omega. \quad (11.42)$$

Note that one can choose  $\omega = 1$  without loss of generality since if  $\varphi$  is a solution for  $\omega = 1$  then

$$\varphi_\omega(x) = \omega^{\frac{1}{\alpha-1}}\varphi(\omega^{1/2}x) \quad (11.43)$$

is a solution for  $\omega > 0$ . Moreover, if  $\varphi$  is a solution, so is  $e^{i\theta}\varphi(\cdot - a)$  for any  $\theta \in \mathbb{R}$  and  $a \in \mathbb{R}^n$ .

If one multiplies (11.42) with a test function  $v \in H^1(\mathbb{R}^n)$  and integrates over  $\mathbb{R}^n$  one obtains the weak formulation

$$\int_{\mathbb{R}^n} (\partial\varphi \cdot \partial v + \varphi v - |\varphi|^{\alpha-1}\varphi v) d^n x = 0, \quad v \in H^1(\mathbb{R}^n). \quad (11.44)$$

In particular, choosing  $v = \varphi^*$  we obtain

$$\int_{\mathbb{R}^n} (|\partial\varphi|^2 + |\varphi|^2 - |\varphi|^{\alpha+1}) d^n x = 0, \quad (11.45)$$

which shows that, if we flip the sign in front of the nonlinearity (defocusing case), there is only the trivial solution.

In one-dimension one has the explicit solution

$$\varphi(x) = \left( \frac{\sqrt{1+\beta}}{\cosh(\beta x)} \right)^{1/\beta}, \quad \beta = \frac{\alpha-1}{2}. \quad (11.46)$$

In higher dimensions we can apply Theorem 9.26 to get existence of solutions:

**Theorem 11.20.** Suppose  $n \geq 2$  and  $1 < \alpha < \frac{n+2}{n-2}$ . Then the nonlinear elliptic problem (11.42) has a weak positive radial solution in  $H^1(\mathbb{R}^n)$ .

**Proof.** To apply Theorem 9.26 we choose  $X = H_{\text{rad}}^1(\mathbb{R}^n, \mathbb{R})$  and  $Y = L_{\text{rad}}^{\alpha+1}(\mathbb{R}^n, \mathbb{R})$  and note that the Strauss lemma (Problem 7.33 from [48]) implies compactness of the embedding  $X \hookrightarrow Y$  for the range of  $\alpha$  under consideration. Hence minimizing

$$F(u) = \frac{1}{2} \int_{\mathbb{R}^n} (|\partial u|^2 + |u|^2) d^n x$$

under the constraint (cf. Example 9.5)

$$N(u) = \frac{1}{\alpha + 1} \int_{\mathbb{R}^n} |u|^{\alpha+1} d^n x = 1$$

gives a weak radial solution  $u_0$  of the problem

$$-\Delta u + u = \lambda |u|^{\alpha-1} u.$$

In particular, choosing  $u_0$  as a test function for the weak formulation shows  $\lambda > 0$ . Moreover, by Lemma 7.12 from [48] we have  $|u_0| \in H_{\text{rad}}^1(\mathbb{R}^n)$  with  $F(|u_0|) = F(u_0)$  and hence  $|u_0|$  is also a minimizer. Rescaling this solution according to  $\varphi(x) = \lambda^{1/(\alpha-1)} |u_0(x)|$  establishes the claim.  $\square$

Note that for  $\alpha \leq \frac{n}{n-2}$  we have  $|u|^{\alpha-1} u \in L^2(\mathbb{R}^n)$  for  $u \in H^1(\mathbb{R}^n)$  and hence  $(-\Delta + 1)\varphi \in L^2(\mathbb{R}^n)$  implying  $\varphi \in H^2(\mathbb{R}^n)$ .

**Problem 11.10.** Let  $1 < \alpha < \frac{n+2}{n-2}$  be an odd integer (i.e.  $n = 2$  and  $\alpha = 3, 5, 6, \dots$  or  $n = 3$  and  $\alpha = 3$ ). Show that  $\varphi \in H^k(\mathbb{R}^n)$  for any  $k \in \mathbb{N}$ . (Hint: As already pointed out we have  $\varphi \in H^2$ .)



# The Brouwer mapping degree

## 12.1. Introduction

Many applications lead to the problem of finding zeros of a mapping  $f : U \subseteq X \rightarrow X$ , where  $X$  is some (real) Banach space. That is, we are interested in the solutions of

$$f(x) = 0, \quad x \in U. \quad (12.1)$$

In most cases it turns out that this is too much to ask for, since determining the zeros analytically is in general impossible.

Hence one has to ask some weaker questions and hope to find answers for them. One such question would be “Are there any solutions, respectively, how many are there?”. Luckily, these questions allow some progress.

To see how, let's consider the case  $f \in \mathcal{H}(\mathbb{C})$ , where  $\mathcal{H}(U)$  denotes the set of **holomorphic functions** on a domain  $U \subset \mathbb{C}$ . Recall the concept of the **winding number** from complex analysis. The winding number of a path  $\gamma : [0, 1] \rightarrow \mathbb{C} \setminus \{z_0\}$  around a point  $z_0 \in \mathbb{C}$  is defined by

$$n(\gamma, z_0) := \frac{1}{2\pi i} \int_{\gamma} \frac{dz}{z - z_0} \in \mathbb{Z}. \quad (12.2)$$

It gives the number of times  $\gamma$  encircles  $z_0$  taking orientation into account. That is, encirclings in opposite directions are counted with opposite signs.

In particular, if we pick  $f \in \mathcal{H}(\mathbb{C})$  one computes (assuming  $0 \notin f(\gamma)$ )

$$n(f(\gamma), 0) = \frac{1}{2\pi i} \int_{\gamma} \frac{f'(z)}{f(z)} dz = \sum_k n(\gamma, z_k) \alpha_k, \quad (12.3)$$

where  $z_k$  denote the zeros of  $f$  and  $\alpha_k$  their respective multiplicity. Moreover, if  $\gamma$  is a Jordan curve encircling a simply connected domain  $U \subset \mathbb{C}$ , then  $n(\gamma, z_k) = 0$  if  $z_k \notin \overline{U}$  and  $n(\gamma, z_k) = 1$  if  $z_k \in U$ . Hence  $n(f(\gamma), 0)$  counts the number of zeros inside  $U$ .

However, this result is useless unless we have an efficient way of computing  $n(f(\gamma), 0)$  (which does not involve the knowledge of the zeros  $z_k$ ). This is our next task.

Now, let's recall how one would compute complex integrals along complicated paths. Clearly, one would use homotopy invariance and look for a simpler path along which the integral can be computed and which is homotopic to the original one. In particular, if  $f : \gamma \rightarrow \mathbb{C} \setminus \{0\}$  and  $g : \gamma \rightarrow \mathbb{C} \setminus \{0\}$  are homotopic, we have  $n(f(\gamma), 0) = n(g(\gamma), 0)$  (which is known as Rouché's theorem).

More explicitly, we need to find a mapping  $g$  for which  $n(g(\gamma), 0)$  can be computed and a **homotopy**  $H : [0, 1] \times \gamma \rightarrow \mathbb{C} \setminus \{0\}$  such that  $H(0, z) = f(z)$  and  $H(1, z) = g(z)$  for  $z \in \gamma$ . For example, how many zeros of  $f(z) = \frac{1}{2}z^6 + z - \frac{1}{3}$  lie inside the unit circle? Consider  $g(z) = z$ , then  $H(t, z) = (1-t)f(z) + tg(z)$  is the required homotopy since  $|f(z) - g(z)| < |g(z)|$ ,  $|z| = 1$ , implying  $H(t, z) \neq 0$  on  $[0, 1] \times \gamma$ . Hence  $f(z)$  has one zero inside the unit circle.

Summarizing, given a (sufficiently smooth) domain  $U$  with enclosing Jordan curve  $\partial U$ , we have defined a degree  $\deg(f, U, z_0) = n(f(\partial U), z_0) = n(f(\partial U) - z_0, 0) \in \mathbb{Z}$  which counts the number of solutions of  $f(z) = z_0$  inside  $U$ . The invariance of this degree with respect to certain deformations of  $f$  allowed us to explicitly compute  $\deg(f, U, z_0)$  even in nontrivial cases.

Our ultimate goal is to extend this approach to continuous functions  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . However, such a generalization runs into several problems. First of all, it is unclear how one should define the multiplicity of a zero. But even more severe is the fact, that the number of zeros is unstable with respect to small perturbations. For example, consider  $f_\varepsilon : [-1, 2] \rightarrow \mathbb{R}$ ,  $x \mapsto x^2 - \varepsilon$ . Then  $f_\varepsilon$  has no zeros for  $\varepsilon < 0$ , one zero for  $\varepsilon = 0$ , two zeros for  $0 < \varepsilon \leq 1$ , one for  $1 < \varepsilon \leq \sqrt{2}$ , and none for  $\varepsilon > \sqrt{2}$ . This shows the following facts.

- (i) Zeros with  $f' \neq 0$  are stable under small perturbations.
- (ii) The number of zeros can change if two zeros with opposite sign change (i.e., opposite signs of  $f'$ ) run into each other.
- (iii) The number of zeros can change if a zero drops over the boundary.

Hence we see that we cannot expect too much from our degree. In addition, since it is unclear how it should be defined, we will first require some basic properties a degree should have and then we will look for functions satisfying these properties.

## 12.2. Definition of the mapping degree and the determinant formula

To begin with, let us introduce some useful notation. Throughout this section  $U$  will be a bounded open subset of  $\mathbb{R}^n$ . For  $f \in C^1(U, \mathbb{R}^n)$  the Jacobi matrix of  $f$  at  $x \in U$  is  $df(x) = (\partial_{x_j} f_i(x))_{1 \leq i, j \leq n}$  and the Jacobi determinant of  $f$  at  $x \in U$  is

$$J_f(x) := \det df(x). \quad (12.4)$$

The set of **regular values** is

$$\text{RV}(f) := \{y \in \mathbb{R}^n \mid \forall x \in f^{-1}(y) : J_f(x) \neq 0\}. \quad (12.5)$$

Its complement  $\text{CV}(f) := \mathbb{R}^n \setminus \text{RV}(f)$  is called the set of **critical values**. We will also need the spaces

$$\bar{C}^k(U, \mathbb{R}^n) := C^k(U, \mathbb{R}^n) \cap C(\bar{U}, \mathbb{R}^n) \quad (12.6)$$

and regard them as subspaces of the Banach space  $C(\bar{U}, \mathbb{R}^n)$  (cf. Section 9.2). Note that  $\bar{C}^\infty(U, \mathbb{R}^n)$  is dense in  $C(\bar{U}, \mathbb{R}^n)$ . To see this you can either apply Stone–Weierstraß or use the Tietze extension theorem to extend  $f \in C(\bar{U}, \mathbb{R}^n)$  to all of  $\mathbb{R}^n$  and then mollify. If you use mollification and  $f \in \bar{C}^k(U, \mathbb{R}^n)$  then all derivatives up to order  $k$  will converge uniformly on compact subsets of  $U$ . Finally, for  $y \in \mathbb{R}^n$  we set

$$\bar{C}_y^k(\bar{U}, \mathbb{R}^n) := \{f \in \bar{C}^k(U, \mathbb{R}^n) \mid y \notin f(\partial U)\} \quad (12.7)$$

and  $\bar{C}_y(U, \mathbb{R}^n) := \bar{C}_y^0(U, \mathbb{R}^n)$ .

Note that, since  $U$  is bounded,  $\partial U$  is compact and so is  $f(\partial U)$  if  $f \in C(\bar{U}, \mathbb{R}^n)$ . In particular,

$$\text{dist}(y, f(\partial U)) = \min_{x \in \partial U} |y - f(x)| \quad (12.8)$$

is positive for  $f \in \bar{C}_y(U, \mathbb{R}^n)$  and thus  $\bar{C}_y(U, \mathbb{R}^n)$  is an open subset of  $C(\bar{U}, \mathbb{R}^n)$ .

Now that these things are out of the way, we come to the formulation of the requirements for our degree.

A function  $\deg$  which assigns each  $f \in \bar{C}_y(U, \mathbb{R}^n)$ ,  $y \in \mathbb{R}^n$ , a real number  $\deg(f, U, y)$  will be called degree if it satisfies the following conditions.

- (D1).  $\deg(f, U, y) = \deg(f - y, U, 0)$  (*translation invariance*).
- (D2).  $\deg(\mathbb{I}, U, y) = 1$  if  $y \in U$  (*normalization*).
- (D3). If  $U_{1,2}$  are open, disjoint subsets of  $U$  such that  $y \notin f(\bar{U} \setminus (U_1 \cup U_2))$ , then  $\deg(f, U, y) = \deg(f, U_1, y) + \deg(f, U_2, y)$  (*additivity*).
- (D4). If  $H(t) = (1 - t)f + tg \in \bar{C}_y(U, \mathbb{R}^n)$ ,  $t \in [0, 1]$ , then  $\deg(f, U, y) = \deg(g, U, y)$  (*homotopy invariance*).

Before we draw some first conclusions from this definition, let us discuss the properties (D1)–(D4) first. (D1) is natural since  $\deg(f, U, y)$  should have something to do with the solutions of  $f(x) = y$ ,  $x \in U$ , which is the same as the solutions of  $f(x) - y = 0$ ,  $x \in U$ . (D2) is a normalization since any multiple of  $\deg$  would also satisfy the other requirements. (D3) is also quite natural since it requires  $\deg$  to be additive with respect to components. In addition, it implies that sets where  $f \neq y$  do not contribute. (D4) is not that natural since it already rules out the case where  $\deg$  is the cardinality of  $f^{-1}(\{y\})$ . On the other hand it will give us the ability to compute  $\deg(f, U, y)$  in several cases.

**Theorem 12.1.** *Suppose  $\deg$  satisfies (D1)–(D4) and let  $f, g \in \bar{C}_y(U, \mathbb{R}^n)$ , then the following statements hold.*

- (i). *We have  $\deg(f, \emptyset, y) = 0$ . Moreover, if  $U_i$ ,  $1 \leq i \leq N$ , are disjoint open subsets of  $U$  such that  $y \notin f(\bar{U} \setminus \bigcup_{i=1}^N U_i)$ , then  $\deg(f, U, y) = \sum_{i=1}^N \deg(f, U_i, y)$ .*
- (ii). *If  $y \notin f(U)$ , then  $\deg(f, U, y) = 0$  (but not the other way round). Equivalently, if  $\deg(f, U, y) \neq 0$ , then  $y \in f(U)$ .*
- (iii). *If  $|f(x) - g(x)| < |f(x) - y|$ ,  $x \in \partial U$ , then  $\deg(f, U, y) = \deg(g, U, y)$ . In particular, this is true if  $f(x) = g(x)$  for  $x \in \partial U$ .*

**Proof.** For the first part of (i) use (D3) with  $U_1 = U$  and  $U_2 = \emptyset$ . For the second part use  $U_2 = \emptyset$  in (D3) if  $N = 1$  and the rest follows from induction. For (ii) use  $N = 1$  and  $U_1 = \emptyset$  in (i). For (iii) note that  $H(t, x) = (1 - t)f(x) + tg(x)$  satisfies  $|H(t, x) - y| \geq \text{dist}(y, f(\partial U)) - |f(x) - g(x)|$  for  $x$  on the boundary.  $\square$

Item (iii) is a version of **Rouché's theorem** for our degree. Next we show that (D4) implies several at first sight stronger looking facts.

**Theorem 12.2.** *We have that  $\deg(\cdot, U, y)$  and  $\deg(f, U, \cdot)$  are both continuous. In fact, we even have*

- (i).  $\deg(\cdot, U, y)$  *is constant on each component of  $\bar{C}_y(U, \mathbb{R}^n)$ .*
- (ii).  $\deg(f, U, \cdot)$  *is constant on each component of  $\mathbb{R}^n \setminus f(\partial U)$ .*

*Moreover, if  $H : [0, 1] \times \bar{U} \rightarrow \mathbb{R}^n$  and  $y : [0, 1] \rightarrow \mathbb{R}^n$  are both continuous such that  $H(t) \in \bar{C}_{y(t)}(U, \mathbb{R}^n)$ ,  $t \in [0, 1]$ , then  $\deg(H(0), U, y(0)) = \deg(H(1), U, y(1))$ .*

**Proof.** For (i) it suffices to show that  $\deg(\cdot, U, y)$  is locally constant. But if  $|g - f| < \text{dist}(y, f(\partial U))$ , then  $\deg(f, U, y) = \deg(g, U, y)$  by (D4) since  $|H(t) - y| \geq |f - y| - |g - f| > 0$ ,  $H(t) = (1 - t)f + tg$ . The proof of (ii) is similar.

For the remaining part observe, that if  $H : [0, 1] \times \bar{U} \rightarrow \mathbb{R}^n$ ,  $(t, x) \mapsto H(t, x)$ , is continuous, then so is  $H : [0, 1] \rightarrow C(\bar{U}, \mathbb{R}^n)$ ,  $t \mapsto H(t)$ , since  $\bar{U}$  is compact. Hence, if in addition  $H(t) \in \bar{C}_y(U, \mathbb{R}^n)$ , then  $\deg(H(t), U, y)$  is independent of  $t$  and if  $y = y(t)$  we can use  $\deg(H(0), U, y(0)) = \deg(H(t) - y(t), U, 0) = \deg(H(1), U, y(1))$ .  $\square$

In this context note that a Banach space  $X$  is locally path-connected and hence the components of any open subset are open (in the topology of  $X$ ) and path-connected (see Lemma B.33 (vi)).

Moreover, note that this result also shows why  $\deg(f, U, y)$  cannot be defined meaningful for  $y \in f(\partial U)$ . Indeed, approaching  $y$  from within different components of  $\mathbb{R}^n \setminus f(\partial U)$  will result in different limits in general!

Now let us try to compute  $\deg$  using its properties. If you are not interested in how to derive the determinant formula for the degree from its properties you can of course take it as a definition and skip to the next section.

Let's start with a simple case and suppose  $f \in \bar{C}_y^1(U, \mathbb{R}^n)$  and  $y \notin CV(f)$ . Without restriction we consider  $y = 0$ . In addition, we avoid the trivial case  $f^{-1}(\{0\}) = \emptyset$ . Since the points of  $f^{-1}(\{0\})$  inside  $U$  are isolated (use  $J_f(x) \neq 0$  and the inverse function theorem) they can only cluster at the boundary  $\partial U$ . But this is also impossible since  $f$  would equal 0 at the limit point on the boundary by continuity. Hence  $f^{-1}(\{0\}) = \{x^i\}_{i=1}^N$ . Picking sufficiently small neighborhoods  $U(x^i)$  around  $x^i$  we consequently get

$$\deg(f, U, 0) = \sum_{i=1}^N \deg(f, U(x^i), 0). \quad (12.9)$$

It suffices to consider one of the zeros, say  $x^1$ . Moreover, we can even assume  $x^1 = 0$  and  $U(x^1) = B_\delta(0)$ . Next we replace  $f$  by its linear approximation around 0. By the definition of the derivative we have

$$f(x) = df(0)x + |x|r(x), \quad r \in C(B_\delta(0), \mathbb{R}^n), \quad r(0) = 0. \quad (12.10)$$

Now consider the homotopy  $H(t, x) = df(0)x + (1 - t)|x|r(x)$ . In order to conclude  $\deg(f, B_\delta(0), 0) = \deg(df(0), B_\delta(0), 0)$  we need to show  $0 \notin H(t, \partial B_\delta(0))$ . Since  $J_f(0) \neq 0$  we can find a constant  $\lambda$  such that  $|df(0)x| \geq \lambda|x|$  and since  $r(0) = 0$  we can decrease  $\delta$  such that  $|r| < \lambda$ . This implies  $|H(t, x)| \geq |df(0)x| - (1 - t)|x||r(x)| \geq \lambda\delta - \delta|r| > 0$  for  $x \in \partial B_\delta(0)$  as desired.

In summary we have

$$\deg(f, U, 0) = \sum_{i=1}^N \deg(df(x^i), B_\delta(0), 0) \quad (12.11)$$



and it remains to compute the degree of a nonsingular matrix. To this end we need the following lemma.

**Lemma 12.3.** *Two nonsingular matrices  $M_{1,2} \in \text{GL}(n)$  are homotopic in  $\text{GL}(n)$  if and only if  $\text{sign det } M_1 = \text{sign det } M_2$ .*

**Proof.** We will show that any given nonsingular matrix  $M$  is homotopic to  $\text{diag}(\text{sign det } M, 1, \dots, 1)$ , where  $\text{diag}(m_1, \dots, m_n)$  denotes a diagonal matrix with diagonal entries  $m_i$ .

In fact, note that adding one row to another and multiplying a row by a positive constant can be realized by continuous deformations such that all intermediate matrices are nonsingular. Hence we can reduce  $M$  to a diagonal matrix  $\text{diag}(m_1, \dots, m_n)$  with  $(m_i)^2 = 1$ . Next,

$$\begin{pmatrix} \pm \cos(\pi t) & \mp \sin(\pi t) \\ \sin(\pi t) & \cos(\pi t) \end{pmatrix},$$

shows that  $\text{diag}(\pm 1, 1)$  and  $\text{diag}(\mp 1, -1)$  are homotopic. Now we apply this result to all two by two subblocks as follows. For each  $i$  starting from  $n$  and going down to 2 transform the subblock  $\text{diag}(m_{i-1}, m_i)$  into  $\text{diag}(1, 1)$  respectively  $\text{diag}(-1, 1)$ . The result is the desired form for  $M$ .

To conclude the proof note that a continuous deformation within  $\text{GL}(n)$  cannot change the sign of the determinant since otherwise the determinant would have to vanish somewhere in between (i.e., we would leave  $\text{GL}(n)$ ).  $\square$

Using this lemma we can now show the main result of this section.

**Theorem 12.4.** *Suppose  $f \in \bar{C}_y^1(U, \mathbb{R}^n)$  and  $y \notin \text{CV}(f)$ , then a degree satisfying (D1)–(D4) satisfies*

$$\deg(f, U, y) = \sum_{x \in f^{-1}(\{y\})} \text{sign } J_f(x), \quad (12.12)$$

where the sum is finite and we agree to set  $\sum_{x \in \emptyset} = 0$ .

**Proof.** By the previous lemma we obtain

$$\deg(df(0), B_\delta(0), 0) = \deg(\text{diag}(\text{sign } J_f(0), 1, \dots, 1), B_\delta(0), 0)$$

since  $\det M \neq 0$  is equivalent to  $Mx \neq 0$  for  $x \in \partial B_\delta(0)$ . Hence it remains to show  $\deg(M_\pm, B_\delta(0), 0) = \pm 1$ , where  $M_\pm := \text{diag}(\pm 1, 1, \dots, 1)$ . For  $M_+$  this is true by (D2) and for  $M_-$  we note that we can replace  $B_\delta(0)$  by any neighborhood  $U$  of 0. Now abbreviate  $U_1 := \{x \in \mathbb{R}^n \mid |x_i| < 1, 1 \leq i \leq n\}$ ,  $U_2 := \{x \in \mathbb{R}^n \mid 1 < x_1 < 3, |x_i| < 1, 2 \leq i \leq n\}$ ,  $U := \{x \in \mathbb{R}^n \mid -1 < x_1 < 3, |x_i| < 1, 2 \leq i \leq n\}$ , and  $g(r) = 2 - |r - 1|$ ,  $h(r) = 1 - r^2$ . Consider the two functions  $f_1(x) = (1 - g(x_1)h(x_2) \cdots h(x_n), x_2, \dots, x_n)$  and  $f_2(x) = (1, x_2, \dots, x_n)$ . Clearly  $f_1^{-1}(\{0\}) = \{x^1, x^2\}$  with  $x^1 = 0$ ,  $x^2 = (2, 0, \dots, 0)$  and  $f_2^{-1}(0) = \emptyset$ . Since  $f_1(x) = f_2(x)$  for  $x \in \partial U$  we infer  $\deg(f_1, U, 0) = \deg(f_2, U, 0) = 0$ .

Moreover, we have  $\deg(f_1, U, 0) = \deg(f_1, U_1, 0) + \deg(f_1, U_2, 0)$  and hence  $\deg(M_-, U_1, 0) = \deg(df_1(x^1), U_1, 0) = \deg(f_1, U_1, 0) = -\deg(f_1, U_2, 0) = -\deg(df_1(x^2), U_1, 0) = -\deg(\mathbb{I}, U_1, 0) = -1$  as claimed.  $\square$

Up to this point we have only shown that a degree (provided there is one at all) necessarily satisfies (12.12). Once we have shown that regular values are dense, it will follow that the degree is uniquely determined by (12.12) since the remaining values follow from point (iii) of Theorem 12.1. On the other hand, we don't even know whether a degree exists since it is unclear whether (12.12) satisfies (D4). Hence we need to show that (12.12) can be extended to  $f \in \bar{C}_y(U, \mathbb{R}^n)$  and that this extension satisfies our requirements (D1)–(D4).

### 12.3. Extension of the determinant formula

Our present objective is to show that the determinant formula (12.12) can be extended to all  $f \in \bar{C}_y(U, \mathbb{R}^n)$ . As a preparation we prove that the set of regular values is dense. This is a consequence of a special case of Sard's theorem which says that  $\text{CV}(f)$  has zero measure.

**Lemma 12.5** (Sard). *Suppose  $f \in C^1(U, \mathbb{R}^n)$ , then the Lebesgue measure of  $\text{CV}(f)$  is zero.*

**Proof.** Since the claim is easy for linear mappings our strategy is as follows. We divide  $U$  into sufficiently small subsets. Then we replace  $f$  by its linear approximation in each subset and estimate the error.

Let  $\text{CP}(f) := \{x \in U \mid J_f(x) = 0\}$  be the set of critical points of  $f$ . We first pass to cubes which are easier to divide. Let  $\{Q_i\}_{i \in \mathbb{N}}$  be a countable cover for  $U$  consisting of open cubes such that  $\bar{Q}_i \subset U$ . Then it suffices to prove that  $f(\text{CP}(f) \cap Q_i)$  has zero measure since  $\text{CV}(f) = f(\text{CP}(f)) = \bigcup_i f(\text{CP}(f) \cap Q_i)$  (the  $Q_i$ 's are a cover).

Let  $Q$  be anyone of these cubes and denote by  $\rho$  the length of its edges. Fix  $\varepsilon > 0$  and divide  $Q$  into  $N^n$  cubes  $Q_i$  of length  $\rho/N$ . These cubes don't have to be open and hence we can assume that they cover  $Q$ . Since  $df(x)$  is uniformly continuous on  $\bar{Q}$  we can find an  $N$  (independent of  $i$ ) such that

$$|f(x) - f(\tilde{x}) - df(\tilde{x})(x - \tilde{x})| \leq \int_0^1 |df(\tilde{x} + t(x - \tilde{x})) - df(\tilde{x})| |\tilde{x} - x| dt \leq \frac{\varepsilon \rho}{N} \quad (12.13)$$

for  $\tilde{x}, x \in Q_i$ . Now pick a  $Q_i$  which contains a critical point  $\tilde{x}_i \in \text{CP}(f)$ . Without restriction we assume  $\tilde{x}_i = 0$ ,  $f(\tilde{x}_i) = 0$  and set  $M := df(\tilde{x}_i)$ . By  $\det M = 0$  there is an orthonormal basis  $\{b^i\}_{1 \leq i \leq n}$  of  $\mathbb{R}^n$  such that  $b^n$  is

orthogonal to the image of  $M$ . In addition,

$$Q_i \subseteq \left\{ \sum_{i=1}^n \lambda_i b^i \mid \sqrt{\sum_{i=1}^n |\lambda_i|^2} \leq \sqrt{n} \frac{\rho}{N} \right\}$$

and hence there is a constant (again independent of  $i$ ) such that

$$MQ_i \subseteq \left\{ \sum_{i=1}^{n-1} \lambda_i b^i \mid |\lambda_i| \leq C \sqrt{n} \frac{\rho}{N} \right\}$$

(e.g.,  $C := \max_{x \in \overline{Q}} |df(x)|$ ). Next, by our estimate (12.13) we even have

$$f(Q_i) \subseteq \left\{ \sum_{i=1}^n \lambda_i b^i \mid |\lambda_i| \leq (C + \varepsilon) \sqrt{n} \frac{\rho}{N}, |\lambda_n| \leq \varepsilon \sqrt{n} \frac{\rho}{N} \right\}$$

and hence the measure of  $f(Q_i)$  is smaller than  $\frac{\tilde{C}\varepsilon}{N^n}$ . Since there are at most  $N^n$  such  $Q_i$ 's, we see that the measure of  $f(\text{CP}(f) \cap Q)$  is smaller than  $\tilde{C}\varepsilon$ .  $\square$

By (ii) of Theorem 12.2,  $\deg(f, U, y)$  should be constant on each component of  $\mathbb{R}^n \setminus f(\partial U)$ . Unfortunately, if we connect  $y$  and a nearby regular value  $\tilde{y}$  by a path, then there might be some critical values in between.

**Example 12.1.** The function  $f(x) := x^2 \sin(\frac{\pi}{2x})$  is in  $\bar{C}_0^1([-1, 1], \mathbb{R})$ . It has 0 as a critical value and the critical values accumulate at 0.  $\diamond$

To overcome this problem we need a definition for  $\deg$  which works for critical values as well. Let us try to look for an integral representation. Formally (12.12) can be written as  $\deg(f, U, y) = \int_U \delta_y(f(x)) J_f(x) d^n x$ , where  $\delta_y(\cdot)$  is the Dirac distribution at  $y$ . But since we don't want to mess with distributions, we replace  $\delta_y(\cdot)$  by  $\phi_\varepsilon(\cdot - y)$ , where  $\{\phi_\varepsilon\}_{\varepsilon>0}$  is a family of functions such that  $\phi_\varepsilon$  is supported on the ball  $B_\varepsilon(0)$  of radius  $\varepsilon$  around 0 and satisfies  $\int_{\mathbb{R}^n} \phi_\varepsilon(x) d^n x = 1$ .

**Lemma 12.6** (Heinz). *Suppose  $f \in \bar{C}_y^1(U, \mathbb{R}^n)$  and  $y \notin \text{CV}(f)$ . Then the degree defined as in (12.12) satisfies*

$$\deg(f, U, y) = \int_U \phi_\varepsilon(f(x) - y) J_f(x) d^n x \quad (12.14)$$

for all positive  $\varepsilon$  smaller than a certain  $\varepsilon_0$  depending on  $f$  and  $y$ . Moreover,  $\text{supp}(\phi_\varepsilon(f(\cdot) - y)) \subset U$  for  $\varepsilon < \text{dist}(y, f(\partial U))$ .

**Proof.** If  $f^{-1}(\{y\}) = \emptyset$ , we can set  $\varepsilon_0 = \text{dist}(y, f(\overline{U}))$ , implying  $\phi_\varepsilon(f(x) - y) = 0$  for  $x \in \overline{U}$ .

If  $f^{-1}(\{y\}) = \{x^i\}_{1 \leq i \leq N}$ , the inverse function theorem ensures that we can find an  $\varepsilon_0 > 0$  such that  $f^{-1}(B_{\varepsilon_0}(y))$  is a union of disjoint neighborhoods

$U(x^i)$  of  $x^i$  with  $f|_{U(x^i)}$  a bijection and  $J_f(x)$  nonzero on  $U(x^i)$ . Again  $\phi_\varepsilon(f(x) - y) = 0$  for  $x \in \overline{U} \setminus \bigcup_{i=1}^N U(x^i)$  and hence

$$\begin{aligned} \int_U \phi_\varepsilon(f(x) - y) J_f(x) d^n x &= \sum_{i=1}^N \int_{U(x^i)} \phi_\varepsilon(f(x) - y) J_f(x) d^n x \\ &= \sum_{i=1}^N \text{sign}(J_f(x^i)) \int_{B_{\varepsilon_0}(0)} \phi_\varepsilon(\tilde{x}) d^n \tilde{x} = \deg(f, U, y), \end{aligned}$$

where we have used the change of variables  $\tilde{x} = f(x) - y$  in the second step.  $\square$

Our new integral representation makes sense even for critical values. But since  $\varepsilon_0$  depends on  $f$  and  $y$ , continuity is not clear. This will be tackled next.

The key idea is to show that the integral representation is independent of  $\varepsilon$  as long as  $\varepsilon < \text{dist}(y, f(\partial U))$ . To this end we will rewrite the difference as an integral over a divergence supported in  $U$  and then apply the Gauss–Green theorem. For this purpose the following result will be used.

**Lemma 12.7.** *Suppose  $f \in C^2(U, \mathbb{R}^n)$  and  $u \in C^1(\mathbb{R}^n, \mathbb{R}^n)$ , then*

$$(\text{div } u)(f) J_f = \text{div } D_f(u), \quad (12.15)$$

where  $D_f(u)_j$  is the determinant of the matrix obtained from  $df$  by replacing the  $j$ -th column by  $u(f)$ . Here  $\text{div } u = \sum_{j=1}^n \partial_j u_j$  is the **divergence** of a vector field.

**Proof.** We compute

$$\text{div } D_f(u) = \sum_{j=1}^n \partial_{x_j} D_f(u)_j = \sum_{j,k=1}^n D_f(u)_{j,k},$$

where  $D_f(u)_{j,k}$  is the determinant of the matrix obtained from the matrix associated with  $D_f(u)_j$  by applying  $\partial_{x_j}$  to the  $k$ -th column. Since  $\partial_{x_j} \partial_{x_k} f = \partial_{x_k} \partial_{x_j} f$  we infer  $D_f(u)_{j,k} = -D_f(u)_{k,j}$ ,  $j \neq k$ , by exchanging the  $k$ -th and the  $j$ -th column. Hence

$$\text{div } D_f(u) = \sum_{i=1}^n D_f(u)_{i,i}.$$

Now let  $J_f^{(i,j)}(x)$  denote the  $(i, j)$  cofactor of  $df(x)$  and recall the cofactor expansion of the determinant  $\sum_{i=1}^n J_f^{(i,j)} \partial_{x_i} f_k = \delta_{j,k} J_f$ . Using this to expand

the determinant  $D_f(u)_{i,i}$  along the  $i$ -th column shows

$$\begin{aligned} \operatorname{div} D_f(u) &= \sum_{i,j=1}^n J_f^{(i,j)} \partial_{x_i} u_j(f) = \sum_{i,j=1}^n J_f^{(i,j)} \sum_{k=1}^n (\partial_{x_k} u_j)(f) \partial_{x_i} f_k \\ &= \sum_{j,k=1}^n (\partial_{x_k} u_j)(f) \sum_{i=1}^n J_f^{(i,j)} \partial_{x_i} f_k = \sum_{j=1}^n (\partial_{x_j} u_j)(f) J_f \end{aligned}$$

as required.  $\square$

Now we can prove

**Theorem 12.8.** *There is a unique degree  $\deg$  satisfying (D1)–(D4). Moreover,  $\deg(\cdot, U, y) : \bar{C}_y(U, \mathbb{R}^n) \rightarrow \mathbb{Z}$  is constant on each component and given  $f \in \bar{C}_y(U, \mathbb{R}^n)$  we have*

$$\deg(f, U, y) = \sum_{x \in \tilde{f}^{-1}(y)} \operatorname{sign} J_{\tilde{f}}(x), \quad (12.16)$$

where  $\tilde{f} \in \bar{C}_y^1(U, \mathbb{R}^n)$  is in the same component of  $\bar{C}_y(U, \mathbb{R}^n)$ , say  $\|f - \tilde{f}\|_\infty < \operatorname{dist}(y, f(\partial U))$ , such that  $y \in \operatorname{RV}(\tilde{f})$ .

**Proof.** We will first show that our integral formula works in fact for all  $\varepsilon < \rho := \operatorname{dist}(y, f(\partial U))$ . For this we will make some additional assumptions: Let  $f \in \bar{C}^2(U, \mathbb{R}^n)$  and choose a family of functions  $\phi_\varepsilon \in C^\infty((0, \infty))$  with  $\operatorname{supp}(\phi_\varepsilon) \subset (0, \varepsilon)$  such that  $S_n \int_0^\varepsilon \phi(r) r^{n-1} dr = 1$ . Consider

$$I_\varepsilon(f, U, y) := \int_U \phi_\varepsilon(|f(x) - y|) J_f(x) d^n x.$$

Then  $I := I_{\varepsilon_1} - I_{\varepsilon_2}$  will be of the same form but with  $\phi_\varepsilon$  replaced by  $\varphi := \phi_{\varepsilon_1} - \phi_{\varepsilon_2}$ , where  $\varphi \in C^\infty((0, \infty))$  with  $\operatorname{supp}(\varphi) \subset (0, \rho)$  and  $\int_0^\rho \varphi(r) r^{n-1} dr = 0$ . To show that  $I = 0$  we will use our previous lemma with  $u$  chosen such that  $\operatorname{div}(u(x)) = \varphi(|x|)$ . To this end we make the ansatz  $u(x) = \psi(|x|)x$  such that  $\operatorname{div}(u(x)) = |x|\psi'(|x|) + n\psi(|x|)$ . Our requirement now leads to an ordinary differential equation whose solution is

$$\psi(r) = \frac{1}{r^n} \int_0^r s^{n-1} \varphi(s) ds.$$

Moreover, one checks  $\psi \in C^\infty((0, \infty))$  with  $\operatorname{supp}(\psi) \subset (0, \rho)$ . Thus our lemma shows

$$I = \int_U \operatorname{div} D_{f-y}(u) d^n x$$

and since the integrand vanishes in a neighborhood of  $\partial U$  we can extend it to all of  $\mathbb{R}^n$  by setting it zero outside  $U$  and choose a cube  $Q \supset U$ . Then elementary coordinatewise integration gives  $I = \int_Q \operatorname{div} D_{f-y}(u) d^n x = 0$ .

Now fix  $\delta < \rho$  and look at  $I_\varepsilon(f + g, U, y)$  for  $g \in B_\delta(f) \cap \bar{C}^2(U, \mathbb{R}^n) \subset C(\bar{U}, \mathbb{R}^n)$  and  $\varepsilon < \rho - \delta < \text{dist}((f + g)(\partial U), y)$  fixed. Then  $t \mapsto I_\varepsilon(f + tg, U, y)$ ,  $t \in [0, 1]$ , is continuous and it is integer valued (since it is equal to our determinant formula) on a dense set. Consequently it must be constant and we can extend  $I_\varepsilon$  to a function  $\bar{I}_\varepsilon$  on all of  $B_\delta(f)$  and hence on all of  $\bar{C}_y(U, \mathbb{R}^n)$ . Note that by mollifying  $f \in \bar{C}^1(U, \mathbb{R}^n)$  we get a sequence of smooth functions for which both  $f$  and  $df$  converge uniformly on compact subsets of  $U$  and hence  $I_\varepsilon$  converges, such that for such  $f$  we still have  $\bar{I}_\varepsilon(f, U, y) = I_\varepsilon(f, U, y)$ .

Now we set

$$\deg(f, U, y) := I_\varepsilon(\tilde{f}, U, y),$$

where  $\tilde{f} \in \bar{C}^1(U, \mathbb{R}^n)$  with  $\varepsilon < \rho$  and  $|\tilde{f} - f| < \rho - \varepsilon$ . Then (D1) holds since it holds for  $I_\varepsilon$ , (D2) holds since  $I_\varepsilon$  extends the determinant formula, (D3) holds since the integrand of  $I_\varepsilon$  vanishes on  $\bar{U} \setminus (U_1 \cup U_2)$ , and (D4) holds since we can choose  $\varepsilon < \min_{t \in [0, 1]} \text{dist}(H(t)(\partial U), y)$  such that  $I_\varepsilon(H(t), U, y)$  is continuous and hence constant for  $t \in [0, 1]$ .  $\square$

To conclude this section, let us give a few simple examples illustrating the use of the Brouwer degree.

**Example 12.2.** First, let's investigate the zeros of

$$f(x_1, x_2) := (x_1 - 2x_2 + \cos(x_1 + x_2), x_2 + 2x_1 + \sin(x_1 + x_2)).$$

Denote the linear part by

$$g(x_1, x_2) := (x_1 - 2x_2, x_2 + 2x_1).$$

Then we have  $|g(x)| = \sqrt{5}|x|$  and  $|f(x) - g(x)| = 1$  and hence  $h(t) = (1 - t)g + tf = g + t(f - g)$  satisfies  $|h(t)| \geq |g| - t|f - g| > 0$  for  $|x| > 1/\sqrt{5}$  implying

$$\deg(f, B_r(0), 0) = \deg(g, B_r(0), 0) = 1, \quad r > 1/\sqrt{5}.$$

Moreover, since  $J_f(x) = 5 + 3\cos(x_1 + x_2) + \sin(x_1 + x_2) > 1$  the determinant formula (12.12) for the degree implies that  $f(x) = 0$  has a unique solution in  $\mathbb{R}^2$ . This solution even has to lie on the circle  $|x| = 1/\sqrt{5}$  since  $f(x) = 0$  implies  $1 = |f(x) - g(x)| = |g(x)| = \sqrt{5}|x|$ .  $\diamond$

Next let us prove the following result which implies the **hairy ball (or hedgehog) theorem**.

**Theorem 12.9.** *Suppose  $U$  is open, bounded and contains the origin and let  $f : \partial U \rightarrow \mathbb{R}^n \setminus \{0\}$  be continuous. If  $n$  is odd, then there exists an  $x \in \partial U$  and a  $\lambda \neq 0$  such that  $f(x) = \lambda x$ .*

**Proof.** By Theorem 13.10 we can assume  $f \in C(\bar{U}, \mathbb{R}^n)$  and since  $n$  is odd we have  $\deg(-\mathbb{I}, U, 0) = -1$ . Now if  $\deg(f, U, 0) \neq -1$ , then  $H(t, x) =$

$(1-t)f(x) - tx$  must have a zero  $(t_0, x_0) \in (0, 1) \times \partial U$  and hence  $f(x_0) = \frac{t_0}{1-t_0}x_0$ . Otherwise, if  $\deg(f, U, 0) = -1$  we can apply the same argument to  $H(t, x) = (1-t)f(x) + tx$ .  $\square$

In particular, this result implies that a continuous tangent vector field on the unit sphere  $f : S^{n-1} \rightarrow \mathbb{R}^n$  (with  $f(x)x = 0$  for all  $x \in S^{n-1}$ ) must vanish somewhere if  $n$  is odd. Or, for  $n = 3$ , you cannot smoothly comb a hedgehog without leaving a bald spot or making a parting. It is however possible to comb the hair smoothly on a torus and that is why the magnetic containers in nuclear fusion are toroidal.

**Example 12.3.** The result fails in even dimensions as the example  $n = 2$ ,  $U = B_1(0)$ ,  $f(x_1, x_2) = (-x_2, x_1)$  shows.  $\diamond$

Another illustration is the fact that a vector field on  $\mathbb{R}^n$ , which points outwards (or inwards) on a sphere, must vanish somewhere inside the sphere (Problem 12.2).

One more useful observation is that odd functions have odd degree:

**Theorem 12.10** (Borsuk). *Let  $0 \in U \subseteq \mathbb{R}^n$  be open, bounded and symmetric with respect to the origin (i.e.,  $U = -U$ ). Let  $f \in \bar{C}_0(U, \mathbb{R}^n)$  be odd (i.e.,  $f(-x) = -f(x)$ ). Then  $\deg(f, U, 0)$  is odd.*

**Proof.** If  $f \in \bar{C}_0^1(U)$  and  $0 \in \text{RV}(f)$ , then the claim is straightforward since

$$\deg(f, U, 0) = \text{sign } J_f(0) + \sum_{x \in f^{-1}(0) \setminus \{0\}} \text{sign } J_f(x),$$

where the sum is even since for every  $x \in f^{-1}(0) \setminus \{0\}$  we also have  $-x \in f^{-1}(0) \setminus \{0\}$  as well as  $J_f(x) = J_f(-x)$ .

Hence we need to reduce the general case to this one. Clearly if  $f \in \bar{C}_0(U, \mathbb{R}^n)$  we can choose an approximating  $f_0 \in \bar{C}_0^1(U, \mathbb{R}^n)$  and replacing  $f_0$  by its odd part  $\frac{1}{2}(f_0(x) - f_0(-x))$  we can assume  $f_0$  to be odd. Moreover, if  $J_{f_0}(0) = 0$  we can replace  $f_0$  by  $f_0(x) + \delta x$  such that 0 is regular. However, if we choose a nearby regular value  $y$  and consider  $f_0(x) - y$  we have the problem that constant functions are even. Hence we will try the next best thing and perturb by a function which is constant in all except one direction. To this end we choose an odd function  $\varphi \in C^1(\mathbb{R})$  such that  $\varphi'(0) = 0$  (since we don't want to alter the behavior at 0) and  $\varphi(t) \neq 0$  for  $t \neq 0$ . Now we consider  $f_1(x) = f_0(x) - \varphi(x_1)y^1$  and note

$$df_1(x) = df_0(x) - d\varphi(x_1)y^1 = df_0(x) - d\varphi(x_1)\frac{f_0(x)}{\varphi(x_1)} = \varphi(x_1)d\left(\frac{f_0(x)}{\varphi(x_1)}\right)$$

for every  $x \in U_1 := \{x \in U \mid x_1 \neq 0\}$  with  $f_1(x) = 0$ . Hence if  $y^1$  is chosen such that  $y^1 \in \text{RV}(h_1)$ , where  $h_1 : U_1 \rightarrow \mathbb{R}^n$ ,  $x \mapsto \frac{f_0(x)}{\varphi(x_1)}$ , then 0 will be

a regular value for  $f_1$  when restricted to  $V_1 := U_1$ . Now we repeat this procedure and consider  $f_2(x) = f_1(x) - \varphi(x_2)y^2$  with  $y^2 \in \text{RV}(h_2)$  as before. Then every point  $x \in V_2 := U_1 \cup U_2$  with  $f_2(x) = 0$  either satisfies  $x_2 \neq 0$  and thus is regular by our choice of  $y^2$  or satisfies  $x_2 = 0$  and thus is regular since it is in  $V_1$  and  $df_2(x) = df_1(x)$  by our assumption  $\phi'(0) = 0$ . After  $n$  steps we reach  $V_n = U \setminus \{0\}$  and  $f_n$  is the approximation we are looking for.  $\square$

At first sight the obvious conclusion that an odd function has a zero does not seem too spectacular since the fact that  $f$  is odd already implies  $f(0) = 0$ . However, the result gets more interesting upon observing that it suffices when the boundary values are odd. Moreover, local constancy of the degree implies that  $f$  does not only attain 0 but also any  $y$  in a neighborhood of 0. The next two important consequences are based on this observation:

**Theorem 12.11** (Borsuk–Ulam). *Let  $0 \in U \subseteq \mathbb{R}^n$  be open, bounded and symmetric with respect to the origin. Let  $f \in C(\partial U, \mathbb{R}^m)$  with  $m < n$ . Then there is some  $x \in \partial U$  with  $f(x) = f(-x)$ .*

**Proof.** Consider  $g(x) = f(x) - f(-x)$  and extend it to a continuous odd function  $\bar{U} \rightarrow \mathbb{R}^m$  (extend the domain by Tietze and then take the odd part, finally fill up the missing coordinates by setting them equal to 0). If  $g$  does not vanish on  $\partial U$ , we get that  $\deg(g, U, y) = \deg(g, U, 0) \neq 0$  for  $y$  in a neighborhood of 0 and thus the image of  $g$  contains a neighborhood of 0 (in  $\mathbb{R}^m$ ), which contradicts the fact that the image is in  $\mathbb{R}^m \times \{0\} \subset \mathbb{R}^n$ .  $\square$

This theorem is often illustrated by the fact that there are always two opposite points on the earth which have the same weather (in the sense that they have the same temperature and the same pressure). In a similar manner one can also derive the **invariance of domain theorem**.

**Theorem 12.12** (Brouwer). *Let  $U \subseteq \mathbb{R}^n$  be open and let  $f : U \rightarrow \mathbb{R}^n$  be continuous and locally injective. Then  $f(U)$  is also open.*

**Proof.** It suffices to show that every point  $x \in U$  contains a neighborhood  $B_r(x)$  such that the image  $f(B_r(x))$  contains a ball centered at  $f(x)$ . By simple translations we can assume  $x = 0$  as well as  $f(x) = 0$ . Now choose  $r$  sufficiently small such that  $f$  restricted to  $\bar{B}_r(0)$  is injective and consider  $H(t, x) := f(\frac{1}{1+t}x) - f(-\frac{t}{1+t}x)$  for  $t \in [0, 1]$  and  $x \in \bar{B}_r(0)$ . Moreover, if  $H(t, x) = 0$  then by injectivity  $\frac{1}{1+t}x = -\frac{t}{1+t}x$ , that is,  $x = 0$ . Thus  $\deg(f, B_r(0), 0) = \deg(H(1), B_r(0), 0) \neq 0$  since  $H(1) = f(\frac{1}{2}x) - f(-\frac{1}{2}x)$  is odd. But then we also have  $\deg(f, B_r(0), y) \neq 0$  for  $y \in B_\varepsilon(0)$  and thus  $B_\varepsilon(0) \subseteq f(B_r(0))$ .  $\square$



An easy consequence worth while noting is the topological invariance of dimension:

**Corollary 12.13.** *If  $m < n$  and  $U$  is a nonempty open subset of  $\mathbb{R}^n$ , then there is no continuous injective mapping from  $U$  to  $\mathbb{R}^m$ .*

**Proof.** Suppose there were such a map and extend it to a map from  $U$  to  $\mathbb{R}^n$  by setting the additional coordinates equal to zero. The resulting map contradicts the invariance of domain theorem.  $\square$

In particular,  $\mathbb{R}^m$  and  $\mathbb{R}^n$  are not homeomorphic for  $m \neq n$ .

**Problem 12.1.** *Suppose  $U = (a, b) \subset \mathbb{R}^1$ . Show*

$$\deg(f, (a, b), y) = \frac{1}{2}(\text{sign}(f(b) - y) - \text{sign}(f(a) - y)).$$

*In particular, our degree reduces to the intermediate value theorem in this case.*

**Problem\* 12.2.** *Suppose  $f : \bar{B}_r(0) \rightarrow \mathbb{R}^n$  is continuous and satisfies*

$$f(x)x > 0, \quad |x| = r.$$

*Then  $f(x)$  vanishes somewhere inside  $B_r(0)$ .*

**Problem 12.3.** *Show that in Borsuk's theorem the condition  $f$  is odd can be replaced by  $f(x) \neq tf(-x)$  for all  $x \in \partial U$  and  $t \in (0, 1]$ . Note that this condition will hold if  $\text{sign}(f(x)) \neq \text{sign}(f(-x))$ ,  $x \in \partial U$  (where  $\text{sign}(f(x)) := \frac{f(x)}{|f(x)|}$ ).*

## 12.4. The Brouwer fixed point theorem

Now we can show that the famous Brouwer fixed point theorem is a simple consequence of the properties of our degree.

**Theorem 12.14** (Brouwer fixed point). *Let  $K$  be a topological space homeomorphic to a compact, convex subset of  $\mathbb{R}^n$  and let  $f \in C(K, K)$ , then  $f$  has at least one fixed point.*

**Proof.** Clearly we can assume  $K \subset \mathbb{R}^n$  since homeomorphisms preserve fixed points. Now let's assume  $K = \bar{B}_r(0)$ . If there is a fixed point on the boundary  $\partial B_r(0)$  we are done. Otherwise  $H(t, x) = x - tf(x)$  satisfies  $0 \notin H(t, \partial B_r(0))$  since  $|H(t, x)| \geq |x| - t|f(x)| \geq (1 - t)r > 0$ ,  $0 \leq t < 1$ . And the claim follows from  $\deg(x - f(x), B_r(0), 0) = \deg(x, B_r(0), 0) = 1$ .

Now let  $K$  be convex. Then  $K \subseteq B_\rho(0)$  and, by the Hilbert projection theorem (Theorem 2.11) (or alternatively by the Tietze extension theorem or its variant Theorem 13.10 below), we can find a continuous retraction  $R : \mathbb{R}^n \rightarrow K$  (i.e.,  $R(x) = x$  for  $x \in K$ ) and consider  $\tilde{f} =$

$f \circ R \in C(\bar{B}_\rho(0), \bar{B}_\rho(0))$ . By our previous analysis, there is a fixed point  $x = \tilde{f}(x) \in \text{conv}(f(K)) \subseteq K$ .  $\square$

Note that any compact, convex subset of a finite dimensional Banach space (complex or real) is isomorphic to a compact, convex subset of  $\mathbb{R}^n$  since linear transformations preserve both properties. In addition, observe that all assumptions are needed. For example, the map  $f : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto x+1$ , has no fixed point ( $\mathbb{R}$  is homeomorphic to a bounded set but not to a compact one). The same is true for the map  $f : \partial B_1(0) \rightarrow \partial B_1(0), x \mapsto -x$  ( $\partial B_1(0) \subset \mathbb{R}^n$  is simply connected for  $n \geq 3$  but not homeomorphic to a convex set).

As an easy example of how to use the Brouwer fixed point theorem we show the famous **Perron–Frobenius theorem**.

**Theorem 12.15** (Perron–Frobenius). *Let  $A$  be an  $n \times n$  matrix all whose entries are nonnegative and there is an  $m$  such the entries of  $A^m$  are all positive. Then  $A$  has a positive eigenvalue and the corresponding eigenvector can be chosen to have positive components.*

**Proof.** We equip  $\mathbb{R}^n$  with the norm  $|x|_1 := \sum_{j=1}^n |x_j|$  and set  $\Delta := \{x \in \mathbb{R}^n | x_j \geq 0, |x|_1 = 1\}$ . For  $x \in \Delta$  we have  $Ax \neq 0$  (since  $A^m x \neq 0$ ) and hence

$$f : \Delta \rightarrow \Delta, \quad x \mapsto \frac{Ax}{|Ax|_1}$$

has a fixed point  $x_0$  by the Brouwer fixed point theorem. Then  $Ax_0 = |Ax_0|_1 x_0$  and  $x_0$  has positive components since  $A^m x_0 = |Ax_0|_1^m x_0$  has.  $\square$

Let me remark that the Brouwer fixed point theorem is equivalent to the fact that there is no continuous retraction  $R : \bar{B}_1(0) \rightarrow \partial B_1(0)$  (with  $R(x) = x$  for  $x \in \partial B_1(0)$ ) from the unit ball to the unit sphere in  $\mathbb{R}^n$ .

In fact, if  $R$  would be such a retraction,  $-R$  would have a fixed point  $x_0 \in \partial B_1(0)$  by Brouwer's theorem. But then  $x_0 = -R(x_0) = -x_0$  which is impossible. Conversely, if a continuous function  $f : \bar{B}_1(0) \rightarrow \bar{B}_1(0)$  has no fixed point we can define a retraction  $R(x) = f(x) + t(x)(x - f(x))$ , where  $t(x) \geq 0$  is chosen such that  $|R(x)|^2 = 1$  (i.e.,  $R(x)$  lies on the intersection of the line spanned by  $x, f(x)$  with the unit sphere).

Using this equivalence the Brouwer fixed point theorem can also be derived easily by showing that the homology groups of the unit ball  $\bar{B}_1(0)$  and its boundary (the unit sphere) differ (see, e.g., [38] for details).

### 12.5. Kakutani's fixed point theorem and applications to game theory

In this section we want to apply Brouwer's fixed point theorem to show the existence of Nash equilibria for  $n$ -person games. As a preparation we extend Brouwer's fixed point theorem to set-valued functions.

Denote by  $\text{CS}(K)$  the set of all nonempty convex subsets of  $K$ .

**Theorem 12.16** (Kakutani). *Suppose  $K$  is a compact convex subset of  $\mathbb{R}^n$  and  $f : K \rightarrow \text{CS}(K)$ . If the set*

$$\Gamma := \{(x, y) | y \in f(x)\} \subseteq K^2 \quad (12.17)$$

*is closed, then there is a point  $x \in K$  such that  $x \in f(x)$ .*

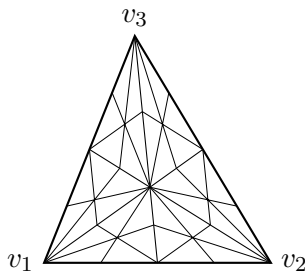
**Proof.** Our strategy is to apply Brouwer's theorem, hence we need a function related to  $f$ . For this purpose it is convenient to assume that  $K$  is a simplex

$$K = \text{conv}(v_1, \dots, v_m), \quad m \leq n + 1,$$

where  $v_i$  are the vertices. Recall that each point  $x \in K$  can be uniquely represented by its barycentric coordinates  $\lambda_i(x)$  (i.e.,  $\lambda_i \geq 0$ ,  $\sum_{i=1}^m \lambda_i(x) = 1$  and  $x = \sum_{i=1}^m \lambda_i v_i$ ). Now if we pick  $y_i \in f(v_i)$  we could set

$$f^1(x) = \sum_{i=1}^m \lambda_i(x) y_i.$$

By construction,  $f^1 \in C(K, K)$  and there is a fixed point  $x^1$ . But unless  $x^1$  is one of the vertices, this doesn't help us too much. So let's choose a better function as follows. Consider the  $k$ -th barycentric subdivision, that is, for every permutation  $v_{\sigma_1}, \dots, v_{\sigma_m}$  of the vertices you consider the simplex  $\text{conv}(v_{\sigma_1}, \frac{1}{2}(v_{\sigma_1} + v_{\sigma_2}), \dots, \frac{1}{m}(v_{\sigma_1} + \dots + v_{\sigma_m}))$ . This gives you  $m!$  smaller simplices (note that the maximal distance between vertices of the subsimplices decreases by a factor  $\frac{m-1}{m}$  during the subdivision) whose union is the simplex you have started with. Now repeat this construction  $k$  times.



For each vertex  $v_i$  in this subdivision pick an element  $y_i \in f(v_i)$ . Now define  $f^k(v_i) = y_i$  and extend  $f^k$  to the interior of each subsimplex as before.

Hence  $f^k \in C(K, K)$  and there is a fixed point  $x^k$  in one of the subsimplices. Denote this subsimplex by  $\text{conv}(v_1^k, \dots, v_m^k)$  such that

$$x^k = \sum_{i=1}^m \lambda_i^k v_i^k = \sum_{i=1}^m \lambda_i^k y_i^k, \quad y_i^k = f^k(v_i^k). \quad (12.18)$$

Since  $(x^k, \lambda_1^k, \dots, \lambda_m^k, y_1^k, \dots, y_m^k) \in K \times [0, 1]^m \times K^m$  we can assume that this sequence converges to some limit  $(x^0, \lambda_1^0, \dots, \lambda_m^0, y_1^0, \dots, y_m^0)$  after passing to a subsequence. Since the subsimplices shrink to a point, this implies  $v_i^k \rightarrow x^0$  and hence  $y_i^0 \in f(x^0)$  since  $(v_i^k, y_i^k) \in \Gamma \rightarrow (v_i^0, y_i^0) \in \Gamma$  by the closedness assumption. Now (12.18) tells us

$$x^0 = \sum_{i=1}^m \lambda_i^0 y_i^0 \in f(x^0)$$

since  $f(x^0)$  is convex and the claim holds if  $K$  is a simplex.

If  $K$  is not a simplex, we can pick a simplex  $S$  containing  $K$  and proceed as in the proof of the Brouwer theorem.  $\square$

If  $f(x)$  contains precisely one point for all  $x$ , then Kakutani's theorem reduces to the Brouwer's theorem (show that the closedness of  $\Gamma$  is equivalent to continuity of  $f$ ).

Now we want to see how this applies to game theory.

An  **$n$ -person game** consists of  $n$  players who have  $m_i$  possible actions to choose from. The set of all possible actions for the  $i$ -th player will be denoted by  $\Phi_i = \{1, \dots, m_i\}$ . An element  $\varphi_i \in \Phi_i$  is also called a pure strategy for reasons to become clear in a moment. Once all players have chosen their move  $\varphi_i$ , the payoff for each player is given by the **payoff** function

$$R_i(\varphi) \in \mathbb{R}, \quad \varphi = (\varphi_1, \dots, \varphi_n) \in \Phi = \bigtimes_{i=1}^n \Phi_i \quad (12.19)$$

of the  $i$ -th player. We will consider the case where the game is repeated a large number of times and where in each step the players choose their action according to a fixed strategy. Here a **strategy**  $s_i$  for the  $i$ -th player is a probability distribution on  $\Phi_i$ , that is,  $s_i = (s_i^1, \dots, s_i^{m_i})$  such that  $s_i^k \geq 0$  and  $\sum_{k=1}^{m_i} s_i^k = 1$ . The set of all possible strategies for the  $i$ -th player is denoted by  $S_i$ . The number  $s_i^k$  is the probability for the  $k$ -th pure strategy to be chosen. Consequently, if  $s = (s_1, \dots, s_n) \in S = \bigtimes_{i=1}^n S_i$  is a collection of strategies, then the probability that a given collection of pure strategies gets chosen is

$$s(\varphi) = \prod_{i=1}^n s_i(\varphi), \quad s_i(\varphi) = s_i^{k_i}, \quad \varphi = (k_1, \dots, k_n) \in \Phi \quad (12.20)$$

(assuming all players make their choice independently) and the expected payoff for player  $i$  is

$$R_i(s) = \sum_{\varphi \in \Phi} s(\varphi) R_i(\varphi). \quad (12.21)$$

By construction,  $R_i : S \rightarrow \mathbb{R}$  is polynomial and hence in particular continuous.

The question is of course, what is an optimal strategy for a player? If the other strategies are known, a **best reply** of player  $i$  against  $s$  would be a strategy  $\bar{s}_i$  satisfying

$$R_i(s \setminus \bar{s}_i) = \max_{\tilde{s}_i \in S_i} R_i(s \setminus \tilde{s}_i) \quad (12.22)$$

Here  $s \setminus \tilde{s}_i$  denotes the strategy combination obtained from  $s$  by replacing  $s_i$  by  $\tilde{s}_i$ . The set of all best replies against  $s$  for the  $i$ -th player is denoted by  $B_i(s)$ . Since

$$R_i(s) = \sum_{k=1}^{m_i} s_i^k R_i(s/k) \quad (12.23)$$

we have  $\bar{s}_i \in B_i(s)$  if and only if  $\bar{s}_i^k = 0$  whenever  $R_i(s \setminus k) < \max_{1 \leq l \leq m_i} R_i(s \setminus l)$ . In particular, since there are no restrictions on the other entries,  $B_i(s)$  is a nonempty convex set.

Let  $s, \bar{s} \in S$ , we call  $\bar{s}$  a best reply against  $s$  if  $\bar{s}_i$  is a best reply against  $s$  for all  $i$ . The set of all best replies against  $s$  is  $B(s) = \times_{i=1}^n B_i(s)$ .

A strategy combination  $\bar{s} \in S$  is a **Nash equilibrium** for the game if it is a best reply against itself, that is,

$$\bar{s} \in B(\bar{s}). \quad (12.24)$$

Or, put differently,  $\bar{s}$  is a Nash equilibrium if no player can increase his payoff by changing his strategy as long as all others stick to their respective strategies. In addition, if a player sticks to his equilibrium strategy, he is assured that his payoff will not decrease no matter what the others do.

To illustrate these concepts, let us consider the famous *prisoner's dilemma*. Here we have two players which can choose to defect or to cooperate. The payoff is symmetric for both players and given by the following diagram

$$\begin{array}{c|cc} R_1 & d_2 & c_2 \\ \hline d_1 & 0 & 2 \\ c_1 & -1 & 1 \end{array} \quad \begin{array}{c|cc} R_2 & d_2 & c_2 \\ \hline d_1 & 0 & -1 \\ c_1 & 2 & 1 \end{array} \quad (12.25)$$

where  $c_i$  or  $d_i$  means that player  $i$  cooperates or defects, respectively. You should think of two prisoners who are offered a reduced sentence if they testify against the other.

It is easy to see that the (pure) strategy pair  $(d_1, d_2)$  is the only Nash equilibrium for this game and that the expected payoff is 0 for both players.

Of course, both players could get the payoff 1 if they both agree to cooperate. But if one would break this agreement in order to increase his payoff, the other one would get less. Hence it might be safer to defect.

Now that we have seen that Nash equilibria are a useful concept, we want to know when such an equilibrium exists. Luckily we have the following result.

**Theorem 12.17** (Nash). *Every  $n$ -person game has at least one Nash equilibrium.*

**Proof.** The definition of a Nash equilibrium begs us to apply Kakutani's theorem to the set-valued function  $s \mapsto B(s)$ . First of all,  $S$  is compact and convex and so are the sets  $B(s)$ . Next, observe that the closedness condition of Kakutani's theorem is satisfied since if  $s^m \in S$  and  $\bar{s}^m \in B(s^m)$  both converge to  $s$  and  $\bar{s}$ , respectively, then (12.22) for  $s^m, \bar{s}^m$

$$R_i(s^m \setminus \tilde{s}_i) \leq R_i(s^m \setminus \bar{s}_i^m), \quad \tilde{s}_i \in S_i, 1 \leq i \leq n,$$

implies (12.22) for the limits  $s, \bar{s}$

$$R_i(s \setminus \tilde{s}_i) \leq R_i(s \setminus \bar{s}_i), \quad \tilde{s}_i \in S_i, 1 \leq i \leq n,$$

by continuity of  $R_i(s)$ . □

## 12.6. Further properties of the degree

We now prove some additional properties of the mapping degree. The first one will relate the degree in  $\mathbb{R}^n$  with the degree in  $\mathbb{R}^m$ . It will be needed later on to extend the definition of degree to infinite dimensional spaces. By virtue of the canonical embedding  $\mathbb{R}^m \hookrightarrow \mathbb{R}^m \times \{0\} \subset \mathbb{R}^n$  we can consider  $\mathbb{R}^m$  as a subspace of  $\mathbb{R}^n$ . We can project  $\mathbb{R}^n$  to  $\mathbb{R}^m$  by setting the last  $n - m$  coordinates equal to zero.

**Theorem 12.18** (Reduction property). *Let  $U \subseteq \mathbb{R}^n$  be open and bounded,  $f \in C(\bar{U}, \mathbb{R}^m)$  and  $y \in \mathbb{R}^m \setminus (\mathbb{I} + f)(\partial U)$ , then*

$$\deg(\mathbb{I} + f, U, y) = \deg(\mathbb{I} + f_m, U_m, y), \quad (12.26)$$

where  $f_m = f|_{U_m}$ , where  $U_m$  is the projection of  $U$  to  $\mathbb{R}^m$ .

**Proof.** After perturbing  $f$  a little, we can assume  $f \in C^1(U, \mathbb{R}^m)$  without loss of generality. Let  $x \in (\mathbb{I} + f)^{-1}(\{y\})$ , then  $x = y - f(x) \in \mathbb{R}^m$  implies  $(\mathbb{I} + f)^{-1}(\{y\}) = (\mathbb{I} + f_m)^{-1}(\{y\})$ . Moreover,

$$\begin{aligned} J_{\mathbb{I}+f}(x) &= \det(\mathbb{I} + df)(x) = \det \begin{pmatrix} \delta_{ij} + \partial_j f_i(x) & \partial_j f_j(x) \\ 0 & \delta_{ij} \end{pmatrix} \\ &= \det(\delta_{ij} + \partial_j f_i) = J_{\mathbb{I}+f_m}(x) \end{aligned}$$

So if  $y \in \text{RV}(\mathbb{I} + f_m)$  we immediately get  $\deg(\mathbb{I} + f, U, y) = \deg(\mathbb{I} + f_m, U_m, y)$  as desired. Otherwise, if  $y \in \text{CV}(\mathbb{I} + f_m)$  we can choose some  $\tilde{y} \in \text{RV}(\mathbb{I} + f_m)$  (and hence also  $\tilde{y} \in \text{RV}(\mathbb{I} + f)$  by the first part) with  $|y - \tilde{y}| < \min(\text{dist}(y, f(\partial U)), \text{dist}(y, f_m(\partial U_m)))$  and the claim follows from the regular case.  $\square$

Let  $U \subseteq \mathbb{R}^n$  and  $f \in C(\overline{U}, \mathbb{R}^n)$  be as usual. By Theorem 12.2 we know that  $\deg(f, U, y)$  is the same for every  $y$  in a connected component of  $\mathbb{R}^n \setminus f(\partial U)$ . Since  $\mathbb{R}^n \setminus f(\partial U)$  is open and locally path connected, these components are open. We will denote these components by  $G_j$  and write  $\deg(f, U, G_j) := \deg(f, U, y)$  if  $y \in G_j$ . In this context observe that since  $f(\partial U)$  is compact any unbounded component (there will be two for  $n = 1$  and one for  $n > 1$ ) will have degree zero.

**Theorem 12.19** (Product formula). *Let  $U \subseteq \mathbb{R}^n$  be a bounded and open set and denote by  $G_j$  the connected components of  $\mathbb{R}^n \setminus f(\partial U)$ . If  $g \circ f \in \bar{C}_y(U, \mathbb{R}^n)$ , then*

$$\deg(g \circ f, U, y) = \sum_j \deg(f, U, G_j) \deg(g, G_j, y), \quad (12.27)$$

where only finitely many terms in the sum are nonzero (and in particular, summands corresponding to unbounded components are considered to be zero).

**Proof.** Since  $y \notin (g \circ f)(\partial U)$  we have  $g^{-1}(\{y\}) \cap f(\partial U) = \emptyset$ , that is,  $g^{-1}(\{y\}) \subset \bigcup_j G_j$ . Moreover, since  $f(\overline{U})$  is compact, we can find an  $r > 0$  such that  $f(\overline{U}) \subseteq B_r(0)$ . Moreover, since  $g^{-1}(\{y\})$  is closed,  $g^{-1}(\{y\}) \cap B_r(0)$  is compact and hence can be covered by finitely many components, say  $g^{-1}(\{y\}) \subset \bigcup_{j=1}^m G_j$ . In particular, the others will be either unbounded or have  $\deg(f, G_k, y) = 0$  and hence only finitely many terms in the above sum are nonzero.

We begin by computing  $\deg(g \circ f, U, y)$  in the case where  $f, g \in C^1$  and  $y \notin \text{CV}(g \circ f)$ . Since  $d(g \circ f)(x) = dg(f(x)) \circ df(x)$  the claim is a

straightforward calculation

$$\begin{aligned}
\deg(g \circ f, U, y) &= \sum_{x \in (g \circ f)^{-1}(\{y\})} \text{sign}(J_{g \circ f}(x)) \\
&= \sum_{x \in (g \circ f)^{-1}(\{y\})} \text{sign}(J_g(f(x))) \text{sign}(J_f(x)) \\
&= \sum_{z \in g^{-1}(\{y\})} \text{sign}(J_g(z)) \sum_{x \in f^{-1}(\{z\})} \text{sign}(J_f(x)) \\
&= \sum_{z \in g^{-1}(\{y\})} \text{sign}(J_g(z)) \deg(f, U, z)
\end{aligned}$$

and, using our cover  $\{G_j\}_{j=1}^m$ ,

$$\begin{aligned}
\deg(g \circ f, U, y) &= \sum_{j=1}^m \sum_{z \in g^{-1}(\{y\}) \cap G_j} \text{sign}(J_g(z)) \deg(f, U, z) \\
&= \sum_{j=1}^m \deg(f, U, G_j) \sum_{z \in g^{-1}(\{y\}) \cap G_j} \text{sign}(J_g(z)) \\
&= \sum_{j=1}^m \deg(f, U, G_j) \deg(g, G_j, y).
\end{aligned}$$

Moreover, this formula still holds for  $y \in \text{CV}(g \circ f)$  and for  $g \in C$  by construction of the Brouwer degree. However, the case  $f \in C$  will need a closer investigation since the components  $G_j$  depend on  $f$ . To overcome this problem we will introduce the sets

$$L_l := \{z \in \mathbb{R}^n \setminus f(\partial U) \mid \deg(f, U, z) = l\}.$$

Observe that  $L_l$ ,  $l \neq 0$ , must be a union of some sets from  $\{G_j\}_{j=1}^m$ , that is,  $L_l = \bigcup_{k=1}^{m_l} G_{j_k^l}$  and  $\bigcup_{l \neq 0} L_l = \bigcup_{j=1}^m G_j$ .

Now choose  $\tilde{f} \in C^1$  such that  $|f(x) - \tilde{f}(x)| < 2^{-1} \text{dist}(g^{-1}(\{y\}), f(\partial U))$  for  $x \in \bar{U}$  and define  $\tilde{G}_j$ ,  $\tilde{L}_l$  accordingly. Then we have  $L_l \cap g^{-1}(\{y\}) = \tilde{L}_l \cap g^{-1}(\{y\})$  by Theorem 12.1 (iii) and hence  $\deg(g, \tilde{L}_l, y) = \deg(g, L_l, y)$  by Theorem 12.1 (i) implying

$$\begin{aligned}
\deg(g \circ f, U, y) &= \deg(g \circ \tilde{f}, U, y) = \sum_{j=1}^{\tilde{m}} \deg(\tilde{f}, U, \tilde{G}_j) \deg(g, \tilde{G}_j, y) \\
&= \sum_{l \neq 0} l \deg(g, \tilde{L}_l, y) = \sum_{l \neq 0} l \deg(g, L_l, y) \\
&= \sum_{l \neq 0} \sum_{k=1}^{m_l} l \deg(g, G_{j_k^l}, y) = \sum_{j=1}^m \deg(f, U, G_j) \deg(g, G_j, y)
\end{aligned}$$



which proves the claim.  $\square$

### 12.7. The Jordan curve theorem

In this section we want to show how the product formula (12.27) for the Brouwer degree can be used to prove the famous **Jordan curve theorem** which states that a homeomorphic image of the circle dissects  $\mathbb{R}^2$  into two components (which necessarily have the image of the circle as common boundary). In fact, we will even prove a slightly more general result.

**Theorem 12.20.** *Let  $C_j \subset \mathbb{R}^n$ ,  $j = 1, 2$ , be homeomorphic compact sets. Then  $\mathbb{R}^n \setminus C_1$  and  $\mathbb{R}^n \setminus C_2$  have the same number of connected components.*

**Proof.** Denote the components of  $\mathbb{R}^n \setminus C_1$  by  $H_j$  and those of  $\mathbb{R}^n \setminus C_2$  by  $K_j$ . Since our sets are closed these components are open. Moreover,  $\partial H_j \subseteq C_1$  since a sequence from  $H_j$  cannot converge to a (necessarily interior) point of  $H_k$  for some  $k \neq j$ . Let  $h : C_1 \rightarrow C_2$  be a homeomorphism with inverse  $k : C_2 \rightarrow C_1$ . By Theorem 13.10 we can extend both to  $\mathbb{R}^n$ . Then Theorem 12.1 (iii) and the product formula imply

$$1 = \deg(k \circ h, H_j, y) = \sum_l \deg(h, H_j, G_l) \deg(k, G_l, y)$$

for any  $y \in H_j$ , where  $G_l$  are the components of  $\mathbb{R}^n \setminus h(\partial H_j)$ . Now we have

$$\bigcup_i K_i = \mathbb{R}^n \setminus C_2 \subseteq \mathbb{R}^n \setminus h(\partial H_j) = \bigcup_l G_l$$

and hence for every  $i$  we have  $K_i \subseteq G_l$  for some  $l$  since components are connected. Let  $N_l := \{i | K_i \subseteq G_l\}$  and observe that we have  $\deg(k, G_l, y) = \sum_{i \in N_l} \deg(k, K_i, y)$  by Theorem 12.1 (i) since  $k^{-1}(\{y\}) \subseteq \bigcup_j K_j$  and of course  $\deg(h, H_j, K_i) = \deg(h, H_j, G_l)$  for every  $i \in N_l$ . Therefore,

$$1 = \sum_l \sum_{i \in N_l} \deg(h, H_j, K_i) \deg(k, K_i, y) = \sum_i \deg(h, H_j, K_i) \deg(k, K_i, H_j)$$

By reversing the role of  $C_1$  and  $C_2$ , the same formula holds with  $H_j$  and  $K_i$  interchanged.

Hence

$$\sum_i 1 = \sum_i \sum_j \deg(h, H_j, K_i) \deg(k, K_i, H_j) = \sum_j 1$$

shows that if either the number of components of  $\mathbb{R}^n \setminus C_1$  or the number of components of  $\mathbb{R}^n \setminus C_2$  is finite, then so is the other and both are equal. Otherwise there is nothing to prove.  $\square$

# The Leray–Schauder mapping degree

## 13.1. The mapping degree on finite dimensional Banach spaces

The objective of this section is to extend the mapping degree from  $\mathbb{R}^n$  to general Banach spaces. Naturally, we will first consider the finite dimensional case.

Let  $X$  be a (real) Banach space of dimension  $n$  and let  $\phi$  be any isomorphism between  $X$  and  $\mathbb{R}^n$ . Then, for  $f \in \bar{C}_y(U, X)$ ,  $U \subset X$  open,  $y \in X$ , we can define

$$\deg(f, U, y) := \deg(\phi \circ f \circ \phi^{-1}, \phi(U), \phi(y)) \quad (13.1)$$

provided this definition is independent of the isomorphism chosen. To see this let  $\psi$  be a second isomorphism. Then  $A = \psi \circ \phi^{-1} \in \text{GL}(n)$ . Abbreviate  $f^* = \phi \circ f \circ \phi^{-1}$ ,  $y^* = \phi(y)$  and pick  $\tilde{f}^* \in \bar{C}_{y^*}^1(\phi(U), \mathbb{R}^n)$  in the same component of  $\bar{C}_{y^*}(\phi(U), \mathbb{R}^n)$  as  $f^*$  such that  $y^* \in \text{RV}(\tilde{f}^*)$ . Then  $A \circ \tilde{f}^* \circ A^{-1} \in \bar{C}_y^1(\psi(U), \mathbb{R}^n)$  is in the same component of  $\bar{C}_y(\psi(U), \mathbb{R}^n)$  as  $A \circ f^* \circ A^{-1} = \psi \circ f \circ \psi^{-1}$  (since  $A$  is also a homeomorphism) and

$$J_{A \circ \tilde{f}^* \circ A^{-1}}(Ay^*) = \det(A) J_{\tilde{f}^*}(y^*) \det(A^{-1}) = J_{\tilde{f}^*}(y^*) \quad (13.2)$$

by the chain rule. Thus we have  $\deg(\psi \circ f \circ \psi^{-1}, \psi(U), \psi(y)) = \deg(\phi \circ f \circ \phi^{-1}, \phi(U), \phi(y))$  and our definition is independent of the basis chosen. In addition, it inherits all properties from the mapping degree in  $\mathbb{R}^n$ . Note also that the reduction property holds if  $\mathbb{R}^m$  is replaced by an arbitrary subspace  $X_1$  since we can always choose  $\phi : X \rightarrow \mathbb{R}^n$  such that  $\phi(X_1) = \mathbb{R}^m$ .

Our next aim is to tackle the infinite dimensional case. The following example due to Kakutani shows that the Brouwer fixed point theorem (and hence also the Brouwer degree) does not generalize to infinite dimensions directly.

**Example 13.1.** Let  $X$  be the Hilbert space  $\ell^2(\mathbb{N})$  and let  $R$  be the right shift given by  $Rx := (0, x_1, x_2, \dots)$ . Define  $f : \bar{B}_1(0) \rightarrow \bar{B}_1(0)$ ,  $x \mapsto \sqrt{1 - \|x\|^2}\delta_1 + Rx = (\sqrt{1 - \|x\|^2}, x_1, x_2, \dots)$ . Then a short calculation shows  $\|f(x)\|^2 = (1 - \|x\|^2) + \|x\|^2 = 1$  and any fixed point must satisfy  $\|x\| = 1$ ,  $x_1 = \sqrt{1 - \|x\|^2} = 0$  and  $x_{j+1} = x_j$ ,  $j \in \mathbb{N}$  giving the contradiction  $x_j = 0$ ,  $j \in \mathbb{N}$ .  $\diamond$

However, by the reduction property we expect that the degree should hold for functions of the type  $\mathbb{I} + F$ , where  $F$  has finite dimensional range. In fact, it should work for functions which can be approximated by such functions. Hence as a preparation we will investigate this class of functions.

### 13.2. Compact maps

Let  $X, Y$  be Banach spaces and  $U \subseteq X$ . A map  $F : U \subset X \rightarrow Y$  is called **finite dimensional** if its range is finite dimensional. In addition, it is called **compact** if it is continuous and maps bounded sets into relatively compact ones. The set of all compact maps is denoted by  $\mathcal{K}(U, Y)$  and the set of all compact, finite dimensional maps is denoted by  $\mathcal{F}(U, Y)$ . Both sets are normed linear spaces and we have  $\mathcal{K}(U, Y) \subseteq C_b(U, Y)$  if  $U$  is bounded (recall that compact sets are automatically bounded).

If  $K$  is compact, then  $\mathcal{K}(K, Y) = C(K, Y)$  (since the continuous image of a compact set is compact) and if  $\dim(Y) < \infty$ , then  $\mathcal{F}(U, Y) = \mathcal{K}(U, Y)$ . In particular, if  $U \subset \mathbb{R}^n$  is bounded, then  $\mathcal{F}(\bar{U}, \mathbb{R}^n) = \mathcal{K}(\bar{U}, \mathbb{R}^n) = C(\bar{U}, \mathbb{R}^n)$ .

**Example 13.2.** Note that for nonlinear functions it is important to include continuity in the definition of compactness. Indeed, if  $X$  is a Hilbert space with an orthonormal basis  $\{x_j\}_{j \in \mathbb{N}}$ , then

$$\phi(x) = \begin{cases} j(1 - 2|x - x_j|), & x \in B_{1/2}(x_j), \\ 0, & \text{else,} \end{cases}$$

is in  $C(B_1(0), \mathbb{R})$  but not bounded. Hence  $F(x) = \phi(x)x_1$  is one-dimensional but not compact. Choosing  $\phi(x) = 1$  for  $x \in B_{1/2}(0)$  and  $\phi(x) = 0$  else gives a map  $F$  which maps bounded sets to relatively compact ones but which is not continuous.  $\diamond$

Now let us collect some results needed in the sequel.

**Lemma 13.1.** *If  $K \subset X$  is compact, then for every  $\varepsilon > 0$  there is a finite dimensional subspace  $X_\varepsilon \subseteq X$  and a continuous map  $P_\varepsilon : K \rightarrow X_\varepsilon$  such that  $|P_\varepsilon(x) - x| \leq \varepsilon$  for all  $x \in K$ .*

**Proof.** Pick  $\{x_i\}_{i=1}^n \subseteq K$  such that  $\bigcup_{i=1}^n B_\varepsilon(x_i)$  covers  $K$ . Let  $\{\phi_i\}_{i=1}^n$  be a partition of unity (restricted to  $K$ ) subordinate to  $\{B_\varepsilon(x_i)\}_{i=1}^n$ , that is,  $\phi_i \in C(K, [0, 1])$  with  $\text{supp}(\phi_i) \subset B_\varepsilon(x_i)$  and  $\sum_{i=1}^n \phi_i(x) = 1$ ,  $x \in K$ . Set

$$P_\varepsilon(x) = \sum_{i=1}^n \phi_i(x)x_i,$$

then

$$|P_\varepsilon(x) - x| = \left| \sum_{i=1}^n \phi_i(x)x - \sum_{i=1}^n \phi_i(x)x_i \right| \leq \sum_{i=1}^n \phi_i(x)|x - x_i| \leq \varepsilon. \quad \square$$

This lemma enables us to prove the following important result.

**Theorem 13.2.** *Let  $U$  be bounded, then the closure of  $\mathcal{F}(U, Y)$  in  $C_b(U, Y)$  is  $\mathcal{K}(U, Y)$ .*

**Proof.** Suppose  $F_N \in \mathcal{K}(U, Y)$  converges to  $F$ . If  $F \notin \mathcal{K}(U, Y)$  then we can find a sequence  $x_n \in U$  such that  $|F(x_n) - F(x_m)| \geq \rho > 0$  for  $n \neq m$ . If  $N$  is so large that  $|F - F_N| \leq \rho/4$ , then

$$\begin{aligned} |F_N(x_n) - F_N(x_m)| &\geq |F(x_n) - F(x_m)| - |F_N(x_n) - F(x_n)| \\ &\quad - |F_N(x_m) - F(x_m)| \\ &\geq \rho - 2\frac{\rho}{4} = \frac{\rho}{2} \end{aligned}$$

This contradiction shows  $\overline{\mathcal{F}(U, Y)} \subseteq \mathcal{K}(U, Y)$ . Conversely, let  $F \in \mathcal{K}(U, Y)$ , set  $K := \overline{F(U)}$  and choose  $P_\varepsilon$  according to Lemma 13.1. Then  $F_\varepsilon = P_\varepsilon \circ F \in \mathcal{F}(U, Y)$  converges to  $F$ . Hence  $\mathcal{K}(U, Y) \subseteq \overline{\mathcal{F}(U, Y)}$  and we are done.  $\square$

Finally, let us show some interesting properties of mappings  $\mathbb{I} + F$ , where  $F \in \mathcal{K}(U, Y)$ .

**Lemma 13.3.** *Let  $\overline{U} \subseteq X$  be bounded and closed. Suppose  $F \in \mathcal{K}(\overline{U}, X)$ , then  $\mathbb{I} + F$  is **proper** (i.e., inverse images of compact sets are compact) and maps closed subsets to closed subsets.*

**Proof.** Let  $A \subseteq \overline{U}$  be closed and suppose  $y_n = (\mathbb{I} + F)(x_n) \in (\mathbb{I} + F)(A)$  converges to some point  $y$ . Since  $y_n - x_n = F(x_n) \in F(U)$  we can assume that  $y_n - x_n \rightarrow z$  after passing to a subsequence and hence  $x_n \rightarrow x = y - z \in A$ . Since  $y = x + F(x) \in (\mathbb{I} + F)(A)$ ,  $(\mathbb{I} + F)(A)$  is closed.

Next, let  $\overline{U}$  be closed and  $K \subset X$  be compact. Let  $\{x_n\} \subseteq (\mathbb{I} + F)^{-1}(K)$ . Then  $y_n := x_n + F(x_n) \in K$  and we can pass to a subsequence  $y_{n_m} =$

$x_{n_m} + F(x_{n_m})$  such that  $y_{n_m} \rightarrow y$ . As before this implies  $x_{n_m} \rightarrow x$  and thus  $(\mathbb{I} + F)^{-1}(K)$  is compact.  $\square$

Finally note that if  $F \in \mathcal{K}(\overline{U}, Y)$  and  $G \in C(Y, Z)$ , then  $G \circ F \in \mathcal{K}(\overline{U}, Z)$  and similarly, if  $G \in C_b(\overline{V}, \overline{U})$ , then  $F \circ G \in \mathcal{K}(\overline{V}, Y)$ .

Now we are all set for the definition of the Leray–Schauder degree, that is, for the extension of our degree to infinite dimensional Banach spaces.

### 13.3. The Leray–Schauder mapping degree

For an open set  $U \subset X$  we set

$$\bar{\mathcal{K}}_y(U, X) := \{F \in \mathcal{K}(\overline{U}, X) \mid y \notin (\mathbb{I} + F)(\partial U)\} \quad (13.3)$$

and  $\bar{\mathcal{F}}_y(U, X) := \{F \in \mathcal{F}(\overline{U}, X) \mid y \notin (\mathbb{I} + F)(\partial U)\}$ . Note that for  $F \in \bar{\mathcal{K}}_y(U, X)$  we have  $\text{dist}(y, (\mathbb{I} + F)(\partial U)) > 0$  since  $\mathbb{I} + F$  maps closed sets to closed sets (cf. Problem B.28).

Abbreviate  $\rho := \text{dist}(y, (\mathbb{I} + F)(\partial U))$  and pick  $F_1 \in \mathcal{F}(\overline{U}, X)$  such that  $|F - F_1| < \rho$  implying  $F_1 \in \bar{\mathcal{F}}_y(U, X)$ . Next, let  $X_1$  be a finite dimensional subspace of  $X$  such that  $F_1(U) \subset X_1$ ,  $y \in X_1$  and set  $U_1 := U \cap X_1$ . Then we have  $F_1 \in \bar{\mathcal{F}}_y(U_1, X_1)$  and might define

$$\deg(\mathbb{I} + F, U, y) := \deg(\mathbb{I} + F_1, U_1, y) \quad (13.4)$$

provided we show that this definition is independent of  $F_1$  and  $X_1$  (as above). Pick another map  $F_2 \in \mathcal{F}(\overline{U}, X)$  such that  $|F - F_2| < \rho$  and let  $X_2$  be a corresponding finite dimensional subspace as above. Consider  $X_0 := X_1 + X_2$ ,  $U_0 = U \cap X_0$ , then  $F_i \in \bar{\mathcal{F}}_y(U_0, X_0)$ ,  $i = 1, 2$ , and

$$\deg(\mathbb{I} + F_i, U_0, y) = \deg(\mathbb{I} + F_i, U_i, y), \quad i = 1, 2, \quad (13.5)$$

by the reduction property. Moreover, set  $H(t) = \mathbb{I} + (1 - t)F_1 + tF_2$  implying  $H(t) \in \bar{\mathcal{K}}_y(U_0, X_0)$ ,  $t \in [0, 1]$ , since  $|H(t) - (\mathbb{I} + F)| < \rho$  for  $t \in [0, 1]$ . Hence homotopy invariance

$$\deg(\mathbb{I} + F_1, U_0, y) = \deg(\mathbb{I} + F_2, U_0, y) \quad (13.6)$$

shows that (13.4) is independent of  $F_1$ ,  $X_1$ .

**Theorem 13.4.** *Let  $U$  be a bounded open subset of a (real) Banach space  $X$  and let  $F \in \bar{\mathcal{K}}_y(U, X)$ ,  $y \in X$ . Then the following hold true.*

- (i).  $\deg(\mathbb{I} + F, U, y) = \deg(\mathbb{I} + F - y, U, 0)$ .
- (ii).  $\deg(\mathbb{I}, U, y) = 1$  if  $y \in U$ .
- (iii). If  $U_{1,2}$  are open, disjoint subsets of  $U$  such that  $y \notin f(\overline{U} \setminus (U_1 \cup U_2))$ , then  $\deg(\mathbb{I} + F, U, y) = \deg(\mathbb{I} + F, U_1, y) + \deg(\mathbb{I} + F, U_2, y)$ .

- (iv). If  $H : [0, 1] \times \bar{U} \rightarrow X$  and  $y : [0, 1] \rightarrow X$  are both continuous such that  $H(t) \in \bar{\mathcal{K}}_{y(t)}(U, X)$ ,  $t \in [0, 1]$ , then  $\deg(\mathbb{I} + H(0), U, y(0)) = \deg(\mathbb{I} + H(1), U, y(1))$ .

**Proof.** Except for (iv) all statements follow easily from the definition of the degree and the corresponding property for the degree in finite dimensional spaces. Considering  $H(t, x) - y(t)$ , we can assume  $y(t) = 0$  by (i). Since  $H([0, 1], \partial U)$  is compact, we have  $\rho = \text{dist}(y, H([0, 1], \partial U)) > 0$ . By Theorem 13.2 we can pick  $H_1 \in \mathcal{F}([0, 1] \times U, X)$  such that  $|H(t) - H_1(t)| < \rho$ ,  $t \in [0, 1]$ . This implies  $\deg(\mathbb{I} + H(t), U, 0) = \deg(\mathbb{I} + H_1(t), U, 0)$  and the rest follows from Theorem 12.2.  $\square$

In addition, Theorem 12.1 and Theorem 12.2 hold for the new situation as well (no changes are needed in the proofs).

**Theorem 13.5.** Let  $F, G \in \bar{\mathcal{K}}_y(U, X)$ , then the following statements hold.

- (i). We have  $\deg(\mathbb{I} + F, \emptyset, y) = 0$ . Moreover, if  $U_i$ ,  $1 \leq i \leq N$ , are disjoint open subsets of  $U$  such that  $y \notin (\mathbb{I} + F)(\bar{U} \setminus \bigcup_{i=1}^N U_i)$ , then  $\deg(\mathbb{I} + F, U, y) = \sum_{i=1}^N \deg(\mathbb{I} + F, U_i, y)$ .
- (ii). If  $y \notin (\mathbb{I} + F)(U)$ , then  $\deg(\mathbb{I} + F, U, y) = 0$  (but not the other way round). Equivalently, if  $\deg(\mathbb{I} + F, U, y) \neq 0$ , then  $y \in (\mathbb{I} + F)(U)$ .
- (iii). If  $|F(x) - G(x)| < \text{dist}(y, (\mathbb{I} + F)(\partial U))$ ,  $x \in \partial U$ , then  $\deg(\mathbb{I} + F, U, y) = \deg(\mathbb{I} + G, U, y)$ . In particular, this is true if  $F(x) = G(x)$  for  $x \in \partial U$ .
- (iv).  $\deg(\mathbb{I} + \cdot, U, y)$  is constant on each component of  $\bar{\mathcal{K}}_y(U, X)$ .
- (v).  $\deg(\mathbb{I} + F, U, \cdot)$  is constant on each component of  $X \setminus (\mathbb{I} + F)(\partial U)$ .

Note that it is easy to generalize Borsuk's theorem.

**Theorem 13.6** (Borsuk). Let  $U \subseteq X$  be open, bounded and symmetric with respect to the origin (i.e.,  $U = -U$ ). Let  $F \in \bar{\mathcal{K}}_0(U, X)$  be odd (i.e.,  $F(-x) = -F(x)$ ). Then  $\deg(\mathbb{I} + F, U, 0)$  is odd.

**Proof.** Choose  $F_1$  and  $U_1$  as in the definition of the degree. Then  $U_1$  is symmetric and  $F_1$  can be chosen to be odd by replacing it by its odd part. Hence the claim follows from the finite dimensional version.  $\square$

In the same way as in the finite dimensional case we also obtain the **invariance of domain theorem**.

**Theorem 13.7** (Brouwer). Let  $U \subseteq X$  be open and let  $F \in \mathcal{K}(U, X)$  be compact with  $\mathbb{I} + F$  locally injective. Then  $\mathbb{I} + F$  is also open.

### 13.4. The Leray–Schauder principle and the Schauder fixed point theorem

As a first consequence we note the Leray–Schauder principle which says that a priori estimates yield existence.

**Theorem 13.8** (Schaefer fixed point or Leray–Schauder principle). *Suppose  $F \in \mathcal{K}(X, X)$  and any solution  $x$  of  $x = tF(x)$ ,  $t \in [0, 1]$  satisfies the a priori bound  $|x| \leq M$  for some  $M > 0$ , then  $F$  has a fixed point.*

**Proof.** Pick  $\rho > M$  and observe  $\deg(\mathbb{I} - F, B_\rho(0), 0) = \deg(\mathbb{I}, B_\rho(0), 0) = 1$  using the compact homotopy  $H(t, x) := -tF(x)$ . Here  $H(t) \in \tilde{\mathcal{K}}_0(B_\rho(0), X)$  due to the a priori bound.  $\square$

Now we can extend the Brouwer fixed point theorem to infinite dimensional spaces as well.

**Theorem 13.9** (Schauder fixed point). *Let  $K$  be a closed, convex, and bounded subset of a Banach space  $X$ . If  $F \in \mathcal{K}(K, K)$ , then  $F$  has at least one fixed point. The result remains valid if  $K$  is only homeomorphic to a closed, convex, and bounded subset.*

**Proof.** Since  $K$  is bounded, there is a  $\rho > 0$  such that  $K \subseteq B_\rho(0)$ . By Theorem 13.10 below we can find a continuous retraction  $R : X \rightarrow K$  (i.e.,  $R(x) = x$  for  $x \in K$ ) and consider  $\tilde{F} = F \circ R \in \mathcal{K}(\overline{B_\rho(0)}, \overline{B_\rho(0)})$ . Now either  $t\tilde{F}$  has a fixed point on the boundary  $\partial B_\rho(0)$  or the compact homotopy  $H(t, x) := -t\tilde{F}(x)$  satisfies  $0 \notin (\mathbb{I} - t\tilde{F})(\partial B_\rho(0))$  and thus  $\deg(\mathbb{I} - \tilde{F}, B_\rho(0), 0) = \deg(\mathbb{I}, B_\rho(0), 0) = 1$ . Hence there is a point  $x_0 = \tilde{F}(x_0) \in K$ . Since  $\tilde{F}(x_0) = F(x_0)$  for  $x_0 \in K$  we are done.  $\square$

It remains to prove the following variant of the **Tietze extension theorem** needed in the proof.

**Theorem 13.10.** *Let  $X$  be a metric space,  $Y$  a normed space and let  $K$  be a closed subset of  $X$ . Then  $F \in C(K, Y)$  has a continuous extension  $\bar{F} \in C(X, Y)$  such that  $\bar{F}(X) \subseteq \text{conv}(F(K))$ .*

**Proof.** Consider the open cover  $\{B_{\rho(x)}(x)\}_{x \in X \setminus K}$  for  $X \setminus K$ , where  $\rho(x) = \text{dist}(x, K)/2$ . Choose a locally finite refinement  $\{O_\lambda\}_{\lambda \in \Lambda}$  of this cover (see Lemma B.14) and define

$$\phi_\lambda(x) := \frac{\text{dist}(x, X \setminus O_\lambda)}{\sum_{\mu \in \Lambda} \text{dist}(x, X \setminus O_\mu)}.$$

Set

$$\bar{F}(x) := \sum_{\lambda \in \Lambda} \phi_\lambda(x) F(x_\lambda) \text{ for } x \in X \setminus K,$$

where  $x_\lambda \in K$  satisfies  $\text{dist}(x_\lambda, O_\lambda) \leq 2 \text{dist}(K, O_\lambda)$ . By construction,  $\bar{F}$  is continuous except for possibly at the boundary of  $K$ . Fix  $x_0 \in \partial K$ ,  $\varepsilon > 0$  and choose  $\delta > 0$  such that  $|F(x) - F(x_0)| \leq \varepsilon$  for all  $x \in K$  with  $|x - x_0| < 4\delta$ . We will show that  $|\bar{F}(x) - F(x_0)| \leq \varepsilon$  for all  $x \in X$  with  $|x - x_0| < \delta$ . Suppose  $x \notin K$ , then  $|\bar{F}(x) - F(x_0)| \leq \sum_{\lambda \in \Lambda} \phi_\lambda(x) |F(x_\lambda) - F(x_0)|$ . By our construction,  $x_\lambda$  should be close to  $x$  for all  $\lambda$  with  $x \in \overline{O_\lambda}$  since  $x$  is close to  $K$ . In fact, if  $x \in \overline{O_\lambda}$  we have

$$\begin{aligned} |x - x_\lambda| &\leq \text{dist}(x_\lambda, O_\lambda) + \text{diam}(O_\lambda) \\ &\leq 2 \text{dist}(K, O_\lambda) + \text{diam}(O_\lambda), \end{aligned}$$

where  $\text{diam}(O_\lambda) := \sup_{x, y \in O_\lambda} |x - y|$ . Since our partition of unity is subordinate to the cover  $\{B_{\rho(x)}(x)\}_{x \in X \setminus K}$  we can find a  $\tilde{x} \in X \setminus K$  such that  $O_\lambda \subset B_{\rho(\tilde{x})}(\tilde{x})$  and hence  $\text{diam}(O_\lambda) \leq 2\rho(\tilde{x}) \leq \text{dist}(K, B_{\rho(\tilde{x})}(\tilde{x})) \leq \text{dist}(K, O_\lambda)$ . Putting it all together implies that we have  $|x - x_\lambda| \leq 3 \text{dist}(K, O_\lambda) \leq 3|x_0 - x|$  whenever  $x \in \overline{O_\lambda}$  and thus

$$|x_0 - x_\lambda| \leq |x_0 - x| + |x - x_\lambda| \leq 4|x_0 - x| \leq 4\delta$$

as expected. By our choice of  $\delta$  we have  $|F(x_\lambda) - F(x_0)| \leq \varepsilon$  for all  $\lambda$  with  $\phi_\lambda(x) \neq 0$ . Hence  $|F(x) - F(x_0)| \leq \varepsilon$  whenever  $|x - x_0| \leq \delta$  and we are done.  $\square$

**Example 13.3.** Consider the nonlinear integral equation

$$x = F(x), \quad F(x)(t) := \int_0^1 e^{-ts} \cos(\lambda x(s)) ds$$

in  $X := C[0, 1]$  with  $\lambda > 0$ . Then one checks that  $F \in C(X, X)$  since

$$\begin{aligned} |F(x)(t) - F(y)(t)| &\leq \int_0^1 e^{-ts} |\cos(\lambda x(s)) - \cos(\lambda y(s))| ds \\ &\leq \int_0^1 e^{-ts} \lambda |x(s) - y(s)| ds \leq \lambda \|x - y\|_\infty. \end{aligned}$$

In particular, for  $\lambda < 1$  we have a contraction and the contraction principle gives us existence of a unique fixed point. Moreover, proceeding similarly, one obtains estimates for the norm of  $F(x)$  and its derivative:

$$\|F(x)\|_\infty \leq 1, \quad \|F(x)'\|_\infty \leq 1.$$

Hence the Arzelà–Ascoli theorem (Theorem B.40) implies that the image of  $F$  is a compact subset of the unit ball and hence  $F \in \mathcal{K}(\bar{B}_1(0), \bar{B}_1(0))$ . Thus the Schauder fixed point theorem guarantees a fixed point for all  $\lambda > 0$ .  $\diamond$

Finally, let us prove another fixed point theorem which covers several others as special cases.



**Theorem 13.11.** *Let  $U \subset X$  be open and bounded and let  $F \in \mathcal{K}(\bar{U}, X)$ . Suppose there is an  $x_0 \in U$  such that*

$$F(x) - x_0 \neq \alpha(x - x_0), \quad x \in \partial U, \alpha \in (1, \infty). \quad (13.7)$$

*Then  $F$  has a fixed point.*

**Proof.** Consider  $H(t, x) := x - x_0 - t(F(x) - x_0)$ , then we have  $H(t, x) \neq 0$  for  $x \in \partial U$  and  $t \in [0, 1]$  by assumption. If  $H(1, x) = 0$  for some  $x \in \partial U$ , then  $x$  is a fixed point and we are done. Otherwise we have  $\deg(\mathbb{I} - F, U, 0) = \deg(\mathbb{I} - x_0, U, 0) = \deg(\mathbb{I}, U, x_0) = 1$  and hence  $F$  has a fixed point.  $\square$

Now we come to the anticipated corollaries.

**Corollary 13.12.** *Let  $F \in \mathcal{K}(\bar{B}_\rho(0), X)$ . Then  $F$  has a fixed point if one of the following conditions holds.*

- (i)  $F(\partial B_\rho(0)) \subseteq \bar{B}_\rho(0)$  (Rothe).
- (ii)  $|F(x) - x|^2 \geq |F(x)|^2 - |x|^2$  for  $x \in \partial B_\rho(0)$  (Altman).
- (iii)  $X$  is a Hilbert space and  $\langle F(x), x \rangle \leq |x|^2$  for  $x \in \partial B_\rho(0)$  (Krasnosel'skii).

**Proof.** Our strategy is to verify (13.7) with  $x_0 = 0$ . (i).  $F(\partial B_\rho(0)) \subseteq \bar{B}_\rho(0)$  and  $F(x) = \alpha x$  for  $|x| = \rho$  implies  $|\alpha|\rho \leq \rho$  and hence (13.7) holds. (ii).  $F(x) = \alpha x$  for  $|x| = \rho$  implies  $(\alpha - 1)^2 \rho^2 \geq (\alpha^2 - 1)\rho^2$  and hence  $\alpha \leq 1$ . (iii). Special case of (ii) since  $|F(x) - x|^2 = |F(x)|^2 - 2\langle F(x), x \rangle + |x|^2$ .  $\square$

### 13.5. Applications to integral and differential equations

In this section we want to show how our results can be applied to integral and differential equations. To be able to apply our results we will need to know that certain integral operators are compact.

**Lemma 13.13.** *Suppose  $I = [a, b] \subset \mathbb{R}$  and  $f \in C(I \times I \times \mathbb{R}^n, \mathbb{R}^n)$ ,  $\tau \in C(I, I)$ , then*

$$\begin{aligned} F : C(I, \mathbb{R}^n) &\rightarrow C(I, \mathbb{R}^n) \\ x(t) &\mapsto F(x)(t) = \int_a^{\tau(t)} f(t, s, x(s)) ds \end{aligned} \quad (13.8)$$

*is compact.*

**Proof.** We first need to prove that  $F$  is continuous. Fix  $x_0 \in C(I, \mathbb{R}^n)$  and  $\varepsilon > 0$ . Set  $\rho := \|x_0\|_\infty + 1$  and abbreviate  $\bar{B} = \bar{B}_\rho(0) \subset \mathbb{R}^n$ . The function  $f$  is uniformly continuous on  $Q := I \times I \times \bar{B}$  since  $Q$  is compact. Hence for

$\varepsilon_1 := \varepsilon/(b-a)$  we can find a  $\delta \in (0, 1]$  such that  $|f(t, s, x) - f(t, s, y)| \leq \varepsilon_1$  for  $|x - y| < \delta$ . But this implies

$$\begin{aligned} \|F(x) - F(x_0)\|_\infty &\leq \sup_{t \in I} \int_a^{\tau(t)} |f(t, s, x(s)) - f(t, s, x_0(s))| ds \\ &\leq \sup_{t \in I} (b-a) \varepsilon_1 = \varepsilon, \end{aligned}$$

for  $\|x - x_0\|_\infty < \delta$ . In other words,  $F$  is continuous. Next we note that if  $U \subset C(I, \mathbb{R}^n)$  is bounded, say  $U \subset \bar{B}_\rho(0)$ , then

$$\|F(x)\|_\infty \leq \sup_{x \in U} \left| \int_a^{\tau(t)} f(t, s, x(s)) ds \right| \leq (b-a)M, \quad x \in U,$$

where  $M := \max |f(I, I, \bar{B})|$ . Moreover, the family  $F(U)$  is equicontinuous. Fix  $\varepsilon$  and  $\varepsilon_1 := \varepsilon/(2(b-a))$ ,  $\varepsilon_2 := \varepsilon/(2M)$ . Since  $f$  and  $\tau$  are uniformly continuous on  $I \times I \times \bar{B}$  and  $I$ , respectively, we can find a  $\delta > 0$  such that  $|f(t, s, x) - f(t_0, s, x)| \leq \varepsilon_1$  and  $|\tau(t) - \tau(t_0)| \leq \varepsilon_2$  for  $|t - t_0| < \delta$ . Hence we infer for  $|t - t_0| < \delta$

$$\begin{aligned} |F(x)(t) - F(x)(t_0)| &= \left| \int_a^{\tau(t)} f(t, s, x(s)) ds - \int_a^{\tau(t_0)} f(t_0, s, x(s)) ds \right| \\ &\leq \int_a^{\tau(t_0)} |f(t, s, x(s)) - f(t_0, s, x(s))| ds + \left| \int_{\tau(t_0)}^{\tau(t)} f(t, s, x(s)) ds \right| \\ &\leq (b-a)\varepsilon_1 + \varepsilon_2 M = \varepsilon. \end{aligned}$$

This implies that  $F(U)$  is relatively compact by the Arzelà–Ascoli theorem (Theorem B.40). Thus  $F$  is compact.  $\square$

As a first application we use this result to show existence of solutions to integral equations.

**Theorem 13.14.** *Let  $F$  be as in the previous lemma. Then the integral equation*

$$x - \lambda F(x) = y, \quad \lambda \in \mathbb{R}, y \in C(I, \mathbb{R}^n) \quad (13.9)$$

*has at least one solution  $x \in C(I, \mathbb{R}^n)$  if  $|\lambda| \leq \frac{\rho}{(b-a)M(\rho)}$ , where  $M(\rho) = \max_{(s,t,x) \in I \times I \times \bar{B}_\rho(0)} |f(s, t, x - y(s))|$  and  $\rho > 0$  is arbitrary.*

**Proof.** Note that, by our assumption on  $\lambda$ ,  $\lambda F + y$  maps  $\bar{B}_\rho(y)$  into itself. Now apply the Schauder fixed point theorem.  $\square$

This result immediately gives the Peano theorem for ordinary differential equations.

**Theorem 13.15** (Peano). *Consider the initial value problem*

$$\dot{x} = f(t, x), \quad x(t_0) = x_0, \quad (13.10)$$

where  $f \in C(I \times U, \mathbb{R}^n)$  and  $I \subset \mathbb{R}$  is an interval containing  $t_0$ . Then (13.10) has at least one local solution  $x \in C^1([t_0 - \varepsilon, t_0 + \varepsilon], \mathbb{R}^n)$ ,  $\varepsilon > 0$ . For example, any  $\varepsilon$  satisfying  $\varepsilon M(\varepsilon, \rho) \leq \rho$ ,  $\rho > 0$  with  $M(\varepsilon, \rho) := \max |f|([t_0 - \varepsilon, t_0 + \varepsilon] \times \bar{B}_\rho(x_0))$  works. In addition, if  $M(\varepsilon, \rho) \leq \tilde{M}(\varepsilon)(1 + \rho)$ , then there exists a global solution.

**Proof.** For notational simplicity we make the shift  $t \rightarrow t - t_0$ ,  $x \rightarrow x - x_0$ ,  $f(t, x) \rightarrow f(t + t_0, x + t_0)$  and assume  $t_0 = 0$ ,  $x_0 = 0$ . In addition, it suffices to consider  $t \geq 0$  since  $t \rightarrow -t$  amounts to  $f \rightarrow -f$ .

Now observe, that (13.10) is equivalent to

$$x(t) - \int_0^t f(s, x(s)) ds = 0, \quad x \in C([0, \varepsilon], \mathbb{R}^n)$$

and the first part follows from our previous theorem. To show the second, fix  $\varepsilon > 0$  and assume  $M(\varepsilon, \rho) \leq \tilde{M}(\varepsilon)(1 + \rho)$ . Then

$$|x(t)| \leq \int_0^t |f(s, x(s))| ds \leq \tilde{M}(\varepsilon) \int_0^t (1 + |x(s)|) ds$$

implies  $|x(t)| \leq \exp(\tilde{M}(\varepsilon)\varepsilon)$  by Gronwall's inequality (Problem 13.1). Hence we have an a priori bound which implies existence by the Leray–Schauder principle. Since  $\varepsilon$  was arbitrary we are done.  $\square$

As another example we look at the stationary Navier–Stokes equation. Our goal is to use the Leray–Schauder principle to prove an existence and uniqueness result for solutions.

Let  $U (\neq \emptyset)$  be an open, bounded, and connected subset of  $\mathbb{R}^3$ . We assume that  $U$  is filled with an incompressible fluid described by its velocity field  $v(t, x) \in \mathbb{R}^3$  and its pressure  $p(t, x) \in \mathbb{R}$  at time  $t \in \mathbb{R}$  and at a point  $x \in U$ . The requirement that the fluid is incompressible implies  $\nabla \cdot v = 0$  (here we use a dot to emphasize a scalar product in  $\mathbb{R}^3$ ), which follows from the Gauss theorem since the flux through any closed surface must be zero. Moreover, the outer force density (force per volume) will be denoted by  $K(x) \in \mathbb{R}^3$  and is assumed to be known (e.g. gravity).

Then the **Navier–Stokes equation** governing the motion of the fluid reads

$$\rho \partial_t v = \eta \Delta v - \rho(v \cdot \nabla)v - \nabla p + K, \quad (13.11)$$

where  $\eta > 0$  is the viscosity constant and  $\rho > 0$  is the density of the fluid. In addition to the incompressibility condition  $\nabla \cdot v = 0$  we also require the boundary condition  $v|_{\partial U} = 0$ , which follows from experimental observations.

In what follows we will only consider the stationary Navier–Stokes equation

$$0 = \eta \Delta v - \rho(v \cdot \nabla)v - \nabla p + K. \quad (13.12)$$

Our first step is to switch to a weak formulation and rewrite this equation in integral form, which is more suitable for our further analysis. We pick as underlying Hilbert space  $H_0^1(U, \mathbb{R}^3)$  with scalar product

$$\langle u, v \rangle = \sum_{i,j=1}^3 \int_U (\partial_j u_i)(\partial_j v_i) dx. \quad (13.13)$$

Recall that by the Poincaré inequality (Theorem 7.38 from [48]) the corresponding norm is equivalent to the usual one. In order to take care of the incompressibility condition we will choose

$$\mathcal{X} := \{v \in H_0^1(U, \mathbb{R}^3) \mid \nabla \cdot v = 0\}. \quad (13.14)$$

as our configuration space (check that this is a closed subspace of  $H_0^1(U, \mathbb{R}^3)$ ).

Now we multiply (13.12) by  $w \in \mathcal{X}$  and integrate over  $U$

$$\begin{aligned} \int_U (\eta \Delta v - \rho(v \cdot \nabla)v - K) \cdot w \, d^3x &= \int_U (\nabla p) \cdot w \, d^3x \\ &= \int_U p(\nabla w) \, d^3x = 0, \end{aligned} \quad (13.15)$$

where we have used integration by parts (Lemma 7.11 from [48] (iii)) to conclude that the pressure term drops out of our picture. Using further integration by parts we finally arrive at the weak formulation of the stationary Navier–Stokes equation

$$\eta \langle v, w \rangle - a(v, v, w) - \int_U K \cdot w \, d^3x = 0, \quad \text{for all } w \in \mathcal{X}, \quad (13.16)$$

where

$$a(u, v, w) := \sum_{j,k=1}^3 \int_U u_k v_j (\partial_k w_j) \, d^3x. \quad (13.17)$$

In other words, (13.16) represents a necessary solubility condition for the Navier–Stokes equations and a solution of (13.16) will also be called a **weak solution**. If we can show that a weak solution is in  $H^2(U, \mathbb{R}^3)$ , then we can undo the integration by parts and obtain again (13.15). Since the integral on the left-hand side vanishes for all  $w \in \mathcal{X}$ , one can conclude that the expression in parenthesis must be the gradient of some function  $p \in L^2(U, \mathbb{R})$  and hence one recovers the original equation. In particular, note that  $p$  follows from  $v$  up to a constant if  $U$  is connected.

For later use we note

$$\begin{aligned} a(v, v, v) &= \sum_{j,k} \int_U v_k v_j (\partial_k v_j) d^3x = \frac{1}{2} \sum_{j,k} \int_U v_k \partial_k (v_j v_j) d^3x \\ &= -\frac{1}{2} \sum_{j,k} \int_U (v_j v_j) \partial_k v_k d^3x = 0, \quad v \in \mathcal{X}. \end{aligned} \quad (13.18)$$

We proceed by studying (13.16). Let  $K \in L^2(U, \mathbb{R}^3)$ , then  $\int_U K \cdot w d^3x$  is a bounded linear functional on  $\mathcal{X}$  and hence there is a  $\tilde{K} \in \mathcal{X}$  such that

$$\int_U K \cdot w d^3x = \langle \tilde{K}, w \rangle, \quad w \in \mathcal{X}. \quad (13.19)$$

Moreover, applying the Cauchy–Schwarz inequality twice to each summand in  $a(u, v, w)$  we see

$$\begin{aligned} |a(u, v, w)| &\leq \sum_{j,k} \left( \int_U (u_k v_j)^2 dx \right)^{1/2} \left( \int_U (\partial_k w_j)^2 dx \right)^{1/2} \\ &\leq \|w\| \sum_{j,k} \left( \int_U (u_k)^4 dx \right)^{1/4} \left( \int_U (v_j)^4 dx \right)^{1/4} = \|u\|_4 \|v\|_4 \|w\|. \end{aligned} \quad (13.20)$$

Since by the Gagliardo–Nirenberg–Sobolev inequality (Theorem 7.26 from [48]) there is a continuous embedding  $H^1(U, \mathbb{R}^3) \hookrightarrow L^4(U, \mathbb{R}^3)$  (which in this context also known as Ladyzhenskaya inequality), the map  $a(u, v, \cdot)$  is a bounded linear functional in  $\mathcal{X}$  whenever  $u, v \in \mathcal{X}$ , and hence there is an element  $B(u, v) \in \mathcal{X}$  such that

$$a(u, v, w) = \langle B(u, v), w \rangle, \quad w \in \mathcal{X}. \quad (13.21)$$

In addition, the map  $B : \mathcal{X}^2 \rightarrow \mathcal{X}$  is bilinear and bounded  $\|B(u, v)\| \leq \|u\|_4 \|v\|_4$ . In summary we obtain

$$\langle \eta v - B(v, v) - \tilde{K}, w \rangle = 0, \quad w \in \mathcal{X}, \quad (13.22)$$

and hence

$$\eta v - B(v, v) = \tilde{K}. \quad (13.23)$$

So in order to apply the theory from our previous chapter, we choose the Banach space  $Y := L^4(U, \mathbb{R}^3)$  such that  $\mathcal{X} \hookrightarrow Y$  is compact by the Rellich–Kondrachov theorem (Theorem 7.35 from [48]).

Motivated by this analysis we formulate the following theorem which implies existence of weak solutions and uniqueness for sufficiently small outer forces.

**Theorem 13.16.** *Let  $\mathcal{X}$  be a Hilbert space,  $Y$  a Banach space, and suppose there is a compact embedding  $\mathcal{X} \hookrightarrow Y$ . In particular,  $\|u\|_Y \leq \beta\|u\|$ . Let  $a : \mathcal{X}^3 \rightarrow \mathbb{R}$  be a multilinear form such that*

$$|a(u, v, w)| \leq \alpha\|u\|_Y\|v\|_Y\|w\| \quad (13.24)$$

*and  $a(v, v, v) = 0$ . Then for any  $\tilde{K} \in \mathcal{X}$ ,  $\eta > 0$  we have a solution  $v \in \mathcal{X}$  to the problem*

$$\eta\langle v, w \rangle - a(v, v, w) = \langle \tilde{K}, w \rangle, \quad w \in \mathcal{X}. \quad (13.25)$$

*Moreover, if  $2\alpha\beta\|\tilde{K}\| < \eta^2$  this solution is unique.*

**Proof.** It is no loss to set  $\eta = 1$ . Arguing as before we see that our equation is equivalent to

$$v - B(v, v) + \tilde{K} = 0,$$

where our assumption (13.24) implies

$$\|B(u, v)\| \leq \alpha\|u\|_Y\|v\|_Y \leq \alpha\beta^2\|u\|\|v\|$$

Here the second equality follows since the embedding  $\mathcal{X} \hookrightarrow Y$  is continuous.

Abbreviate  $F(v) = B(v, v)$ . Observe that  $F$  is locally Lipschitz continuous since if  $\|u\|, \|v\| \leq \rho$  we have

$$\begin{aligned} \|F(u) - F(v)\| &= \|B(u - v, u) + B(v, u - v)\| \leq 2\alpha\beta\rho\|u - v\|_Y \\ &\leq 2\alpha\beta^2\rho\|u - v\|. \end{aligned}$$

Moreover, let  $v_n$  be a bounded sequence in  $\mathcal{X}$ . After passing to a subsequence we can assume that  $v_n$  is Cauchy in  $Y$  and hence  $F(v_n)$  is Cauchy in  $\mathcal{X}$  by  $\|F(u) - F(v)\| \leq 2\alpha\beta\rho\|u - v\|_Y$ . Thus  $F : \mathcal{X} \rightarrow \mathcal{X}$  is compact.

Hence all we need to apply the Leray–Schauder principle is an a priori estimate. Suppose  $v$  solves  $v = tF(v) + t\tilde{K}$ ,  $t \in [0, 1]$ , then

$$\langle v, v \rangle = t a(v, v, v) + t \langle \tilde{K}, v \rangle = t \langle \tilde{K}, v \rangle.$$

Hence  $\|v\| \leq \|\tilde{K}\|$  is the desired estimate and the Leray–Schauder principle yields existence of a solution.

Now suppose there are two solutions  $v_i$ ,  $i = 1, 2$ . By our estimate they satisfy  $\|v_i\| \leq \|\tilde{K}\|$  and hence  $\|v_1 - v_2\| = \|F(v_1) - F(v_2)\| \leq 2\alpha\beta^2\|\tilde{K}\|\|v_1 - v_2\|$  which is a contradiction if  $2\alpha\beta^2\|\tilde{K}\| < 1$ .  $\square$

**Problem\* 13.1** (Gronwall's inequality). *Let  $\alpha \geq 0$  and  $\beta, \varphi : [0, T] \rightarrow [0, \infty)$  be integrable functions satisfying*

$$\varphi(t) \leq \alpha + \int_0^t \beta(s)\varphi(s)ds.$$

*Then  $\varphi(t) \leq \alpha e^{\int_0^t \beta(s)ds}$ . (Hint: Differentiate  $\log(\alpha + \int_0^t \beta(s)\varphi(s)ds)$ .)*



# Monotone maps

## 14.1. Monotone maps

The Leray–Schauder theory can only be applied to compact perturbations of the identity. If  $F$  is not compact, we need different tools. In this section we briefly present another class of maps, namely monotone ones, which allow some progress.

If  $F : \mathbb{R} \rightarrow \mathbb{R}$  is continuous and we want  $F(x) = y$  to have a unique solution for every  $y \in \mathbb{R}$ , then  $f$  should clearly be strictly monotone increasing (or decreasing) and satisfy  $\lim_{x \rightarrow \pm\infty} F(x) = \pm\infty$ . Rewriting these conditions slightly such that they make sense for vector valued functions the analogous result holds.

**Lemma 14.1.** *Suppose  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is continuous and satisfies*

$$\lim_{|x| \rightarrow \infty} \frac{F(x)x}{|x|} = \infty. \quad (14.1)$$

*Then the equation*

$$F(x) = y \quad (14.2)$$

*has a solution for every  $y \in \mathbb{R}^n$ . If  $F$  is strictly monotone*

$$(F(x) - F(y))(x - y) > 0, \quad x \neq y, \quad (14.3)$$

*then this solution is unique.*

**Proof.** Our first assumption implies that  $G(x) = F(x) - y$  satisfies  $G(x)x = F(x)x - yx > 0$  for  $|x|$  sufficiently large. Hence the first claim follows from Problem 12.2. The second claim is trivial.  $\square$



Now we want to generalize this result to infinite dimensional spaces. Throughout this chapter,  $\mathfrak{H}$  will be a real Hilbert space with scalar product  $\langle \cdot, \cdot \rangle$ . A map  $F : \mathfrak{H} \rightarrow \mathfrak{H}$  is called **monotone** if

$$\langle F(x) - F(y), x - y \rangle \geq 0, \quad x, y \in \mathfrak{H}, \quad (14.4)$$

**strictly monotone** if

$$\langle F(x) - F(y), x - y \rangle > 0, \quad x \neq y \in \mathfrak{H}, \quad (14.5)$$

and finally **strongly monotone** if there is a constant  $C > 0$  such that

$$\langle F(x) - F(y), x - y \rangle \geq C\|x - y\|^2, \quad x, y \in \mathfrak{H}. \quad (14.6)$$

Note that the same definitions can be made for a Banach space  $X$  and mappings  $F : X \rightarrow X^*$ .

Observe that if  $F$  is strongly monotone, then it is **coercive**

$$\lim_{\|x\| \rightarrow \infty} \frac{\langle F(x), x \rangle}{\|x\|} = \infty. \quad (14.7)$$

(Just take  $y = 0$  in the definition of strong monotonicity.) Hence the following result is not surprising.

**Theorem 14.2** (Zarantonello). *Suppose  $F \in C(\mathfrak{H}, \mathfrak{H})$  is (globally) Lipschitz continuous and strongly monotone. Then, for each  $y \in \mathfrak{H}$  the equation*

$$F(x) = y \quad (14.8)$$

*has a unique solution  $x(y) \in \mathfrak{H}$  which depends continuously on  $y$ .*

**Proof.** Set

$$G(x) := x - t(F(x) - y), \quad t > 0,$$

then  $F(x) = y$  is equivalent to the fixed point equation

$$G(x) = x.$$

It remains to show that  $G$  is a contraction. We compute

$$\begin{aligned} \|G(x) - G(\tilde{x})\|^2 &= \|x - \tilde{x}\|^2 - 2t\langle F(x) - F(\tilde{x}), x - \tilde{x} \rangle + t^2\|F(x) - F(\tilde{x})\|^2 \\ &\leq (1 - 2\frac{C}{L}(Lt) + (Lt)^2)\|x - \tilde{x}\|^2, \end{aligned}$$

where  $L$  is a Lipschitz constant for  $F$  (i.e.,  $\|F(x) - F(\tilde{x})\| \leq L\|x - \tilde{x}\|$ ). Thus, if  $t \in (0, \frac{2C}{L^2})$ ,  $G$  is a uniform contraction and the rest follows from the uniform contraction principle.  $\square$

Again observe that our proof is constructive. In fact, the best choice for  $t$  is clearly  $t = \frac{C}{L^2}$  such that the contraction constant  $\theta = 1 - (\frac{C}{L})^2$  is minimal. Then the sequence

$$x_{n+1} = x_n - \frac{C}{L^2}(F(x_n) - y), \quad x_0 = y, \quad (14.9)$$

converges to the solution.

**Example 14.1.** Let  $A \in \mathcal{L}(\mathfrak{H})$  and consider  $F(x) = Ax$ . Then the condition

$$\langle Ax, x \rangle \geq C\|x\|^2$$

implies that  $A$  has a bounded inverse  $A^{-1} : \mathfrak{H} \rightarrow \mathfrak{H}$  with  $\|A^{-1}\| \leq C^{-1}$  (cf. Problem 1.37).  $\diamond$

## 14.2. The nonlinear Lax–Milgram theorem

As a consequence of the last theorem we obtain a nonlinear version of the Lax–Milgram theorem. We want to investigate the following problem:

$$a(x, y) = b(y), \quad \text{for all } y \in \mathfrak{H}, \quad (14.10)$$

where  $a : \mathfrak{H}^2 \rightarrow \mathbb{R}$  and  $b : \mathfrak{H} \rightarrow \mathbb{R}$ . For this equation the following result holds.

**Theorem 14.3** (Nonlinear Lax–Milgram theorem). *Suppose  $b \in \mathcal{L}(\mathfrak{H}, \mathbb{R})$  and  $a(x, \cdot) \in \mathcal{L}(\mathfrak{H}, \mathbb{R})$ ,  $x \in \mathfrak{H}$ , are linear functionals such that there are positive constants  $L$  and  $C$  such that for all  $x, y, z \in \mathfrak{H}$  we have*

$$a(x, x - y) - a(y, x - y) \geq C\|x - y\|^2 \quad (14.11)$$

and

$$|a(x, z) - a(y, z)| \leq L\|z\|\|x - y\|. \quad (14.12)$$

Then there is a unique  $x \in \mathfrak{H}$  such that (14.10) holds.

**Proof.** By the Riesz lemma (Theorem 2.10) there are elements  $F(x) \in \mathfrak{H}$  and  $z \in \mathfrak{H}$  such that  $a(x, y) = b(y)$  is equivalent to  $\langle F(x) - z, y \rangle = 0$ ,  $y \in \mathfrak{H}$ , and hence to

$$F(x) = z.$$

By (14.11) the map  $F$  is strongly monotone. Moreover, by (14.12) we infer

$$\|F(x) - F(y)\| = \sup_{\tilde{x} \in \mathfrak{H}, \|\tilde{x}\|=1} |\langle F(x) - F(y), \tilde{x} \rangle| \leq L\|x - y\|$$

that  $F$  is Lipschitz continuous. Now apply Theorem 14.2.  $\square$

The special case where  $a \in \mathcal{L}^2(\mathfrak{H}, \mathbb{R})$  is a bounded bilinear form which is strongly coercive, that is,

$$a(x, x) \geq C\|x\|^2, \quad x \in \mathfrak{H}, \quad (14.13)$$

is usually known as (linear) Lax–Milgram theorem (Theorem 2.16).

**Example 14.2.** For example, let  $U \subset \mathbb{R}^n$  be a bounded domain and consider the Dirichlet problem for the second order nonlinear **elliptic problem**

$$-\sum_{i,j=1}^n \partial_i A_{ij}(x) \partial_j u(x) + F(x, u(x)) = w(x)$$

with  $A_{i,j} \in L^\infty(U, \mathbb{R})$  and  $F : U \times \mathbb{R} \rightarrow \mathbb{R}$ . If we impose Dirichlet boundary conditions, then it is natural to look for solutions from the class  $\mathfrak{D} = \{u \in H_0^1(U, \mathbb{R}) \mid A_{ij} \partial_j u \in H^1(U, \mathbb{R}), 1 \leq i, j \leq n\}$ . Multiplying this equation with a function  $v \in H_0^1(U, \mathbb{R})$  and integrating over  $U$  gives the associated weak formulation

$$a(u, v) = w(v),$$

where

$$a(u, v) := \sum_{i,j=1}^n \int_U (A_{ij}(x) (\partial_j u(x)) (\partial_i v(x)) + F(x, u(x)) v(x)) d^n x,$$

$$w(v) := \int_U w(x) v(x) d^n x.$$

Here we have assumed  $w \in L^2(U, \mathbb{R})$  but, somewhat more general,  $w \in H^1(U, \mathbb{R})^*$  would also suffice.

If we require

$$C := \inf_{e \in S^n, x \in U} e_i A_{ij}(x) e_j > 0$$

as well as

$$|F(x, u_1) - F(x, u_2)| \leq L|u_1 - u_2| \text{ and } (F(x, u_1) - F(x, u_2))(u_1 - u_2) \geq 0,$$

then the assumption of the nonlinear Lax–Milgram theorem are satisfied on  $H_0^1(U, \mathbb{R})$ .  $\diamond$

### 14.3. The main theorem of monotone maps

Now we return to the investigation of  $F(x) = y$  and weaken the conditions of Theorem 14.2. We will assume that  $\mathfrak{H}$  is a separable Hilbert space and that  $F : \mathfrak{H} \rightarrow \mathfrak{H}$  is a continuous, coercive monotone map. In fact, it suffices to assume that  $F$  is **demicontinuous**

$$\lim_{n \rightarrow \infty} \langle F(x_n), y \rangle = \langle F(x), y \rangle, \quad \text{for all } y \in \mathfrak{H} \quad (14.14)$$

whenever  $x_n \rightarrow x$ .

The idea is as follows: Start with a finite dimensional subspace  $\mathfrak{H}_n \subset \mathfrak{H}$  and project the equation  $F(x) = y$  to  $\mathfrak{H}_n$  resulting in an equation

$$F_n(x_n) = y_n, \quad x_n, y_n \in \mathfrak{H}_n. \quad (14.15)$$

More precisely, let  $P_n$  be the (linear) projection onto  $\mathfrak{H}_n$  and set  $F_n(x_n) = P_n F(x_n)$ ,  $y_n = P_n y$  (verify that  $F_n$  is continuous and monotone!).

Now Lemma 14.1 ensures that there exists a solution  $x_n$ . Now choose the subspaces  $\mathfrak{H}_n$  such that  $\mathfrak{H}_n \rightarrow \mathfrak{H}$  (i.e.,  $\mathfrak{H}_n \subset \mathfrak{H}_{n+1}$  and  $\bigcup_{n=1}^{\infty} \mathfrak{H}_n$  is dense). Then our hope is that  $x_n$  converges to a solution  $x$ .

This approach is quite common when solving equations in infinite dimensional spaces and is known as **Galerkin approximation**. It can often be used for numerical computations and the right choice of the spaces  $\mathfrak{H}_n$  will have a significant impact on the quality of the approximation.

So how should we show that  $x_n$  converges? First of all observe that our construction of  $x_n$  shows that  $x_n$  lies in some ball with radius  $R_n$ , which is chosen such that

$$\langle F_n(x), x \rangle > \|y_n\| \|x\|, \quad \|x\| \geq R_n, \quad x \in \mathfrak{H}_n. \quad (14.16)$$

Since  $\langle F_n(x), x \rangle = \langle P_n F(x), x \rangle = \langle F(x), P_n x \rangle = \langle F(x), x \rangle$  for  $x \in \mathfrak{H}_n$  we can drop all  $n$ 's to obtain a constant  $R$  (depending on  $\|y\|$ ) which works for all  $n$ . So the sequence  $x_n$  is uniformly bounded

$$\|x_n\| \leq R. \quad (14.17)$$

Now by Theorem 4.28 there exists a weakly convergent subsequence. That is, after dropping some terms, we can assume that there is some  $x$  such that  $x_n \rightharpoonup x$ , that is,

$$\langle x_n, z \rangle \rightarrow \langle x, z \rangle, \quad \text{for every } z \in \mathfrak{H}. \quad (14.18)$$

And it remains to show that  $x$  is indeed a solution. This follows from

**Lemma 14.4.** *Suppose  $F : \mathfrak{H} \rightarrow \mathfrak{H}$  is demicontinuous and monotone, then*

$$\langle y - F(z), x - z \rangle \geq 0 \quad \text{for every } z \in \mathfrak{H} \quad (14.19)$$

*implies  $F(x) = y$ .*

**Proof.** Choose  $z = x \pm tw$ , then  $\mp \langle y - F(x \pm tw), w \rangle \geq 0$  and by continuity  $\mp \langle y - F(x), w \rangle \geq 0$ . Thus  $\langle y - F(x), w \rangle = 0$  for every  $w$  implying  $y - F(x) = 0$ .  $\square$

Now we can show

**Theorem 14.5** (Browder–Minty). *Let  $\mathfrak{H}$  be a separable Hilbert space. Suppose  $F : \mathfrak{H} \rightarrow \mathfrak{H}$  is demicontinuous, coercive and monotone. Then the equation*

$$F(x) = y \quad (14.20)$$

*has a solution for every  $y \in \mathfrak{H}$ . If  $F$  is strictly monotone then this solution is unique.*

**Proof.** We have  $\langle y - F(z), x_n - z \rangle = \langle y_n - F_n(z), x_n - z \rangle \geq 0$  for  $z \in \mathfrak{H}_n$ . Taking the limit (Problem 4.28) implies  $\langle y - F(z), x - z \rangle \geq 0$  for every  $z \in \mathfrak{H}_\infty = \bigcup_{n=1}^\infty \mathfrak{H}_n$ . Since  $\mathfrak{H}_\infty$  is dense,  $\langle y - F(z), x - z \rangle \geq 0$  for every  $z \in \mathfrak{H}$  by continuity and hence  $F(x) = y$  by our lemma.  $\square$

Note that in the infinite dimensional case we need monotonicity even to show existence. Moreover, this result can be further generalized in two more ways. First of all, the Hilbert space  $\mathfrak{H}$  can be replaced by a reflexive Banach space  $X$  if  $F : X \rightarrow X^*$ . The proof is similar. Secondly, it suffices if

$$t \mapsto \langle F(x + ty), z \rangle \quad (14.21)$$

is continuous for  $t \in [0, 1]$  and all  $x, y, z \in \mathfrak{H}$ , since this condition together with monotonicity can be shown to imply demicontinuity.

## Some set theory

At the beginning of the 20th century Russell showed with his famous paradox "Is  $\{x|x \notin x\}$  a set?" that naive set theory can lead to contradictions. Hence it was replaced by **axiomatic set theory**, more specific we will take the **Zermelo–Fraenkel set theory (ZF)**, which assumes existence of some sets (like the empty set and the integers) and defines what operations are allowed. Somewhat informally (i.e. without writing them using the symbolism of first order logic) they can be stated as follows:

- **Axiom of empty set.** There is a set  $\emptyset$  which contains no elements.
- **Axiom of extensionality.** Two sets  $A$  and  $B$  are equal  $A = B$  if they contain the same elements. If a set  $A$  contains all elements from a set  $B$ , it is called a subset  $A \subseteq B$ . In particular  $A \subseteq B$  and  $B \subseteq A$  if and only if  $A = B$ .

The last axiom implies that the empty set is unique and that any set which is not equal to the empty set has at least one element.

- **Axiom of pairing.** If  $A$  and  $B$  are sets, then there exists a set  $\{A, B\}$  which contains  $A$  and  $B$  as elements. One writes  $\{A, A\} = \{A\}$ . By the axiom of extensionality we have  $\{A, B\} = \{B, A\}$ .
- **Axiom of union.** Given a set  $\mathcal{F}$  whose elements are again sets, there is a set  $A = \bigcup \mathcal{F}$  containing every element that is a member of some member of  $\mathcal{F}$ . In particular, given two sets  $A, B$  there exists a set  $A \cup B = \bigcup \{A, B\}$  consisting of the elements of both sets. Note that this definition ensures that the union is commutative  $A \cup B = B \cup A$  and associative  $(A \cup B) \cup C = A \cup (B \cup C)$ . Note also  $\bigcup \{A\} = A$ .

- **Axiom schema of specification.** Given a set  $A$  and a logical statement  $\phi(x)$  depending on  $x \in A$  we can form the set  $B = \{x \in A | \phi(x)\}$  of all elements from  $A$  obeying  $\phi$ . For example, given two sets  $A$  and  $B$  we can define their intersection as  $A \cap B = \{x \in A \cup B | (x \in A) \wedge (x \in B)\}$  and their complement as  $A \setminus B = \{x \in A | x \notin B\}$ . Or the intersection of a family of sets  $\mathcal{F}$  as  $\bigcap \mathcal{F} = \{x \in \bigcup \mathcal{F} | \forall F \in \mathcal{F} : x \in F\}$ .
- **Axiom of power set.** For any set  $A$ , there is a **power set**  $\mathfrak{P}(A)$  that contains every subset of  $A$ .

From these axioms one can define ordered pairs as  $(x, y) = \{\{x\}, \{x, y\}\}$  and the Cartesian product as  $A \times B = \{z \in \mathfrak{P}(A \cup \mathfrak{P}(A \cup B)) | \exists x \in A, y \in B : z = (x, y)\}$ . Functions  $f : A \rightarrow B$  are defined as single valued relations, that is  $f \subseteq A \times B$  such that  $(x, y) \in f$  and  $(x, \tilde{y}) \in f$  implies  $y = \tilde{y}$ .

- **Axiom schema of replacement.** For every function  $f$  the image of a set  $A$  is again a set  $B = \{f(x) | x \in A\}$ .

So far the previous axioms were concerned with ensuring that the usual set operations required in mathematics yield again sets. In particular, we can start constructing sets with any given finite number of elements starting from the empty set:  $\emptyset$  (no elements),  $\{\emptyset\}$  (one element),  $\{\emptyset, \{\emptyset\}\}$  (two elements), etc. However, while existence of infinite sets (like e.g. the integers) might seem *obvious* at this point, it cannot be deduced from the axioms we have so far. Hence it has to be added as well.

- **Axiom of infinity.** There exists a set  $A$  which contains the empty set and for every element  $x \in A$  we also have  $x \cup \{x\} \in A$ . The smallest such set  $\{\emptyset, \{\emptyset\}, \{\emptyset, \{\emptyset\}\}, \dots\}$  can be identified with the integers via  $0 = \emptyset$ ,  $1 = \{\emptyset\}$ ,  $2 = \{\emptyset, \{\emptyset\}\}$ ,  $\dots$

Now we finally have the integers and thus everything we need to start constructing the rational, real, and complex numbers in the usual way. Hence we only add one more axiom to exclude some pathological objects which will lead to contradictions.

- **Axiom of Regularity.** Every nonempty set  $A$  contains an element  $x$  with  $x \cap A = \emptyset$ . This excludes for example the possibility that a set contains itself as an element (apply the axiom to  $\{A\}$ ). Similarly, we can only have  $A \in B$  or  $B \in A$  but not both (apply it to the set  $\{A, B\}$ ).

Hence a set is something which can be constructed from the above axioms. Of course this raises the question if these axioms are consistent but as has been shown by Gödel this question cannot be answered: If ZF contains a statement of its own consistency then ZF is inconsistent. In fact, the same

holds for any other sufficiently rich (such that one can do basic math) system of axioms. In particular, it also holds for ZFC defined below. So we have to live with the fact that someday someone might come and prove that ZFC is inconsistent.

Starting from ZF one can develop basic analysis (including the construction of the real numbers). However, it turns out that several fundamental results require yet another construction for their proof:

Given an index set  $A$  and for every  $\alpha \in A$  some set  $M_\alpha$  the product  $\times_{\alpha \in A} M_\alpha$  is defined to be the set of all functions  $\varphi : A \rightarrow \bigcup_{\alpha \in A} M_\alpha$  which assign each element  $\alpha \in A$  some element  $m_\alpha \in M_\alpha$ . If all sets  $M_\alpha$  are nonempty it seems quite reasonable that there should be such a *choice function* which chooses an element from  $M_\alpha$  for every  $\alpha \in A$ . However, no matter how obvious this might seem, it cannot be deduced from the ZF axioms alone and hence has to be added:

- **Axiom of Choice:** Given an index set  $A$  and nonempty sets  $\{M_\alpha\}_{\alpha \in A}$  their product  $\times_{\alpha \in A} M_\alpha$  is nonempty.

ZF augmented by the axiom of choice is known as **ZFC** and we accept it as the fundament upon which our functional analytic house is built.

Note that the axiom of choice is not only used to ensure that infinite products are nonempty but also in many proofs! For example, suppose you start with a set  $M_1$  and recursively construct some sets  $M_n$  such that in every step you have a nonempty set. Then the axiom of choice guarantees the existence of a sequence  $x = (x_n)_{n \in \mathbb{N}}$  with  $x_n \in M_n$ .

The axiom of choice has many important consequences (many of which are in fact equivalent to the axiom of choice and it is hence a matter of taste which one to choose as axiom).

A **partial order** is a binary relation " $\preceq$ " over a set  $\mathcal{P}$  such that for all  $A, B, C \in \mathcal{P}$ :

- $A \preceq A$  (reflexivity),
- if  $A \preceq B$  and  $B \preceq A$  then  $A = B$  (antisymmetry),
- if  $A \preceq B$  and  $B \preceq C$  then  $A \preceq C$  (transitivity).

It is custom to write  $A \prec B$  if  $A \preceq B$  and  $A \neq B$ .

**Example A.1.** Let  $\mathcal{P}(X)$  be the collections of all subsets of a set  $X$ . Then  $\mathcal{P}$  is partially ordered by inclusion  $\subseteq$ .  $\diamond$

It is important to emphasize that two elements of  $\mathcal{P}$  need not be comparable, that is, in general neither  $A \preceq B$  nor  $B \preceq A$  might hold. However, if any two elements are comparable,  $\mathcal{P}$  will be called **totally ordered**. A set with a total order is called **well-ordered** if every nonempty subset has



a **least element**, that is some  $A \in \mathcal{P}$  with  $A \preceq B$  for every  $B \in \mathcal{P}$ . Note that the least element is unique by antisymmetry.

**Example A.2.**  $\mathbb{R}$  with  $\leq$  is totally ordered and  $\mathbb{N}$  with  $\leq$  is well-ordered.  $\diamond$

On every well-ordered set we have the

**Theorem A.1** (Induction principle). *Let  $K$  be well ordered and let  $S(k)$  be a statement for arbitrary  $k \in K$ . Then, if  $A(l)$  true for all  $l \prec k$  implies  $A(k)$  true, then  $A(k)$  is true for all  $k \in K$ .*

**Proof.** Otherwise the set of all  $k$  for which  $A(k)$  is false had a least element  $k_0$ . But by our choice of  $k_0$ ,  $A(l)$  holds for all  $l \prec k_0$  and thus for  $k_0$  contradicting our assumption.  $\square$

The induction principle also shows that in a well-ordered set functions  $f$  can be defined recursively, that is, by a function  $\varphi$  which computes the value of  $f(k)$  from the values  $f(l)$  for all  $l \prec k$ . Indeed, the induction principle implies that on the set  $M_k = \{l \in K \mid l \prec k\}$  there is at most one such function  $f_k$ . Since  $k$  is arbitrary,  $f$  is unique. In case of the integers existence of  $f_k$  is also clear provided  $f(1)$  is given. In general, one can prove existence provided  $f_k$  is given for some  $k$  but we will not need this.

If  $\mathcal{P}$  is partially ordered, then every totally ordered subset is also called a **chain**. If  $\mathcal{Q} \subseteq \mathcal{P}$ , then an element  $M \in \mathcal{P}$  satisfying  $A \preceq M$  for all  $A \in \mathcal{Q}$  is called an **upper bound**.

**Example A.3.** Let  $\mathcal{P}(X)$  as before. Then a collection of subsets  $\{A_n\}_{n \in \mathbb{N}} \subseteq \mathcal{P}(X)$  satisfying  $A_n \subseteq A_{n+1}$  is a chain. The set  $\bigcup_n A_n$  is an upper bound.  $\diamond$

An element  $M \in \mathcal{P}$  for which  $M \preceq A$  for some  $A \in \mathcal{P}$  is only possible if  $M = A$  is called a **maximal element**.

**Theorem A.2** (Zorn's lemma). *Every partially ordered set in which every chain has an upper bound contains at least one maximal element.*

**Proof.** Suppose it were false. Then to every chain  $C$  we can assign an element  $m(C)$  such that  $m(C) \succ x$  for all  $x \in C$  (here we use the axiom of choice). We call a chain  $C$  distinguished if it is well-ordered and if for every segment  $C_x = \{y \in C \mid y \prec x\}$  we have  $m(C_x) = x$ . We will also regard  $C$  as a segment of itself.

Then (since for the least element of  $C$  we have  $C_x = \emptyset$ ) every distinguished chain must start like  $m(\emptyset) \prec m(m(\emptyset)) \prec \dots$  and given two segments  $C, D$  we expect that always one must be a segment of the other.

So let us first prove this claim. Suppose  $D$  is not a segment of  $C$ . Then we need to show  $C = D_z$  for some  $z$ . We start by showing that  $x \in C$  implies  $x \in D$  and  $C_x = D_x$ . To see this suppose it were wrong and let  $x$  be the

least  $x \in C$  for which it fails. Then  $y \in K_x$  implies  $y \in L$  and hence  $K_x \subset L$ . Then, since  $C_x \neq D$  by assumption, we can find a least  $z \in D \setminus C_x$ . In fact we must even have  $z \succ C_x$  since otherwise we could find a  $y \in C_x$  such that  $x \succ y \succ z$ . But then, using that it holds for  $y$ ,  $y \in D$  and  $C_y = D_y$  so we get the contradiction  $z \in D_y = C_y \subset C_x$ . So  $z \succ C_x$  and thus also  $C_x = D_z$  which in turn shows  $x = m(C_x) = m(D_z) = z$  and proves that  $x \in C$  implies  $x \in D$  and  $C_x = D_x$ . In particular  $C \subset D$  and as before  $C = D_z$  for the least  $z \in D \setminus C$ . This proves the claim.

Now using this claim we see that we can take the union over all distinguished chains to get a maximal distinguished chain  $C_{max}$ . But then we could add  $m(C_{max}) \notin C_{max}$  to  $C_{max}$  to get a larger distinguished chain  $C_{max} \cup \{m(C_{max})\}$  contradicting maximality.  $\square$

We will also frequently use the **cardinality** of sets: Two sets  $A$  and  $B$  have the same cardinality, written as  $|A| = |B|$ , if there is a bijection  $\varphi : A \rightarrow B$ . We write  $|A| \leq |B|$  if there is an injective map  $\varphi : A \rightarrow B$ . Note that  $|A| \leq |B|$  and  $|B| \leq |C|$  implies  $|A| \leq |C|$ . A set  $A$  is called infinite if  $|A| \geq |\mathbb{N}|$ , countable if  $|A| \leq |\mathbb{N}|$ , and countably infinite if  $|A| = |\mathbb{N}|$ .

**Theorem A.3** (Schröder–Bernstein).  $|A| \leq |B|$  and  $|B| \leq |A|$  implies  $|A| = |B|$ .

**Proof.** Suppose  $\varphi : A \rightarrow B$  and  $\psi : B \rightarrow A$  are two injective maps. Now consider sequences  $x_n$  defined recursively via  $x_{2n+1} = \varphi(x_{2n})$  and  $x_{2n+1} = \psi(x_{2n})$ . Given a start value  $x_0 \in A$  the sequence is uniquely defined but might terminate at a negative integer since our maps are not surjective. In any case, if an element appears in two sequences, the elements to the left and to the right must also be equal (use induction) and hence the two sequences differ only by an index shift. So the ranges of such sequences form a partition for  $A \cup B$  and it suffices to find a bijection between elements in one partition. If the sequence stops at an element in  $A$  we can take  $\varphi$ . If the sequence stops at an element in  $B$  we can take  $\psi^{-1}$ . If the sequence is doubly infinite either of the previous choices will do.  $\square$

**Theorem A.4** (Zermelo). Either  $|A| \leq |B|$  or  $|B| \leq |A|$ .

**Proof.** Consider the set of all bijective functions  $\varphi_\alpha : A_\alpha \rightarrow B$  with  $A_\alpha \subseteq A$ . Then we can define a partial ordering via  $\varphi_\alpha \preceq \varphi_\beta$  if  $A_\alpha \subseteq A_\beta$  and  $\varphi_\beta|_{A_\alpha} = \varphi_\alpha$ . Then every chain has an upper bound (the unique function defined on the union of all domains) and by Zorn's lemma there is a maximal element  $\varphi_m$ . For  $\varphi_m$  we have either  $A_m = A$  or  $\varphi_m(A_m) = B$  since otherwise there is some  $x \in A \setminus A_m$  and some  $y \in B \setminus \varphi_m(A_m)$  which could be used to extend  $\varphi_m$  to  $A_m \cup \{x\}$  by setting  $\varphi(x) = y$ . But if  $A_m = A$  we have  $|A| \leq |B|$  and if  $\varphi_m(A_m) = B$  we have  $|B| \leq |A|$ .  $\square$

The cardinality of the power set  $\mathfrak{P}(A)$  is strictly larger than the cardinality of  $A$ .

**Theorem A.5** (Cantor).  $|A| < |\mathfrak{P}(A)|$ .

**Proof.** Suppose there were a bijection  $\varphi : A \rightarrow \mathfrak{P}(A)$ . Then, for  $B = \{x \in A \mid x \notin \varphi(x)\}$  there must be some  $y$  such that  $B = \varphi(y)$ . But  $y \in B$  if and only if  $y \notin \varphi(y) = B$ , a contradiction.  $\square$

This innocent looking result also caused some grief when announced by Cantor as it clearly gives a contradiction when applied to the *set of all sets* (which is fortunately not a legal object in ZFC).

The following result and its corollary will be used to determine the cardinality of unions and products.

**Lemma A.6.** *Any infinite set can be written as a disjoint union of countably infinite sets.*

**Proof.** Consider collections of disjoint countably infinite subsets. Such collections can be partially ordered by inclusion and hence there is a maximal collection by Zorn's lemma. If the union of such a maximal collection falls short of the whole set the complement must be finite. Since this finite remainder can be added to a set of the collection we are done.  $\square$

**Corollary A.7.** *Any infinite set can be written as a disjoint union of two disjoint subsets having the same cardinality as the original set.*

**Proof.** By the lemma we can write  $A = \bigcup A_\alpha$ , where all  $A_\alpha$  are countably infinite. Now split  $A_\alpha = B_\alpha \cup C_\alpha$  into two disjoint countably infinite sets (map  $A_\alpha$  bijective to  $\mathbb{N}$  and then split into even and odd elements). Then the desired splitting is  $A = B \cup C$  with  $B = \bigcup B_\alpha$  and  $C = \bigcup C_\alpha$ .  $\square$

**Theorem A.8.** *Suppose  $A$  or  $B$  is infinite. Then  $|A \cup B| = \max\{|A|, |B|\}$ .*

**Proof.** Suppose  $A$  is infinite and  $|B| \leq |A|$ . Then  $|A| \leq |A \cup B| \leq |A \cup B| \leq |A \cup A| = |A|$  by the previous corollary. Here  $\cup$  denotes the disjoint union.  $\square$

A standard theorem proven in every introductory course is that  $\mathbb{N} \times \mathbb{N}$  is countable. The generalization of this result is also true.

**Theorem A.9** (Hessenberg). *Suppose  $A$  is infinite and  $B \neq \emptyset$ . Then  $|A \times B| = \max\{|A|, |B|\}$ .*

**Proof.** Without loss of generality we can assume  $|B| \leq |A|$  (otherwise exchange both sets). Then  $|A| \leq |A \times B| \leq |A \times A|$  and it suffices to show  $|A \times A| = |A|$ .

We proceed as before and consider the set of all bijective functions  $\varphi_\alpha : A_\alpha \rightarrow A_\alpha \times A_\alpha$  with  $A_\alpha \subseteq A$  with the same partial ordering as before. By Zorn's lemma there is a maximal element  $\varphi_m$ . Let  $A_m$  be its domain and let  $A'_m = A \setminus A_m$ . We claim that  $|A'_m| < |A_m|$ . If not,  $A'_m$  had a subset  $A''_m$  with the same cardinality of  $A_m$  and hence we had a bijection from  $A''_m \rightarrow A''_m \times A''_m$  which could be used to extend  $\varphi$ . So  $|A'_m| < |A_m|$  and thus  $|A| = |A_m \cup A'_m| = |A_m|$ . Since we have shown  $|A_m \times A_m| = |A_m|$  the claim follows.  $\square$

**Example A.4.** Note that for  $A = \mathbb{N}$  we have  $|\mathfrak{P}(\mathbb{N})| = |\mathbb{R}|$ . Indeed, since  $|\mathbb{R}| = |\mathbb{Z} \times [0, 1)| = |[0, 1)|$  it suffices to show  $|\mathfrak{P}(\mathbb{N})| = |[0, 1)|$ . To this end note that  $\mathfrak{P}(\mathbb{N})$  can be identified with the set of all sequences with values in  $\{0, 1\}$  (the value of the sequence at a point tells us whether it is in the corresponding subset). Now every point in  $[0, 1)$  can be mapped to such a sequence via its binary expansion. This map is injective but not surjective since a point can have different binary expansions:  $|[0, 1)| \leq |\mathfrak{P}(\mathbb{N})|$ . Conversely, given a sequence  $a_n \in \{0, 1\}$  we can map it to the number  $\sum_{n=1}^{\infty} a_n 4^{-n}$ . Since this map is again injective (note that we avoid expansions which are eventually 1) we get  $|\mathfrak{P}(\mathbb{N})| \leq |[0, 1)|$ .  $\diamond$

Hence we have

$$|\mathbb{N}| < |\mathfrak{P}(\mathbb{N})| = |\mathbb{R}| \quad (\text{A.1})$$

and the **continuum hypothesis** states that there are no sets whose cardinality lie in between. It was shown by Gödel and Cohen that it, as well as its negation, is consistent with ZFC and hence cannot be decided within this framework.

**Problem A.1.** Show that Zorn's lemma implies the axiom of choice. (Hint: Consider the set of all partial choice functions defined on a subset.)

**Problem A.2.** Show  $|\mathbb{R}^{\mathbb{N}}| = |\mathbb{R}|$ . (Hint: Without loss we can replace  $\mathbb{R}$  by  $(0, 1)$  and identify each  $x \in (0, 1)$  with its decimal expansion. Now the digits in a given sequence are indexed by two countable parameters.)



# Metric and topological spaces

This chapter collects some basic facts from metric and topological spaces as a reference for the main text. I presume that you are familiar with most of these topics from your calculus course. As a general reference I can warmly recommend Kelley's classical book [26] or the nice book by Kaplansky [24]. As always such a brief compilation introduces a zoo of properties. While sometimes the connection between these properties are straightforward, other times they might be quite tricky. So if at some point you are wondering if there exists an infinite multi-variable sub-polynomial Woffle which does not satisfy the lower regular  $Q$ -property, start searching in the book by Steen and Seebach [44].

## B.1. Basics

One of the key concepts in analysis is convergence. To define convergence requires the notion of distance. Motivated by the Euclidean distance one is lead to the following definition:

A **metric space** is a space  $X$  together with a distance function  $d : X \times X \rightarrow [0, \infty)$  such that for arbitrary points  $x, y, z \in X$  we have

- (i)  $d(x, y) = 0$  if and only if  $x = y$ ,
- (ii)  $d(x, y) = d(y, x)$ ,
- (iii)  $d(x, z) \leq d(x, y) + d(y, z)$  (**triangle inequality**).

If (i) does not hold,  $d$  is called a **pseudometric**. As a straightforward consequence we record the **inverse triangle inequality** (Problem B.1)

$$|d(x, y) - d(z, y)| \leq d(x, z). \quad (\text{B.1})$$

**Example B.1.** The role model for a metric space is of course Euclidean space  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) together with  $d(x, y) := (\sum_{k=1}^n |x_k - y_k|^2)^{1/2}$ .  $\diamond$

Several notions from  $\mathbb{R}^n$  carry over to metric spaces in a straightforward way. The set

$$B_r(x) := \{y \in X \mid d(x, y) < r\} \quad (\text{B.2})$$

is called an **open ball** around  $x$  with radius  $r > 0$ . We will write  $B_r^X(x)$  in case we want to emphasize the corresponding space. A point  $x$  of some set  $U \subseteq X$  is called an **interior point** of  $U$  if  $U$  contains some open ball around  $x$ . If  $x$  is an interior point of  $U$ , then  $U$  is also called a **neighborhood** of  $x$ . A point  $x$  is called a **limit point** of  $U$  (also **accumulation** or **cluster point**) if for any open ball  $B_r(x)$ , there exists at least one point in  $B_r(x) \cap U$  distinct from  $x$ . Note that a limit point  $x$  need not lie in  $U$ , but  $U$  must contain points arbitrarily close to  $x$ . A point  $x$  is called an **isolated point** of  $U$  if there exists a neighborhood of  $x$  not containing any other points of  $U$ . A set which consists only of isolated points is called a **discrete set**. If any neighborhood of  $x$  contains at least one point in  $U$  and at least one point not in  $U$ , then  $x$  is called a **boundary point** of  $U$ . The set of all boundary points of  $U$  is called the boundary of  $U$  and denoted by  $\partial U$ .

**Example B.2.** Consider  $\mathbb{R}$  with the usual metric and let  $U := (-1, 1)$ . Then every point  $x \in U$  is an interior point of  $U$ . The points  $[-1, 1]$  are limit points of  $U$ , and the points  $\{-1, +1\}$  are boundary points of  $U$ .

Let  $U := \mathbb{Q}$ , the set of rational numbers. Then  $U$  has no interior points and  $\partial U = \mathbb{R}$ .  $\diamond$

A set all of whose points are interior points is called **open**. The family of open sets  $\mathcal{O}$  satisfies the properties

- (i)  $\emptyset, X \in \mathcal{O}$ ,
- (ii)  $O_1, O_2 \in \mathcal{O}$  implies  $O_1 \cap O_2 \in \mathcal{O}$ ,
- (iii)  $\{O_\alpha\} \subseteq \mathcal{O}$  implies  $\bigcup_\alpha O_\alpha \in \mathcal{O}$ .

That is,  $\mathcal{O}$  is closed under finite intersections and arbitrary unions. Indeed, (i) is obvious, (ii) follows since the intersection of two open balls centered at  $x$  is again an open ball centered at  $x$  (explicitly  $B_{r_1}(x) \cap B_{r_2}(x) = B_{\min(r_1, r_2)}(x)$ ), and (iii) follows since every ball contained in one of the sets is also contained in the union.

**Example B.3.** Of course every open ball  $B_r(x)$  is an open set since  $y \in B_r(x)$  implies  $B_s(y) \subseteq B_r(x)$  for  $s < r - d(x, y)$ .  $\diamond$

Now it turns out that for defining convergence, a distance is slightly more than what is actually needed. In fact, it suffices to know when a point is in the neighborhood of another point. And if we adapt the definition of a neighborhood by requiring it to contain an open set around  $x$ , then we see that it suffices to know when a set is open. This motivates the following definition:

A space  $X$  together with a family of sets  $\mathcal{O}$ , the open sets, satisfying (i)–(iii), is called a **topological space**. The notions of interior point, limit point, and neighborhood carry over to topological spaces if we replace open ball around  $x$  by open set containing  $x$ .

There are usually different choices for the topology. Two not too interesting examples are the **trivial topology**  $\mathcal{O} = \{\emptyset, X\}$  and the **discrete topology**  $\mathcal{O} = \mathfrak{P}(X)$  (the power set of  $X$ ). Given two topologies  $\mathcal{O}_1$  and  $\mathcal{O}_2$  on  $X$ ,  $\mathcal{O}_1$  is called **weaker** (or **coarser**) than  $\mathcal{O}_2$  if  $\mathcal{O}_1 \subseteq \mathcal{O}_2$ . Conversely,  $\mathcal{O}_1$  is called **stronger** (or **finer**) than  $\mathcal{O}_2$  if  $\mathcal{O}_2 \subseteq \mathcal{O}_1$ .

Given two topologies on  $X$  their intersection will again be a topology on  $X$ . In fact, the intersection of an arbitrary collection of topologies is again a topology and one can define the weakest topology with a certain property to be the intersection of all topologies with this property (provided there is one at all).

**Example B.4.** Note that different metrics can give rise to the same topology. For example, we can equip  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) with the Euclidean distance  $d(x, y)$  as before or we could also use

$$\tilde{d}(x, y) := \sum_{k=1}^n |x_k - y_k|. \quad (\text{B.3})$$

Then

$$\frac{1}{\sqrt{n}} \sum_{k=1}^n |x_k| \leq \sqrt{\sum_{k=1}^n |x_k|^2} \leq \sum_{k=1}^n |x_k| \quad (\text{B.4})$$

shows  $B_{r/\sqrt{n}}(x) \subseteq \tilde{B}_r(x) \subseteq B_r(x)$ , where  $B, \tilde{B}$  are balls computed using  $d, \tilde{d}$ , respectively. In particular, both distances will lead to the same notion of convergence.  $\diamond$

**Example B.5.** We can always replace a metric  $d$  by the bounded metric

$$\tilde{d}(x, y) := \frac{d(x, y)}{1 + d(x, y)} \quad (\text{B.5})$$

without changing the topology (since the family of open balls does not change:  $B_\delta(x) = \tilde{B}_{\delta/(1+\delta)}(x)$ ). To see that  $\tilde{d}$  is again a metric, observe that  $f(r) = \frac{r}{1+r}$  is monotone as well as concave and hence subadditive,  $f(r+s) \leq f(r) + f(s)$  (cf. Problem B.3).  $\diamond$



Every subspace  $Y$  of a topological space  $X$  becomes a topological space of its own if we call  $O \subseteq Y$  open if there is some open set  $\tilde{O} \subseteq X$  such that  $O = \tilde{O} \cap Y$ . This natural topology  $\mathcal{O} \cap Y$  is known as the **relative topology** (also **subspace**, **trace** or **induced topology**).

**Example B.6.** The set  $(0, 1] \subseteq \mathbb{R}$  is not open in the topology of  $X := \mathbb{R}$ , but it is open in the relative topology when considered as a subset of  $Y := [-1, 1]$ .  $\diamond$

A family of open sets  $\mathcal{B} \subseteq \mathcal{O}$  is called a **base** for the topology if for each  $x$  and each neighborhood  $U(x)$ , there is some set  $O \in \mathcal{B}$  with  $x \in O \subseteq U(x)$ . Since an open set  $O$  is a neighborhood of every one of its points, it can be written as  $O = \bigcup_{O \supseteq \tilde{O} \in \mathcal{B}} \tilde{O}$  and we have

**Lemma B.1.** *A family of open sets  $\mathcal{B} \subseteq \mathcal{O}$  is a base for the topology if and only if every open set can be written as a union of elements from  $\mathcal{B}$ .*

**Proof.** To see the converse let  $x$  and  $U(x)$  be given. Then  $U(x)$  contains an open set  $O$  containing  $x$  which can be written as a union of elements from  $\mathcal{B}$ . One of the elements in this union must contain  $x$  and this is the set we are looking for.  $\square$

A family of open sets  $\mathcal{B} \subseteq \mathcal{O}$  is called a **subbase** for the topology if every open set can be written as a union of finite intersections of elements from  $\mathcal{B}$ .

**Example B.7.** The intervals form a base for the topology on  $\mathbb{R}$ . Slightly more general, the open balls are a base for the topology in a metric space. Intervals of the form  $(\alpha, \infty)$  or  $(-\infty, \alpha)$  with  $\alpha \in \mathbb{R}$  are a subbase for topology of  $\mathbb{R}$ .  $\diamond$

Note that a subbase  $\mathcal{B}$  generates the topology in the sense that the corresponding topology is the weakest topology for which all of the sets from  $\mathcal{B}$  are open.

**Example B.8.** The **extended real numbers**  $\bar{\mathbb{R}} = \mathbb{R} \cup \{\infty, -\infty\}$  have a topology generated by the extended intervals of the form  $(\alpha, \infty]$  or  $[-\infty, \alpha)$  with  $\alpha \in \mathbb{R}$ . Note that the map  $f(x) := \frac{x}{1+|x|}$  maps  $\mathbb{R} \rightarrow [-1, 1]$ . This map becomes isometric if we choose  $d(x, y) = |f(x) - f(y)|$  as a metric on  $\bar{\mathbb{R}}$ . It is not hard to verify that this metric generates our topology and hence we can think of  $\bar{\mathbb{R}}$  as  $[-1, 1]$ .  $\diamond$

There is also a local version of the previous notions. A **neighborhood base** for a point  $x$  is a collection of neighborhoods  $\mathcal{B}(x)$  of  $x$  such that for each neighborhood  $U(x)$ , there is some set  $B \in \mathcal{B}(x)$  with  $B \subseteq U(x)$ . Note that the sets in a neighborhood base are not required to be open.

If every point has a countable neighborhood base, then  $X$  is called **first countable**. If there exists a countable base, then  $X$  is called **second countable**. Note that a second countable space is in particular first countable since for every base  $\mathcal{B}$  the subset  $\mathcal{B}(x) := \{O \in \mathcal{B} | x \in O\}$  is a neighborhood base for  $x$ .

**Example B.9.** In a metric space the open balls  $\{B_{1/m}(x)\}_{m \in \mathbb{N}}$  are a neighborhood base for  $x$ . Hence every metric space is first countable. Taking the union over all  $x$ , we obtain a base. In the case of  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) it even suffices to take balls with rational center, and hence  $\mathbb{R}^n$  (as well as  $\mathbb{C}^n$ ) is second countable.  $\diamond$

**Example B.10.** Let  $X = \mathbb{R}$  (or any other uncountable set) and call  $O \subseteq X$  open if either  $O$  is empty or  $X \setminus O$  is countable. This defines a topology (check this), known as the **cofinite topology**. Then  $X$  is not first countable. Indeed, suppose there were a countable neighborhood base  $B_j$  for  $0 \in X$ . Without loss of generality we can assume  $B_j \neq X$  for all  $j$  and hence there is some  $x_j \in X \setminus B_j$  for every  $j$ . Now  $N := X \setminus \bigcup_j \{x_j\}$  is a neighborhood of  $0$  for which  $B_j \not\subseteq N$  for all  $j$ , a contradiction.  $\diamond$

Given a collection  $\mathcal{M}$  of subsets of  $X$  we can define the topology generated by  $\mathcal{M}$  as the weakest topology (i.e., the intersection of all topologies) containing  $\mathcal{M}$ . Note that if  $\mathcal{M}$  is closed under finite intersections and  $\emptyset, X \in \mathcal{M}$ , then it will be a base for the topology generated by  $\mathcal{M}$  (Problem B.9).

Given two bases we can use them to check if the corresponding topologies are equal.

**Lemma B.2.** Let  $\mathcal{O}_j$ ,  $j = 1, 2$  be two topologies for  $X$  with corresponding bases  $\mathcal{B}_j$ . Then  $\mathcal{O}_1 \subseteq \mathcal{O}_2$  if and only if for every  $x \in X$  and every  $B_1 \in \mathcal{B}_1$  with  $x \in B_1$  there is some  $B_2 \in \mathcal{B}_2$  such that  $x \in B_2 \subseteq B_1$ .

**Proof.** Suppose  $\mathcal{O}_1 \subseteq \mathcal{O}_2$ , then  $B_1 \in \mathcal{O}_2$  and there is a corresponding  $B_2$  by the very definition of a base. Conversely, let  $O_1 \in \mathcal{O}_1$  and pick some  $x \in O_1$ . Then there is some  $B_1 \in \mathcal{B}_1$  with  $x \in B_1 \subseteq O_1$  and by assumption some  $B_2 \in \mathcal{B}_2$  such that  $x \in B_2 \subseteq B_1 \subseteq O_1$  which shows that  $x$  is an interior point with respect to  $\mathcal{O}_2$ . Since  $x$  was arbitrary we conclude  $O_1 \in \mathcal{O}_2$ .  $\square$

The next definition will ensure that limits are unique: A topological space is called a **Hausdorff space** if for any two different points there are always two disjoint neighborhoods.

**Example B.11.** Any metric space is a Hausdorff space: Given two different points  $x$  and  $y$ , the balls  $B_{d/2}(x)$  and  $B_{d/2}(y)$ , where  $d := d(x, y) > 0$ , are disjoint neighborhoods. A pseudometric space will in general not be

Hausdorff since two points of distance 0 cannot be separated by open balls.  $\diamond$

The complement of an open set is called a **closed set**. It follows from **De Morgan's laws**

$$X \setminus \left( \bigcup_{\alpha} U_{\alpha} \right) = \bigcap_{\alpha} (X \setminus U_{\alpha}), \quad X \setminus \left( \bigcap_{\alpha} U_{\alpha} \right) = \bigcup_{\alpha} (X \setminus U_{\alpha}) \quad (\text{B.6})$$

that the family of closed sets  $\mathcal{C}$  satisfies

- (i)  $\emptyset, X \in \mathcal{C}$ ,
- (ii)  $C_1, C_2 \in \mathcal{C}$  implies  $C_1 \cup C_2 \in \mathcal{C}$ ,
- (iii)  $\{C_{\alpha}\} \subseteq \mathcal{C}$  implies  $\bigcap_{\alpha} C_{\alpha} \in \mathcal{C}$ .

That is, closed sets are closed under finite unions and arbitrary intersections.

The smallest closed set containing a given set  $U$  is called the **closure**

$$\overline{U} := \bigcap_{C \in \mathcal{C}, U \subseteq C} C, \quad (\text{B.7})$$

and the largest open set contained in a given set  $U$  is called the **interior**

$$U^{\circ} := \bigcup_{O \in \mathcal{O}, O \subseteq U} O. \quad (\text{B.8})$$

It is not hard to see that the closure satisfies the following axioms (**Kuratowski closure axioms**):

- (i)  $\overline{\emptyset} = \emptyset$ ,
- (ii)  $U \subseteq \overline{U}$ ,
- (iii)  $\overline{\overline{U}} = \overline{U}$ ,
- (iv)  $\overline{U \cup V} = \overline{U} \cup \overline{V}$ .

In fact, one can show that these axioms can equivalently be used to define the topology by observing that the closed sets are precisely those which satisfy  $\overline{U} = U$ . Similarly, the open sets are precisely those which satisfy  $U^{\circ} = U$ .

**Lemma B.3.** *Let  $X$  be a topological space. Then the interior of  $U$  is the set of all interior points of  $U$ , and the closure of  $U$  is the union of  $U$  with all limit points of  $U$ . Moreover,  $\partial U = \overline{U} \setminus U^{\circ}$ .*

**Proof.** The first claim is straightforward. For the second claim observe that by Problem B.7 we have that  $\overline{U} = (X \setminus (X \setminus U)^{\circ})$ , that is, the closure is the set of all points which are not interior points of the complement. That is,  $x \notin \overline{U}$  iff there is some open set  $O$  containing  $x$  with  $O \subseteq X \setminus U$ . Hence,  $x \in \overline{U}$  iff for all open sets  $O$  containing  $x$  we have  $O \not\subseteq X \setminus U$ , that is,  $O \cap U \neq \emptyset$ . Hence,  $x \in \overline{U}$  iff  $x \in U$  or if  $x$  is a limit point of  $U$ . The last claim is left as Problem B.8.  $\square$

**Example B.12.** For any  $x \in X$  in a metric space  $X$  the **closed ball**

$$\bar{B}_r(x) := \{y \in X \mid d(x, y) \leq r\} \quad (\text{B.9})$$

is a closed set (check that its complement is open). But in general we have only

$$\overline{B_r(x)} \subseteq \bar{B}_r(x) \quad (\text{B.10})$$

since an isolated point  $y$  with  $d(x, y) = r$  will not be a limit point. In  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) we have of course equality.  $\diamond$

Note that in a Hausdorff space sets consisting of one point are closed (Problem B.5).

**Problem B.1.** Show that  $|d(x, y) - d(z, y)| \leq d(x, z)$ .

**Problem B.2.** Show the **quadrangle inequality**  $|d(x, y) - d(x', y')| \leq d(x, x') + d(y, y')$ .

**Problem B.3.** Show that if  $f : [0, \infty) \rightarrow \mathbb{R}$  is concave,  $f(\lambda x + (1 - \lambda)y) \geq \lambda f(x) + (1 - \lambda)f(y)$  for  $\lambda \in [0, 1]$ , and satisfies  $f(0) = 0$ , then it is subadditive,  $f(x + y) \leq f(x) + f(y)$ . If in addition  $f$  is increasing and  $d$  is a pseudometric, then so is  $f(d)$ . (Hint: Begin by showing  $f(\lambda x) \geq \lambda f(x)$ .)

**Problem B.4.** Show De Morgan's laws.

**Problem\* B.5.** Show that in a (nonempty) Hausdorff space  $X$  singleton sets  $\{x\}$  (with  $x \in X$ ) are closed.

**Problem B.6.** Show that the closure satisfies the Kuratowski closure axioms.

**Problem B.7.** Show that the closure and interior operators are dual in the sense that

$$X \setminus \bar{U} = (X \setminus U)^\circ \quad \text{and} \quad X \setminus U^\circ = \overline{(X \setminus U)}.$$

In particular, the closure is the set of all points which are not interior points of the complement. (Hint: De Morgan's laws.)

**Problem B.8.** Show that the boundary of  $U$  is given by  $\partial U = \bar{U} \setminus U^\circ$ .

**Problem B.9.** Suppose  $\mathcal{M}$  is a collection of sets closed under finite intersections containing  $\emptyset$  and  $X$ . Then the topology generated by  $\mathcal{M}$  is given by  $\mathcal{O}(\mathcal{M}) := \{\bigcup_\alpha M_\alpha \mid M_\alpha \in \mathcal{M}\}$ .

## B.2. Convergence and completeness

A sequence  $(x_n)_{n=1}^\infty \in X^\mathbb{N}$  is said to **converge** to some point  $x \in X$  if  $\lim_{n \rightarrow \infty} d(x, x_n) = 0$ . We write  $\lim_{n \rightarrow \infty} x_n = x$  or  $x_n \rightarrow x$  as usual in this case. Clearly the **limit**  $x$  is unique if it exists (this is not true for a pseudometric). We will also frequently identify the sequence with its values  $x_n$  for simplicity of notation.

Note that convergent sequences are bounded. Here a set  $U \subseteq X$  is called **bounded** if it is contained within a ball, that is,  $U \subseteq B_r(x)$  for some  $x \in X$  and  $r > 0$ .

Note that convergence can also be equivalently formulated in topological terms: A sequence  $(x_n)_{n=1}^\infty$  converges to  $x$  if and only if for every neighborhood  $U(x)$  of  $x$  there is some  $N \in \mathbb{N}$  such that  $x_n \in U(x)$  for  $n \geq N$ . In a Hausdorff space the limit is unique. However, sequences usually do not suffice to describe a topology and, in general, definitions in terms of sequences are weaker (see the example below). This could be avoided by using generalized sequences, so-called nets, where the index set  $\mathbb{N}$  is replaced by arbitrary directed sets. We will not pursue this here.

**Example B.13.** For example, we can call a set  $U \subseteq X$  **sequentially closed** if every convergent sequence from  $U$  also has its limit in  $U$ . If  $U$  is closed, then every point in the complement is an interior point of the complement, thus no sequence from  $U$  can converge to such a point. Hence every closed set is sequentially closed. In a metric space (or more generally in a first countable space) we can find a sequence for every limit point  $x$  by choosing a point (different from  $x$ ) from every set in a neighborhood base. Hence the converse is also true in this case.  $\diamond$

Note that the argument from the previous example shows that in a first countable space sequentially closed is the same as closed. In particular, in this case the family of closed sets is uniquely determined by the convergent sequences:

**Lemma B.4.** *Two first countable topologies agree if and only if they give rise to the same convergent sequences.*

Of course every subsequence of a convergent sequence will converge to the same limit and we have the following converse:

**Lemma B.5.** *Let  $X$  be a topological space,  $(x_n)_{n=1}^\infty \in X^\mathbb{N}$  a sequence and  $x \in X$ . If every subsequence has a further subsequence which converges to  $x$ , then  $x_n$  converges to  $x$ .*

**Proof.** We argue by contradiction: If  $x_n \not\rightarrow x$  we can find a neighborhood  $U(x)$  and a subsequence  $x_{n_k} \notin U(x)$ . But then no subsequence of  $x_{n_k}$  can converge to  $x$ .  $\square$

This innocent observation is often useful to establish convergence in situations where one knows that the limit of a subsequence solves a given problem together with uniqueness of solutions for this problem. It can also be used to show that a notion of convergence does not stem from a topology (cf. Problem 5.11 from [48]).

In summary: A metric induces a natural topology and a topology induces a natural notion of convergence. However, a notion of convergence might not stem from a topology (or different topologies might give rise to the same notion of convergence) and a topology might not stem from a metric.

A sequence  $(x_n)_{n=1}^{\infty} \in X^{\mathbb{N}}$  is called a **Cauchy sequence** if for every  $\varepsilon > 0$  there is some  $N \in \mathbb{N}$  such that

$$d(x_n, x_m) \leq \varepsilon, \quad n, m \geq N. \quad (\text{B.11})$$

Every convergent sequence is a Cauchy sequence. If the converse is also true, that is, if every Cauchy sequence has a limit, then  $X$  is called **complete**. It is easy to see that a Cauchy sequence converges if and only if it has a convergent subsequence.

**Example B.14.** Both  $\mathbb{R}^n$  and  $\mathbb{C}^n$  are complete metric spaces.  $\diamond$

**Example B.15.** The metric

$$d(x, y) := |\arctan(x) - \arctan(y)| \quad (\text{B.12})$$

gives rise to the standard topology on  $\mathbb{R}$  (since  $\arctan$  is bi-Lipschitz on every compact interval). However,  $x_n = n$  is a Cauchy sequence with respect to this metric but not with respect to the usual metric. Moreover, any sequence with  $x_n \rightarrow \infty$  or  $x_n \rightarrow -\infty$  will be Cauchy with respect to this metric and hence (show this) for the completion of  $\mathbb{R}$  precisely the two new points  $-\infty$  and  $+\infty$  have to be added (cf. Example B.8).  $\diamond$

As noted before, in a metric space  $x$  is a limit point of  $U$  if and only if there exists a sequence  $(x_n)_{n=1}^{\infty} \subseteq U \setminus \{x\}$  with  $\lim_{n \rightarrow \infty} x_n = x$ . Hence  $U$  is closed if and only if for every convergent sequence the limit is in  $U$ . In particular,

**Lemma B.6.** *A subset of a complete metric space is again a complete metric space if and only if it is closed.*

A set  $U \subseteq X$  is called **dense** if its closure is all of  $X$ , that is, if  $\overline{U} = X$ . A space is called **separable** if it contains a countable dense set.

**Lemma B.7.** *A metric space is separable if and only if it is second countable as a topological space.*

**Proof.** From every dense set we get a countable base by considering open balls with rational radii and centers in the dense set. Conversely, from every countable base we obtain a dense set by choosing an element from each set in the base.  $\square$

**Lemma B.8.** *Let  $X$  be a separable metric space. Every subset  $Y$  of  $X$  is again separable.*

**Proof.** Let  $A = \{x_n\}_{n \in \mathbb{N}}$  be a dense set in  $X$ . The only problem is that  $A \cap Y$  might contain no elements at all. However, some elements of  $A$  must be at least arbitrarily close to this intersection: Let  $J \subseteq \mathbb{N}^2$  be the set of all pairs  $(n, m)$  for which  $B_{1/m}(x_n) \cap Y \neq \emptyset$  and choose some  $y_{n,m} \in B_{1/m}(x_n) \cap Y$  for all  $(n, m) \in J$ . Then  $B = \{y_{n,m}\}_{(n,m) \in J} \subseteq Y$  is countable. To see that  $B$  is dense, choose  $y \in Y$ . Then there is some sequence  $x_{n_k}$  with  $d(x_{n_k}, y) < 1/k$ . Hence  $(n_k, k) \in J$  and  $d(y_{n_k, k}, y) \leq d(y_{n_k, k}, x_{n_k}) + d(x_{n_k}, y) \leq 2/k \rightarrow 0$ .  $\square$

If  $X$  is an (incomplete) metric space, consider the set of all Cauchy sequences  $\mathcal{X}$  from  $X$ . Call two Cauchy sequences  $x = (x_n)_{n \in \mathbb{N}}$  and  $y = (y_n)_{n \in \mathbb{N}}$  equivalent if  $d_X(x_n, y_n) \rightarrow 0$  and denote by  $\bar{X}$  the set of all equivalence classes  $[x]$ . Moreover, the quadrangle inequality (Problem B.2) shows that if  $x = (x_n)_{n \in \mathbb{N}}$  and  $y = (y_n)_{n \in \mathbb{N}}$  are Cauchy sequences, so is  $d(x_n, y_n)$  and hence we can define a metric on  $\bar{X}$  via

$$d_{\bar{X}}([x], [y]) = \lim_{n \rightarrow \infty} d_X(x_n, y_n). \quad (\text{B.13})$$

Indeed, it is straightforward to check that  $d_{\bar{X}}$  is well defined (independent of the representative) and inherits all properties of a metric from  $d_X$ . Moreover,  $d_{\bar{X}}$  agrees with  $d_X$  of the limits whenever both Cauchy sequences converge in  $X$ .

**Theorem B.9.** *The space  $\bar{X}$  is a metric space containing  $X$  as a dense subspace if we identify  $x \in X$  with the equivalence class of all sequences converging to  $x$ . Moreover, this embedding is isometric.*

**Proof.** The map  $J : X \rightarrow \bar{X}$ ,  $x_0 \mapsto [(x_0, x_0, \dots)]$  is an isometric embedding (i.e., it is injective and preserves the metric). Moreover, for a Cauchy sequence  $x = (x_n)_{n \in \mathbb{N}}$  the sequence  $J(x_n)$  converges to  $[x]$  and hence  $J(X)$  is dense in  $\bar{X}$ . It remains to show that  $\bar{X}$  is complete. Let  $\xi_n = [(x_{n,j})_{j \in \mathbb{N}}]$  be a Cauchy sequence in  $\bar{X}$ . Without loss of generality (by dropping terms) we can choose the representatives  $x_{n,j}$  such that  $d(x_{n,j}, x_{n,k}) \leq \frac{1}{n}$  for  $j, k \geq n$ . Then it is not hard to see that  $\xi = [(x_{j,j})_{j \in \mathbb{N}}]$  is its limit.  $\square$

Notice that the completion  $\bar{X}$  is unique. More precisely, suppose  $\tilde{X}$  is another complete metric space which contains  $X$  as a dense subset such that the embedding  $\tilde{J} : X \hookrightarrow \tilde{X}$  is isometric. Then  $I = \tilde{J} \circ J^{-1} : J(X) \rightarrow \tilde{J}(X)$  has a unique isometric extension  $\bar{I} : \bar{X} \rightarrow \tilde{X}$  (compare Theorem B.39 below). In particular, it is no restriction to assume that a metric space is complete.

**Problem B.10.** *Let  $X$  be some nonempty set and define  $d(x, y) = 0$  if  $x = y$  and  $d(x, y) = 1$  if  $x \neq y$ . Show that  $(X, d)$  is a metric space. When are sequences convergent? When is  $X$  to be separable?*

**Problem B.11.** *Let  $U \subseteq V$  be subsets of a metric space  $X$ . Show that if  $U$  is dense in  $V$  and  $V$  is dense in  $X$ , then  $U$  is dense in  $X$ .*

**Problem B.12.** Let  $X$  be a metric space and denote by  $B(X)$  the set of all bounded functions  $X \rightarrow \mathbb{C}$ . Introduce the metric

$$d(f, g) = \sup_{x \in X} |f(x) - g(x)|.$$

Show that  $B(X)$  is complete.

**Problem B.13.** Let  $X$  be a metric space and  $B(X)$  as in the previous problem. Consider the embedding  $J : X \hookrightarrow B(X)$  defined via

$$y \mapsto J(x)(y) = d(x, y) - d(x_0, y)$$

for some fixed  $x_0 \in X$ . Show that this embedding is isometric. Hence  $\overline{J(X)}$  is another (equivalent) completion of  $X$ .

### B.3. Functions

Next, we come to functions  $f : X \rightarrow Y$ ,  $x \mapsto f(x)$ . We use the usual conventions  $f(U) := \{f(x) | x \in U\}$  for  $U \subseteq X$  and  $f^{-1}(V) := \{x | f(x) \in V\}$  for  $V \subseteq Y$ . Note

$$U \subseteq f^{-1}(f(U)), \quad f(f^{-1}(V)) \subseteq V. \quad (\text{B.14})$$

The set  $\text{Ran}(f) := f(X)$  is called the **range** of  $f$ , and  $X$  is called the **domain** of  $f$ . A function  $f$  is called **injective** or **one-to-one** if for each  $y \in Y$  there is at most one  $x \in X$  with  $f(x) = y$  (i.e.,  $f^{-1}(\{y\})$  contains at most one point) and **surjective** or **onto** if  $\text{Ran}(f) = Y$ . A function  $f$  which is both injective and surjective is called **bijective**.

Recall that we always have

$$\begin{aligned} f^{-1}\left(\bigcup_{\alpha} V_{\alpha}\right) &= \bigcup_{\alpha} f^{-1}(V_{\alpha}), & f^{-1}\left(\bigcap_{\alpha} V_{\alpha}\right) &= \bigcap_{\alpha} f^{-1}(V_{\alpha}), \\ f^{-1}(Y \setminus V) &= X \setminus f^{-1}(V) \end{aligned} \quad (\text{B.15})$$

as well as

$$\begin{aligned} f\left(\bigcup_{\alpha} U_{\alpha}\right) &= \bigcup_{\alpha} f(U_{\alpha}), & f\left(\bigcap_{\alpha} U_{\alpha}\right) &\subseteq \bigcap_{\alpha} f(U_{\alpha}), \\ f(X) \setminus f(U) &\subseteq f(X \setminus U) \end{aligned} \quad (\text{B.16})$$

with equality if  $f$  is injective.

A function  $f$  between metric spaces  $X$  and  $Y$  is called continuous at a point  $x \in X$  if for every  $\varepsilon > 0$  we can find a  $\delta > 0$  such that

$$d_Y(f(x), f(y)) \leq \varepsilon \quad \text{if} \quad d_X(x, y) < \delta. \quad (\text{B.17})$$

If  $f$  is continuous at every point, it is called **continuous**. In the case  $d_Y(f(x), f(y)) = d_X(x, y)$  we call  $f$  **isometric** and every isometry is of course continuous.



**Lemma B.10.** *Let  $f : X \rightarrow Y$  be a function between metric spaces  $X, Y$ . The following are equivalent:*

- (i)  $f$  is continuous at  $x$  (i.e., (B.17) holds).
- (ii)  $f(x_n) \rightarrow f(x)$  whenever  $x_n \rightarrow x$ .
- (iii) For every neighborhood  $V$  of  $f(x)$  the preimage  $f^{-1}(V)$  is a neighborhood of  $x$ .

**Proof.** (i)  $\Rightarrow$  (ii) is obvious. (ii)  $\Rightarrow$  (iii): If (iii) does not hold, there is a neighborhood  $V$  of  $f(x)$  such that  $B_\delta(x) \not\subseteq f^{-1}(V)$  for every  $\delta$ . Hence we can choose a sequence  $x_n \in B_{1/n}(x)$  such that  $x_n \notin f^{-1}(V)$ . Thus  $x_n \rightarrow x$  but  $f(x_n) \not\rightarrow f(x)$ . (iii)  $\Rightarrow$  (i): Choose  $V = B_\varepsilon(f(x))$  and observe that by (iii),  $B_\delta(x) \subseteq f^{-1}(V)$  for some  $\delta$ .  $\square$

In a general topological space we use (iii) as the definition of continuity and (ii) is called **sequential continuity**. Then continuity will imply sequential continuity but the converse will not be true unless we assume (e.g.) that  $X$  is first countable (Problem B.14).

In particular, (iii) implies that  $f$  is continuous if and only if the preimage of every open set is again open (equivalently, the inverse image of every closed set is closed). Note that by (B.15) it suffices to check this for sets in a subbase. If the image of every open set is open, then  $f$  is called **open**. A bijection  $f$  is called a **homeomorphism** if both  $f$  and its inverse  $f^{-1}$  are continuous. Note that if  $f$  is a bijection, then  $f^{-1}$  is continuous if and only if  $f$  is open. Two topological spaces are called **homeomorphic** if there is a homeomorphism between them.

In a topological space  $X$  a function  $f : X \rightarrow \overline{\mathbb{R}}$  is **lower semicontinuous** if the set  $f^{-1}((a, \infty])$  is open for every  $a \in \mathbb{R}$ . Similarly,  $f$  is **upper semicontinuous** if the set  $f^{-1}([-\infty, a))$  is open for every  $a \in \mathbb{R}$ . Clearly  $f$  is lower semicontinuous if and only if  $-f$  is upper semicontinuous.

Finally, the **support** of a function  $f : X \rightarrow \mathbb{C}^n$  is the closure of all points  $x$  for which  $f(x)$  does not vanish; that is,

$$\text{supp}(f) := \overline{\{x \in X \mid f(x) \neq 0\}}. \quad (\text{B.18})$$

**Problem B.14.** *Let  $X, Y$  be topological spaces. Show that if  $f : X \rightarrow Y$  is continuous at  $x \in X$  then it is also sequential continuous. Show that the converse holds if  $X$  is first countable.*

**Problem B.15.** *Let  $f : X \rightarrow Y$  be continuous. Then  $f(\overline{A}) \subseteq \overline{f(A)}$ .*

**Problem B.16.** *Let  $X, Y$  be topological spaces and let  $f : X \rightarrow Y$  be continuous. Show that if  $X$  is separable, then so is  $f(X)$ .*

**Problem B.17.** Let  $X$  be a topological space and  $f : X \rightarrow \mathbb{R}$ . Let  $x_0 \in X$  and let  $\mathcal{B}(x_0)$  be a neighborhood base for  $x_0$ . Define

$$\liminf_{x \rightarrow x_0} f(x) := \sup_{U \in \mathcal{B}(x_0)} \inf_{U(x_0)} f, \quad \limsup_{x \rightarrow x_0} f(x) := \inf_{U \in \mathcal{B}(x_0)} \sup_{U(x_0)} f.$$

Show that both are independent of the neighborhood base and satisfy

- (i)  $\liminf_{x \rightarrow x_0} (-f(x)) = -\limsup_{x \rightarrow x_0} f(x)$ .
- (ii)  $\liminf_{x \rightarrow x_0} (\alpha f(x)) = \alpha \liminf_{x \rightarrow x_0} f(x)$ ,  $\alpha \geq 0$ .
- (iii)  $\liminf_{x \rightarrow x_0} (f(x) + g(x)) \geq \liminf_{x \rightarrow x_0} f(x) + \liminf_{x \rightarrow x_0} g(x)$ .

Moreover, show that

$$\liminf_{n \rightarrow \infty} f(x_n) \geq \liminf_{x \rightarrow x_0} f(x), \quad \limsup_{n \rightarrow \infty} f(x_n) \leq \limsup_{x \rightarrow x_0} f(x)$$

for every sequence  $x_n \rightarrow x_0$  and there exists a sequence attaining equality if  $X$  is a metric space.

**Problem\* B.18.** Show that the supremum over lower semicontinuous functions is again lower semicontinuous.

**Problem\* B.19.** Let  $X$  be a topological space and  $f : X \rightarrow \overline{\mathbb{R}}$ . Show that  $f$  is lower semicontinuous if and only if

$$\liminf_{x \rightarrow x_0} f(x) \geq f(x_0), \quad x_0 \in X.$$

Similarly,  $f$  is upper semicontinuous if and only if

$$\limsup_{x \rightarrow x_0} f(x) \leq f(x_0), \quad x_0 \in X.$$

Show that a lower semicontinuous function is also sequentially lower semicontinuous

$$\liminf_{n \rightarrow \infty} f(x_n) \geq f(x_0), \quad x_n \rightarrow x_0, x_0 \in X.$$

Show the converse if  $X$  is a metric space. (Hint: Problem B.17.)

## B.4. Product topologies

If  $X$  and  $Y$  are metric spaces, then  $X \times Y$  together with

$$d((x_1, y_1), (x_2, y_2)) := d_X(x_1, x_2) + d_Y(y_1, y_2) \quad (\text{B.19})$$

is a metric space. A sequence  $(x_n, y_n)$  converges to  $(x, y)$  if and only if  $x_n \rightarrow x$  and  $y_n \rightarrow y$ . In particular, the projections onto the first  $(x, y) \mapsto x$ , respectively, onto the second  $(x, y) \mapsto y$ , coordinate are continuous. Moreover, if  $X$  and  $Y$  are complete, so is  $X \times Y$ .

In particular, by the inverse triangle inequality (B.1),

$$|d(x_n, y_n) - d(x, y)| \leq d(x_n, x) + d(y_n, y), \quad (\text{B.20})$$

we see that  $d : X \times X \rightarrow \mathbb{R}$  is continuous.

**Example B.16.** If we consider  $\mathbb{R} \times \mathbb{R}$ , we do not get the Euclidean distance of  $\mathbb{R}^2$  unless we modify (B.19) as follows:

$$\tilde{d}((x_1, y_1), (x_2, y_2)) := \sqrt{d_X(x_1, x_2)^2 + d_Y(y_1, y_2)^2}. \quad (\text{B.21})$$

As noted in our previous example, the topology (and thus also convergence/continuity) is independent of this choice.  $\diamond$

If  $X$  and  $Y$  are just topological spaces, the **product topology** is defined by calling  $O \subseteq X \times Y$  open if for every point  $(x, y) \in O$  there are open neighborhoods  $U$  of  $x$  and  $V$  of  $y$  such that  $U \times V \subseteq O$ . In other words, the products of open sets form a base of the product topology. Again the projections onto the first and second component are continuous. In the case of metric spaces this clearly agrees with the topology defined via the product metric (B.19). There is also another way of constructing the product topology, namely, as the weakest topology which makes the projections continuous. In fact, this topology must contain all sets which are inverse images of open sets  $U \subseteq X$ , that is all sets of the form  $U \times Y$  as well as all inverse images of open sets  $V \subseteq Y$ , that is all sets of the form  $X \times V$ . Adding finite intersections we obtain all sets of the form  $U \times V$  and hence the same base as before. In particular, a sequence  $(x_n, y_n)$  will converge if and only if both components converge.

Note that the product topology immediately extends to the product of an arbitrary number of spaces  $X := \bigtimes_{\alpha \in A} X_\alpha$  by defining it as the weakest topology which makes all projections  $\pi_\alpha : X \rightarrow X_\alpha$  continuous.

**Example B.17.** Let  $X$  be a topological space and  $A$  an index set. Then  $X^A = \bigtimes_A X$  is the set of all functions  $x : A \rightarrow X$  and a neighborhood base at  $x$  are sets of functions which coincide with  $x$  at a given finite number of points. Convergence with respect to the product topology corresponds to pointwise convergence (note that the projection  $\pi_\alpha$  is the point evaluation at  $\alpha$ :  $\pi_\alpha(x) = x(\alpha)$ ). If  $A$  is uncountable (and  $X$  is not equipped with the trivial topology), then there is no countable neighborhood base (if there were such a base, it would involve only a countable number of points, now choose a point from the complement ...). In particular, there is no corresponding metric even if  $X$  has one. Moreover, this topology cannot be characterized with sequences alone. For example, let  $X = \{0, 1\}$  (with the discrete topology) and  $A = \mathbb{R}$ . Then the set  $F = \{x | x^{-1}(1) \text{ is countable}\}$  is sequentially closed but its closure is all of  $\{0, 1\}^{\mathbb{R}}$  (every set from our neighborhood base contains an element which vanishes except at finitely many points).  $\diamond$

In fact this is a special case of a more general construction which is often used. Let  $\{f_\alpha\}_{\alpha \in A}$  be a collection of functions  $f_\alpha : X \rightarrow Y_\alpha$ , where  $Y_\alpha$  are some topological spaces. Then we can equip  $X$  with the weakest topology (known as the **initial topology**) which makes all  $f_\alpha$  continuous. That is, we

take the topology generated by sets of the forms  $f_\alpha^{-1}(O_\alpha)$ , where  $O_\alpha \subseteq Y_\alpha$  is open, known as open **cylinders**. Finite intersections of such sets, known as open **cylinder sets**, are hence a base for the topology and a sequence  $x_n$  will converge to  $x$  if and only if  $f_\alpha(x_n) \rightarrow f_\alpha(x)$  for all  $\alpha \in A$ . In particular, if the collection is countable, then  $X$  will be first (or second) countable if all  $Y_\alpha$  are.

The initial topology has the following characteristic property:

**Lemma B.11.** *Let  $X$  have the initial topology from a collection of functions  $\{f_\alpha : X \rightarrow Y_\alpha\}_{\alpha \in A}$  and let  $Z$  be another topological space. A function  $f : Z \rightarrow X$  is continuous (at  $z$ ) if and only if  $f_\alpha \circ f$  is continuous (at  $z$ ) for all  $\alpha \in A$ .*

**Proof.** If  $f$  is continuous at  $z$ , then so is the composition  $f_\alpha \circ f$ . Conversely, let  $U \subseteq X$  be a neighborhood of  $f(z)$ . Then  $\bigcap_{j=1}^n f_{\alpha_j}^{-1}(O_{\alpha_j}) \subseteq U$  for some  $\alpha_j$  and some open neighborhoods  $O_{\alpha_j}$  of  $f_{\alpha_j}(f(z))$ . But then  $f^{-1}(U)$  contains the neighborhood  $f^{-1}(\bigcap_{j=1}^n f_{\alpha_j}^{-1}(O_{\alpha_j})) = \bigcap_{j=1}^n (f_{\alpha_j} \circ f)^{-1}(O_{\alpha_j})$  of  $z$ .  $\square$

If all  $Y_\alpha$  are Hausdorff and if the collection  $\{f_\alpha\}_{\alpha \in A}$  **separates points**, that is for every  $x \neq y$  there is some  $\alpha$  with  $f_\alpha(x) \neq f_\alpha(y)$ , then  $X$  will again be Hausdorff. Indeed for  $x \neq y$  choose  $\alpha$  such that  $f_\alpha(x) \neq f_\alpha(y)$  and let  $U_\alpha, V_\alpha$  be two disjoint neighborhoods separating  $f_\alpha(x), f_\alpha(y)$ . Then  $f_\alpha^{-1}(U_\alpha), f_\alpha^{-1}(V_\alpha)$  are two disjoint neighborhoods separating  $x, y$ . In particular,  $X = \bigtimes_{\alpha \in A} X_\alpha$  is Hausdorff if all  $X_\alpha$  are.

Note that a similar construction works in the other direction. Let  $\{f_\alpha\}_{\alpha \in A}$  be a collection of functions  $f_\alpha : X_\alpha \rightarrow Y$ , where  $X_\alpha$  are some topological spaces. Then we can equip  $Y$  with the strongest topology (known as the **final topology**) which makes all  $f_\alpha$  continuous. That is, we take as open sets those for which  $f_\alpha^{-1}(O)$  is open for all  $\alpha \in A$ .

**Example B.18.** Let  $\sim$  be an equivalence relation on  $X$  with equivalence classes  $[x] = \{y \in X \mid x \sim y\}$ . Then the **quotient topology** on the set of equivalence classes  $X/\sim$  is the final topology of the projection map  $\pi : X \rightarrow X/\sim$ .  $\diamond$

**Example B.19.** Let  $X_\alpha$  be a collection of topological spaces. The **disjoint union**

$$X := \bigsqcup_{\alpha \in A} X_\alpha$$

is usually given the final topology from the canonical injections  $i_\alpha : X_\alpha \hookrightarrow X$  such that  $O \subseteq X$  is open if and only if  $O \cap X_\alpha$  is open for all  $\alpha \in A$ .  $\diamond$

**Lemma B.12.** *Let  $Y$  have the final topology from a collection of functions  $\{f_\alpha : X_\alpha \rightarrow Y\}_{\alpha \in A}$  and let  $Z$  be another topological space. A function  $f : Y \rightarrow Z$  is continuous if and only if  $f \circ f_\alpha$  is continuous for all  $\alpha \in A$ .*

**Proof.** If  $f$  is continuous, then so is the composition  $f \circ f_\alpha$ . Conversely, let  $V \subseteq Z$  be open. Then  $f \circ f_\alpha$  implies  $(f \circ f_\alpha)^{-1}(V) = f_\alpha^{-1}(f^{-1}(V))$  open for all  $\alpha$  and hence  $f^{-1}(V)$  open.  $\square$

**Problem B.20.** Show that  $X$  is Hausdorff if and only if the diagonal  $\Delta := \{(x, x) | x \in X\} \subseteq X \times X$  is closed.

**Problem B.21.** Let  $X = \prod_{\alpha \in A} X_\alpha$  with the product topology. Show that the projection maps are open.

**Problem B.22.** Let  $X = \prod_{\alpha \in A} X_\alpha$  with the product topology. Show that the product  $\prod_{\alpha \in A} C_\alpha$  of closed sets  $C_\alpha \subseteq X_\alpha$  is closed.

**Problem B.23** (Gluing lemma). Suppose  $X, Y$  are topological spaces and  $f_\alpha : A_\alpha \rightarrow Y$  are continuous functions defined on  $A_\alpha \subseteq X$ . Suppose  $f_\alpha = f_\beta$  on  $A_\alpha \cap A_\beta$  such that  $f : A := \bigcup_\alpha A_\alpha \rightarrow Y$  is well defined by  $f(x) = f_\alpha(x)$  if  $x \in A_\alpha$ . Show that  $f$  is continuous if either all sets  $A_\alpha$  are open or if the collection  $A_\alpha$  is finite and all are closed.

**Problem B.24.** Let  $\{(X_j, d_j)\}_{j \in \mathbb{N}}$  be a sequence of metric spaces. Show that

$$d(x, y) = \sum_{j \in \mathbb{N}} \frac{1}{2^j} \frac{d_j(x_j, y_j)}{1 + d_j(x_j, y_j)} \quad \text{or} \quad d(x, y) = \max_{j \in \mathbb{N}} \frac{1}{2^j} \frac{d_j(x_j, y_j)}{1 + d_j(x_j, y_j)}$$

is a metric on  $X = \prod_{n \in \mathbb{N}} X_n$  which generates the product topology. Show that  $X$  is complete if all  $X_n$  are.

## B.5. Compactness

A **cover** of a set  $Y \subseteq X$  is a family of sets  $\{U_\alpha\}$  such that  $Y \subseteq \bigcup_\alpha U_\alpha$ . A cover is called open if all  $U_\alpha$  are open. Any subset of  $\{U_\alpha\}$  which still covers  $Y$  is called a **subcover**.

**Lemma B.13** (Lindelöf). If  $X$  is second countable, then every open cover has a countable subcover.

**Proof.** Let  $\{U_\alpha\}$  be an open cover for  $Y$ , and let  $\mathcal{B}$  be a countable base. Since every  $U_\alpha$  can be written as a union of elements from  $\mathcal{B}$ , the set of all  $B \in \mathcal{B}$  which satisfy  $B \subseteq U_\alpha$  for some  $\alpha$  form a countable open cover for  $Y$ . Moreover, for every  $B_n$  in this set we can find an  $\alpha_n$  such that  $B_n \subseteq U_{\alpha_n}$ . By construction,  $\{U_{\alpha_n}\}$  is a countable subcover.  $\square$

A **refinement**  $\{V_\beta\}$  of a cover  $\{U_\alpha\}$  is a cover such that for every  $\beta$  there is some  $\alpha$  with  $V_\beta \subseteq U_\alpha$ . A cover is called **locally finite** if every point has a neighborhood that intersects only finitely many sets in the cover.

**Lemma B.14** (Stone). In a metric space every countable open cover has a locally finite open refinement.

**Proof.** Denote the cover by  $\{O_j\}_{j \in \mathbb{N}}$  and introduce the sets

$$\hat{O}_{j,n} := \bigcup_{x \in A_{j,n}} B_{2^{-n}}(x), \text{ where}$$

$$A_{j,n} := \{x \in O_j \setminus (O_1 \cup \cdots \cup O_{j-1}) \mid x \notin \bigcup_{k \in \mathbb{N}, 1 \leq l < n} \hat{O}_{k,l} \text{ and } B_{3 \cdot 2^{-n}}(x) \subseteq O_j\}.$$

Then, by construction,  $\hat{O}_{j,n}$  is open,  $\hat{O}_{j,n} \subseteq O_j$ , and it is a cover since for every  $x$  there is a smallest  $j$  such that  $x \in O_j$  and a smallest  $n$  such that  $B_{3 \cdot 2^{-n}}(x) \subseteq O_j$  implying  $x \in \hat{O}_{j,n}$  for some  $n$ .

To show that  $\hat{O}_{j,n}$  is locally finite fix some  $x$  and let  $j$  be the smallest integer such that  $x \in \hat{O}_{j,n}$  for some  $n$ . Moreover, choose  $m$  such that  $B_{2^{-m}}(x) \subseteq \hat{O}_{j,n}$ . It suffices to show that:

- (i) If  $i \geq n + m$  then  $B_{2^{-n-m}}(x)$  is disjoint from  $\hat{O}_{k,i}$  for all  $k$ .
- (ii) If  $i < n + m$  then  $B_{2^{-n-m}}(x)$  intersects  $\hat{O}_{k,i}$  for at most one  $k$ .

To show (i) observe that since  $i > n$  every ball  $B_{2^{-i}}(y)$  used in the definition of  $\hat{O}_{k,i}$  has its center outside of  $\hat{O}_{j,n}$ . Hence  $d(x, y) \geq 2^{-m}$  and  $B_{2^{-n-m}}(x) \cap B_{2^{-i}}(y) = \emptyset$  since  $i \geq m + 1$  as well as  $n + m \geq m + 1$ .

To show (ii) let  $y \in \hat{O}_{j,i}$  and  $z \in \hat{O}_{k,i}$  with  $j < k$ . We will show  $d(y, z) > 2^{-n-m+1}$ . There are points  $r$  and  $s$  such that  $y \in B_{2^{-i}}(r) \subseteq \hat{O}_{j,i}$  and  $z \in B_{2^{-i}}(s) \subseteq \hat{O}_{k,i}$ . Then by definition  $B_{3 \cdot 2^{-i}}(r) \subseteq O_j$  but  $s \notin O_j$ . So  $d(r, s) \geq 3 \cdot 2^{-i}$  and  $d(y, z) > 2^{-i} \geq 2^{-n-m+1}$ .  $\square$

A subset  $K \subset X$  is called **compact** if every open cover of  $K$  has a finite subcover. A set is called **relatively compact** if its closure is compact.

**Lemma B.15.** *A topological space is compact if and only if it has the **finite intersection property**: The intersection of a family of closed sets is empty if and only if the intersection of some finite subfamily is empty.*

**Proof.** By taking complements, to every family of open sets there corresponds a family of closed sets and vice versa. Moreover, the open sets are a cover if and only if the corresponding closed sets have empty intersection.  $\square$

**Lemma B.16.** *Let  $X$  be a topological space.*

- (i) *The continuous image of a compact set is compact.*
- (ii) *Every closed subset of a compact set is compact.*
- (iii) *If  $X$  is Hausdorff, every compact set is closed.*
- (iv) *The finite union of compact sets is compact.*
- (v) *If  $X$  is Hausdorff, any intersection of compact sets is compact.*

**Proof.** (i) Observe that if  $\{O_\alpha\}$  is an open cover for  $f(Y)$ , then  $\{f^{-1}(O_\alpha)\}$  is one for  $Y$ .

(ii) Let  $\{O_\alpha\}$  be an open cover for the closed subset  $Y$  (in the induced topology). Then there are open sets  $\tilde{O}_\alpha$  with  $O_\alpha = \tilde{O}_\alpha \cap Y$  and  $\{\tilde{O}_\alpha\} \cup \{X \setminus Y\}$  is an open cover for  $X$  which has a finite subcover. This subcover induces a finite subcover for  $Y$ .

(iii) Let  $Y \subseteq X$  be compact. We show that  $X \setminus Y$  is open. Fix  $x \in X \setminus Y$  (if  $Y = X$  there is nothing to do). By the definition of Hausdorff, for every  $y \in Y$  there are disjoint neighborhoods  $V(y)$  of  $y$  and  $U_y(x)$  of  $x$ . By compactness of  $Y$ , there are  $y_1, \dots, y_n$  such that the  $V(y_j)$  cover  $Y$ . But then  $\bigcap_{j=1}^n U_{y_j}(x)$  is a neighborhood of  $x$  which does not intersect  $Y$ .

(iv) Note that a cover of the union is a cover for each individual set and the union of the individual subcovers is the subcover we are looking for.

(v) Follows from (ii) and (iii) since an intersection of closed sets is closed.  $\square$

As a consequence we obtain a simple criterion when a continuous function is a homeomorphism.

**Corollary B.17.** *Let  $X$  and  $Y$  be topological spaces with  $X$  compact and  $Y$  Hausdorff. Then every continuous bijection  $f : X \rightarrow Y$  is a homeomorphism.*

**Proof.** It suffices to show that  $f$  maps closed sets to closed sets. By (ii) every closed set is compact, by (i) its image is also compact, and by (iii) it is also closed.  $\square$

Concerning products of compact sets we have

**Theorem B.18** (Tychonoff). *The product  $\prod_{\alpha \in A} K_\alpha$  of an arbitrary collection of compact topological spaces  $\{K_\alpha\}_{\alpha \in A}$  is compact with respect to the product topology.*

**Proof.** We say that a family  $\mathcal{F}$  of closed subsets of  $K$  has the finite intersection property if the intersection of every finite subfamily has nonempty intersection. The collection of all such families which contain  $\mathcal{F}$  is partially ordered by inclusion and every chain has an upper bound (the union of all sets in the chain). Hence, by Zorn's lemma, there is a maximal family  $\mathcal{F}_M$  (note that this family is closed under finite intersections).

Denote by  $\pi_\alpha : K \rightarrow K_\alpha$  the projection onto the  $\alpha$  component. Then the closed sets  $\{\pi_\alpha(F)\}_{F \in \mathcal{F}_M}$  also have the finite intersection property and since  $K_\alpha$  is compact, there is some  $x_\alpha \in \bigcap_{F \in \mathcal{F}_M} \overline{\pi_\alpha(F)}$ . Consequently, if  $F_\alpha$  is a closed neighborhood of  $x_\alpha$ , then  $\pi_\alpha^{-1}(F_\alpha) \in \mathcal{F}_M$  (otherwise there would be some  $F \in \mathcal{F}_M$  with  $F \cap \pi_\alpha^{-1}(F_\alpha) = \emptyset$  contradicting  $F_\alpha \cap \pi_\alpha(F) \neq \emptyset$ ).

Furthermore, for every finite subset  $A_0 \subseteq A$  we have  $\bigcap_{\alpha \in A_0} \pi_\alpha^{-1}(F_\alpha) \in \mathcal{F}_M$  and so every neighborhood of  $x = (x_\alpha)_{\alpha \in A}$  intersects  $F$ . Since  $F$  is closed,  $x \in F$  and hence  $x \in \bigcap_{\mathcal{F}_M} F$ .  $\square$

A subset  $K \subseteq X$  is called **sequentially compact** if every sequence from  $K$  has a convergent subsequence whose limit is in  $K$ . In a metric space, compactness and sequentially compactness are equivalent.

**Lemma B.19.** *Let  $X$  be a metric space. Then a subset is compact if and only if it is sequentially compact.*

**Proof.** Without loss of generality we can assume the subset to be all of  $X$ . Suppose  $X$  is compact and let  $x_n$  be a sequence which has no convergent subsequence. Then  $K := \{x_n\}$  has no limit points and is hence compact by Lemma B.16 (ii). For every  $n$  there is a ball  $B_{\varepsilon_n}(x_n)$  which contains only finitely many elements of  $K$ . However, finitely many balls suffice to cover  $K$ , a contradiction.

Conversely, suppose  $X$  is sequentially compact and let  $\{O_\alpha\}$  be some open cover which has no finite subcover. For every  $x \in X$  we can choose some  $\alpha(x)$  such that if  $B_r(x)$  is the largest ball contained in  $O_{\alpha(x)}$ , then there is no  $\beta$  with  $B_{2r}(x) \subset O_\beta$  (show that this is possible). Now choose a sequence  $x_n$  such that  $x_n \notin \bigcup_{m < n} O_{\alpha(x_m)}$ . Note that by construction the distance  $d = d(x_m, x_n)$  to every successor of  $x_m$  satisfies that the ball  $B_{2d}(x_m)$  does not fit into any of the  $O_\alpha$ .

Now let  $y$  be the limit of some convergent subsequence and fix some  $r > 0$  such that  $B_r(y) \subseteq O_{\alpha(y)}$ . Then this subsequence must eventually be in  $B_{r/6}(y)$ , but this is impossible since if  $d := d(x_{n_1}, x_{n_2}) < r/3$  is the distance between two consecutive elements of this subsequence within  $B_{r/6}(y)$ , then  $B_{2d}(x_{n_1})$  cannot fit into  $O_{\alpha(y)}$  by construction whereas on the other hand  $B_{2d}(x_{n_1}) \subseteq B_r(y) \subseteq O_{\alpha(y)}$ .  $\square$

If we drop the requirement that the limit must be in  $K$ , we obtain relatively compact sets:

**Corollary B.20.** *Let  $X$  be a metric space and  $K \subset X$ . Then  $K$  is relatively compact if and only if every sequence from  $K$  has a convergent subsequence (the limit need not be in  $K$ ).*

**Proof.** For any sequence  $x_n \in \overline{K}$  we can find a nearby sequence  $y_n \in K$  with  $x_n - y_n \rightarrow 0$ . If we can find a convergent subsequence of  $y_n$  then the corresponding subsequence of  $x_n$  will also converge (to the same limit) and  $\overline{K}$  is (sequentially) compact in this case. The converse is trivial.  $\square$

As another simple consequence observe that



**Corollary B.21.** *A compact metric space  $X$  is complete and separable.*

**Proof.** Completeness is immediate from the previous lemma. To see that  $X$  is separable note that, by compactness, for every  $n \in \mathbb{N}$  there is a finite set  $S_n \subseteq X$  such that the balls  $\{B_{1/n}(x)\}_{x \in S_n}$  cover  $X$ . Then  $\bigcup_{n \in \mathbb{N}} S_n$  is a countable dense set.  $\square$

Recall that a set in a metric space is called bounded if it is contained inside some ball. Clearly the union of two bounded sets is bounded. Moreover, compact sets are always bounded since they can be covered by finitely many balls. In  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) the converse also holds.

**Theorem B.22** (Heine–Borel). *In  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) a set is compact if and only if it is bounded and closed.*

**Proof.** By Lemma B.16 (ii), (iii), and Tychonoff's theorem it suffices to show that a closed interval in  $I \subseteq \mathbb{R}$  is compact. Moreover, by Lemma B.19, it suffices to show that every sequence in  $I = [a, b]$  has a convergent subsequence. Let  $x_n$  be our sequence and divide  $I = [a, \frac{a+b}{2}] \cup [\frac{a+b}{2}, b]$ . Then at least one of these two intervals, call it  $I_1$ , contains infinitely many elements of our sequence. Let  $y_1 = x_{n_1}$  be the first one. Subdivide  $I_1$  and pick  $y_2 = x_{n_2}$ , with  $n_2 > n_1$  as before. Proceeding like this, we obtain a Cauchy sequence  $y_n$  (note that by construction  $I_{n+1} \subseteq I_n$  and hence  $|y_n - y_m| \leq \frac{b-a}{2^n}$  for  $m \geq n$ ).  $\square$

By Lemma B.19 this is equivalent to

**Theorem B.23** (Bolzano–Weierstraß). *Every bounded infinite subset of  $\mathbb{R}^n$  (or  $\mathbb{C}^n$ ) has at least one limit point.*

Combining Theorem B.22 with Lemma B.16 (i) we also obtain the **extreme value theorem**.

**Theorem B.24** (Weierstraß). *Let  $X$  be compact. Every continuous function  $f : X \rightarrow \mathbb{R}$  attains its maximum and minimum.*

A metric space  $X$  for which the Heine–Borel theorem holds is called **proper**. Lemma B.16 (ii) shows that  $X$  is proper if and only if every closed ball is compact. Note that a proper metric space must be complete (since every Cauchy sequence is bounded). A topological space is called **locally compact** if every point has a compact neighborhood. Clearly a proper metric space is locally compact. A topological space is called  **$\sigma$ -compact**, if it can be written as a countable union of compact sets. Again a proper space is  $\sigma$ -compact.

**Lemma B.25.** *For a metric space  $X$  the following are equivalent:*

- (i)  $X$  is separable and locally compact.
- (ii)  $X$  contains a countable base consisting of relatively compact sets.
- (iii)  $X$  is locally compact and  $\sigma$ -compact.
- (iv)  $X$  can be written as the union of an increasing sequence  $U_n$  of relatively compact open sets satisfying  $\overline{U_n} \subseteq U_{n+1}$  for all  $n$ .

**Proof.** (i)  $\Rightarrow$  (ii): Let  $\{x_n\}$  be a dense set. Then the balls  $B_{n,m} = B_{1/m}(x_n)$  form a base. Moreover, for every  $n$  there is some  $m_n$  such that  $B_{n,m}$  is relatively compact for  $m \leq m_n$ . Since those balls are still a base we are done. (ii)  $\Rightarrow$  (iii): Take the union over the closures of all sets in the base. (iii)  $\Rightarrow$  (vi): Let  $X = \bigcup_n K_n$  with  $K_n$  compact. Without loss  $K_n \subseteq K_{n+1}$ . For a given compact set  $K$  we can find a relatively compact open set  $V(K)$  such that  $K \subseteq V(K)$  (cover  $K$  by relatively compact open balls and choose a finite subcover). Now define  $U_n = V(\overline{U_n})$ . (vi)  $\Rightarrow$  (i): Each of the sets  $\overline{U_n}$  has a countable dense subset by Corollary B.21. The union gives a countable dense set for  $X$ . Since every  $x \in U_n$  for some  $n$ ,  $X$  is also locally compact.  $\square$

A sequence of compact sets  $\overline{U_n}$  as in (iv) is known as a **compact exhaustion** of  $X$ . Since every point of  $X$  is an interior point of some  $U_n$ , every compact set  $K \subseteq X$  is contained in  $U_n$  for  $n$  sufficiently large (cover  $K$  by taking an open neighborhood contained in some  $U_n$  for every point).

**Example B.20.**  $X = (0, 1)$  with the usual metric is locally compact and  $\sigma$ -compact but not proper.  $\diamond$

**Example B.21.** Consider  $\ell^2(\mathbb{N})$  with the standard basis  $\delta^j$ . Let  $X_j := \{\lambda \delta^j \mid \lambda \in [0, 1]\}$  and note that the metric on  $X_j$  inherited from  $\ell^2(\mathbb{Z})$  is the same as the usual metric from  $\mathbb{R}$ . Then  $X := \bigcup_{j \in \mathbb{N}} X_j$  is a complete separable  $\sigma$ -compact space, which is not locally compact. In fact, consider a ball of radius  $\varepsilon$  around zero. Then  $(\varepsilon/2)\delta^j \in B_\varepsilon(0)$  is a bounded sequence which has no convergent subsequence since  $d((\varepsilon/2)\delta^j, (\varepsilon/2)\delta^k) = \varepsilon/\sqrt{2}$  for  $k \neq j$ .  $\diamond$

However, under the assumptions of the previous lemma we can always switch to a new metric which generates the same topology and for which  $X$  is proper. To this end recall that a function  $f : X \rightarrow Y$  between topological spaces is called **proper** if the inverse image of a compact set is again compact. Now given a proper function (Problem B.32) there is a new metric with the claimed properties (Problem B.27).

A subset  $U$  of a complete metric space  $X$  is called **totally bounded** if for every  $\varepsilon > 0$  it can be covered with a finite number of balls of radius  $\varepsilon$ . We will call such a cover an  $\varepsilon$ -cover. Clearly every totally bounded set is bounded.

**Example B.22.** Of course in  $\mathbb{R}^n$  the totally bounded sets are precisely the bounded sets. This is in fact true for every proper metric space since the closure of a bounded set is compact and hence has a finite cover.  $\diamond$

**Lemma B.26.** *Let  $X$  be a complete metric space. Then a set is relatively compact if and only if it is totally bounded.*

**Proof.** Without loss of generality we can assume our set to be closed. Clearly a compact set  $K$  is closed and totally bounded (consider the cover by all balls of radius  $\varepsilon$  with center in the set and choose a finite subcover). Conversely, we will show that  $K$  is sequentially compact. So start with  $\varepsilon_1 = 1$  and choose a finite cover of balls with radius  $\varepsilon_1$ . One of these balls contains an infinite number of elements  $x_n^1$  from our sequence  $x_n$ . Choose  $\varepsilon_2 = \frac{1}{2}$  and repeat the process with the sequence  $x_n^1$ . The resulting diagonal sequence  $x_n^n$  gives a subsequence which is Cauchy and hence converges by completeness.  $\square$

**Problem\* B.25** (Alexandroff one-point compactification). *Suppose  $X$  is a locally compact Hausdorff space which is not compact. Introduce a new point  $\infty$ , set  $\hat{X} = X \cup \{\infty\}$  and make it into a topological space by calling  $O \subseteq \hat{X}$  open if either  $\infty \notin O$  and  $O$  is open in  $X$  or if  $\infty \in O$  and  $\hat{X} \setminus O$  is compact. Show that  $\hat{X}$  is a compact Hausdorff space which contains  $X$  as a dense subset.*

**Problem B.26.** *Show that every open set  $O \subseteq \mathbb{R}$  can be written as a countable union of disjoint intervals. (Hint: Consider the set  $\{I_\alpha\}$  of all maximal open subintervals of  $O$ ; that is,  $I_\alpha \subseteq O$  and there is no other subinterval of  $O$  which contains  $I_\alpha$ .)*

**Problem B.27.** *Let  $(X, d)$  be a metric space. Show that if there is a proper function  $f : X \rightarrow \mathbb{R}$ , then*

$$\tilde{d}(x, y) = d(x, y) + |f(x) - f(y)|$$

*is a metric which generates the same topology and for which  $(X, \tilde{d})$  is proper.*

## B.6. Separation

The **distance** between a point  $x \in X$  and a subset  $Y \subseteq X$  is

$$\text{dist}(x, Y) := \inf_{y \in Y} d(x, y). \quad (\text{B.22})$$

Note that  $x \in \overline{Y}$  if and only if  $\text{dist}(x, Y) = 0$  (Problem B.28).

**Lemma B.27.** *Let  $X$  be a metric space and  $Y \subseteq X$  nonempty. Then*

$$|\text{dist}(x, Y) - \text{dist}(z, Y)| \leq d(x, z). \quad (\text{B.23})$$

*In particular,  $x \mapsto \text{dist}(x, Y)$  is continuous.*

**Proof.** Taking the infimum in the triangle inequality  $d(x, y) \leq d(x, z) + d(z, y)$  shows  $\text{dist}(x, Y) \leq d(x, z) + \text{dist}(z, Y)$ . Hence  $\text{dist}(x, Y) - \text{dist}(z, Y) \leq d(x, z)$ . Interchanging  $x$  and  $z$  shows  $\text{dist}(z, Y) - \text{dist}(x, Y) \leq d(x, z)$ .  $\square$

A topological space is called **normal** if for any two disjoint closed sets  $C_1$  and  $C_2$ , there are disjoint open sets  $O_1$  and  $O_2$  such that  $C_j \subseteq O_j$ ,  $j = 1, 2$ .

**Lemma B.28** (Urysohn). *Let  $X$  be a topological space. Then  $X$  is normal if and only if for every pair of disjoint closed sets  $C_1$  and  $C_2$ , there exists a continuous function  $f : X \rightarrow [0, 1]$  which is one on  $C_1$  and zero on  $C_2$ .*

*If in addition  $X$  is locally compact and  $C_1$  is compact, then  $f$  can be chosen to have compact support.*

**Proof.** To construct  $f$  we choose an open neighborhood  $O_0$  of  $C_1$  (e.g.  $O_0 := X \setminus C_2$ ). Now we could set  $f$  equal to one on  $C_1$ , equal to zero outside  $O_0$  and equal to  $\frac{1}{2}$  on the layer  $O_0 \setminus C_1$  in between. Clearly this function is not continuous, but we can successively improve the situation by introducing additional layers in between.

To this end observe that  $X$  is normal if and only if for every closed set  $C$  and every open neighborhood  $U$  of  $C$ , there exists an open set  $O_1$  and a closed set  $C_1$  such that  $C \subset O_1 \subset C_1 \subset U$  (just use the identification  $C_1 \leftrightarrow C$ ,  $C_2 \leftrightarrow X \setminus O$  and  $O_1 \leftrightarrow O_1$ ,  $O_2 \leftrightarrow X \setminus C$ ).

Using our observation we can find an open set  $O_{1/2}$  and a closed set  $C_{1/2}$  such that  $C_1 \subseteq O_{1/2} \subseteq C_{1/2} \subseteq O_0$ . Repeating this argument we get two more open and two more closed sets such that

$$C_1 \subseteq O_{3/4} \subseteq C_{3/4} \subseteq O_{1/2} \subseteq C_{1/2} \subseteq O_{1/4} \subseteq C_{1/4} \subseteq O_0.$$

Iterating this construction we get open sets  $O_q$  and closed sets  $C_q$  for every dyadic rational  $q = k/2^n \in [0, 1]$  such that  $O_q \subseteq C_q$  and  $C_p \subseteq O_q$  for  $p < q$ . Now set  $f_n^s(x) := \max\{q = k/2^n | x \in O_q, 0 \leq k < 2^n\}$  for  $x \in O_0$  and  $f_n^s(x) := 0$  else as well as  $f_n^i(x) := \min\{q = k/2^n | x \notin C_q, 0 \leq k < 2^n\}$  for  $x \notin C_1$  and  $f_n^i(x) := 1$  else. Then  $f_n^s(x) \nearrow f^s(x) := \sup\{q | x \in O_q\}$  and  $f_n^i(x) \searrow f^i(x) := \inf\{q | x \notin C_q\}$ . Moreover, if  $f_n^s(x) = q$  we have  $x \in O_q \setminus O_{q+2^{-n}}$  and depending on  $x \in C_{q+2^{-n}}$  or  $x \notin C_{q+2^{-n}}$  we have  $f_n^i(x) = q + 2^{-n+1}$  or  $f_n^i(x) = q + 2^{-n}$ , respectively. In particular,  $f^s(x) = f^i(x)$ . Finally, since  $(f^s)^{-1}((r, 1]) = \bigcup_{q > r} O_q$  and  $(f^i)^{-1}([0, r)) = \bigcap_{q < r} X \setminus C_q$  are open we see that  $f := f^s = f^i$  is continuous.

Conversely, given  $f$  choose  $O_1 := f^{-1}([0, 1/2))$  and  $O_2 := f^{-1}((1/2, 1])$ .

For the second claim, observe that there is an open set  $O_0$  such that  $\overline{O_0}$  is compact and  $C_1 \subset O_0 \subset \overline{O_0} \subset X \setminus C_2$ . In fact, for every  $x \in C_1$ , there is a ball  $B_\varepsilon(x)$  such that  $\overline{B_\varepsilon(x)}$  is compact and  $\overline{B_\varepsilon(x)} \subset X \setminus C_2$ . Since  $C_1$

is compact, finitely many of them cover  $C_1$  and we can choose the union of those balls to be  $O_0$ .  $\square$

**Example B.23.** In a metric space we can choose  $f(x) := \frac{\text{dist}(x, C_2)}{\text{dist}(x, C_1) + \text{dist}(x, C_2)}$  and hence every metric space is normal.  $\diamond$

an important class of normal spaces are locally compact Hausdorff spaces.

**Lemma B.29.** *Every locally compact Hausdorff space is normal.*

**Proof.** By replacing  $X$  by its one point compactification (Problem B.25) we can assume that  $X$  is compact without loss of generality.

We first show that we can separate a point  $x \in X$  and a closed subset  $C \subset X$  provided  $x \notin C$ . Indeed, since  $X$  is Hausdorff, for every  $y \in C$  there are disjoint open sets  $U_y$  containing  $x$  and  $V_y$  containing  $y$ . Moreover, since  $X$  is compact so is  $C$  (Lemma B.16 (ii)) and hence finitely many sets  $V_{y_1}, \dots, V_{y_n}$  cover  $C$ . Hence  $U := \bigcap U_{y_j}$  is an open set containing  $x$  which is disjoint from the open set  $V := \bigcup V_{y_j}$  containing  $C$ .

Now given two disjoint closed sets  $C_1, C_2 \subset X$ , the first part implies that for every  $y \in C_1$  there are disjoint open sets  $U_y$  containing  $C_1$  and  $V_y$  containing  $y$ . Hence we can proceed as before.  $\square$

Another important result is the **Tietze extension theorem**:

**Theorem B.30** (Tietze). *Suppose  $C$  is a closed subset of a normal topological space  $X$ . For every continuous function  $f : C \rightarrow [-1, 1]$  there is a continuous extension  $\tilde{f} : X \rightarrow [-1, 1]$ . If in addition  $X$  is locally compact and  $C$  is compact, then  $f$  can be chosen to have compact support.*

**Proof.** The idea is to construct a rough approximation using Urysohn's lemma and then iteratively improve this approximation. To this end we set  $C_1 := f^{-1}([\frac{1}{3}, 1])$  and  $C_2 := f^{-1}([-1, -\frac{1}{3}])$  and let  $g$  be the function from Urysohn's lemma. Then  $f_1 := \frac{2g-1}{3}$  satisfies  $|f(x) - f_1(x)| \leq \frac{2}{3}$  for  $x \in C$  as well as  $|f_1(x)| \leq \frac{1}{3}$  for all  $x \in X$ . Applying this same procedure to  $f - f_1$  we obtain a function  $f_2$  such that  $|f(x) - f_1(x) - f_2(x)| \leq (\frac{2}{3})^2$  for  $x \in C$  and  $|f_2(x)| \leq \frac{1}{3} (\frac{2}{3})$ . Continuing this process we arrive at a sequence of functions  $f_n$  such that  $|f(x) - \sum_{j=1}^n f_j(x)| \leq (\frac{2}{3})^n$  for  $x \in C$  and  $|f_n(x)| \leq \frac{1}{3} (\frac{2}{3})^{n-1}$ . By construction the corresponding series converges uniformly to the desired extension  $\tilde{f} := \sum_{j=1}^{\infty} f_j$ .

To see the last claim multiply  $f$  with a continuous function which equals 1 on  $C$  and has compact support (which exists by Urysohn's lemma).  $\square$

Note that by extending each component we can also handle functions with values in  $\mathbb{R}^n$ .

There is also a corresponding result for Lipschitz continuous functions.

**Lemma B.31** (McShane). *Let  $X$  be a metric space and  $Y \subseteq X$  nonempty. Then a Lipschitz continuous functions  $f : Y \rightarrow \mathbb{R}$*

$$|f(x) - f(y)| \leq L d(x, y), \quad x, y \in Y, \quad (\text{B.24})$$

*has a continuous extensions*

$$\bar{f}(x) = \sup_{y \in Y} (f(y) - L d(x, y)) \quad (\text{B.25})$$

*such that  $|\bar{f}(x) - \bar{f}(y)| \leq L d(x, y)$  for all  $x, y \in X$ .*

**Proof.** Since  $f(y) - L d(x, y) \leq \bar{f}(x)$  with equality for  $x = y$  in case  $x \in Y$  we see that  $\bar{f}(x) = f(x)$  in this case. To see the Lipschitz condition choose  $x, x' \in X$  and suppose  $\bar{f}(x') \leq \bar{f}(x)$  without loss of generality. Then, if  $\bar{f}(x, ) < \infty$  we have

$$\begin{aligned} \bar{f}(x) - \bar{f}(x') &= \sup_{y \in Y} (f(y) - L d(x, y)) - \sup_{y \in Y} (f(y) - L d(x', y)) \\ &\leq \sup_{y \in Y} ((f(y) - L d(x, y)) - (f(y) - L d(x', y))) \\ &= L \sup_{y \in Y} (d(x', y) - d(x, y)) \leq L d(x', x) \end{aligned}$$

In particular, for  $x' \in Y$  we conclude that  $\bar{f}(x)$  is finite for all  $x \in X$ .  $\square$

Note that the proof easily extends to (e.g.) Hölder continuous functions by replacing  $d$  by  $d^\alpha$  in the definition of  $\bar{f}$ .

A **partition of unity** is a collection of functions  $h_\alpha : X \rightarrow [0, 1]$  such that  $\sum_\alpha h_\alpha(x) = 1$ . We will only consider the case when the partition of unity is **locally finite**, that is, when every  $x$  has a neighborhood where all but a finite number of the functions  $h_\alpha$  vanish. Moreover, given a cover  $\{O_\beta\}$  of  $X$  it is called **subordinate** to this cover if every  $h_\alpha$  has support contained in some set  $O_\beta$  from this cover.

In the case of subsets of  $\mathbb{R}^n$  we are interested in the existence of smooth partitions of unity. To this end recall that for every point  $x \in \mathbb{R}^n$  there is a smooth bump function with values in  $[0, 1]$  which is positive at  $x$  and supported in a given neighborhood of  $x$ .

**Example B.24.** The standard bump function is  $\phi(x) := \exp(\frac{1}{|x|^2 - 1})$  for  $|x| < 1$  and  $\phi(x) = 0$  otherwise. To show that this function is indeed smooth it suffices to show that all left derivatives of  $f(r) = \exp(\frac{1}{r^2 - 1})$  at  $r = 1$  vanish, which can be done using l'Hôpital's rule. By scaling and translation  $\phi(\frac{x - x_0}{r})$  we get a bump function which is supported in  $B_r(x_0)$  and satisfies  $\phi(\frac{x - x_0}{r})|_{x=x_0} = \phi(0) = e^{-1}$ .  $\diamond$

**Lemma B.32.** *Let  $X \subseteq \mathbb{R}^n$  be open and  $\{O_j\}$  a countable open cover. Then there is a locally finite partition of unity of functions from  $C_c^\infty(X)$  subordinate to this cover.*

**Proof.** Let  $U_j$  be as in Lemma B.25 (iv). For the compact set  $\overline{U}_j$  choose finitely many bump functions  $\tilde{h}_{j,k}$  such that  $\tilde{h}_{j,1}(x) + \cdots + \tilde{h}_{j,k_j}(x) > 0$  for every  $x \in \overline{U}_j \setminus U_{j-1}$  and such that  $\text{supp}(\tilde{h}_{j,k})$  is contained in one of the  $O_k$  and in  $U_{j+1} \setminus U_{j-1}$ . Then  $\{\tilde{h}_{j,k}\}_{j,k}$  is locally finite and hence  $h := \sum_{j,k} \tilde{h}_{j,k}$  is a smooth function which is everywhere positive. Finally,  $\{\tilde{h}_{j,k}/h\}_{j,k}$  is a partition of unity of the required type.  $\square$

**Problem B.28.** *Show  $\text{dist}(x, Y) = \text{dist}(x, \overline{Y})$ . Moreover, show  $x \in \overline{Y}$  if and only if  $\text{dist}(x, Y) = 0$ .*

**Problem B.29.** *Let  $Y, Z \subseteq X$  and define*

$$\text{dist}(Y, Z) := \inf_{y \in Y, z \in Z} d(y, z).$$

*Show  $\text{dist}(Y, Z) = \text{dist}(\overline{Y}, \overline{Z})$ . Moreover, show that if  $K$  is compact, then  $\text{dist}(K, Y) > 0$  if and only if  $K \cap \overline{Y} = \emptyset$ .*

**Problem B.30.** *Let  $X$  be some normed vector space and  $Y \subseteq Z \subseteq X$  with  $Z$  open. Show  $\text{dist}(Y, \partial Z) = \text{dist}(Y, X \setminus Z)$ .*

**Problem B.31.** *Let  $K \subseteq U$  with  $K$  compact and  $U$  open. Show that there is some  $\varepsilon > 0$  such that  $K_\varepsilon := \{x \in X \mid \text{dist}(x, K) < \varepsilon\} \subseteq U$ .*

**Problem B.32.** *Let  $(X, d)$  be a locally compact metric space. Then  $X$  is  $\sigma$ -compact if and only if there exists a proper function  $f : X \rightarrow [0, \infty)$ . (Hint: Let  $U_n$  be as in item (iv) of Lemma B.25 and use Uryson's lemma to find functions  $f_n : X \rightarrow [0, 1]$  such that  $f_n(x) = 0$  for  $x \in \overline{U}_n$  and  $f_n(x) = 1$  for  $x \in X \setminus U_{n+1}$ . Now consider  $f = \sum_{n=1}^\infty f_n$ .)*

## B.7. Connectedness

Roughly speaking a topological space  $X$  is disconnected if it can be split into two (nonempty) separated sets. This of course raises the question what should be meant by separated. Evidently it should be more than just disjoint since otherwise we could split any space containing more than one point. Hence we will consider two sets separated if each is disjoint from the closure of the other. Note that if we can split  $X$  into two separated sets  $X = U \cup V$  then  $\overline{U} \cap V = \emptyset$  implies  $\overline{U} = U$  (and similarly  $\overline{V} = V$ ). Hence both sets must be closed and thus also open (being complements of each other). This brings us to the following definition:

A topological space  $X$  is called **disconnected** if one of the following equivalent conditions holds

- $X$  is the union of two nonempty separated sets.
- $X$  is the union of two nonempty disjoint open sets.
- $X$  is the union of two nonempty disjoint closed sets.

In this case the sets from the splitting are both open and closed. A topological space  $X$  is called **connected** if it cannot be split as above. That is, in a connected space  $X$  the only sets which are both open and closed are  $\emptyset$  and  $X$ . This last observation is frequently used in proofs: If the set where a property holds is both open and closed it must either hold nowhere or everywhere. In particular, any continuous mapping from a connected to a discrete space must be constant since the inverse image of a point is both open and closed.

A subset of  $X$  is called (dis-)connected if it is (dis-)connected with respect to the relative topology. In other words, a subset  $A \subseteq X$  is disconnected if there are disjoint nonempty open sets  $U$  and  $V$  which split  $A$  according to  $A = (U \cap A) \cup (V \cap A)$ .

**Example B.25.** In  $\mathbb{R}$  the nonempty connected sets are precisely the intervals (Problem B.33). Consequently  $A = [0, 1] \cup [2, 3]$  is disconnected with  $[0, 1]$  and  $[2, 3]$  being its components (to be defined precisely below). While you might be reluctant to consider the closed interval  $[0, 1]$  as open, it is important to observe that it is the relative topology which is relevant here.  $\diamond$

The maximal connected subsets (ordered by inclusion) of a nonempty topological space  $X$  are called the **connected components** of  $X$ .

**Example B.26.** Consider  $\mathbb{Q} \subseteq \mathbb{R}$ . Then every rational point is its own component (if a set of rational points contains more than one point there would be an irrational point in between which can be used to split the set).  $\diamond$

In many applications one also needs the following stronger concept. A space  $X$  is called **path-connected** if any two points  $x, y \in X$  can be joined by a **path**, that is a continuous map  $\gamma : [0, 1] \rightarrow X$  with  $\gamma(0) = x$  and  $\gamma(1) = y$ . A space is called **locally (path-)connected** if for every given point and every open set containing that point there is a smaller open set which is (path-)connected.

**Example B.27.** Every normed vector space is (locally) path-connected since every ball is path-connected (consider straight lines). In fact this also holds for locally convex spaces. Every open subset of a locally (path-)connected space is locally (path-)connected.  $\diamond$

Every path-connected space is connected. In fact, if  $X = U \cup V$  were disconnected but path-connected we could choose  $x \in U$  and  $y \in V$  plus a path  $\gamma$  joining them. But this would give a splitting  $[0, 1] = \gamma^{-1}(U) \cup \gamma^{-1}(V)$



contradicting our assumption. The converse however is not true in general as a space might be impassable (an example will follow).

**Example B.28.** The spaces  $\mathbb{R}$  and  $\mathbb{R}^n$ ,  $n > 1$ , are not homeomorphic. In fact, removing any point from  $\mathbb{R}$  gives a disconnected space while removing a point from  $\mathbb{R}^n$  still leaves it (path-)connected.  $\diamond$

We collect a few simple but useful properties below.

**Lemma B.33.** *Suppose  $X$  and  $Y$  are topological spaces.*

- (i) *Suppose  $f : X \rightarrow Y$  is continuous. Then if  $X$  is (path-)connected so is the image  $f(X)$ .*
- (ii) *Suppose  $A_\alpha \subseteq X$  are (path-)connected and  $\bigcap_\alpha A_\alpha \neq \emptyset$ . Then  $\bigcup_\alpha A_\alpha$  is (path-)connected*
- (iii)  *$A \subseteq X$  is (path-)connected if and only if any two points  $x, y \in A$  are contained in a (path-)connected set  $B \subseteq A$*
- (iv) *Suppose  $X_1, \dots, X_n$  are (path-)connected then so is  $\times_{j=1}^n X_j$ .*
- (v) *Suppose  $A \subseteq X$  is connected, then  $\bar{A}$  is connected.*
- (vi) *A locally path-connected space is path-connected if and only if it is connected.*

**Proof.** (i). Suppose we have a splitting  $f(X) = U \cup V$  into nonempty disjoint sets which are open in the relative topology. Hence, there are open sets  $U'$  and  $V'$  such that  $U = U' \cap f(X)$  and  $V = V' \cap f(X)$  implying that the sets  $f^{-1}(U) = f^{-1}(U')$  and  $f^{-1}(V) = f^{-1}(V')$  are open. Thus we get a corresponding splitting  $X = f^{-1}(U) \cup f^{-1}(V)$  into nonempty disjoint open sets contradicting connectedness of  $X$ .

If  $X$  is path connected, let  $y_1 = f(x_1)$  and  $y_2 = f(x_2)$  be given. If  $\gamma$  is a path connecting  $x_1$  and  $x_2$ , then  $f \circ \gamma$  is a path connecting  $y_1$  and  $y_2$ .

(ii). Let  $A = \bigcup_\alpha A_\alpha$  and suppose there is a splitting  $A = (U \cap A) \cup (V \cap A)$ . Since there is some  $x \in \bigcap_\alpha A_\alpha$  we can assume  $x \in U$  w.l.o.g. Hence there is a splitting  $A_\alpha = (U \cap A_\alpha) \cup (V \cap A_\alpha)$  and since  $A_\alpha$  is connected and  $U \cap A_\alpha$  is nonempty we must have  $V \cap A_\alpha = \emptyset$ . Hence  $V \cap A = \emptyset$  and  $A$  is connected.

If the  $x \in A_\alpha$  and  $y \in A_\beta$  then choose a point  $z \in A_\alpha \cap A_\beta$  and paths  $\gamma_\alpha$  from  $x$  to  $z$  and  $\gamma_\beta$  from  $z$  to  $y$ , then  $\gamma_\alpha \odot \gamma_\beta$  is a path from  $x$  to  $y$ , where  $\gamma_\alpha \odot \gamma_\beta(t) = \gamma_\alpha(2t)$  for  $0 \leq t \leq \frac{1}{2}$  and  $\gamma_\alpha \odot \gamma_\beta(t) = \gamma_\beta(2t - 1)$  for  $\frac{1}{2} \leq t \leq 1$  (cf. Problem B.23).

(iii). If  $X$  is connected we can choose  $B = A$ . Conversely, fix some  $x \in A$  and let  $B_y$  be the corresponding set for the pair  $x, y$ . Then  $A = \bigcup_{y \in A} B_y$  is (path-)connected by the previous item.

(iv). We first consider two spaces  $X = X_1 \times X_2$ . Let  $x, y \in X$ . Then  $\{x_1\} \times X_2$  is homeomorphic to  $X_2$  and hence (path-)connected. Similarly

$X_1 \times \{y_2\}$  is (path-)connected as well as  $\{x_1\} \times X_2 \cup X_1 \times \{y_2\}$  by (ii) since both sets contain  $(x_1, y_2) \in X$ . But this last set contains both  $x, y$  and hence the claim follows from (iii). The general case follows by iterating this result.

(v). Let  $x \in \bar{A}$ . Then  $\{x\}$  and  $A$  cannot be separated and hence  $\{x\} \cup A$  is connected. The rest follows from (ii).

(vi). Consider the set  $U(x)$  of all points connected to a fixed point  $x$  (via paths). If  $y \in U(x)$  then so is any path-connected neighborhood of  $y$  by gluing paths (as in item (ii)). Hence  $U(x)$  is open. Similarly, if  $y \in \overline{U(x)}$  then any path-connected neighborhood of  $y$  will intersect  $U(y)$  and hence  $y \in U(x)$ . Thus  $U(x)$  is also closed and hence must be all of  $X$  by connectedness. The converse is trivial.  $\square$

A few simple consequences are also worth while noting: If two different components contain a common point, their union is again connected contradicting maximality. Hence two different components are always disjoint. Moreover, every point is contained in a component, namely the union of all connected sets containing this point. In other words, the components of any topological space  $X$  form a partition of  $X$  (i.e., they are disjoint, nonempty, and their union is  $X$ ). Moreover, every component is a closed subset of the original space  $X$ . In the case where their number is finite we can take complements and each component is also an open subset (the rational numbers from our first example show that components are not open in general). In a locally (path-)connected space, components are open and (path-)connected by (vi) of the last lemma. Note also that in a second countable space an open set can have at most countably many components (take those sets from a countable base which are contained in some component, then we have a surjective map from these sets to the components).

**Example B.29.** Consider the graph of the function  $f : (0, 1] \rightarrow \mathbb{R}$ ,  $x \mapsto \sin(\frac{1}{x})$ . Then  $\Gamma(f) \subseteq \mathbb{R}^2$  is path-connected and its closure  $\overline{\Gamma(f)} = \Gamma(f) \cup \{0\} \times [-1, 1]$  is connected. However,  $\overline{\Gamma(f)}$  is not path-connected as there is no path from  $(1, 0)$  to  $(0, 0)$ . Indeed, suppose  $\gamma$  were such a path. Then, since  $\gamma_1$  covers  $[0, 1]$  by the intermediate value theorem (see below), there is a sequence  $t_n \rightarrow 1$  such that  $\gamma_1(t_n) = \frac{2}{(2n+1)\pi}$ . But then  $\gamma_2(t_n) = (-1)^n \not\rightarrow 0$  contradicting continuity.  $\diamond$

**Theorem B.34** (Intermediate Value Theorem). *Let  $X$  be a connected topological space and  $f : X \rightarrow \mathbb{R}$  be continuous. For any  $x, y \in X$  the function  $f$  attains every value between  $f(x)$  and  $f(y)$ .*

**Proof.** The image  $f(X)$  is connected and hence an interval.  $\square$

**Problem B.33.** *A nonempty subset of  $\mathbb{R}$  is connected if and only if it is an interval.*

**Problem B.34.** Let  $U = \bigcup_j U_j \subseteq \mathbb{R}^n$  be an open set with  $U_j$  its connected components. Show  $\partial U_j \subseteq \partial U$ . Show that in the case of finitely many components we have  $\bigcup_{j=1}^n \partial U_j = \partial U$ . Show that this fails for infinitely many components in general.

## B.8. Continuous functions on metric spaces

Let  $X, Y$  be topological spaces and let  $C(X, Y)$  be the set of all continuous functions  $f : X \rightarrow Y$ . Set  $C(X) := C(X, \mathbb{C})$ . Moreover, if  $Y$  is a metric space then  $C_b(X, Y)$  will denote the set of all bounded continuous functions, that is, those continuous functions for which  $\sup_{x \in X} d_Y(f(x), y)$  is finite for some (and hence for all)  $y \in Y$ . Note that by the extreme value theorem  $C_b(X, Y) = C(X, Y)$  if  $X$  is compact. For these functions we can introduce a metric via

$$d(f, g) := \sup_{x \in X} d_Y(f(x), g(x)). \quad (\text{B.26})$$

In fact, the requirements for a metric are readily checked. Of course convergence with respect to this metric implies pointwise convergence but not the other way round.

**Example B.30.** Consider  $X := [0, 1]$ , then  $f_n(x) := \max(1 - |nx - 1|, 0)$  converges pointwise to 0 (in fact,  $f_n(0) = 0$  and  $f_n(x) = 0$  on  $[\frac{2}{n}, 1]$ ) but not with respect to the above metric since  $f_n(\frac{1}{n}) = 1$ .  $\diamond$

This kind of convergence is known as **uniform convergence** since for every positive  $\varepsilon$  there is some index  $N$  (independent of  $x$ ) such that  $d_Y(f_n(x), f(x)) < \varepsilon$  for  $n \geq N$ . In contradistinction, in the case of pointwise convergence,  $N$  is allowed to depend on  $x$ . One advantage is that continuity of the limit function comes for free.

**Theorem B.35.** Let  $X$  be a topological space and  $Y$  a metric space. Suppose  $f_n \in C(X, Y)$  converges uniformly to some function  $f : X \rightarrow Y$ . Then  $f$  is continuous.

**Proof.** Let  $x \in X$  be given and write  $y := f(x)$ . We need to show that  $f^{-1}(B_\varepsilon(y))$  is a neighborhood of  $x$  for every  $\varepsilon > 0$ . So fix  $\varepsilon$ . Then we can find an  $N$  such that  $d(f_n, f) < \frac{\varepsilon}{2}$  for  $n \geq N$  implying  $f_N^{-1}(B_{\varepsilon/2}(y)) \subseteq f^{-1}(B_\varepsilon(y))$  since  $d(f_n(z), y) < \frac{\varepsilon}{2}$  implies  $d(f(z), y) \leq d(f(z), f_n(z)) + d(f_n(z), y) \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon$  for  $n \geq N$ .  $\square$

**Corollary B.36.** Let  $X$  be a topological space and  $Y$  a complete metric space. The space  $C_b(X, Y)$  together with the metric  $d$  is complete.

**Proof.** Suppose  $f_n$  is a Cauchy sequence with respect to  $d$ , then  $f_n(x)$  is a Cauchy sequence for fixed  $x$  and has a limit since  $Y$  is complete. Call this limit  $f(x)$ . Then  $d_Y(f(x), f_n(x)) = \lim_{m \rightarrow \infty} d_Y(f_m(x), f_n(x)) \leq$

$\sup_{m \geq n} d(f_m, f_n)$  and since this last expression goes to 0 as  $n \rightarrow \infty$ , we see that  $f_n$  converges uniformly to  $f$ . Moreover,  $f \in C(X, Y)$  by the previous theorem so we are done.  $\square$

Let  $Y$  be a vector space. By  $C_c(X, Y) \subseteq C_b(X, Y)$  we will denote the set of continuous functions with compact support. Its closure will be denoted by  $C_0(X, Y) := \overline{C_c(X, Y)} \subseteq C_b(X, Y)$ . Of course if  $X$  is compact all these spaces agree  $C_c(X, Y) = C_0(X, Y) = C_b(X, Y) = C(X, Y)$ . In the general case one at least assumes  $X$  to be locally compact since if we take a closed neighborhood  $V$  of  $f(x) \neq 0$  which does not contain 0, then  $f^{-1}(V)$  will be a compact neighborhood of  $x$ . Hence without this assumption  $f$  must vanish on every point which does not have a compact neighborhood and  $C_c(X, Y)$  will not be sufficiently rich.

**Example B.31.** Let  $X$  be a separable and locally compact metric space and  $Y = \mathbb{C}^n$ . Then

$$C_0(X, \mathbb{C}^n) = \{f \in C_b(X, \mathbb{C}^n) \mid \forall \varepsilon > 0, \exists K \subseteq X \text{ compact} : \\ |f(x)| < \varepsilon, x \in X \setminus K\}. \quad (\text{B.27})$$

To see this denote the set on the right-hand side by  $C$ . Let  $K_m$  be an increasing sequence of compact sets with  $K_m \nearrow X$  (Lemma B.25) and let  $\varphi_m$  be a corresponding sequence as in Urysohn's lemma (Lemma B.28). Then for  $f \in C$  the sequence  $f_m = \varphi_m f \in C_c(X, \mathbb{C}^n)$  will converge to  $f$ . Conversely, if  $f_n \in C_c(X, \mathbb{C}^n)$  converges to  $f \in C_b(X, \mathbb{C}^n)$ , then given  $\varepsilon > 0$  choose  $K = \text{supp}(f_m)$  for some  $m$  with  $d(f_m, f) < \varepsilon$ .

In the case where  $X$  is an open subset of  $\mathbb{R}^n$  this says that  $C_0(X, Y)$  are those which vanish at the boundary (including the case as  $|x| \rightarrow \infty$  if  $X$  is unbounded).  $\diamond$

**Lemma B.37.** *If  $X$  is a separable and locally compact space then  $C_0(X, \mathbb{C}^n)$  is separable.*

**Proof.** Choose a countable base  $\mathcal{B}$  for  $X$  and let  $\mathcal{I}$  the collection of all balls in  $\mathbb{C}^n$  with rational radius and center. Given  $O_1, \dots, O_m \in \mathcal{B}$  and  $I_1, \dots, I_m \in \mathcal{I}$  we say that  $f \in C_c(X, \mathbb{C}^n)$  is adapted to these sets if  $\text{supp}(f) \subseteq \bigcup_{j=1}^m O_j$  and  $f(O_j) \subseteq I_j$ . The set of all tuples  $(O_j, I_j)_{1 \leq j \leq m}$  is countable and for each tuple we choose a corresponding adapted function (if there exists one at all). Then the set of these functions  $\mathcal{F}$  is dense. It suffices to show that the closure of  $\mathcal{F}$  contains  $C_c(X, \mathbb{C}^n)$ . So let  $f \in C_c(X, \mathbb{C}^n)$  and let  $\varepsilon > 0$  be given. Then for every  $x \in X$  there is some neighborhood  $O(x) \in \mathcal{B}$  such that  $|f(x) - f(y)| < \varepsilon$  for  $y \in O(x)$ . Since  $\text{supp}(f)$  is compact, it can be covered by  $O(x_1), \dots, O(x_m)$ . In particular  $f(O(x_j)) \subseteq B_\varepsilon(f(x_j))$  and we can find a ball  $I_j$  of radius at most  $2\varepsilon$  with  $f(O(x_j)) \subseteq I_j$ . Now let  $g$  be

the function from  $\mathcal{F}$  which is adapted to  $(O(x_j), I_j)_{1 \leq j \leq m}$  and observe that  $|f(x) - g(x)| < 4\varepsilon$  since  $x \in O(x_j)$  implies  $f(x), g(x) \in I_j$ .  $\square$

Let  $X, Y$  be metric spaces. A function  $f \in C(X, Y)$  is called **uniformly continuous** if for every  $\varepsilon > 0$  there is a  $\delta > 0$  such that

$$d_Y(f(x), f(y)) \leq \varepsilon \quad \text{whenever} \quad d_X(x, y) < \delta. \quad (\text{B.28})$$

Note that with the usual definition of continuity one fixes  $x$  and then chooses  $\delta$  depending on  $x$ . Here  $\delta$  has to be independent of  $x$ . Note that the uniform limit of uniformly continuous functions is again uniformly continuous (Problem B.35) and hence  $C_{buc}(X, Y) \subseteq C_b(X, Y)$  is a closed subspace. If the domain is compact, this extra condition comes for free.

**Theorem B.38.** *Let  $X$  be a compact metric space and  $Y$  a metric space. Then every  $f \in C(X, Y)$  is uniformly continuous.*

**Proof.** Suppose the claim were wrong. Fix  $\varepsilon > 0$ . Then for every  $\delta_n = \frac{1}{n}$  we can find  $x_n, y_n$  with  $d_X(x_n, y_n) < \delta_n$  but  $d_Y(f(x_n), f(y_n)) \geq \varepsilon$ . Since  $X$  is compact we can assume that  $x_n$  converges to some  $x \in X$  (after passing to a subsequence if necessary). Then we also have  $y_n \rightarrow x$  implying  $d_Y(f(x_n), f(y_n)) \rightarrow 0$ , a contradiction.  $\square$

**Example B.32.** If  $X$  is not compact, there are bounded continuous functions which are not uniformly continuous. For example,  $f(x) := \sin(x^2)$  is in  $C_b(\mathbb{R})$  but not in  $C_{buc}(\mathbb{R})$ .  $\diamond$

Note that a uniformly continuous function maps Cauchy sequences to Cauchy sequences. This fact can be used to extend a uniformly continuous function to boundary points.

**Theorem B.39.** *Let  $X$  be a metric space and  $Y$  a complete metric space. A uniformly continuous function  $f : A \subseteq X \rightarrow Y$  has a unique continuous extension  $\bar{f} : \bar{A} \rightarrow Y$ . This extension is again uniformly continuous.*

**Proof.** If there is an extension it must be  $\bar{f}(x) := \lim_{n \rightarrow \infty} f(x_n)$ , where  $x_n \in A$  is some sequence converging to  $x \in \bar{A}$ . Indeed, since  $x_n$  converges,  $f(x_n)$  is Cauchy and hence has a limit since  $Y$  is assumed complete. Moreover, uniqueness of limits shows that  $\bar{f}(x)$  is independent of the sequence chosen. Also  $\bar{f}(x) = f(x)$  for  $x \in A$  by continuity. To see that  $\bar{f}$  is uniformly continuous, let  $\varepsilon > 0$  be given and choose a  $\delta$  which works for  $f$ . Then for given  $x, y$  with  $d_X(x, y) < \frac{\delta}{3}$  we can find  $\tilde{x}, \tilde{y} \in A$  with  $d_X(\tilde{x}, x) < \frac{\delta}{3}$  and  $d_Y(f(\tilde{x}), \bar{f}(x)) \leq \varepsilon$  as well as  $d_X(\tilde{y}, y) < \frac{\delta}{3}$  and  $d_Y(f(\tilde{y}), \bar{f}(y)) \leq \varepsilon$ . Hence  $d_Y(\bar{f}(x), \bar{f}(y)) \leq d_Y(\bar{f}(x), f(\tilde{x})) + d_Y(f(\tilde{x}), f(\tilde{y})) + d_Y(f(\tilde{y}), \bar{f}(y)) \leq 3\varepsilon$ .  $\square$

Next we want to identify relatively compact subsets in  $C(X, Y)$ . A family of functions  $F \subset C(X, Y)$  is called (pointwise) **equicontinuous** if for every  $\varepsilon > 0$  and every  $x \in X$  there is a neighborhood  $U(x)$  of  $x$  such that

$$d_Y(f(x), f(y)) \leq \varepsilon \quad \text{whenever } y \in U(x), \quad \forall f \in F. \quad (\text{B.29})$$

**Theorem B.40** (Arzelà–Ascoli). *Let  $X$  be a compact space and  $Y$  a proper metric space. Let  $F \subset C(X, Y)$  be a family of continuous functions. Then every sequence from  $F$  has a uniformly convergent subsequence if and only if  $F$  is equicontinuous and the set  $\{f(x) | f \in F\}$  is bounded for every  $x \in X$ . In this case  $F$  is even bounded.*

**Proof.** Suppose  $F$  is equicontinuous and pointwise bounded. Fix  $\varepsilon > 0$ . By compactness of  $X$  there are finitely many points  $x_1, \dots, x_n \in X$  such that the neighborhoods  $U(x_j)$  (from the definition of equicontinuity) cover  $X$ . Now first of all note that,  $F$  is bounded since  $d_Y(f(x), y) \leq \max_j \sup_{f \in F} d_Y(f(x_j), y) + \varepsilon$  for every  $x \in X$  and every  $f \in F$ .

Next consider  $P : C(X, Y) \rightarrow Y^n$ ,  $P(f) = (f(x_1), \dots, f(x_n))$ . Then  $P(F)$  is bounded and  $d(f, g) \leq 3\varepsilon$  whenever  $d_Y(P(f), P(g)) < \varepsilon$ . Indeed, just note that for every  $x$  there is some  $j$  such that  $x \in U(x_j)$  and thus  $d_Y(f(x), g(x)) \leq d_Y(f(x), f(x_j)) + d_Y(f(x_j), g(x_j)) + d_Y(g(x_j), g(x)) \leq 3\varepsilon$ . Hence  $F$  is relatively compact by Lemma 1.11.

Conversely, suppose  $F$  is relatively compact. Then  $F$  is totally bounded and hence bounded. To see equicontinuity fix  $x \in X$ ,  $\varepsilon > 0$  and choose a corresponding  $\varepsilon$ -cover  $\{B_\varepsilon(f_j)\}_{j=1}^n$  for  $F$ . Pick a neighborhood  $U(x)$  such that  $y \in U(x)$  implies  $d_Y(f_j(y), f_j(x)) < \varepsilon$  for all  $1 \leq j \leq n$ . Then  $f \in B_\varepsilon(f_j)$  for some  $j$  and hence  $d_Y(f(y), f(x)) \leq d_Y(f(y), f_j(y)) + d_Y(f_j(y), f_j(x)) + d_Y(f_j(x), f(x)) \leq 3\varepsilon$ , proving equicontinuity.  $\square$

In many situations a certain property can be seen for a class of *nice* functions and then extended to a more general class of functions by approximation. In this respect it is important to identify classes of functions which allow to approximate *all* functions. That is, in our present situation we are looking for functions which are dense in  $C(X, Y)$ . For example, the classical Weierstraß approximation theorem (Theorem 1.3) says that the polynomials are dense in  $C([a, b])$  for any compact interval. Here we will present a generalization of this result. For its formulation observe that  $C(X)$  is not only a vector space but also comes with a natural product, given by pointwise multiplication of functions, which turns it into an algebra over  $\mathbb{C}$ . By a subalgebra we will mean a subspace which is closed under multiplication and by a  $*$ -subalgebra we will mean a subalgebra which is also closed under complex conjugation. The  $(*)$ -subalgebra generated by a set is of course the smallest  $(*)$ -subalgebra containing this set.

The proof will use the fact that the absolute value can be approximated by polynomials on  $[-1, 1]$ . This of course follows from the Weierstraß approximation theorem but can also be seen directly by defining the sequence of polynomials  $p_n$  via

$$p_1(t) := 0, \quad p_{n+1}(t) := p_n(t) + \frac{t^2 - p_n(t)^2}{2}. \quad (\text{B.30})$$

Then this sequence of polynomials satisfies  $p_n(t) \leq p_{n+1}(t) \leq |t|$  and converges pointwise to  $|t|$  for  $t \in [-1, 1]$ . Hence by Dini's theorem (Problem B.38) it converges uniformly. By scaling we get the corresponding result for arbitrary compact subsets of the real line.

**Theorem B.41** (Stone–Weierstraß, real version). *Suppose  $K$  is a compact topological space and consider  $C(K, \mathbb{R})$ . If  $F \subset C(K, \mathbb{R})$  contains the identity 1 and separates points (i.e., for every  $x_1 \neq x_2$  there is some function  $f \in F$  such that  $f(x_1) \neq f(x_2)$ ), then the subalgebra generated by  $F$  is dense.*

**Proof.** Denote by  $A$  the subalgebra generated by  $F$ . Note that if  $f \in \overline{A}$ , we have  $|f| \in \overline{A}$ : Choose a polynomial  $p_n(t)$  such that  $||t| - p_n(t)| < \frac{1}{n}$  for  $t \in f(K)$  and hence  $p_n(f) \rightarrow |f|$ .

In particular, if  $f, g \in \overline{A}$ , we also have

$$\max\{f, g\} = \frac{(f+g) + |f-g|}{2} \in \overline{A}, \quad \min\{f, g\} = \frac{(f+g) - |f-g|}{2} \in \overline{A}.$$

Now fix  $f \in C(K, \mathbb{R})$ . We need to find some  $f^\varepsilon \in \overline{A}$  with  $\|f - f^\varepsilon\|_\infty < \varepsilon$ .

First of all, since  $A$  separates points, observe that for given  $y, z \in K$  there is a function  $f_{y,z} \in A$  such that  $f_{y,z}(y) = f(y)$  and  $f_{y,z}(z) = f(z)$  (show this). Next, for every  $y \in K$  there is a neighborhood  $U(y)$  such that

$$f_{y,z}(x) > f(x) - \varepsilon, \quad x \in U(y),$$

and since  $K$  is compact, finitely many, say  $U(y_1), \dots, U(y_j)$ , cover  $K$ . Then

$$f_z = \max\{f_{y_1,z}, \dots, f_{y_j,z}\} \in \overline{A}$$

and satisfies  $f_z > f - \varepsilon$  by construction. Since  $f_z(z) = f(z)$  for every  $z \in K$ , there is a neighborhood  $V(z)$  such that

$$f_z(x) < f(x) + \varepsilon, \quad x \in V(z),$$

and a corresponding finite cover  $V(z_1), \dots, V(z_k)$ . Now

$$f^\varepsilon = \min\{f_{z_1}, \dots, f_{z_k}\} \in \overline{A}$$

satisfies  $f^\varepsilon < f + \varepsilon$ . Since  $f - \varepsilon < f_{z_l}$  for all  $1 \leq l \leq k$  we have  $f - \varepsilon < f^\varepsilon$  and we have found a required function.  $\square$

**Example B.33.** The set  $\{f \in C(K, \mathbb{R}) \mid f(x_0) = 0\}$  for some  $x_0 \in K$  is a closed algebra which, in particular, is not dense. The same is true for the set  $\{f \in C(K, \mathbb{R}) \mid f(x_1) = f(x_2)\}$  for some  $x_1, x_2 \in K$ . These examples show that the above two conditions are also necessary.  $\diamond$

**Theorem B.42** (Stone–Weierstraß). *Suppose  $K$  is a compact topological space and consider  $C(K)$ . If  $F \subset C(K)$  contains the identity 1 and separates points, then the  $*$ -subalgebra generated by  $F$  is dense.*

**Proof.** Just observe that  $\tilde{F} = \{\operatorname{Re}(f), \operatorname{Im}(f) \mid f \in F\}$  satisfies the assumption of the real version. Hence every real-valued continuous function can be approximated by elements from the subalgebra generated by  $\tilde{F}$ ; in particular, this holds for the real and imaginary parts for every given complex-valued function. Finally, note that the subalgebra spanned by  $\tilde{F}$  is contained in the  $*$ -subalgebra spanned by  $F$ .  $\square$

Note that the additional requirement of being closed under complex conjugation is crucial: The functions holomorphic on the unit disc and continuous on the boundary separate points, but they are not dense (since the uniform limit of holomorphic functions is again holomorphic).

**Corollary B.43.** *Suppose  $K$  is a compact topological space and consider  $C(K)$ . If  $F \subset C(K)$  separates points, then the closure of the  $*$ -subalgebra generated by  $F$  is either  $C(K)$  or  $\{f \in C(K) \mid f(t_0) = 0\}$  for some  $t_0 \in K$ .*

**Proof.** There are two possibilities: either all  $f \in F$  vanish at one point  $t_0 \in K$  (there can be at most one such point since  $F$  separates points) or there is no such point.

If there is no such point, then the identity can be approximated by elements in  $\overline{A}$ : First of all note that  $|f| \in \overline{A}$  if  $f \in \overline{A}$ , since the polynomials  $p_n(t)$  used to prove this fact can be replaced by  $p_n(t) - p_n(0)$  which contain no constant term. Hence for every point  $y$  we can find a nonnegative function in  $\overline{A}$  which is positive at  $y$  and by compactness we can find a finite sum of such functions which is positive everywhere, say  $m \leq f(t) \leq M$ . Now approximate  $\min(m^{-1}t, t^{-1})$  by polynomials  $q_n(t)$  (again a constant term is not needed) to conclude that  $q_n(f) \rightarrow f^{-1} \in \overline{A}$ . Hence  $1 = f \cdot f^{-1} \in \overline{A}$  as claimed and so  $\overline{A} = C(K)$  by the Stone–Weierstraß theorem.

If there is such a  $t_0$  we have  $\overline{A} \subseteq \{f \in C(K) \mid f(t_0) = 0\}$  and the identity is clearly missing from  $\overline{A}$ . However, adding the identity to  $\overline{A}$  we get  $\overline{A} + \mathbb{C} = C(K)$  by the Stone–Weierstraß theorem. Moreover, if  $f \in C(K)$  with  $f(t_0) = 0$  we get  $f = \tilde{f} + \alpha$  with  $\tilde{f} \in \overline{A}$  and  $\alpha \in \mathbb{C}$ . But  $0 = f(t_0) = \tilde{f}(t_0) + \alpha = \alpha$  implies  $f = \tilde{f} \in \overline{A}$ , that is,  $\overline{A} = \{f \in C(K) \mid f(t_0) = 0\}$ .  $\square$

**Problem B.35.** *Show that the uniform limit of uniformly continuous functions is again uniformly continuous.*



**Problem B.36.** Suppose  $X$  is compact and connected and let  $F \subset C(X, Y)$  be a family of equicontinuous functions. Then  $\{f(x_0) | f \in F\}$  bounded for one  $x_0$  implies  $F$  bounded.

**Problem B.37.** Let  $X, Y$  be metric spaces. A family of functions  $F \subset C(X, Y)$  is called **uniformly equicontinuous** if for every  $\varepsilon > 0$  there is a  $\delta > 0$  such that

$$d_Y(f(x), f(y)) \leq \varepsilon \quad \text{whenever} \quad d_X(x, y) < \delta, \quad \forall f \in F. \quad (\text{B.31})$$

Show that if  $X$  is compact, then a family  $F$  is pointwise equicontinuous if and only if it is uniformly equicontinuous.

**Problem\* B.38** (Dini's theorem). Suppose  $X$  is compact and let  $f_n \in C(X)$  be a sequence of decreasing (or increasing) functions converging pointwise  $f_n(x) \searrow f(x)$  to some function  $f \in C(X)$ . Then  $f_n \rightarrow f$  uniformly. (Hint: Reduce it to the case  $f_n \searrow 0$  and apply the finite intersection property to  $f_n^{-1}([\varepsilon, \infty))$ .)

**Problem B.39.** Let  $k \in \mathbb{N}$  and  $I \subseteq \mathbb{R}$ . Show that the  $*$ -subalgebra generated by  $f_{z_0}(t) = \frac{1}{(t-z_0)^k}$  for one  $z_0 \in \mathbb{C}$  is dense in the set  $C_0(I)$  of continuous functions vanishing at infinity:

- for  $I = \mathbb{R}$  if  $z_0 \in \mathbb{C} \setminus \mathbb{R}$  and  $k = 1$  or  $k = 2$ ,
- for  $I = [a, \infty)$  if  $z_0 \in (-\infty, a)$  and  $k$  arbitrary,
- for  $I = (-\infty, a] \cup [b, \infty)$  if  $z_0 \in (a, b)$  and  $k$  odd.

(Hint: Add  $\infty$  to  $\mathbb{R}$  to make it compact.)

**Problem B.40.** Let  $U \subseteq \mathbb{C} \setminus \mathbb{R}$  be a set which has a limit point and is symmetric under complex conjugation. Show that the span of  $\{(t-z)^{-1} | z \in U\}$  is dense in the set  $C_0(\mathbb{R})$  of continuous functions vanishing at infinity. (Hint: The product of two such functions is in the span provided they are different.)

**Problem B.41.** Let  $K \subseteq \mathbb{C}$  be a compact set. Show that the set of all functions  $f(z) = p(x, y)$ , where  $p: \mathbb{R}^2 \rightarrow \mathbb{C}$  is polynomial and  $z = x + iy$ , is dense in  $C(K)$ .

---

# Bibliography

- [1] H. W. Alt, *Lineare Funktionalanalysis*, 4th ed., Springer, Berlin, 2002.
- [2] H. Bauer, *Measure and Integration Theory*, de Gruyter, Berlin, 2001.
- [3] M. Berger and M. Berger, *Perspectives in Nonlinearity*, Benjamin, New York, 1968.
- [4] A. Bowers and N. Kalton, *An Introductory Course in Functional Analysis*, Springer, New York, 2014.
- [5] H. Brezis, *Functional Analysis, Sobolev Spaces and Partial Differential Equations*, Springer, New York, 2011.
- [6] S.-N. Chow and J. K. Hale, *Methods of Bifurcation Theory*, Springer, New York, 1982.
- [7] J. B. Conway, *A Course in Functional Analysis*, 2nd ed., Springer, New York, 1994.
- [8] K. Deimling, *Nichtlineare Gleichungen und Abbildungsgrade*, Springer, Berlin, 1974.
- [9] K. Deimling, *Nonlinear Functional Analysis*, Springer, Berlin, 1985.
- [10] E. DiBenedetto, *Real Analysis*, Birkhäuser, Boston, 2002.
- [11] K.-J. Engel and R. Nagel, *One-Parameter Semigroups for Linear Evolution Equations*, Springer, Berlin, 2000.
- [12] L. C. Evans, *Weak Convergence Methods for nonlinear Partial Differential Equations*, CBMS 74, American Mathematical Society, Providence, 1990.
- [13] L. C. Evans, *Partial Differential Equations*, 2nd ed., American Mathematical Society, Providence, 2010.
- [14] G. B. Folland, *Real Analysis: Modern Techniques and Their Applications*, 2nd ed., Wiley, Hoboken NJ, 1999.
- [15] J. Franklin, *Methods of Mathematical Economics*, Springer, New York, 1980.
- [16] I. Gohberg, S. Goldberg, and M.A. Kaashoek, *Basic Classes of Linear Operators*, Springer, Basel, 2003.
- [17] J. Goldstein, *Semigroups of Linear Operators and Applications*, Oxford University Press, New York, 1985.

- [18] L. Grafakos, *Classical Fourier Analysis*, 2nd ed., Springer, New York, 2008.
- [19] L. Grafakos, *Modern Fourier Analysis*, 2nd ed., Springer, New York, 2009.
- [20] G. Grubb, *Distributions and Operators*, Springer, New York, 2009.
- [21] E. Hewitt and K. Stromberg, *Real and Abstract Analysis*, Springer, Berlin, 1965.
- [22] D. Hundertmark, M. Meyries, L. Machinek, and R. Schnaubelt, *Operator Semigroups and Dispersive Equations*, Lecture Notes (16th Internet Seminar on Evolution Equations), 2013. [https://isem.math.kit.edu/images/b/b3/Isem16\\_final.pdf](https://isem.math.kit.edu/images/b/b3/Isem16_final.pdf)
- [23] K. Jänich, *Topology*, Springer, New York, 1995.
- [24] I. Kaplansky, *Set Theory and Metric Spaces*, AMS Chelsea, Providence, 1972.
- [25] T. Kato, *Perturbation Theory for Linear Operators*, Springer, New York, 1966.
- [26] J. L. Kelley, *General Topology*, Springer, New York, 1955.
- [27] O. A. Ladyzhenskaya, *The Boundary Value Problems of Mathematical Physics*, Springer, New York, 1985.
- [28] P. D. Lax, *Functional Analysis*, Wiley, New York, 2002.
- [29] E. Lieb and M. Loss, *Analysis*, 2nd ed., Amer. Math. Soc., Providence, 2000.
- [30] F. Linares and G. Ponce, *Introduction to Nonlinear Dispersive Equations*, 2nd ed., Springer, New York, 2015.
- [31] G. Leoni, *A First Course in Sobolev Spaces*, Amer. Math. Soc., Providence, 2009.
- [32] N. Lloyd, *Degree Theory*, Cambridge University Press, London, 1978.
- [33] R. Meise and D. Vogt, *Introduction to Functional Analysis*, Oxford University Press, Oxford, 2007.
- [34] F. W. J. Olver et al., *NIST Handbook of Mathematical Functions*, Cambridge University Press, Cambridge, 2010.
- [35] I. K. Rana, *An Introduction to Measure and Integration*, 2nd ed., Amer. Math. Soc., Providence, 2002.
- [36] M. Reed and B. Simon, *Methods of Modern Mathematical Physics I. Functional Analysis*, rev. and enl. edition, Academic Press, San Diego, 1980.
- [37] J. R. Retherford, *Hilbert Space: Compact Operators and the Trace Theorem*, Cambridge University Press, Cambridge, 1993.
- [38] J.J. Rotman, *Introduction to Algebraic Topology*, Springer, New York, 1988.
- [39] H. Royden, *Real Analysis*, Prentice Hall, New Jersey, 1988.
- [40] W. Rudin, *Real and Complex Analysis*, 3rd edition, McGraw-Hill, New York, 1987.
- [41] M. Růžička, *Nichtlineare Funktionalanalysis*, Springer, Berlin, 2004.
- [42] H. Schröder, *Funktionalanalysis*, 2nd ed., Harri Deutsch Verlag, Frankfurt am Main 2000.
- [43] B. Simon, *A Comprehensive Course in Analysis*, Amer. Math. Soc., Providence, 2015.
- [44] L. A. Steen and J. A. Seebach, Jr., *Counterexamples in Topology*, Springer, New York, 1978.
- [45] T. Tao, *Nonlinear Dispersive Equations: Local and Global Analysis*, Amer. Math. Soc., Providence, 2006.
- [46] M. E. Taylor, *Measure Theory and Integration*, Amer. Math. Soc., Providence, 2006.

- 
- [47] G. Teschl, *Mathematical Methods in Quantum Mechanics; With Applications to Schrödinger Operators*, Amer. Math. Soc., Providence, 2009.
  - [48] G. Teschl, *Topics in Real Analysis*, Amer. Math. Soc., Providence, to appear.
  - [49] J. Weidmann, *Lineare Operatoren in Hilberträumen I: Grundlagen*, B.G.Teubner, Stuttgart, 2000.
  - [50] D. Werner, *Funktionalanalysis*, 7th edition, Springer, Berlin, 2011.
  - [51] M. Willem, *Functional Analysis*, Birkhäuser, Basel, 2013.
  - [52] E. Zeidler, *Applied Functional Analysis: Applications to Mathematical Physics*, Springer, New York 1995.
  - [53] E. Zeidler, *Applied Functional Analysis: Main Principles and Their Applications*, Springer, New York 1995.



---

## Glossary of notation

$\arg(z)$	... argument of $z \in \mathbb{C}$ ; $\arg(z) \in (-\pi, \pi]$ , $\arg(0) = 0$
$B_r(x)$	... open ball of radius $r$ around $x$ , 370
$B(X)$	... Banach space of bounded measurable functions
$\mathfrak{B}$	$= \mathfrak{B}^1$
$\mathfrak{B}^n$	... Borel $\sigma$ -algebra of $\mathbb{R}^n$ , see [48]
$\mathbb{C}$	... the set of complex numbers
$C(U)$	... set of continuous functions from $U$ to $\mathbb{C}$
$C_0(U)$	... set of continuous functions vanishing on the boundary $\partial U$ , 399
$C_c(U)$	... set of compactly supported continuous functions
$C_{per}[a, b]$	... set of periodic continuous functions (i.e. $f(a) = f(b)$ )
$C^k(U)$	... set of $k$ times continuously differentiable functions
$C_c^\infty(U)$	... set of compactly supported smooth functions
$C(U, Y)$	... set of continuous functions from $U$ to $Y$ , 232
$C^r(U, Y)$	... set of $r$ times continuously differentiable functions, 238
$C_b^r(U, Y)$	... functions in $C^r$ with derivatives bounded, 37, 243
$C_c^r(U, Y)$	... functions in $C^r$ with compact support
$c_0(\mathbb{N})$	... set of sequences converging to zero, 10
$\mathcal{C}(X, Y)$	... set of closed linear operators from $X$ to $Y$ , 199
$\text{CP}(f)$	... critical points of $f$ , 321
$\text{CS}(K)$	... nonempty convex subsets of $K$ , 334
$\text{CV}(f)$	... critical values of $f$ , 321

$\chi_\Omega(\cdot)$	... characteristic function of the set $\Omega$
$\mathfrak{D}(\cdot)$	... domain of an operator
$\delta_{n,m}$	... Kronecker delta, 12
$\deg(D, f, y)$	... mapping degree, 321, 328
$\det$	... determinant
$\dim$	... dimension of a linear space
$\operatorname{div}$	... divergence of a vector field, 327
$\operatorname{diam}(U)$	$= \sup_{(x,y) \in U^2} d(x, y)$ diameter of a set
$\operatorname{dist}(U, V)$	$= \inf_{(x,y) \in U \times V} d(x, y)$ distance of two sets
$D_y^r(\overline{U}, Y)$	... functions in $C^r(\overline{U}, Y)$ which do not attain $y$ on the boundary, 321
$e$	... Napier's constant, $e^z = \exp(z)$
$\operatorname{epi} F$	... epigraph of $F$ , 257
$dF$	... derivative of $F$ , 232
$\mathcal{F}(X, Y)$	... set of compact finite dimensional functions, 342
$\Phi(X, Y)$	... set of all linear Fredholm operators from $X$ to $Y$ , 185
$\Phi_0(X, Y)$	... set of all linear Fredholm operators of index 0, 185
$\operatorname{GL}(n)$	... general linear group in $n$ dimensions
$\Gamma(z)$	... gamma function, see [48]
$\Gamma(A)$	... graph of an operator, 103
$\Gamma(f_1, \dots, f_n)$	... Gram determinant, 47
$\mathfrak{H}$	... a Hilbert space
$\operatorname{conv}(\cdot)$	... convex hull
$\mathcal{H}(U)$	... set of holomorphic functions on a domain $U \subseteq \mathbb{C}$ , 319
$H^k(U)$	$= W^{k,2}(U)$ , Sobolev space
$H_0^k(U)$	$= W_0^{k,2}(U)$ , Sobolev space
$i$	... complex unity, $i^2 = -1$
$\operatorname{Im}(\cdot)$	... imaginary part of a complex number
$\inf$	... infimum
$J_f(x)$	$= \det df(x)$ Jacobi determinant of $f$ at $x$ , 321
$\operatorname{Ker}(A)$	... kernel of an operator $A$ , 28
$\mathcal{K}(X, Y)$	... set of compact linear operators from $X$ to $Y$ , 65
$\mathcal{K}(U, Y)$	... set of compact maps from $U$ to $Y$ , 342
$\bar{\mathcal{K}}_y(U, Y)$	... functions in $\mathcal{K}(\overline{U}, Y)$ which do not attain $y$ on the boundary, 344
$\lambda^n$	... Lebesgue measure in $\mathbb{R}^n$ , see [48]
$\mathcal{L}(X, Y)$	... set of all bounded linear operators from $X$ to $Y$ , 30
$\mathcal{L}(X)$	$= \mathcal{L}(X, X)$
$L^p(X, d\mu)$	... Lebesgue space of $p$ integrable functions, see [48]
$L^\infty(X, d\mu)$	... Lebesgue space of bounded functions, see [48]
$L_{loc}^p(X, d\mu)$	... locally $p$ integrable functions, see [48]
$\mathcal{L}_{cont}^2(I)$	... space of continuous square integrable functions, 20
$\ell^p(\mathbb{N})$	... Banach space of $p$ summable sequences, 9
$\ell^2(\mathbb{N})$	... Hilbert space of square summable sequences, 18
$\ell^\infty(\mathbb{N})$	... Banach space of bounded sequences, 10
$\max$	... maximum

---

$\mathbb{N}$	... the set of positive integers
$\mathbb{N}_0$	$= \mathbb{N} \cup \{0\}$
$n(\gamma, z_0)$	... winding number
$O(\cdot)$	... Landau symbol, $f = O(g)$ iff $\limsup_{x \rightarrow x_0}  f(x)/g(x)  < \infty$
$o(\cdot)$	... Landau symbol, $f = o(g)$ iff $\lim_{x \rightarrow x_0}  f(x)/g(x)  = 0$
$\mathbb{Q}$	... the set of rational numbers
$\mathbb{R}$	... the set of real numbers
$\rho(A)$	... resolvent set of an operator $A$ , 132
$\text{RV}(f)$	... regular values of $f$ , 321
$\text{Ran}(A)$	... range of an operator $A$ , 28
$\text{Re}(\cdot)$	... real part of a complex number
$R(I, X)$	... set of regulated functions, 230
$\sigma(A)$	... spectrum of an operator $A$ , 132, 177
$S^{n-1}$	$= \{x \in \mathbb{R}^n \mid  x  = 1\}$ unit sphere in $\mathbb{R}^n$
$\text{sign}(z)$	$= z/ z $ for $z \neq 0$ and 1 for $z = 0$ ; complex sign function
$S(I, X)$	... step functions $f : I \rightarrow X$ , 230
$\sup$	... supremum
$\text{supp}(f)$	... support of a function $f$ , 380
$\text{span}(M)$	... set of finite linear combinations from $M$ , 12
$W^{k,p}(U)$	... Sobolev space, see [48]
$W_0^{k,p}(U)$	... Sobolev space, see [48]
$\mathbb{Z}$	... the set of integers
$\mathbb{I}$	... identity operator
$\sqrt{z}$	... square root of $z$ with branch cut along $(-\infty, 0)$
$z^*$	... complex conjugation
$A^*$	... adjoint of operators $A$ , 52
$\overline{A}$	... closure of operators $A$ , 200
$\hat{f}$	$= \mathcal{F}f$ , Fourier coefficients/transform of $f$ , 58
$\check{f}$	$= \mathcal{F}^{-1}f$ , inverse Fourier transform of $f$
$ x $	$= \sqrt{\sum_{j=1}^n  x_j ^2}$ Euclidean norm in $\mathbb{R}^n$ or $\mathbb{C}^n$
$ \Omega $	... Lebesgue measure of a Borel set $\Omega$
$\ \cdot\ $	... norm, 18
$\ \cdot\ _p$	... norm in the Banach space $\ell^p$ and $L^p$ , 9, 24
$\langle \cdot, \cdot \rangle$	... scalar product in $\mathfrak{H}$ , 17



---

$\oplus$	... direct/orthogonal sum of vector spaces or operators, 33, 56
$\hat{\oplus}$	... direct sum of operators with the same image space, 34
$\otimes$	... tensor product, 57
$\sqcup$	... union of disjoint sets
$\lfloor x \rfloor$	$= \max\{n \in \mathbb{Z}   n \leq x\}$ , floor function
$\lceil x \rceil$	$= \min\{n \in \mathbb{Z}   n \geq x\}$ , ceiling function
$\partial$	$= (\partial_1 f, \dots, \partial_m f)$ gradient in $\mathbb{R}^m$
$\partial_\alpha$	... partial derivative in multi-index notation
$\partial_x F(x, y)$	... partial derivative with respect to $x$ , 237
$\partial U$	$= \overline{U} \setminus U^\circ$ boundary of the set $U$ , 370
$\overline{U}$	... closure of the set $U$ , 374
$U^\circ$	... interior of the set $U$ , 374
$M^\perp$	... orthogonal complement, 48
$(\lambda_1, \lambda_2)$	$= \{\lambda \in \mathbb{R}   \lambda_1 < \lambda < \lambda_2\}$ , open interval
$[\lambda_1, \lambda_2]$	$= \{\lambda \in \mathbb{R}   \lambda_1 \leq \lambda \leq \lambda_2\}$ , closed interval
$x_n \rightarrow x$	... norm convergence, 8
$x_n \rightharpoonup x$	... weak convergence, 119
$x_n \xrightarrow{*} x$	... weak-* convergence, 126
$A_n \rightarrow A$	... norm convergence of operators
$A_n \xrightarrow{s} A$	... strong convergence of operators, 124
$A_n \rightharpoonup A$	... weak convergence of operators, 124

---

# Index

sigma-comapct, 388

absolute convergence, 16

absolutely convex, 154, 167

absorbing set, 151

accumulation point, 370

adjoint operator, 52, 114, 203

adjugate, 133

Alexandroff extension, 390

almost periodic, 46

analytic, 133

annihilator, 111

approximate point spectrum, 210

ascent, 180

Axiom of Choice, 363

axiomatic set theory, 361

Baire category theorem, 97

balanced set, 170

ball

closed, 375

open, 370

Banach algebra, 32, 130

Banach limit, 114

Banach space, 9

Banach–Steinhaus theorem, 99

base, 372

Basel problem, 82

basis, 12

orthonormal, 43

Bernoulli numbers, 82

Bessel inequality, 42

best reply, 336

bidual space, 109

bifurcation point, 267

bijjective, 379

biorthogonal system, 12, 108

Bolzano–Weierstraß theorem, 388

boundary condition, 6

boundary point, 370

boundary value problem, 6

bounded

operator, 28

sesquilinear form, 23

set, 376

Brouwer fixed point theorem, 332

calculus of variations, 244

direct method, 250

Calkin algebra, 193

Cauchy sequence, 377

weak, 119

Cauchy–Bunyakovsky–Schwarz inequality,

*see* Cauchy–Schwarz inequality

Cauchy–Schwarz inequality, 19

Cayley transform, 143

Cesàro mean, 114

chain rule, 238

character, 191

Chebyshev polynomials, 70

closed

ball, 375

set, 374

closure, 374

cluster point, 370

codimension, 35

coercive, 54, 356

weakly, 250

cofinite topology, 373

cokernel, 35

commutes, 284

commuting operators, 208

- compact, 385
  - locally, 388
  - sequentially, 387
- compact map, 342
- complemented subspace, 34
- complete, 9, 377
- completion, 24
- complexification, 36
- component, 395
- conjugate linear, 17
- connected, 395
- continuous, 379
- contraction principle, 257
  - uniform, 258
- contraction semigroup, 289
- convergence, 375
  - strong, 124
  - weak, 119
  - weak-\*, 126
- convex, 8
  - absolutely, 154, 167
- core, 200
- cover, 384
  - locally finite, 384
  - open, 384
  - refinement, 384
- critical value, 321
- $C^*$  algebra, 140
- cylinder, 383
  - set, 383
- De Morgan's laws, 374
- decent, 180
- delay differential equation, 294
- demicontinuous, 358
- dense, 377
- derivative
  - Fréchet, 233
  - Gâteaux, 234
  - partial, 237
  - variational, 234
- diffeomorphism, 238
- differentiable, 237
- differential equations, 261
- diffusion equation, 3
- dimension, 45
- direct sum, 33
- directed, 167
- Dirichlet kernel, 59
- Dirichlet problem, 56
- disconnected, 394
- discrete set, 370
- discrete topology, 371
- disjoint union topology, 383
- dissipative, 291
- distance, 369, 390
- divergence, 327
- domain, 28
- double dual space, 109
- dual basis, 30
- dual space, 30
- duality set, 290
- Duhamel formula, 275, 281
- eigenspace, 68
- eigenvalue, 68
  - algebraic multiplicity, 183
  - geometric multiplicity, 183
  - index, 183
  - simple, 69
- eigenvector, 68
  - order, 183
- elliptic problem, 358
- epigraph, 257
- equicontinuous, 26, 401
  - uniformly, 404
- equilibrium
  - Nash, 336
- equivalent norms, 22
- exact sequence, 187, 208
- exhaustion, 389
- extended real numbers, 372
- extension, 199
- extension principle, 30
- extremal
  - point, 156
  - subset, 156
- Extreme value theorem, 388
- $F_\sigma$  set, 98
- face, 157
- fat set, 98
- Fejér kernel, 61
- final topology, 383
- finite dimensional map, 342
- finite intersection property, 385
- first category, 98
- first countable, 373
- first resolvent identity, 139, 210
- fixed point theorem
  - Altman, 348
  - Brouwer, 332
  - contraction principle, 257
  - Kakutani, 334
  - Krasnosel'skii, 348
  - Rothe, 348
  - Schauder, 346
  - Weissinger, 258
- form
  - bounded, 23
- Fourier series, 44, 58
  - cosine, 81
  - sine, 80
- FPU lattice, 264

- Fréchet derivative, 233
- Fréchet space, 169
- Fredholm alternative, 182
- Fredholm operator, 185, 222
- Frobenius norm, 92
- from domain, 78
- function
  - open, 380
- fundamental theorem of algebra, 134
- fundamental theorem of calculus, 231
- $G_\delta$  set, 98
- Gâteaux derivative, 234
- Galerkin approximation, 359
- gauge, 151
- Gelfand transform, 194
- global solution, 264
- Gram determinant, 47
- Gram–Schmidt orthogonalization, 44
- graph, 103
- graph norm, 200
- Green function, 76
- Gronwall’s inequality, 353
- group
  - strongly continuous, 276
- growth bound, 276
- half-space, 162
- Hamel basis, 12, 16
- Hankel operator, 95
- Hardy space, 197
- Hausdorff space, 373
- heat equation, 3, 293
- Heine–Borel theorem, 388
- Hermitian form, 17
- Hilbert space, 18
  - dimension, 45
- Hilbert–Schmidt operator, 90
- Hölder continuous, 38
- Hölder’s inequality, 10, 25
- holomorphic function, 319
- homeomorphic, 380
- homeomorphism, 380
- homotopy, 320
- homotopy invariance, 321
- ideal, 192
  - maximal, 192
  - proper, 192
- identity, 32, 130
- implicit function theorem, 259
- index, 185, 222
- induced topology, 372
- Induction Principle, 364
- initial topology, 382
- injective, 379
- inner product, 18
- inner product space, 18
- integral, 230
- interior, 374
- interior point, 370
- inverse function theorem, 260
- involution, 140
- isolated point, 370
- isometric, 379
- isometry, 8
- isomorphic, 8
- Jacobi determinant, 321
- Jacobi matrix, 237
- Jacobi operator, 69, 265, 266, 275, 299
- Jacobi theta function, 7
- Jacobson radical, 195
- Jordan curve theorem, 340
- Kakutani’s fixed point theorem, 334
- kernel, 28
- Kronecker delta, 12
- Kuratowski closure axioms, 374
- Ladyzhenskaya, 352
- Landau inequality, 284
- Landau kernel, 14
- Landau symbols, 232
- Lax–Milgram theorem, 54
  - nonlinear, 357
- Legendre polynomials, 44
- Leray–Schauder principle, 346
- Lidskij trace theorem, 93
- Lie group, 260
- liminf, 381
- limit, 375
- limit point, 370
- limsup, 381
- Lindelöf theorem, 384
- linear
  - functional, 30, 49
  - operator, 28
- linearly independent, 12
- Lipschitz continuous, 38
- locally
  - (path-)connected, 395
- locally convex space, 154, 165
- lower semicontinuous, 380
- Lyapunov–Schmidt reduction, 269
- maximal solution, 264
- maximum norm, 7
- meager set, 98
- mean value theorem, 240
- metric, 369
  - translation invariant, 169
- mild solution, 282
- minimum modulus, 203

- 
- Minkowski functional, 151
  - monotone, 248, 356
    - map, 355
    - strictly, 248, 356
    - strongly, 356
  - Morawetz identity, 315
  - multilinear function, 239
    - symmetric, 239
  - multiplicative linear functional, 191
  - multiplicity
    - algebraic, 183
    - geometric, 183
  - Nash equilibrium, 336
  - Nash theorem, 337
  - Navier–Stokes equation, 350
  - neighborhood, 370
  - neighborhood base, 372
  - Neumann series, 135
  - nilpotent, 136
  - Noether operator, 185
  - nonlinear Schrödinger equation, 266, 299, 301
  - norm, 8
    - operator, 28
    - strictly convex, 16, 172
    - stronger, 20
    - uniformly convex, 172
  - norm-attaining, 112
  - normal
    - operator, 141, 143
    - space, 391
  - normalized, 18
  - normed space, 8
  - nowhere dense, 97
  - $n$ -person game, 335
  - null space, 28
  - one-point compactification, 390
  - one-to-one, 379
  - onto, 379
  - open
    - ball, 370
    - function, 380
    - set, 370
  - operator
    - adjoint, 52, 114
    - bounded, 28
    - closable, 200
    - closed, 103
    - closure, 200
    - compact, 65
    - completely continuous, 123
    - domain, 28
    - finite rank, 87, 116
    - linear, 28
    - nonnegative, 53
    - relatively bounded, 217
    - relatively compact, 220
    - self-adjoint, 68
    - strong convergence, 124
    - symmetric, 68
    - unitary, 46
    - weak convergence, 124
  - order
    - partial, 363
    - total, 363
    - well, 363
  - orthogonal, 18
    - complement, 48
    - projection, 48, 145
    - sum, 56
  - orthonormal
    - basis, 43
    - set, 41
  - parallel, 18
  - parallelogram law, 20
  - parametrix, 189
  - Parseval relation, 44
  - partial order, 363
  - partition of unity, 393
  - path, 395
  - path-connected, 395
  - payoff, 335
  - Peano theorem, 350
  - perpendicular, 18
  - Poisson problem, 253
  - polar decomposition, 86
  - polar set, 155
  - polarization identity, 20
  - power set, 362
  - prisoner's dilemma, 336
  - product rule, 232, 240, 275
  - product topology, 382
  - projection, 139
  - projection-valued measure, 145
  - proper
    - function, 389
    - map, 343
    - metric space, 388
  - pseudometric, 370
  - Pythagorean theorem, 18
  - quadrangle inequality, 375
  - quasiconvex, 251
  - quasinilpotent, 137
  - quasinorm, 17
  - quotient map, 35
  - quotient space, 35
  - quotient topology, 383
  - range, 28
  - rank, 87, 116

- Rayleigh–Ritz method, 83
- reaction-diffusion equation, 5
- reducing subspace, 214
- reduction property, 337
- refinement, 384
- reflexive, 109
- regular value, 321
- regulated function, 230
- relative topology, 372
- relatively bounded, 217
- relatively compact, 220, 385
- reproducing kernel, 81
- residual set, 98
- resolution of the identity, 146
- resolvent, 73, 75, 133, 209
  - formula
    - second, 219
- resolvent identity
  - first, 139, 210
- resolvent set, 132, 209
- Riemann–Lebesgue lemma, 60
- Riesz
  - theorem, 49
- Riesz lemma, 181
- Riesz representation theorem, 49
- Ritz method, 83
- Rouché’s theorem, 320
- Russell’s paradox, 361
  
- Sard’s theorem, 325
- scalar product, 18
- Schatten  $p$ -class, 89
- Schauder basis, 12
- Schrödinger equation, 6, 293, 301
- Schur property, 123
- Schwartz space, 169
- Schwarz’ theorem, 239
- second category, 98
- second countable, 373
- second resolvent formula, 219
- self-adjoint, 52, 141
- semigroup
  - differentiable, 279
  - generator, 277
  - strongly continuous, 276
  - uniform, 273
- seminorm, 8
- separable, 13, 377
- separation, 390
  - of convex sets, 152
  - of points, 383, 402
  - of variables, 4
  - seminorms, 166
- sequentially closed, 376
- sequentially continuous, 380
- series
  - absolutely convergent, 16
  
- sesquilinear form, 17
  - bounded, 23
  - parallelogram law, 23
  - polarization identity, 23
- shift operator, 52, 69
- singular value decomposition, 86
- singular values, 85
- Sobolev space, 201
- span, 12
- spectral bound, 285
- spectral measure, 144
- spectral projections, 145
- spectral radius, 135
- spectral theorem
  - compact operators, 183
  - compact self-adjoint operators, 71
  - normal operators, 197
  - self-adjoint operators, 142, 144
- spectrum, 73, 132, 209
  - approximate point, 210
  - continuous, 177
  - discrete, 190
  - essential, 189
  - Fredholm, 189
  - point, 177
  - residual, 177
- \*-subalgebra, 140
- step function, 230
- Stone–Weierstraß theorem, 403
- strategy, 335
- strictly convex, 16
- strictly convex space, 172
- strong convergence, 124
- strong solution, 280
- Sturm–Liouville problem, 6
- subbase, 372
- subcover, 384
- subspace topology, 372
- support, 380
- support hyperplane, 157
- surjective, 379
- symmetric
  - operator, 68
  - sesquilinear form, 17
  
- Taylor’s theorem, 242
- Taylor’s theorem, 242
- tempered distributions, 169
- tensor product, 57
- theorem
  - Altman, 348
  - Arzelà–Ascoli, 27, 401
  - Atkinson, 187, 189
  - Bair, 97
  - Banach–Alaoglu, 161
  - Banach–Steinhaus, 99, 125
  - Bernstein, 61

- Beurling–Gelfand, 135
- bipolar, 155
- Bolzano–Weierstraß, 388
- Borsuk, 330, 345
- Borsuk–Ulam, 331
- Brouwer, 331, 345
- Browder–Minty, 359
- Carathéodory, 160
- closed graph, 103
- closed range, 118, 206
- Crandall–Rabinowitz, 270
- Dieudonné, 188
- Dini, 404
- Ehrling, 68
- Fejér, 62
- Feller–Miyadera–Phillips, 287
- fundamental thm. of calculus, 231
- Gelfand representation, 194
- Gelfand–Mazur, 134
- Gelfand–Naimark, 196
- Goldstine, 163
- Hahn–Banach, 106
- Hahn–Banach, geometric, 153
- hairy ball, 329
- Heine–Borel, 388
- Hellinger–Toeplitz, 103
- Helly, 127
- Hessenberg, 366
- Hilbert projection, 49
- Hilbert–Schmidt, 71
- Hille–Yosida, 290
- implicit function, 259
- intermediate value, 397
- invariance of domain, 331, 345
- inverse function, 260
- Jordan–von Neumann, 20
- Kakutani, 162
- Kolmogorov, 170
- Krasnosel’skii, 348
- Krein–Milman, 158
- Lax–Milgram, 54, 357
- Leray–Schauder, 346
- Lindelöf, 384
- Lumer–Phillips, 292
- McShane, 393
- mean value, 229
- Milman–Pettis, 174
- Nash, 337
- Omega lemma, 261
- open mapping, 101, 102
- Peano, 350
- Perron–Frobenius, 333
- Pythagorean, 18
- Radon–Riesz, 173
- rank-nullity, 185
- Riesz, 183, 186, 188
- Rothe, 348
- Rouché, 320, 322
- Sard, 325
- Schaefer, 346
- Schauder, 116, 346
- Schröder–Bernstein, 365
- Schwarz, 239
- Šmulian, 121
- spectral, 71, 142, 144, 183, 197
- spectral mapping, 134
- Stone, 293
- Stone–Weierstraß, 403
- Taylor, 242
- Tietze, 346, 392
- Tychonoff, 386
- Urysohn, 391
- Weierstraß, 14, 388
- Weissinger, 258
- Weyl, 189, 224
- Wiener, 196
- Yood, 189
- Zarantonello, 356
- Zermelo, 365
- Zorn, 364
- Toda lattice, 266
- topological space, 371
- topological vector space, 151
- topology
  - base, 372
  - product, 382
  - relative, 372
  - subbase, 372
- total order, 363
- total set, 13, 112
- totally bounded, 389
- trace, 93
  - class, 90
- trace formula, 80
- trace topology, 372
- transcritical bifurcation, 270
- triangle inequality, 8, 369
  - inverse, 8, 370
- trivial topology, 371
- uniform boundedness principle, 99
- uniform contraction principle, 258
- uniform convergence, 398
- uniformly continuous, 400
- uniformly convex space, 172
- unit vector, 18
- unital, 131
- unitarily equivalent, 46
- unitary, 141
- Unitization, 139, 142
- upper semicontinuous, 380
- Urysohn lemma, 391
- Vandermonde determinant, 16

- 
- variational derivative, 234
  - virial identity, 315
  - Volterra integral operator, 136
  
  - wave equation, 5
  - weak convergence, 119
  - weak solution, 351
  - weak topology, 120, 160
  - weak-\* convergence, 126
  - weak-\* topology, 161
  - weakly coercive, 250
  - Weierstraß approximation, 14
  - Weierstraß theorem, 388
  - well-order, 363
  - Weyl asymptotic, 85
  - Weyl sequence, 210
    - singular, 224
  - Wiener algebra, 130
  - winding number, 319
  
  - Young inequality, 10
  
  - Zermelo–Fraenkel set theory, 361
  - ZF, 361
  - ZFC, 363
  - Zorn’s lemma, 364