

Introduction to the
**Modern Theory of
Dynamical Systems**



Anatole Katok

Boris Hasselblatt

This book provides the first self-contained comprehensive exposition of the theory of dynamical systems as a core mathematical discipline closely intertwined with most of the main areas of mathematics. The authors introduce and rigorously develop the theory while providing researchers interested in applications with fundamental tools and paradigms.

The book begins with a discussion of several elementary but fundamental examples. These are used to formulate a program for the general study of asymptotic properties and to introduce the principal theoretical concepts and methods. The main theme of the second part of the book is the interplay between local analysis near individual orbits and the global complexity of the orbit structure. The third and fourth parts develop in depth the theories of low-dimensional dynamical systems and hyperbolic dynamical systems.

The book is aimed at students and researchers in mathematics at all levels from advanced undergraduate up. Scientists and engineers working in applied dynamics, non-linear science, and chaos will also find many fresh insights in this concrete and clear presentation. It contains more than four hundred systematic exercises.

ENCYCLOPEDIA OF MATHEMATICS AND ITS APPLICATIONS

EDITED BY G.-C. ROTA

Editorial Board

R. Doran, M. Ismail, T.-Y. Lam, E. Lutwak, R. Spigler

Volume 54

Introduction to the Modern Theory of Dynamical Systems

ENCYCLOPEDIA OF MATHEMATICS AND ITS APPLICATIONS

- 18 H. O. Fattorini *The Cauchy problem*
- 19 G. G. Lorentz, K. Jetter, and S. D. Riemenschneider *Birkhoff interpolation*
- 21 W. T. Tutte *Graph theory*
- 22 J. R. Bastida *Field extensions and Galois theory*
- 23 J. R. Cannon *The one-dimensional heat equation*
- 25 A. Salomaa *Computation and automata*
- 26 N. White (ed.) *Theory of matroids*
- 27 N. H. Bingham, C. M. Goldie, and J. L. Teugels *Regular variation*
- 28 P. P. Petrushev and V. A. Popov *Rational approximation of real functions*
- 29 N. White (ed.) *Combinatorial geometries*
- 30 M. Pohst and H. Zassenhaus *Algorithmic algebraic number theory*
- 31 J. Aczel and J. Dhombres *Functional equations containing several variables*
- 32 M. Kuczma, B. Chozewski, and R. Ger *Iterative functional equations*
- 33 R. V. Ambartzumian *Factorization calculus and geometric probability*
- 34 G. Gripenberg, S.-O. Londen, and O. Staffans *Volterra integral and functional equations*
- 35 G. Gasper and M. Rahman *Basic hypergeometric series*
- 36 E. Torgersen *Comparison of statistical experiments*
- 37 A. Neumaier *Interval methods for systems of equations*
- 38 N. Korneichuk *Exact constants in approximation theory*
- 39 R. A. Brualdi and H. J. Ryser *Combinatorial matrix theory*
- 40 N. White (ed.) *Matroid applications*
- 41 S. Sakai *Operator algebras in dynamical systems*
- 42 W. Hodges *Model theory*
- 43 H. Stahl and V. Totik *General orthogonal polynomials*
- 44 R. Schneider *Convex bodies*
- 45 G. Da Prato and J. Zabczyk *Stochastic equations in infinite dimensions*
- 46 A. Björner, M. Las Vergnas, B. Sturmfels, N. White, and G. Ziegler *Oriented matroids*
- 47 E. A. Edgar and L. Sucheston *Stopping times and directed processes*
- 48 C. Sims *Computation with finitely presented groups*
- 49 T. Palmer *Banach algebras and the general theory of *-algebras*
- 50 F. Borceux *Handbook of categorical algebra I*
- 51 F. Borceux *Handbook of categorical algebra II*
- 52 F. Borceux *Handbook of categorical algebra III*

ENCYCLOPEDIA OF MATHEMATICS AND ITS APPLICATIONS

*Introduction to the
Modern Theory of Dynamical Systems*

ANATOLE KATOK

Pennsylvania State University

BORIS HASSELBLATT

Tufts University

With a supplement by Anatole Katok and Leonardo Mendoza



CAMBRIDGE UNIVERSITY PRESS
Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore,
São Paulo, Delhi, Dubai, Tokyo, Mexico City

Cambridge University Press
The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by Cambridge University Press, New York

www.cambridge.org
Information on this title: www.cambridge.org/9780521341875

© Cambridge University Press 1995

This publication is in copyright. Subject to statutory exception
and to the provisions of relevant collective licensing agreements,
no reproduction of any part may take place without the written
permission of Cambridge University Press.

First published 1995
Reprinted 1996
First paperback edition 1997
Reprinted 1998, 1999

A catalogue record for this publication is available from the British Library

ISBN 978-0-521-34187-5 Hardback
ISBN 978-0-521-57557-7 Paperback

Cambridge University Press has no responsibility for the persistence or
accuracy of URLs for external or third-party internet websites referred to in
this publication, and does not guarantee that any content on such websites is,
or will remain, accurate or appropriate. Information regarding prices, travel
timetables, and other factual information given in this work is correct at
the time of first printing but Cambridge University Press does not guarantee
the accuracy of such information thereafter.

To Sveta and Kathy

Contents

PREFACE	xiii
O. INTRODUCTION	1
1. Principal branches of dynamics	1
2. Flows, vector fields, differential equations	6
3. Time-one map, section, suspension	8
4. Linearization and localization	10
Part 1 Examples and fundamental concepts	
1. FIRST EXAMPLES	15
1. Maps with stable asymptotic behavior	15
Contracting maps; Stability of contractions; Increasing interval maps	
2. Linear maps	19
3. Rotations of the circle	26
4. Translations on the torus	28
5. Linear flow on the torus and completely integrable systems	32
6. Gradient flows	35
7. Expanding maps	39
8. Hyperbolic toral automorphisms	42
9. Symbolic dynamical systems	47
Sequence spaces; The shift transformation; Topological Markov chains; The Perron-Frobenius operator for positive matrices	
2. EQUIVALENCE, CLASSIFICATION, AND INVARIANTS	57
1. Smooth conjugacy and moduli for maps	57
Equivalence and moduli; Local analytic linearization; Various types of moduli	
2. Smooth conjugacy and time change for flows	64
3. Topological conjugacy, factors, and structural stability	68
4. Topological classification of expanding maps on a circle	71
Expanding maps; Conjugacy via coding; The fixed-point method	
5. Coding, horseshoes, and Markov partitions	79
Markov partitions; Quadratic maps; Horseshoes; Coding of the toral automorphism	
6. Stability of hyperbolic toral automorphisms	87
7. The fast-converging iteration method (Newton method) for the conjugacy problem	90
Methods for finding conjugacies; Construction of the iteration process	
8. The Poincare-Siegel Theorem	94
9. Cocycles and cohomological equations	100
3. PRINCIPAL CLASSES OF ASYMPTOTIC TOPOLOGICAL INVARIANTS	105
1. Growth of orbits	105
Periodic orbits and the λ -function; Topological entropy; Volume growth; Topological complexity; Growth in the fundamental group; Homological growth	

2. Examples of calculation of topological entropy	119
Isometries; Gradient flows; Expanding maps; Shifts and topological Markov chains; The hyperbolic toral automorphism; Finiteness of entropy of Lipschitz maps; Expansive maps	
3. Recurrence properties	128
4. STATISTICAL BEHAVIOR OF ORBITS AND INTRODUCTION TO ERGODIC THEORY	133
1. Asymptotic distribution and statistical behavior of orbits	133
Asymptotic distribution, invariant measures; Existence of invariant measures; The Birkhoff Ergodic Theorem; Existence of asymptotic distribution; Ergodicity and unique ergodicity; Statistical behavior and recurrence; Measure-theoretic isomorphism and factors	
2. Examples of ergodicity; mixing	146
Rotations; Extensions of rotations; Expanding maps; Mixing; Hyperbolic toral automorphisms; Symbolic systems	
3. Measure-theoretic entropy	161
Entropy and conditional entropy of partitions; Entropy of a measure-preserving transformation; Properties of entropy	
4. Examples of calculation of measure-theoretic entropy	173
Rotations and translations; Expanding maps; Bernoulli and Markov measures; Hyperbolic toral automorphisms	
5. The Variational Principle	179
5. SYSTEMS WITH SMOOTH INVARIANT MEASURES AND MORE EXAMPLES	183
1. Existence of smooth invariant measures	183
The smooth measure class; The Perron-Frobenius operator and divergence; Criteria for existence of smooth invariant measures; Absolutely continuous invariant measures for expanding maps; The Moser Theorem	
2. Examples of Newtonian systems	196
The Newton equation; Free particle motion on the torus; The mathematical pendulum; Central forces	
3. Lagrangian mechanics	200
Uniqueness in the configuration space; The Lagrange equation; Lagrangian systems; Geodesic flows; The Legendre transform	
4. Examples of geodesic flows	205
Manifolds with many symmetries; The sphere and the torus; Isometries of the hyperbolic plane; Geodesics of the hyperbolic plane; Compact factors; The dynamics of the geodesic flow on compact hyperbolic surfaces	
5. Hamiltonian systems	219
Symplectic geometry; Cotangent bundles; Hamiltonian vector fields and flows; Poisson brackets; Integrable systems	
6. Contact systems	229
Hamiltonian systems preserving a I-form; Contact forms	
7. Algebraic dynamics: Homogeneous and affine systems	233
Part 2 Local analysis and orbit growth	
6. LOCAL HYPERBOLIC THEORY AND ITS APPLICATIONS	237
1. Introduction	237
2. Stable and unstable manifolds	239
Hyperbolic periodic orbits; Exponential splitting; The Hadamard-Perron Theorem; Proof of the Hadamard-Perron Theorem; The Inclination Lemma	

3. Local stability of a hyperbolic periodic point	260
The Hartman-Grobman Theorem; Local structural stability	
4. Hyperbolic sets	263
Definition and invariant cones; Stable and unstable manifolds; Closing Lemma and periodic orbits; Locally maximal hyperbolic sets	
5. Homoclinic points and horseshoes	273
General horseshoes; Homoclinic points; Horseshoes near homoclinic points	
6. Local smooth linearization and normal forms	278
Jets, formal power series, and smooth equivalence; General formal analysis; The hyperbolic smooth case	
7. TRANSVERSALITY AND GENERICITY	287
1. Generic properties of dynamical systems	287
Residual sets and sets of first category; Hyperbolicity and genericity	
2. Genericity of systems with hyperbolic periodic points	290
Transverse fixed points; The Kupka-Smale Theorem	
3. Nontransversality and bifurcations	298
Structurally stable bifurcations; Hopf bifurcations	
4. The theorem of Artin and Mazur	304
8. ORBIT GROWTH ARISING FROM TOPOLOGY	307
1. Topological and fundamental-group entropies	308
2. A survey of degree theory	310
Motivation; The degree of circle maps; Two definitions of degree for smooth maps; The topological definition of degree	
3. Degree and topological entropy	316
4. Index theory for an isolated fixed point	318
5. The role of smoothness: The Shub-Sullivan Theorem	323
6. The Lefschetz Fixed-Point Formula and applications	326
7. Nielsen theory and periodic points for toral maps	330
9. VARIATIONAL ASPECTS OF DYNAMICS	335
1. Critical points of functions, Morse theory, and dynamics	336
2. The billiard problem	339
3. Twist maps	349
Definition and examples; The generating function; Extensions; Birkhoff periodic orbits; Global minimality of Birkhoff periodic orbits	
4. Variational description of Lagrangian systems	365
5. Local theory and the exponential map	367
6. Minimal geodesics	372
7. Minimal geodesics on compact surfaces	376
Part 3 Low-dimensional phenomena	
10. INTRODUCTION: WHAT IS LOW-DIMENSIONAL DYNAMICS?	381
Motivation; The intermediate value property and conformality; Very low-dimensional and low-dimensional systems; Areas of low-dimensional dynamics	
11. HOMEOMORPHISMS OF THE CIRCLE	387
1. Rotation number	387

2. The Poincare classification	393
Rational rotation number; Irrational rotation number; Orbit types and measurable classification	
12. CIRCLE DIFFEOMORPHISMS	401
1. The Denjoy Theorem	401
2. The Denjoy example	403
3. Local analytic conjugacies for Diophantine rotation number	405
4. Invariant measures and regularity of conjugacies	410
5. An example with singular conjugacy	412
6. Fast-approximation methods	415
Conjugacies of intermediate regularity; Smooth cocycles with wild coboundaries	
7. Ergodicity with respect to Lebesgue measure	419
13. TWIST MAPS	423
1. The Regularity Lemma	424
2. Existence of Aubry-Mather sets and homoclinic orbits	425
Aubry-Mather sets; Invariant circles and regions of instability	
3. Action functionals, minimal and ordered orbits	434
Minimal action; Minimal orbits; Average action and minimal measures; Stable sets for Aubry-Mather sets	
4. Orbits homoclinic to Aubry-Mather sets	441
5. Nonexistence of invariant circles and localization of Aubry-Mather sets	447
14. FLOWS ON SURFACES AND RELATED DYNAMICAL SYSTEMS	451
1. Poincare-Bendixson theory	452
The Poincare-Bendixson Theorem; Existence of transversals	
2. Fixed-point-free flows on the torus	457
Global transversals; Area-preserving flows	
3. Minimal sets	460
4. New phenomena	464
The Cherry flow; Linear flow on the octagon	
5. Interval exchange transformations	470
Definitions and rigid intervals; Coding; Structure of orbit closures; Invariant measures; Minimal nonuniquely ergodic interval exchanges	
6. Application to flows and billiards	479
Classification of orbits; Parallel flows and billiards in polygons	
7. Generalizations of rotation number	483
Rotation vectors for flows on the torus; Asymptotic cycles; Fundamental class and smooth classification of area-preserving flows	
15. CONTINUOUS MAPS OF THE INTERVAL	489
1. Markov covers and partitions	489
2. Entropy, periodic orbits, and horseshoes	493
3. The Sharkovsky Theorem	500
4. Maps with zero topological entropy	505
5. The kneading theory	511
6. The tent model	514

16. SMOOTH MAPS OF THE INTERVAL	519
1. The structure of hyperbolic repellers	519
2. Hyperbolic sets for smooth maps	520
3. Continuity of entropy	525
4. Full families of unimodal maps	526
Part 4 Hyperbolic dynamical systems	
17. SURVEY OF EXAMPLES	531
1. The Smale attractor	532
2. The DA (derived from Anosov) map and the Plykin attractor The DA map; The Plykin attractor	537
3. Expanding maps and Anosov automorphisms of nilmanifolds	541
4. Definitions and basic properties of hyperbolic sets for flows	544
5. Geodesic flows on surfaces of constant negative curvature	549
6. Geodesic flows on compact Riemannian manifolds with negative sectional curvature	551
7. Geodesic flows on rank-one symmetric spaces	555
8. Hyperbolic Julia sets in the complex plane Rational maps of the Riemann sphere; Holomorphic dynamics	559
18. TOPOLOGICAL PROPERTIES OF HYPERBOLIC SETS	565
1. Shadowing of pseudo-orbits	565
2. Stability of hyperbolic sets and Markov approximation	571
3. Spectral decomposition and specification Spectral decomposition for maps; Spectral decomposition for flows; Specifica- tion	574
4. Local product structure	581
5. Density and growth of periodic orbits	583
6. Global classification of Anosov diffeomorphisms on tori	587
7. Markov partitions	591
19. METRIC STRUCTURE OF HYPERBOLIC SETS	597
1. Holder structures The invariant class of Holder-continuous functions; Holder continuity of conju- gacies; Holder continuity of orbit equivalence for flows; Holder continuity and differentiability of the unstable distribution; Holder continuity of the Jacobian	597
2. Cohomological equations over hyperbolic dynamical systems The Livschitz Theorem; Smooth invariant measures for Anosov diffeomor- phisms; Time change and orbit equivalence for hyperbolic flows; Equivalence of torus extensions	-608
20. EQUILIBRIUM STATES AND SMOOTH INVARIANT MEASURES	615
1. Bowen measure	615
2. Pressure and the variational principle	623
3. Uniqueness and classification of equilibrium states Uniqueness of equilibrium states; Classification of equilibrium states	628

4. Smooth invariant measures	637
Properties of smooth invariant measures; Smooth classification of Anosov diffeomorphisms on the torus; Smooth classification of contact Anosov flows on 3-manifolds	
5. Margulis measure	643
6. Multiplicative asymptotic for growth of periodic points	651
Local product flow boxes; The multiplicative asymptotic of orbit growth	
Supplement	
S. DYNAMICAL SYSTEMS WITH NONUNIFORMLY HYPERBOLIC BEHAVIOR BY ANATOLE KATOK AND LEONARDO MENDOZA	659
1. Introduction	659
2. Lyapunov exponents	660
Cocycles over dynamical systems; Examples of cocycles; The Multiplicative Ergodic Theorem; Osedelec-Pesin ϵ -Reduction Theorem; The Ruelle inequality	
3. Regular neighborhoods	672
Existence of regular neighborhoods; Hyperbolic points, admissible manifolds, and the graph transform	
4. Hyperbolic measures	678
Preliminaries; The Closing Lemma; The Shadowing Lemma; Pseudo-Markov covers; The Livschitz Theorem	
5. Entropy and dynamics of hyperbolic measures	693
Hyperbolic measures and hyperbolic periodic points; Continuous measures and transverse homoclinic points; The Spectral Decomposition Theorem; Entropy, horseshoes, and periodic points for hyperbolic measures	
Appendix	
A. BACKGROUND MATERIAL	703
1. Basic topology	703
Topological spaces; Homotopy theory; Metric spaces	
2. Functional analysis	711
3. Differentiable manifolds	715
Differentiable manifolds; Tensor bundles; Exterior calculus; Transversality	
4. Differential geometry	727
5. Topology and geometry of surfaces	730
6. Measure theory	731
Basic notions; Measure and topology	
7. Homology theory	735
8. Locally compact groups and Lie groups	738
NOTES	741
HINTS AND ANSWERS TO THE EXERCISES	765
REFERENCES	781
INDEX	793

Preface

The theory of dynamical systems is a major mathematical discipline closely intertwined with most of the main areas of mathematics. Its mathematical core is the study of the global orbit structure of maps and flows with emphasis on properties invariant under coordinate changes. Its concepts, methods, and paradigms greatly stimulate research in many sciences and have given rise to the vast new area of applied dynamics (also called nonlinear science or chaos theory). The field of dynamical systems comprises several major disciplines, but we are interested mainly in finite-dimensional differentiable dynamics. This theory is inseparably connected with several other areas, primarily ergodic theory, symbolic dynamics, and topological dynamics. So far there has been no account that treats differentiable dynamics from a sufficiently comprehensive point of view encompassing the relations with these areas. This book attempts to fill this gap. It provides a self-contained coherent comprehensive exposition of the fundamentals of the theory of smooth dynamical systems together with the related areas of other fields of dynamics as a core mathematical discipline while providing researchers interested in applications with fundamental tools and paradigms. It introduces and rigorously develops the central concepts and methods in dynamical systems and their applications to a wide variety of topics.

What this book contains. We begin with a detailed discussion of a series of elementary but fundamental examples. These are used to formulate the general program of the study of asymptotic properties as well as to introduce the principal notions (differentiable and topological equivalence, moduli, structural stability, asymptotic orbit growth, entropies, ergodicity, etc.) and, in a simplified way, a number of important methods (fixed-point methods, coding, KAM-type Newton method, local normal forms, homotopy trick, etc.).

The main theme of the second part is the interplay between local analysis near individual (e.g., periodic) orbits and the global complexity of the orbit structure. This is achieved by exploring hyperbolicity, transversality, global topological invariants, and variational methods. The methods include the study of stable and unstable manifolds, bifurcations, index and degree, and construction of orbits as minima and maximaxes of action functionals.

In the third and fourth parts the general program outlined in the first part is carried out to considerable depth for low-dimensional and hyperbolic dynamical systems which are particularly amenable to such analysis. Hyperbolic systems are the prime example of well-understood complexity. This manifests itself in an orbit structure that is rich both from the topological and statistical point of view and stable under perturbation. At the same time the principal features can be described qualitatively and quantitatively with great precision. In low-dimensional dynamical systems on the other hand there are two situations. In the “very low-dimensional” case the orbit structure is simplified and admits only a limited amount of complexity. In the “low-dimensional” case some complexity is possible, yet additional major aspects of the orbit structure can be understood via hyperbolicity or related types of behavior.

Although we develop most themes related to differentiable dynamics in some depth we have not tried to write an encyclopedia of differentiable dynamics. Even if this were possible, the resulting work would be strictly a reference source and not useful as an introduction or a text. Consequently we also do not strive to present the most definitive results available but rather to provide organizing principles for methods and results. This is also not a book on applied dynamics and the examples are not chosen from those models that are widely studied in various disciplines. Instead our examples arise naturally from the internal structure of the subject and contribute to its understanding. The emphasis placed on various areas in the field is not dictated by the relative amount of published work or research activity in those areas, but reflects our understanding of what is basic and fundamental in the subject. An obvious disparity appears in the area of one-dimensional (real and especially complex) dynamics, which witnessed a great surge of activity in the past 15 years producing a number of brilliant results. It plays a relatively modest role in this book. Real one-dimensional dynamics is used mainly as an easy model situation in which various methods can be applied with considerable success. Complex dynamics, which is in our view a fascinating but rather specialized area, appears only as a source of examples of hyperbolic sets. On the other hand we try to point out and emphasize the interactions of dynamics with other areas of mathematics (probability theory, algebraic and differential topology, geometry, calculus of variations, etc.) even in some situations where the current state of knowledge is somewhat tentative.

How to use this book. This book can be used both as a text for a course or for self-study and as a reference book. As a text it would most naturally be used as the primary source for graduate students with background equivalent to one year of graduate study at a major U.S. university who are interested in becoming specialists in dynamical systems or want to acquire solid general knowledge of the field. Some portions of this book do not assume as much background and can be used by advanced undergraduate students or graduate students in science and engineering who want to learn about the subject without becoming experts. Those portions include Chapter 1, most of Chapters 2, 3, and 5, parts of Chapters 4, 6, 8, and 9, Chapters 10 and 11, and most of Chapters 12, 14, 15, and 16. The 472 exercises are a very important part of the book. They fall into several categories. Some of them directly illustrate the use of results or methods from the text; others explore examples that are not discussed in the text or indicate further developments. Sometimes an important side topic is developed in a series of exercises. Those 317 that we do not consider routine have been provided with hints or brief solutions in the back of the book. An asterisk indicates our subjective assessment of higher difficulty, due to the need for either inventiveness or familiarity with material not obviously related to the subject at hand.

Each of the four parts of the book can be the basis of a course roughly at the second-year graduate level running one semester or longer. From this book one can tailor many courses dedicated to more specialized topics, such

as variational methods in classical mechanics, hyperbolic dynamical systems, twist maps and applications, an introduction to ergodic theory and smooth ergodic theory, and the mathematical theory of entropy. In order to assist both students and teachers in selecting material for a course we summarize the principal interrelations between the chapters in Figure F.1. A solid arrow $A \rightarrow B$ indicates that a major portion of the material from Chapter A is used in Chapter B (this relation is transitive). A dashed arrow $A \dashrightarrow B$ indicates that material from Chapter A is used in some parts of Chapter B.

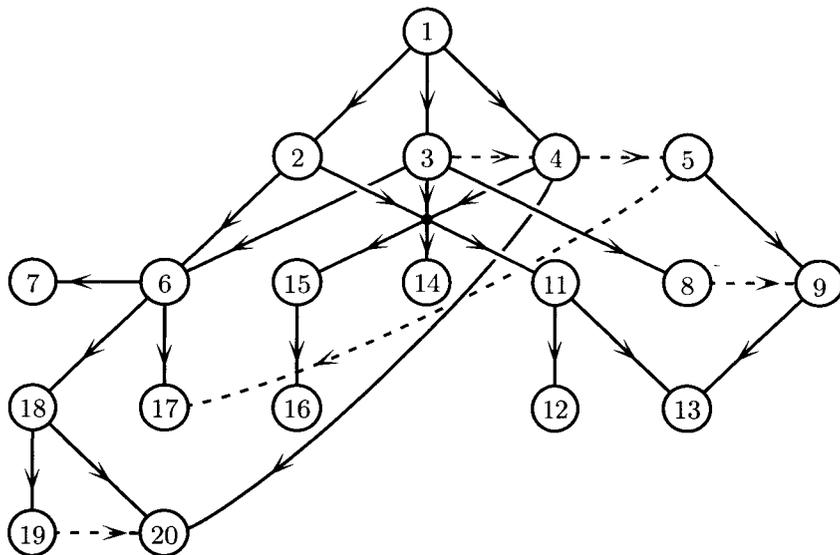


FIGURE F.1.

With the exception of Chapters 1–4, which form a common basis for the rest of the book, generally the material in the left part of the diagram deals with hyperbolic dynamics, that in the middle with low-dimensional dynamics, and that in the right with aspects of differentiable dynamics related to topology and classical mechanics.

There are various kinds of material used in this book. First of all we tacitly use, and assume familiarity with, the results of linear algebra (including Jordan normal forms), calculus of several variables, the basic theory of ordinary differential equations (including systems), elementary complex analysis, basic set theory, elementary Lebesgue integration, basic group theory, and some Fourier series. There is a next higher level of essential background material which is reviewed in the appendix. Most of the material in the appendix is of this nature, namely, the standard theory of topological, metric, and Banach spaces, elementary homotopy theory, the basic theory of differentiable manifolds including

vector fields, bundles and differential forms, and the definition and basic properties of Riemannian metrics. Some topics are used on isolated occasions only. This last level of material includes the basic topology and geometry of surfaces and the general theory of measures, σ -algebras, and Lebesgue spaces, homology theory, material related to Lie groups and symmetric spaces, curvature and connections on manifolds, transversality, and normal families of complex functions. Most, but not all, of this material is also reviewed in the appendix, usually in a less detailed fashion. Either such material can be taken on faith without loss to the application in the text, or otherwise the pertinent portion of the text can be skipped without great loss.

On several occasions we include an important background fact without proof in the text. This happens when a certain result is organically related to a particular section. The Lefschetz Fixed-Point Formula is a good example of such a result.

Sources. Most of the material in this book does not consist of original results. Nevertheless the presentation of most of the material is our own and consists of original or considerably modified proofs of known results, explanations of the structure and interconnectedness of the subject, and so forth. Some portions of the text, roughly a sixth of it, mostly in Parts 3 and 4, closely follow other published sources, the majority of these being original research articles. An outstanding example is the presentation of portions of the hyperbolic theory in Chapters 18 and 20 which was given such a clear treatment by R. Bowen in his articles in the seventies that it could hardly be improved. On several occasions we follow the exposition of a subject in existing books. With the exception of some basic subjects such as Hamiltonian formalism or variational calculus this occurs only in Part 3. The reason for this is that low-dimensional dynamics has a much better developed expository literature than the field as a whole. We acknowledge all borrowings of proofs and presentations of which we are aware in the notes near the end of the book.

Since we aim to present the subject by developing it from first principles and in a self-contained way, rather than to give an exhaustive account of the development and current state of the field, we do not attempt a comprehensive listing of relevant references which would easily increase our bibliography to a thousand items or more. In particular, not all theorems are attributed to the original authors, especially if the results are part of the broader developments in the field, rather than landmark results or of a rather special nature. Most of the attributions are relegated to the notes. These consist of general comments arranged by section and some numbered remarks to particular points in the main text. Furthermore, in order not to interrupt the logical flow of the text all bibliographical references for the main part of the book are also relegated to the notes. Our historical comments, in both the introduction and the notes, do not aim to present a coherent account of the development of the subject, but to subjectively select some of its major moments.

We have included several types of literature in the bibliography. First, we have tried to list all major monographs and representative textbooks and sur-

veys covering the principal branches of dynamics. Next there are landmark papers that introduce and develop various branches of our subject, define principal notions, or contain proofs of major results. We try to list all the sources on which the presentation in various parts of the book is based, or that inspired our presentation in other places and many (but not all) of the original sources for specific results presented in the text. Finally there is a sample of references to important work, both original and surveys, in some areas touched upon but not treated in the text. According to our principle of selecting models for their intrinsic interest rather than their value for concrete scientific problems we omit works by nonmathematicians (even important ones) that are dedicated to the study of models motivated by scientific problems so long as these contain only hypotheses and numerical results. References to such works are widely available in many of the books and surveys that we quote.

History and acknowledgments. The general idea of writing a broad introduction to the theory of dynamical systems first occurred to the first author when he taught a graduate course at the California Institute of Technology in 1984–5. This course resulted in two sets of lecture notes prepared by the second author and by his fellow graduate student John Lindner to whom we are deeply grateful. The key idea of introducing the principal notions and methods via a presentation of a series of basic examples crystallized when the first author was preparing and teaching an intensive four-week course in July 1986 at the Summer Mathematics Institute for graduate students at Fudan University in Shanghai. The summary and notes from that course became the germ for major parts of Chapters 1–4. Further progress was made during another graduate course at the California Institute of Technology in 1986–7 after which it became clear that the original project of a book of 300–350 pages would result in too sketchy and incomplete an account of the subject. In the summer of 1989 we developed a detailed plan of the book which has been carried out with some substantial later modification. Another graduate course during the first author’s first year at the Pennsylvania State University (1990–1) helped to test some existing parts of the book and develop some new material.

We feel deep gratitude to the California Institute of Technology, Tufts University, and the Pennsylvania State University for providing excellent working conditions and supporting several mutual visits. Special thanks are due to the Mathematical Sciences Research Institute in Berkeley, California, where we worked together on major portions of the book in the summer of 1992. During this period our project was transformed from a collection of drafts into an incomplete but coherent product.

We also owe thanks and gratitude to numerous individuals for providing various kinds of help and inspiration during this project. We apologize for any omissions of people whose comments and suggestions may have been incorporated and forgotten.

Jessica Madow, the technical typist at the California Institute of Technology, typed major portions of the then existing manuscript in Exp. Kathy Wyland and Pat Snare at the Pennsylvania State University typed the first drafts of

many chapters in $\text{T}_{\text{E}}\text{X}$. Several people helped with computer support or typesetting advice. David Glaubman at the Mathematical Sciences Research Institute was very helpful, Michael Downes at the Technical Support Department of the American Mathematical Society helped to make the running heads come out right on every single page, and our colleague Uwe Schmock at the ETH wrote the overbar macro for $\text{T}_{\text{E}}\text{X}$ and made other useful comments. Boris Katok made the majority of the illustrations for the book. Bill Schlesinger gave us the initial tutoring that enabled us to make numerous pictures using Matlab. We are deeply grateful to the editors of Cambridge University Press: David Tranah who encouraged and prodded us during the earlier stages of the project and Lauren Cowles who patiently guided us through the process of finishing the book and getting it ready for production. This book was typeset in $\text{T}_{\text{E}}\text{X}$ using $\mathcal{A}\mathcal{M}\mathcal{S}\text{-T}_{\text{E}}\text{X}$, the $\text{T}_{\text{E}}\text{X}$ macro package of the American Mathematical Society.

Viorel Nițică and Alexej Kononenko wrote solutions to the majority of the exercises. Their work helped to correct some flawed exercises, and we used their solutions to write many of our hints.

The following people made numerous suggestions, including pointing out mathematical and stylistical errors, misprints, and the need for better explanations: The greatest amount of this kind of help came from Howie Weiss at the Pennsylvania State University. Further comments were given to us by Luis Barreira, Misha Brin, Mirko Degli-Esposti, David DeLatte, Serge Ferleger, Eugene Gutkin, Moisey Guysinsky, Miaohua Jiang, Tasso Kaper, Alexej Kononenko, Viorel Nițică, Ralf Spatzier, Garrett Stuck, Andrew Török, and Chengbo Yue.

In particular Howie Weiss, Tasso Kaper, Garrett Stuck, Ralf Spatzier, and Misha Brin taught from parts of the book and were very helpful in polishing it.

We had fruitful discussions with Michael Jakobson, Welington de Melo, Mikhael Lyubich, and Zbigniew Nitecki concerning one-dimensional maps and with Eduard Zehnder on variational methods. These were useful in crystallizing the content and presentation of those respective chapters. Gene Wayne helped by providing references concerning infinite-dimensional dynamical systems and Mike Boyle gave some useful guidance for sources in symbolic dynamics.

A number of corrections were made between printings. We would like to thank Luis Barreira, Marlies Gerber, Karl Friedrich Siburg, Garrett Stuck and Andrew Török who pointed out many small errors. Peter Walters found inaccuracies in Lemma 4.5.2 and Lemma 20.2.3. Robert McKay pointed out that some results of Section 14.2 needed recurrence hypotheses and Jonathan Robbins noted problems with the first version of Step 5 in the proof of the Hadamard-Perron Theorem 6.2.8. Tim Hunt corrected the DA construction. Corrections are listed at <http://www.tufts.edu/~bhasselb/thebook.html>.

A serious omission survived three printings: Section 20.6 is entirely due to Charles Toll, an attribution we inadvertently failed to make. Our sincere apologies.

Last, and most importantly, we wish to thank Svetlana Katok and Kathleen Hasselblatt for constant support and inspiration.

Introduction

1. Principal branches of dynamics

The most general and somewhat vague notion of a dynamical system includes the following ingredients:

(i). A “phase space” X , whose elements or “points” represent possible states of the system.

(ii). “Time”, which may be discrete or continuous. It may extend either only into the future (irreversible or noninvertible processes) or into the past as well as the future (reversible or invertible processes). The sequence of time moments for a reversible discrete-time process is in a natural correspondence to the set of all integers; irreversibility corresponds to considering only nonnegative integers. Similarly, for a continuous-time process, time is represented by the set of all real numbers in the reversible case and by the set of nonnegative real numbers for the irreversible case.

(iii). The time-evolution law. In the most general setting this is a rule that allows us to determine the state of the system at each moment of time t from its states at all previous times. Thus, the most general time-evolution law is time dependent and has infinite memory. In the course of this book, however, we will consider only those evolution laws that allow us to define all future (and for reversible systems also past) states given a state at any particular moment. Furthermore we will assume that the law of time evolution itself does not change with time. In other words, the result of time evolution will depend only on the initial position of the system and on the length of the evolution but not on the moment when the state of the system was initially registered. Thus, if our system was initially at a state $x \in X$, it will find itself after time t at a new state, which is uniquely determined by x and t , and thus can be denoted by $F(x, t)$. Fixing t , we obtain a transformation $\varphi^t: x \mapsto F(x, t)$ of the phase space into itself. These transformations for different t are related to each other. Namely, the evolution of the state x for time $s + t$ can be accomplished by first applying

the transformation φ^t to x and then by applying φ^s to the new state $\varphi^t(x)$. Thus, we have $F(x, t+s) = F(\varphi^t(x), s)$ or equivalently, the transformation φ^{t+s} is equal to the composition of φ^t and φ^s . In other words, the transformations φ^t form a semigroup. For a reversible system the transformations φ^t are defined for both positive and negative values of t and each φ^t is invertible. Thus, a reversible discrete-time dynamical system is represented by a cyclic group $\{F^n = (\varphi^1)^n \mid n \in \mathbb{Z}\}$ of one-to-one transformations of the phase space onto itself. Similarly, a reversible continuous-time dynamical system determines a one-parameter group $\{\varphi^t \mid t \in \mathbb{R}\}$ of one-to-one transformations of X onto itself.

The most characteristic feature of dynamical theories, which distinguishes them from other areas of mathematics dealing with groups of automorphisms of various mathematical structures, is the emphasis on asymptotic behavior, especially in the presence of nontrivial recurrence, that is, properties related with the behavior as time goes to infinity. The best way to explain what significant asymptotic properties are is to examine specific examples of dynamical systems and to determine the most characteristic features of their behavior. We will do that in Chapter 1 and then we will summarize some of our findings and present a list of interesting properties in Sections 3.1, 3.3, 4.1, 4.2d, and 4.3. This summary is preceded by an examination of natural equivalence relations for dynamical systems in Chapter 2 which sets the stage for treating asymptotic properties as invariants of those equivalence relations.

Historically, smooth continuous-time dynamical systems appeared first because of Newton's discovery that the motions of mechanical objects can be described by second-order ordinary differential equations. More generally, many other natural and social phenomena, such as radioactive decay, chemical reactions, population growth, or dynamics of prices on the market, may be modeled with various degrees of accuracy by systems of ordinary differential equations. These situations fit into the domain of our investigation if there is no explicit time-dependence in the coefficients and right-hand parts of the equations.

In virtually all situations of interest the phase space of a dynamical system possesses a certain structure which the evolution law respects. Different structures give rise to theories dealing with dynamical systems that preserve those structures. Let us mention the most important of those theories.

1. Ergodic theory. Here the phase space X is a "good" measure space, that is, a Lebesgue space (cf. Section 6 of the Appendix) with a finite or σ -finite measure μ . We can consider as a structure in X either the measure μ itself or its equivalence class which is determined by the collection of all sets of measure zero. Accordingly, ergodic theory concerns groups or semigroups of measurable transformations of X that either preserve μ or transform it into an equivalent measure. In the latter case the measure μ is called *quasi-invariant*. In this book ergodic theory plays an important but auxiliary role. It provides the appropriate paradigms and tools for studying asymptotic distribution and statistical behavior of orbits for smooth dynamical systems. Some central concepts and results of ergodic theory are introduced and discussed in Chapter 4.

The origins of ergodic theory go back to the famous ergodic hypothesis of Boltzmann who postulated equality of time averages and space averages for systems in statistical mechanics. Within mathematics the notions of ergodic theory arose from the study of uniform distributions of sequences. The Kronecker–Weyl Equidistribution Theorem (Proposition 4.2.1) is an early example of such a result. H. Poincaré observed that the preservation of a finite invariant measure forces strong conclusions about recurrence which are encapsulated in his Recurrence Theorem (Theorem 4.1.19). The systematic development of ergodic theory as a mathematical subject started around 1930 by von Neumann who looked at the subject primarily from a functional-analytic viewpoint. Among the early major contributors to the subject were G. D. Birkhoff, E. Hopf, and S. Kakutani. The critical point in the development of ergodic theory which forever changed the emphasis from the functional-analytic to the probabilistic and later geometric and combinatorial viewpoints was the introduction of entropy by A. Kolmogorov around 1958. It built upon C. Shannon’s seminal development of information theory which was given the appropriate mathematical treatment by A. Khinchin. Kolmogorov’s work was quickly followed by the development of an entropy theory based on the probabilistic viewpoint primarily by Y. Sinai and V. Rokhlin which culminated in Sinai’s weak isomorphism theorem. The next crucial juncture was the first proof of the isomorphism of Bernoulli shifts of equal entropy which was obtained by D. Ornstein via combinatorial constructions. This work was followed by the development of the isomorphism theory which in particular gave necessary and sufficient conditions for metric isomorphism to a Bernoulli shift. Among later major developments one should note the Kakutani (monotone) equivalence theory, H. Furstenberg’s theory of multiple recurrence, and the finitary isomorphism theory.

2. Topological dynamics. The phase space in this theory is a good topological space, usually a metrizable compact or locally compact space (see Section 1 of the Appendix). Topological dynamics concerns itself with groups of homeomorphisms and semigroups of continuous transformations of such spaces. Sometimes these objects are called topological dynamical systems. Similarly to the case of ergodic theory we use in this book notions and results from topological dynamics primarily as a framework and a tool for studying smooth dynamical systems. Though we are not making an attempt to provide a comprehensive introduction to the field, a fair amount of material from topological dynamics appears in this book, beginning with our first survey of examples in Chapter 1 and then in Chapter 3. Sections 4.1, 4.5 and later 20.1 and 20.2 provide crucial links between topological dynamics and ergodic theory. Some material in Chapter 8 (for example, Theorem 8.3.1) as well as all of Chapters 11 and 15 deal with particular classes of dynamical systems without any differentiability assumptions and thus belong to topological dynamics.

Topological dynamics was founded by Poincaré when he introduced the idea of qualitatively describing the solutions of differential equations that could not be solved analytically. One of his early achievements was the classification of circle maps (Theorem 11.2.7). M. Morse and G. D. Birkhoff made major

contributions to topological dynamics in the process of trying to understand more classical systems (behavior of geodesics and Hamiltonian systems). Later a more intrinsic approach was developed by G. Hedlund, J. Oxtoby, and others. An important subject in topological dynamics is H. Furstenberg's theory of distal extensions which was further developed by R. Ellis.

3. The theory of smooth dynamical systems or differentiable dynamics. As the name suggests, the phase space here possesses the structure of a smooth manifold, for example, a domain or a closed surface in a Euclidean space (see Section 3 of the Appendix for a more detailed description). This theory, which is the prime subject of this book, concerns diffeomorphisms and flows (smooth one-parameter groups of diffeomorphisms) of such manifolds and iterates of noninvertible differentiable maps. In this book we will deal mostly with finite-dimensional situations. Interest in infinite-dimensional dynamical systems has been growing steadily during the past two decades, to a large extent stimulated by problems in fluid dynamics, statistical mechanics, and other fields of mathematical physics. Several directions in infinite-dimensional dynamics have been developed to a considerable extent starting from analogies with various branches in finite-dimensional dynamics.

Since a finite-dimensional smooth manifold possesses a natural locally compact topology, the theory of smooth dynamical systems naturally draws upon notions and results from topological dynamics. Another deeper reason for these interrelations arises from the fact that in dealing with asymptotic behavior of smooth dynamical systems one is likely to encounter very complicated non-smooth phenomena, which in other contexts would be dismissed as pathological. In particular, some important invariant sets for a smooth system, for example, attractors (Definition 3.3.1), may not have any smooth structure and consequently, such sets should be studied from a different, nonsmooth, point of view. *Symbolic dynamics*, the study of a specific class of topological dynamical systems which occur as closed invariant subsets of the shift transformation in a sequence space (cf. Section 1.9), is particularly important in that respect. For further motivation of the relationships between topological and smooth dynamics see Section 2.3.

Relations with ergodic theory are also intimate, both because invariant measures provide a powerful tool for the study of asymptotic properties of smooth dynamical systems and because the smooth structure on a finite-dimensional manifold determines a natural class of quasi-invariant measures for differentiable dynamical systems (see Section 5.1).

Sometimes the part of the theory of smooth dynamical systems that concerns measure-theoretic properties of such systems is given the separate name *smooth ergodic theory*. One might also say that smooth ergodic theory is the study of automorphisms of a composite structure formed by a smooth manifold and a reasonable measure on it. Chapter 20 and the Supplement are dedicated to this subject. A number of results belonging to smooth ergodic theory are scattered among the earlier chapters.

Poincaré is also the father of differentiable dynamics. His main contribution was to emphasize the qualitative approach as opposed to the traditional emphasis on explicit solutions of differential equations of mechanics. His other achievement was the founding of the local theory of maps and vector fields near fixed and periodic orbits (cf. Sections 2.1, 6.3, 6.6). Other principal figures in the early stages of the field were A. M. Lyapunov and J. Hadamard who introduced various concepts of stability and developed major analytic tools (for example, the Hadamard–Perron Theorem 6.2.8). Part of Poincaré’s program was carried out by G. D. Birkhoff who proved, among other things, Poincaré’s celebrated “Last Geometric Theorem” which gives a mechanism responsible for dynamical complexity in mechanical systems with two degrees of freedom. Another aspect of Poincaré’s program was developed by A. Denjoy who introduced some key new ideas in the process of completing Poincaré’s theory of circle maps and flows on the two-dimensional torus. Symbolic dynamics appeared as a very useful tool beginning with a seminal paper by E. Artin and it was greatly developed by Morse and Hedlund. E. Hopf was the first to realize that hyperbolicity is a key mechanism that produces complicated behavior in nonlinear dynamical systems. His proof of ergodicity of the geodesic flow of surfaces of negative curvature can be viewed as the first major result in smooth ergodic theory.

Another principal root of the modern global approach to the study of smooth dynamical systems was the notion of structural stability which was first introduced by A. Andronov and L. S. Pontryagin in the study of flows on surfaces and later developed in that setting by Peixoto. It was given a second life by Smale who discovered that systems with complicated orbit behavior (the “horseshoe”, Section 2.5) can be structurally stable. Subsequently Smale, Anosov, Sinai, and Bowen developed the core of the theory of hyperbolic dynamical systems. They greatly developed methods from ergodic theory and topological dynamics due to Hopf and Hedlund as well as more classical ideas going back to Hadamard, Perron, and Lyapunov. Identifying a certain hyperbolicity as sufficient (J. Robbin, C. Robinson) and necessary (R. Mañé) for structural stability was one of the crowning achievements of the theory of smooth dynamical systems. A major impetus to smooth ergodic theory was given by D. Ruelle and Y. Sinai who introduced ideas and methods from statistical mechanics to the theory of smooth dynamical systems. The next important step was made by Y. B. Pesin who developed the general structural theory of smooth measure-preserving systems based on the concept of nonuniform hyperbolicity. We should also mention the work of M. Herman and J.-C. Yoccoz on smooth classification of circle diffeomorphisms and the work of D. V. Anosov and A. Katok on constructions of smooth dynamical systems with various often unexpected properties.

4. Hamiltonian or symplectic dynamics. This theory is a natural generalization of a study of differential equations of classical mechanics. The phase space here is an even-dimensional smooth manifold with a nondegenerate closed differential 2-form Ω . One-parameter groups of Ω -preserving diffeomorphisms correspond to Hamiltonian differential equations in classical mechanics. An individual Ω -preserving diffeomorphism generalizes the notion of a canonical

transformation. We first encounter such systems in Section 1.5 and return to this field in a more systematic way in Section 5.5.

The origin of Hamiltonian dynamics as an object of study from the point of view of dynamical systems is largely in the questions of celestial mechanics. Again Poincaré introduced the fundamental approach of the qualitative study of the n -body problem. Later two distinct directions of study emerged: (i) the investigation of dynamical complexity in this problem due to some hyperbolicity (Aleksëev, Conley) and (ii) the study of integrable systems and their perturbations which led to the KAM theory. Though both the hyperbolic and integrable paradigm were available since Poincaré, it was Kolmogorov's profound contribution to realize that many qualitative features of (the very exceptional) integrable systems persist to some extent under perturbations and appear also in generic situations (for example, near an elliptic fixed point). Both of these lines of thought were influenced by the question of the stability of the solar system which was addressed by the hyperbolic approach in terms of the stability of an n -body system and by the KAM approach by considering perturbations, for example, of the (integrable) central force problem without interactions between planets. The work of Conley and Zehnder established a synthesis of topological and variational methods which became the cornerstone of modern global symplectic geometry. A renaissance of the study of completely integrable systems started with a seminal paper by Gardner, Greene, Kruskal, and Miura and the discovery by P. Lax of new mechanisms for producing integrable systems. It led both to a proliferation of new interesting examples of finite-dimensional integrable systems as well as to the theory of infinite-dimensional Hamiltonian systems whose applications to nonlinear partial differential equations were a major breakthrough by providing for the first time means for a complete qualitative analysis in situations other than those with the most simple asymptotic behavior.

2. Flows, vector fields, differential equations

The description of a dynamical system is somewhat easier when time is discrete, because the map generating a discrete-time system often can be given explicitly, usually by means of some formulas. In contrast, a continuous-time dynamical system is usually given infinitesimally (for example, by means of differential equations) and the reconstruction of the dynamics from this infinitesimal description involves some kind of integration process. In this and the next section we will very briefly discuss this local (in time) aspect of the theory of continuous-time dynamical systems and some simple relations between the discrete-time and the continuous-time situations.

We assume that the phase space is a smooth manifold of dimension m which we will usually denote by M , and thus our time evolution is given by a smooth function $F(x, t) = \varphi^t(x)$, $x \in M$, $t \in \mathbb{R}$, which satisfies the group (composition) property $\varphi^t \circ \varphi^s = \varphi^{t+s}$ and may or may not be defined for all x and t . Let us consider first the local aspect of the situation. When we fix $x \in M$ and vary t we obtain a parameterized smooth curve on M . Let $\xi(x)$ be the tangent vector

to this curve at $t = 0$, that is, at the point x . Properly speaking, the vector $\xi(x)$ belongs to the tangent space $T_x M$ which is an m -dimensional linear space “attached” to M at the point x . The map $x \mapsto \xi(x)$ forms a section of the tangent bundle $TM = \bigcup_{x \in M} T_x M$ or a *vector field* on M (see Section 3 of the Appendix for more details). Of course, the local version of this construction is familiar to everybody who completed a standard course of advanced calculus. Namely, let $U \subset M$ be a coordinate neighborhood with coordinates (s_1, \dots, s_m) . Then the tangent bundle TU is simply a direct product $U \times \mathbb{R}^m$ and a vector field is determined by a map from U to \mathbb{R}^m , that is, by m real-valued functions v_1, \dots, v_m , as follows. Denoting by $\frac{\partial}{\partial s_i}$ the basic vector fields which associate to every point the i th vector of the standard basis in \mathbb{R}^m we can represent every vector field locally as $\sum_{i=1}^m v_i(s_1, \dots, s_m) \frac{\partial}{\partial s_i}$. If our initial point x is represented by coordinates s_1^0, \dots, s_m^0 then the evolution of this point is obtained by solving the system of first-order ordinary differential equations

$$\frac{ds_i}{dt} = v_i(s_1, \dots, s_m)$$

with initial conditions $s_i(0) = s_i^0$ for $i = 1, \dots, m$.

We know from the standard theory of ordinary differential equations that under very moderate smoothness assumptions, for example, if the functions v_i are continuously differentiable, the solution for sufficiently small time exists, is unique, and depends smoothly on the initial condition.

Thus, at least for small values of t , the transformation φ^t can be recovered from the vector field. For larger t one should take compositions of maps defined in local coordinates. If solutions exist for all real values of t , the vector field is called *complete*. We should keep in mind that on a manifold we have to work in different local coordinate systems if t is large, but this does not present any difficulties. If the manifold M is compact and has no boundary then it can be covered by a finite number of coordinate charts. Inside any chart the solutions exist for a fixed length of time. Since every point $x \in M$ belongs to a coordinate neighborhood which is not very small, this implies that any C^1 vector field on a closed compact manifold without boundary is complete and thus defines a *smooth flow*, that is, a one-parameter group of diffeomorphisms of M .

This is one of the reasons why we will often prefer to consider dynamical systems on compact manifolds. This preference will not be universal because in many situations such as local and semilocal problems (cf. Section 0.4 and Chapter 6) or systems of differential equations associated to many concrete mechanical and other problems, this assumption would be too restrictive.

Exercise

0.2.1. Show (in detail) that a smooth vector field on a compact manifold is complete.

3. Time-one map, section, suspension

There are several useful relations between continuous-time and discrete-time dynamical systems.

The most obvious way to associate a discrete-time system to a flow $\{\varphi^t\}_{t \in \mathbb{R}}$ is to take the iterates of the map φ^{t_0} for some value of t_0 , say, $t_0 = 1$. However, only very few diffeomorphisms may be obtained that way. For example, let $f = \varphi^{t_0}$ and assume that $f^k(x) = x$, where $k > 1$, but $f(x) \neq x$ so that the orbit of x is periodic, but not fixed. But then for every $t \in \mathbb{R}$

$$f^k \varphi^t(x) = \varphi^{kt_0+t}(x) = \varphi^t(\varphi^{kt_0}(x)) = \varphi^t(f^k(x)) = \varphi^t(x).$$

Hence every point $\varphi^t(x)$ is also a periodic point of period k for f . Thus if f has an isolated periodic point of period greater than one, the map f cannot be obtained as the time- t map of any flow.

Another more local but also more useful method is the construction of a *Poincaré (first-return) map*. Let us take a point $x \in M$ such that $\xi(x) \neq 0$ and an $(m-1)$ -dimensional (codimension-one) submanifold N containing x and transversal to the vector field. The latter property simply means that for every point $y \in N$ the vector $\xi(y)$ is not tangent to N . If we assume that the point x is periodic for the flow, that is, $\varphi^{t_0}(x) = x$ for some $t_0 > 0$, then every nearby orbit of the flow intersects the surface N at a time close to t_0 so we have defined for a neighborhood U of x on N a map $F_N: U \rightarrow N$ such that $F_N(x) = x$. This map is called a *section map* or *first-return map* or *Poincaré map* for the flow. This construction (also called inducing) also works if x is not periodic but comes sufficiently close to itself (see below).

Finally, for any diffeomorphism $f: M \rightarrow M$ one can construct a *suspension flow* on the *suspension manifold* M_f which is obtained from the direct product $M \times [0, 1]$ by identifying pairs of points of the form $(x, 1)$ and $(f(x), 0)$ for $x \in M$. The *suspension flow* σ_f^t is determined by the “vertical” vector field $\frac{\partial}{\partial t}$ on M_f .

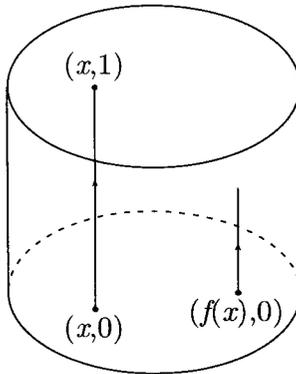


FIGURE 0.3.1. A suspension

This construction is closely related to the solution of a system of ordinary differential equations with periodic coefficients. First recall that a time-dependent system of ordinary differential equations is given by a family of vector fields v_t and thus determines a family of time evolutions $\varphi^{t,s}$ from the moment t to the moment $t + s$ which is not a group. It can, however, be interpreted as a single vector field $w(x, t) = v_t(x) + \frac{\partial}{\partial t}$ in the extended phase space $M \times \mathbb{R}$. The time evolution $\Phi^s(x, t) = (\varphi^{t,s}(x), t + s)$ in $M \times \mathbb{R}$ does have the group property. Of course, the space $M \times \mathbb{R}$ is never compact.

The situation changes, however, if the system of ordinary differential equations is periodic in time with period τ , say. Then $v_{t+\tau} = v_t$ and $\varphi^{t+k\tau,s} = \varphi^{t,s}$ for $k \in \mathbb{Z}$. In this case one can reduce the time evolution in $M \times \mathbb{R}$ to one in a factor space by identifying (x, t) with $(x, t + \tau)$. The factor space thus obtained is compact if M is compact and the projection of the flow Φ^s to that space is diffeomorphic to the suspension flow over the map $\varphi^{0,\tau}$ by the map $h: (\varphi^{0,t}(x), t) \mapsto (\varphi^{0,\tau}(x), t)$ ($0 \leq t \leq \tau$) to $M_{\varphi^{0,\tau}}$.

The suspension construction is generalized to the construction of *the flow under a function* or *special flow*. Namely, add to our data a smooth positive function φ on M and consider the manifold $M_{f,\varphi}$ obtained from the subset $M_\varphi = \{(x, t) \mid x \in M, t \in \mathbb{R}, 0 \leq t \leq \varphi(x)\}$ of the direct product $M \times \mathbb{R}$ by identifying pairs $(x, \varphi(x))$ and $(f(x), 0)$. Of course, topologically $M_{f,\varphi}$ is the same as M_f , but the “vertical” vector field $\frac{\partial}{\partial t}$ on $M_{f,\varphi}$ determines a new flow $\sigma_{f,\varphi}^t$ (the special flow under φ built over f) which differs from the suspension by a time change (see Definition 2.2.3).

Exercises

0.3.1. Let $M = [0, 1]$ and $f(x) = 1 - x$. Show that the manifold M_f is homeomorphic to the Möbius strip. The suspension flow has one orbit of period one and infinitely many orbits of period two. Show that the period-one orbit does not separate M_f and that any period-two orbit except the one that forms the boundary separates it into two pieces, one homeomorphic to the Möbius strip and the other to the cylinder $[0, 1] \times S^1$.

0.3.2. Let $M = S^1 = \{z \in \mathbb{C} \mid |z| = 1\}$ and $f(z) = -z$. Show that the manifold M_f is homeomorphic to the two-torus $\mathbb{T}^2 = S^1 \times S^1$.

0.3.3. Let $M = S^1$ and $f(z) = \bar{z}$. Show that the manifold M_f is homeomorphic to the Klein bottle. The suspension flow has two orbits of period one and infinitely many orbits of period two. Show that period-one orbits do not separate M_f and that any period-two orbit separates it into two pieces homeomorphic to the Möbius strip.

0.3.4. Describe the smooth structure on the suspension manifold M_f and, more generally, on the manifold $M_{f,\varphi}$.

4. Linearization and localization

We will see in the next three chapters that a large number of useful concepts related to the asymptotic behavior of smooth dynamical systems in fact belong to topological dynamics, that is, they are defined only in terms of topology, not the differentiable structure. We already mentioned some reasons for that in Section 0.1. Now we would like to point out some specific features that distinguish the theory of smooth dynamical systems from topological dynamics.

Already in elementary calculus one learns how useful it is to represent a function $\varphi(t)$ of one real variable t near a point t_0 as the main linear part $\varphi(t_0) + \varphi'(t_0)(t - t_0)$ plus an “infinitesimal of higher order”, $o(t - t_0)$. A less elementary version of the same idea plays a central role in the theory of smooth dynamical systems. First, if $U \subset \mathbb{R}^m$ is an open neighborhood of x_0 and $f: U \rightarrow \mathbb{R}^m$ is a differentiable map, we can represent f near the point x_0 as the constant part $f(x_0)$ plus the linear part $Df_{x_0}(x - x_0)$ plus higher-order terms. The differential Df is a linear operator in \mathbb{R}^n that is represented in coordinate form by the matrix of partial derivatives. If $f(t_1, \dots, t_m) = (f_1(t_1, \dots, t_m), \dots, f_m(t_1, \dots, t_m))$, then

$$Df_{x_0}(t_1, \dots, t_m) = \left(\frac{\partial f_i}{\partial t_j} \right)_{i,j=1,\dots,m},$$

where the partial derivatives are calculated at the values of the coordinates corresponding to the point x_0 . If the map is regular at x_0 this operator is invertible.

The picture remains essentially the same for differentiable maps of smooth manifolds with the only difference that instead of the standard coordinate system in \mathbb{R}^m one should use appropriate local coordinate systems near a point and its image. A more invariant way to express the same idea is to describe the differential Df_{x_0} of the map $f: M \rightarrow M$ as a linear map of the tangent space $T_{x_0}M$ into the space $T_{f(x_0)}M$. If f is a diffeomorphism then the differential is invertible. This construction can be globalized by considering the tangent bundle $TM = \bigcup_{x \in M} T_x M$ which can be provided with the structure of a differentiable manifold whose dimension is twice the dimension of M (see Section 3 of the Appendix). Any local coordinate system on M induces a coordinate system in TM which is global in the tangent direction. Namely, tangent vectors to the coordinate curves form a basis in each tangent space and the $2n$ coordinates of a tangent vector include n coordinates of its base point plus the coordinates of the vector with respect to that basis.

When we consider iterates of a map f , the differential $Df_x^n: T_x M \rightarrow T_{f^n(x)} M$ of the n th iterate is a composition of the differentials $Df_{f^i(x)}: T_{f^i(x)} \rightarrow T_{f^{i+1}(x)}$, $i = 0, \dots, n - 1$:

$$T_x M \xrightarrow{Df_x} T_{f(x)} M \xrightarrow{Df_{f(x)}} T_{f^2(x)} M \xrightarrow{Df_{f^2(x)}} \dots \xrightarrow{Df_{f^{n-1}(x)}} T_{f^n(x)} M.$$

$\underbrace{\hspace{15em}}_{Df_x^n}$

In this localized picture the asymptotic properties of f correspond to the properties of products of linear maps thus obtained, when the number of factors goes to infinity. Once the behavior of such products is understood, the question arises as to what extent this behavior reflects the properties of the original nonlinear system. The crucial point here is that the differential at any given point approximates well the behavior of points near to the point at which the differential has been calculated. The quality of this approximation depends on the nonlinear terms, for example, on the size of the second derivatives of the function representing our map in a neighborhood of the original point. When we pass to the iterates of a map, as a rule, the size of second derivatives grows (by the chain rule), so, a priori, the quality of the linear approximation should deteriorate. Under certain conditions the influence of nonlinear terms can be controlled, so that we obtain a picture of the behavior of those orbits that stay near the original orbit for sufficiently long time. Considerations of this kind represent the content of what is usually called the local analysis of smooth dynamical systems. This is the central theme of Chapter 6 and of the first three sections of the Supplement.

An ideal setting for the local approach appears when the original orbit is periodic, say, $f^n(x_0) = x_0$. Then the sequence of differentials is also periodic and the main role in the local analysis is played by the iterates of a single linear operator, $Df_{x_0}^n$, which represents the infinitesimal behavior of nearby orbits for the period. In particular, the eigenvalues of that operator play a crucial role in the local analysis near the point x_0 . See Section 1.2 for an analysis of linear maps and Sections 6.3 and 6.6 for a local analysis of nonlinear maps near a periodic point. For continuous-time dynamical systems the role of the differential is played by the variational equation whose right-hand side represents the infinitesimal generator for the one-parameter group of differentials of the maps forming the flow.

Though the local analysis concerns itself with the relative behavior of nearby orbits or, in the case of a neighborhood of a periodic orbit, with the behavior of orbits or orbit segments as long as they stay near the periodic orbit, the main goal of the theory of smooth dynamical systems is to understand the global behavior of nonlinear maps. Sometimes local analysis plays a crucial role in the global consideration. This happens, for example, if a periodic point represents an attractor, that is, if neighboring orbits approach it asymptotically with time (cf. Sections 1.1 and 3.3). More generally, we may try to localize certain parts of the phase space that play a particularly important role for the asymptotic behavior and to study the orbits inside and nearby this part. It is also possible that the behavior of orbits with certain initial conditions is particularly important due to the nature of a particular scientific problem which is represented by the dynamical system.

All this reasoning leads to a “semilocal” approach which lies between the local analysis and the global study of a system as a whole. Namely, let us assume that M is a smooth manifold (not necessarily compact), $U \subset M$ is an open subset of M , and $\Lambda \subset U$ is a compact set. Let furthermore $f: U \rightarrow M$ be

a smooth map that leaves Λ invariant. We may be interested in the study of orbits of f that lie inside Λ or stay nearby. The local tool of this study is the differential Df localized to the restricted tangent bundle $T_\Lambda M = \bigcup_{x \in \Lambda} T_x M$.

Let us illustrate this approach by an example. Consider the hyperbolic linear map $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$, $f(x, y) = (2x, y/2)$ in a neighborhood of the origin:

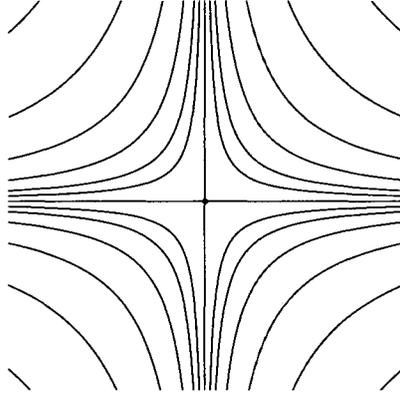


FIGURE 0.4.1. The map $(x, y) \mapsto (2x, y/2)$

The segment of the y -axis consists of points asymptotic to the origin in positive time and the segment of the x -axis consists of points asymptotic to the origin in negative time, while all other points move along hyperbolas $xy = \text{const}$. (cf. Section 1.2 for a more detailed description). Suppose we extend our map nonlinearly to a larger area such that the preimage of the y -axis and the image of the x -axis intersect at a point p and form a nonzero angle (see Figure 6.5.1).

Such a point p is called a *transverse homoclinic point* for the fixed point 0. Obviously $f^n(p) \rightarrow 0$ as $|n| \rightarrow \infty$ so $\Lambda = \{0\} \cup (\bigcup_{n=-\infty}^{\infty} f^n(p))$ is a closed f -invariant set. A theory by G. D. Birkhoff asserts that any neighborhood of Λ contains periodic points of arbitrarily high period. S. Smale gave a complete analysis of the structure of orbits staying in a sufficiently small neighborhood of the set Λ (see Section 6.5). His work played a crucial role in the development of the modern theory of dynamical systems.

Another version of the semilocal analysis involves the study of orbits which stay inside a certain usually open noninvariant set. Of course, there may not be such orbits at all, but under certain conditions the existence of such orbits can be guaranteed. The constructions of an invariant Cantor set in the quadratic map and of Smale's "horseshoe" discussed in Section 2.5 are simple but nontrivial examples of such analysis.

Exercise

0.4.1. Give an example of a diffeomorphism $f: \mathbb{R} \rightarrow \mathbb{R}$ such that $f(0) = (0)$, every orbit of the differential $f'(0)$ is bounded, and every orbit of f different from the origin is unbounded.

Part 1

Examples and fundamental concepts

First examples

This chapter is intended to illustrate the concept of a dynamical system and the notion of asymptotic behavior by an assortment of examples. In the course of our survey we proceed from simple to more complicated types of asymptotic behavior and identify certain properties for a more systematic analysis in the future.

1. Maps with stable asymptotic behavior

a. Contracting maps. The simplest imaginable kind of asymptotic behavior is represented by the convergence of iterates of any given state to a particular state.

Definition 1.1.1. Let (X, d) be a metric space. A map $f: X \rightarrow X$ is called *contracting* if there exists $\lambda < 1$ such that for any $x, y \in X$

$$d(f(x), f(y)) \leq \lambda d(x, y). \quad (1.1.1)$$

The inequality (1.1.1) implies that the map f is continuous and therefore its positive iterates form a discrete-time topological dynamical system.

By iterating (1.1.1), one sees that for any positive integer n

$$d(f^n(x), f^n(y)) \leq \lambda^n d(x, y) \quad (1.1.2)$$

so

$$d(f^n(x), f^n(y)) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

This means that the asymptotic behavior of all points is the same. On the other hand, for any $x \in X$ the sequence $\{f^n(x)\}_{n \in \mathbb{N}}$ is a Cauchy sequence because

for $m \geq n$

$$\begin{aligned} d(f^m(x), f^n(x)) &\leq \sum_{k=0}^{m-n-1} d(f^{n+k+1}(x), f^{n+k}(x)) \\ &\leq \sum_{k=0}^{m-n-1} \lambda^{n+k} d(f(x), x) \leq \frac{\lambda^n}{1-\lambda} d(f(x), x) \xrightarrow{n \rightarrow \infty} 0. \end{aligned} \tag{1.1.3}$$

Thus for any $x \in X$ the limit of $f^n(x)$ as $n \rightarrow \infty$ exists if the space is complete, and by (1.1.2) this limit is the same for all x . Let us denote this limit by p and show that p is a fixed point for f . For any $x \in X$ and any integer n one has

$$\begin{aligned} d(p, f(p)) &\leq d(p, f^n(x)) + d(f^n(x), f^{n+1}(x)) + d(f^{n+1}(x), f(p)) \\ &\leq (1 + \lambda)d(p, f^n(x)) + \lambda^n d(x, f(x)). \end{aligned}$$

Since $d(p, f^n(x)) \rightarrow 0$ as $n \rightarrow \infty$, we have $f(p) = p$. Taking the limit in (1.1.3) as $m \rightarrow \infty$ we obtain that $d(f^n(x), p) \leq \frac{\lambda^n}{1-\lambda} d(f(x), x)$. We will say that two sequences $\{x_n\}_{n \in \mathbb{N}}$ and $\{y_n\}_{n \in \mathbb{N}}$ of points in a metric space *converge exponentially* (or *with exponential speed*) to each other if $d(x_n, y_n) < c\lambda^n$ for some $c > 0$, $\lambda < 1$. In particular, if one of the sequences is constant, that is, $y_n = y$, we will say that x_n *converges exponentially* to y .

The above argument contains the proof of the following fundamental result which gives a complete and very simple picture of asymptotic behavior for the dynamical system generated by a contracting map.

Proposition 1.1.2. (Contraction Mapping Principle) *Let X be a complete metric space. Under the action of iterates of a contracting map $f: X \rightarrow X$ all points converge with exponential speed to the unique fixed point of f .*

Definition 1.1.3. If X is a topological space, $f: X \rightarrow X$, $f(p) = p$, and $f^n(x) \rightarrow p$ as $n \rightarrow \infty$ then we say that x is *(positively) asymptotic* to p . If f is invertible and $f^{-n}(x) \rightarrow p$ as $n \rightarrow \infty$ then we say that x is *negatively asymptotic* to p .

We denote the set of fixed points of any map f by $\text{Fix}(f)$.

Thus for a contracting map all points are asymptotic to a unique fixed point. This result will often be used in the course of our study of dynamical systems with more complicated behavior. Typically, we will be applying it not to the original dynamical system in the phase space but to a certain map in a functional space associated to the dynamical system. Right now, however, we give a straightforward and simple illustration of the use of the *Contraction Mapping Principle* in dynamics where the principle is applied to a certain derived system in the same space.

Proposition 1.1.4. *If p is a periodic point of period m for a C^1 map f and the differential Df_p^m does not have one as an eigenvalue then for every map g sufficiently close to f in the C^1 topology there exists a unique periodic point of period m close to p .*

Proof. Let us introduce local coordinates near p with p as the origin. In these coordinates Df_0^m becomes a matrix. Since 1 is not among its eigenvalues the map $F = f^m - \text{Id}$ defined locally in these coordinates is locally invertible by the Inverse Function Theorem. Now let g be a map C^1 -close to f . Near 0 one can write $g^m = f^m - H$ where H is small together with its first derivatives. A fixed point for g^m can be found from the equation $x = g^m(x) = (f^m - H)(x) = (F + \text{Id} - H)(x)$ or $(F - H)(x) = 0$ or

$$x = F^{-1}H(x).$$

Since F^{-1} has bounded derivatives and H has very small first derivatives one can show that $F^{-1}H$ is a contracting map. More precisely, let $\|\cdot\|_0$ denote the C^0 -norm, $\|dF^{-1}\|_0 = L$, and

$$\max(\|H\|_0, \|dH\|_0) \leq \epsilon.$$

Then, since $F(0) = 0$, we get $\|F^{-1}H(x) - F^{-1}H(y)\| \leq \epsilon L\|x - y\|$ for every x, y close to 0 and $\|F^{-1}H(0)\| \leq L\|H(0)\| \leq \epsilon L$, and hence $\|F^{-1}H(x)\| \leq \|F^{-1}H(x) - F^{-1}H(0)\| + \|F^{-1}H(0)\| \leq \epsilon L\|x\| + \epsilon L$. Thus if $\epsilon \leq \frac{R}{L(1+R)}$

the disc $X = \{x \mid \|x\| \leq R\}$ is mapped by $F^{-1}H$ into itself and the map $F^{-1}H: X \rightarrow X$ is contracting. By the Contraction Mapping Principle it has a unique fixed point in X which is thus a unique fixed point for g^m near 0. \square

Let us illustrate the notion of a contracting map with an elementary example. Consider the real line with the metric induced by the absolute value. Suppose $f: \mathbb{R} \rightarrow \mathbb{R}$ is a continuously differentiable function whose derivative is bounded by some $\lambda < 1$. If $x, y \in \mathbb{R}$ then by the Mean Value Theorem there exists some ξ between x and y such that $f(x) - f(y) = f'(\xi)(x - y)$. Thus $|f(x) - f(y)| = |f'(\xi)||x - y| \leq \lambda|x - y|$ and f is a contracting map according to Definition 1.1.1. Thus any such map has a unique fixed point. Exercise 1.1.3 contains a generalization of this example.

The next section contains a few additional examples of contracting maps.

b. Stability of contractions. We now make an observation about the orbit structure of contraction that is interesting in itself, but also of utility when we apply the Contraction Mapping Principle to operators associated to a dynamical system under study: namely, that changing the contracting map slightly does not move the fixed point very much.

Proposition 1.1.5. *If $f: X \rightarrow X$ is a contraction of a complete metric space X with fixed point x_0 and contraction constant λ as in Definition 1.1.1 then for every $\epsilon > 0$ there exists a $\delta \in (0, 1 - \lambda)$ such that for any map $g: X \rightarrow X$ with*

- (1) $d(f(x), g(x)) < \delta$ for all $x \in X$ and
- (2) $d(g(x), g(y)) \leq (\lambda + \delta)d(x, y)$ for all $x, y \in X$

the fixed point y_0 of g satisfies $d(x_0, y_0) < \epsilon$.

Proof. Take $\delta = \frac{\epsilon(1-\lambda)}{1+\epsilon}$. Since $g^n(x_0) \rightarrow y_0$ we have

$$\begin{aligned} d(x_0, y_0) &\leq \sum_{n=0}^{\infty} d(g^n(x_0), g^{n+1}(x_0)) < d(x_0, g(x_0)) \sum_{n=0}^{\infty} (\lambda + \delta)^n \\ &< \frac{\delta}{1 - \lambda - \delta} = \frac{\epsilon(1-\lambda)}{(1+\epsilon)(1-\lambda - \frac{\epsilon(1-\lambda)}{1+\epsilon})} = \epsilon. \end{aligned}$$

□

c. Increasing interval maps. The next simple asymptotic behavior is convergence of every orbit to a fixed point but in the presence of more than one fixed point. This occurs in the case of increasing functions of a real variable viewed as maps. This example is instructive because it demonstrates an important method in low-dimensional dynamics, the systematic use of the Intermediate Value Theorem.

Proposition 1.1.6. *If $I \subset \mathbb{R}$ is a closed interval and $f: I \rightarrow I$ is a nondecreasing continuous map then all $x \in I$ are asymptotic to a fixed point of f . If f is increasing (hence invertible) then all $x \in I$ are either fixed or positively and negatively asymptotic to adjacent fixed points.*

Proof. Note that the set $\text{Fix}(f)$ is closed by continuity and nonempty by the Intermediate Value Theorem. If $\text{Fix}(f) = I$ then there is nothing to show. Otherwise consider $x \in I \setminus \text{Fix}(f)$ and let (a, b) be the maximal open interval of $I \setminus \text{Fix}(f)$ containing x . Since f is nondecreasing we have $f(a, b) \subset [a, b]$ and by the Intermediate Value Theorem $f - \text{Id}$ does not change sign on (a, b) . To be specific suppose $f(y) > y$ for $y \in (a, b)$ (the other case is similar). Then $x_n := f^n(x)$ defines a nondecreasing sequence bounded by b , hence convergent to some $x_0 \in (a, b]$. But $f(x_0) = f(\lim_{n \rightarrow \infty} x_n) = \lim_{n \rightarrow \infty} f(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = x_0$, so $x_0 \in \text{Fix}(f)$ and in fact $x_0 = b$. Note that for the case $f(x) < x$ on (a, b) we would likewise obtain $f^n(x) \rightarrow a$ for all $x \in (a, b)$ as $n \rightarrow \infty$.

In case f is increasing note that the sign of $f^{-1} - \text{Id}$ on such an interval (a, b) is opposite to that of $f - \text{Id}$, so every $x \in (a, b)$ is positively and negatively asymptotic to opposite ends of $[a, b]$. □

Exercises

Problems 1.1.1 and 1.1.2 investigate the effect of replacing the uniform contraction condition (1.1.1) by weaker assumptions. As in the text X is a complete metric space and $f: X \rightarrow X$ a map of X into itself.

1.1.1. *Construct an example of a map f such that $d(f(x), f(y)) < d(x, y)$ for $x \neq y$, f has no fixed point, and $d(f^n(x), f^n(y))$ does not converge to zero for some x, y .*

1.1.2. Suppose that $d(f(x), f(y)) < d(x, y)$ for $x \neq y$ and in addition the space X is compact. Then the iterates of every point $x \in X$ converge to a single fixed point of f as $n \rightarrow \infty$. Give an example showing that convergence need not be exponential.

1.1.3. Let $f: M \rightarrow M$ be a C^1 map of a complete Riemannian manifold to itself. Then f is a contraction if and only if the norm of the differential is bounded by a constant $\lambda < 1$.

1.1.4. Use Proposition 1.1.5 to show that the fixed point of a contraction depends continuously on the contraction with respect to the C^1 topology.

1.1.5. Prove that no contracting map of a compact metric space with more than one point is invertible.

2. Linear maps

Next let us consider the dynamical system defined by the iterates of a linear map A of the Euclidean space \mathbb{R}^n . If A is invertible, this system is reversible.

Definition 1.2.1. Let $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a linear map. We call the set of eigenvalues the *spectrum* of A and denote it by $\text{sp } A$. We denote the maximal absolute value of an eigenvalue of A by $r(A)$ and call it the *spectral radius* of A .

Given any norm on \mathbb{R}^n we define the norm of a linear map A by $\|A\| := \sup_{\|v\|=1} \|Av\|$. Clearly $\|A\| \geq r(A)$ and with respect to the Euclidean norm we have $\|A\| = r(A)$ whenever A is diagonal. The following fact from linear algebra is useful for the understanding of dynamics of linear maps even if they cannot be diagonalized.

Proposition 1.2.2. For every $\delta > 0$ there exists a norm in \mathbb{R}^n such that $\|A\| < r(A) + \delta$.

Proof. Using the Jordan normal form we can find a basis in \mathbb{R}^n such that the matrix of our map has the block diagonal form

$$\begin{pmatrix} A_1 & & 0 \\ & \ddots & \\ 0 & & A_k \end{pmatrix}$$

where each block is either a Jordan block corresponding to a real eigenvalue λ :

$$\begin{pmatrix} \lambda & 1 & & 0 \\ 0 & \lambda & 1 & 0 \\ & & \ddots & \ddots \\ 0 & & & \lambda & 1 \\ 0 & & & & \lambda \end{pmatrix} \tag{1.2.1}$$

or a combination of two blocks corresponding to a pair of complex conjugate eigenvalues $\lambda = \rho e^{i\varphi}$ and $\bar{\lambda} = \rho e^{-i\varphi}$:

$$\begin{pmatrix} \rho R_\varphi & \text{Id} & \dots & 0 \\ & \rho R_\varphi & \text{Id} & 0 \\ & & \ddots & \\ 0 & & \dots & \rho R_\varphi \end{pmatrix}, \quad (1.2.2)$$

where Id is the 2×2 identity matrix $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ and $R_\varphi := \begin{pmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{pmatrix}$ is the 2×2 matrix corresponding to the rotation of the plane by the angle φ . Let us fix $\delta > 0$. By making an extra diagonal coordinate change of the form

$$\begin{pmatrix} 1 & & & 0 \\ & \delta^{-1} & & \\ & & \ddots & \\ 0 & & & \delta^{-m+1} \end{pmatrix}$$

for an m -block of the form (1.2.1) and of the form

$$\begin{pmatrix} \text{Id} & & & 0 \\ & \delta^{-1} \text{Id} & & \\ & & \ddots & \\ 0 & & & \delta^{-m+1} \text{Id} \end{pmatrix}$$

for a $2m$ -block of the form (1.2.2) one can make the off-diagonal entries in (1.2.1) and (1.2.2) equal to δ . Now for the standard Euclidean norm with respect to this new basis we have (after a slightly messy calculation)

$$\|A\| := \sup_{\|v\|=1} \|Av\| \leq r(A) + \delta. \quad (1.2.3)$$

□

Remark. In fact, since all norms in \mathbb{R}^n are equivalent up to a bounded multiple, one can conclude that for any norm and for every $\epsilon > 0$ there exists C_ϵ such that for any $v \in \mathbb{R}^n$

$$\|A^n v\| \leq C_\epsilon (r(A) + \epsilon)^n \|v\|.$$

We begin our study of asymptotic behavior of linear maps with an important particular case.

Corollary 1.2.3. *Assume all eigenvalues of a linear map A have absolute value less than one. Then there exists a norm in \mathbb{R}^n such that A is a contracting map with respect to the metric generated by that norm.*

Proof. If δ is chosen small enough we can conclude from Proposition 1.2.2 that $\|A\| < 1$ and since $d(x, y) = \|x - y\|$, A is a contracting map. □

The concept of exponential convergence does not depend on a particular choice of a norm. Thus Proposition 1.1.2 and Corollary 1.2.3 immediately imply the following statement.

Corollary 1.2.4. *If all eigenvalues of a linear map $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ have absolute values less than one, then the positive iterates of every point converge to the origin with exponential speed. If in addition A is an invertible map, that is, if zero is not an eigenvalue for A , then negative iterates of every point go to infinity exponentially.*

Let us look at a few examples of linear maps of \mathbb{R}^2 of this kind. The first example that comes to mind is the map given by the matrix $\begin{pmatrix} \lambda & 0 \\ 0 & \lambda \end{pmatrix}$ when $\lambda \in (0, 1)$. Every vector v is contracted by the factor λ , so the iterates of any vector move toward the origin along a line through 0. Since every line through 0 is mapped to itself (while being contracted), we can draw the orbit picture as follows.

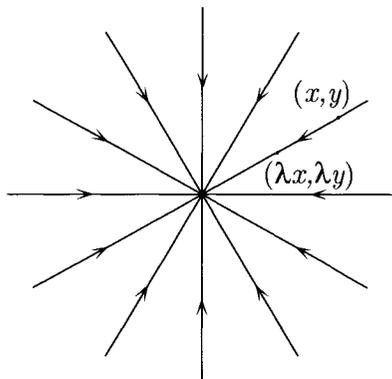


FIGURE 1.2.1. Orbits of contracting homothety

The next example to try is the map $\begin{pmatrix} \lambda & 0 \\ 0 & \mu \end{pmatrix}$ when $\lambda, \mu \in (0, 1)$. Let us suppose that $\mu < \lambda$. Then every point still moves toward the origin, but not on a given straight line. Nevertheless, the orbits move along curves that are invariant under $\begin{pmatrix} \lambda & 0 \\ 0 & \mu \end{pmatrix}$. One can easily verify that these are the axes and the curves given by $x^{\log \mu} = \text{const.} \cdot y^{\log \lambda}$. Thus the corresponding orbit picture is as shown in Figure 1.2.2. A fixed point of a map of this kind is referred to as a *node*.

The next example is that of a map with two complex eigenvalues of absolute value less than one. To be specific consider the map given by $\lambda \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}$. Observe that this is the composition of the rotation by θ and contraction by

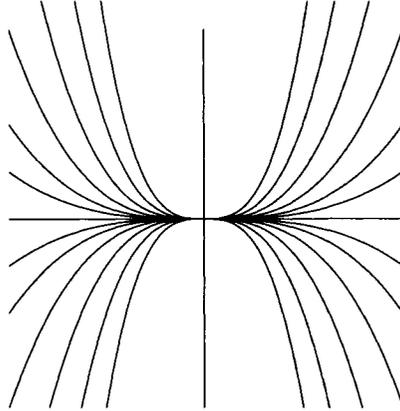


FIGURE 1.2.2. A node

λ and the n th iterate of this map is given by $\lambda^n \begin{pmatrix} \cos n\theta & \sin n\theta \\ -\sin n\theta & \cos n\theta \end{pmatrix}$. Thus, while points still approach the origin with exponential speed, at the same time they rotate around 0. In fact, we still have invariant curves, namely, spirals, which are most easily described in polar coordinates (r, ϕ) . One checks easily that the curves $r = \text{const.} \cdot e^{-(\theta^{-1} \log \lambda)\phi}$ are invariant under this map. Thus here we obtain a portrait as in Figure 1.2.3, an orbit picture called a *focus*.

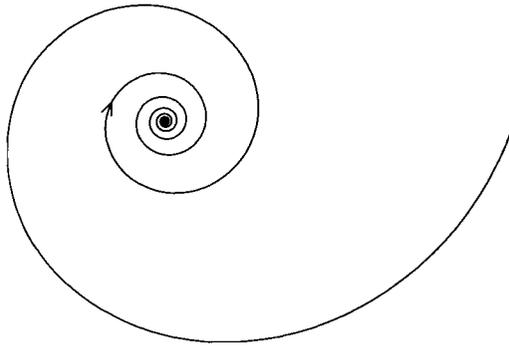


FIGURE 1.2.3. A focus

The appearance of invariant curves in the above examples is not accidental. The linear maps we described above arise from solutions of ordinary differential equations, whose flows interpolate iterates of the maps above. Consider first the ordinary differential equation

$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} \log \lambda & 0 \\ 0 & \log \lambda \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

whose solutions are given by

$$\begin{pmatrix} x(t) \\ y(t) \end{pmatrix} = \begin{pmatrix} x(0)e^{t \log \lambda} \\ y(0)e^{t \log \lambda} \end{pmatrix} = \begin{pmatrix} x(0)\lambda^t \\ y(0)\lambda^t \end{pmatrix}.$$

Thus for $t = 1$ we get the first contracting map. The second one is evidently obtained from the ordinary differential equation with coefficient matrix $\begin{pmatrix} \log \lambda & 0 \\ 0 & \log \mu \end{pmatrix}$. Thus the first two figures above show the orbit structures of two *linear flows*. The second of these is also called a node in the context of ordinary differential equations.

Finally the focus is obtained from the linear ordinary differential equation with coefficient matrix

$$\begin{pmatrix} \log \lambda & \theta \\ -\theta & \log \lambda \end{pmatrix}.$$

Again taking $t = 1$ gives the above map with a focus.

Next we proceed to a somewhat more general case which will play an important role in our future considerations of nonlinear dynamical systems.

Definition 1.2.5. A linear map of \mathbb{R}^n is called *hyperbolic* if all of its eigenvalues have absolute values different from one.

For every linear map $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ and for a real eigenvalue λ of A let us denote by E_λ the root space corresponding to λ , that is, the space of all vectors $v \in \mathbb{R}^n$ such that $(A - \lambda \text{Id})^k v = 0$ for some k .

Similarly, for a pair of complex conjugate eigenvalues $\lambda, \bar{\lambda}$ let $E_{\lambda, \bar{\lambda}}$ be the intersection of \mathbb{R}^n with the sum of root spaces corresponding to E_λ and $E_{\bar{\lambda}}$ for the complexification of A (that is, the extension to the complex space \mathbb{C}^n). For brevity we will call $E_{\lambda, \bar{\lambda}}$ a root space, too. Let

$$E^- = E^-(A) = \bigoplus_{|\lambda| < 1} E_\lambda \oplus \bigoplus_{|\lambda| < 1} E_{\lambda, \bar{\lambda}} \quad (1.2.4)$$

and similarly

$$E^+ = E^+(A) = \bigoplus_{|\lambda| > 1} E_\lambda \oplus \bigoplus_{|\lambda| > 1} E_{\lambda, \bar{\lambda}}. \quad (1.2.5)$$

If the map A is invertible, then $E^+(A) = E^-(A^{-1})$. Finally, let

$$E^0 = E^0(A) = E_1 \oplus E_{-1} \oplus \bigoplus_{|\lambda|=1} E_{\lambda, \bar{\lambda}}. \quad (1.2.6)$$

The spaces E^-, E^+, E^0 are obviously invariant with respect to A and $\mathbb{R}^n = E^- \oplus E^+ \oplus E^0$.

An equivalent way to describe hyperbolic linear maps is to say that A is hyperbolic if $E^0 = \{0\}$ or, equivalently, if $\mathbb{R}^n = E^+ \oplus E^-$

Since the restriction of A to the space $E^-(A)$ is a linear operator with all eigenvalues of absolute value less than one, one obtains immediately from Corollary 1.2.4

Corollary 1.2.6. *There exists a norm such that the restriction of a linear map A to the space $E^-(A)$ is a contracting map. If A is invertible then in addition the restriction of A^{-1} to the space $E^+(A)$ is a contracting map.*

Definition 1.2.7. The space $E^-(A)$ above is called the *contracting* subspace, the space $E^+(A)$ the *expanding* subspace.

Remark. Note that the expanding subspace is not characterized by the fact that vectors in it expand under iterates of the map—all vectors outside the contracting subspace are expanded by a sufficiently large iterate of the map. The characterization of E^+ is given by the description of Corollary 1.2.6, namely that preimages contract.

Our next statement describes the asymptotic behavior of iterates of a hyperbolic linear map.

Proposition 1.2.8. *Let $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a hyperbolic linear map. Then:*

- (1) *For every $v \in E^-$, the positive iterates $A^n v$ converge to the origin with exponential speed as $n \rightarrow \infty$ and if A is invertible then the negative iterates $A^n v$ go to infinity with exponential speed as $n \rightarrow -\infty$.*
- (2) *For every $v \in E^+$ the positive iterates of v go to infinity exponentially and if A is invertible then the negative iterates converge exponentially to the origin.*
- (3) *For every $v \in \mathbb{R}^n \setminus (E^- \cup E^+)$ the iterates $A^n v$ go to infinity exponentially as $n \rightarrow \infty$ and if A is invertible also as $n \rightarrow -\infty$.*

Proof. Statement (1) follows from Corollary 1.2.6 and Proposition 1.1.2; (2) reduces to (1) when A is replaced by A^{-1} . Finally, if $v \in \mathbb{R}^n \setminus (E^- \cup E^+)$ then $v = v^- + v^+$ where $v^- \in E^-$, $v^+ \in E^+$ and $v^-, v^+ \neq 0$.

Then for large positive n one has

$$\|A^n v\| = \|A^n(v^- + v^+)\| \geq \|A^n v^+\| - \|A^n v^-\| \geq \lambda^n c \|v^+\| - \lambda^{-n} c' \|v^-\| \geq \lambda^n c'',$$

where $\lambda > 1$ and $c, c', c'' > 0$ do not depend on n .

The argument for negative iterates is the same with v^+ and v^- exchanged. \square

Let us quickly illustrate the behavior of the hyperbolic linear maps by considering the two-dimensional case. In this case the subspaces E^+ and E^- are of course one-dimensional. In the example where these are the x and y axes (which can be arranged by a coordinate change) we can easily draw a picture. The map is given by $\begin{pmatrix} \lambda & 0 \\ 0 & \mu \end{pmatrix}$ with $0 < \mu < 1 < \lambda$. As for the node one checks that the axes and the curves given by $x^{\log \mu} = \text{const.}$ $y^{\log \lambda}$ are invariant.

Note that the two exponents here now have different sign, so we do not get curves through the origin. Indeed the picture is as in Figure 0.4.1; it is usually referred to as a *saddle*. As in the previous linear examples this one also arises from a linear ordinary differential equation, this time with coefficient matrix

$$\begin{pmatrix} \log \lambda & 0 \\ 0 & \log \mu \end{pmatrix}.$$

An interesting special case is that of a map which is also *area preserving*. In this case we must have $\lambda\mu = 1$, so the invariant curves are the standard hyperbolas $xy = \text{const}$. This is the reason for using the word “hyperbolic” in Definition 1.2.5 and in many other contexts later in this book. There is also an interesting case of a hyperbolic linear map which does *not* come from an ordinary differential equation. It is given by a matrix $\begin{pmatrix} \lambda & 0 \\ 0 & \mu \end{pmatrix}$ where $\lambda < -1 < \mu < 0$, that is, where λ and μ are negative numbers on opposite sides of -1 . These *inverted saddles* play an interesting role in some global problems (cf. Section 8.4 and Exercise 9.2.7).

In order to describe the behavior of iterates for a nonhyperbolic linear map we should first understand what happens inside the subspace E^0 . This subspace splits into root subspaces of E_1 , E_{-1} , and $E_{\lambda, \bar{\lambda}}$ for $|\lambda| = 1$, $\lambda \neq \pm 1$. Inside each of those subspaces there is a corresponding invariant eigenspace which we will denote by \tilde{E}_1 , \tilde{E}_{-1} , and $\tilde{E}_{\lambda, \bar{\lambda}}$, correspondingly. The first two of those spaces yield a fairly trivial behavior; namely, all points inside \tilde{E}_1 are fixed and those in $\tilde{E}_{-1} \setminus \{0\}$ are periodic with period two. A more interesting situation appears in $\tilde{E}_{\lambda, \bar{\lambda}}$ if λ is not real, say $\lambda = e^{2\pi i \varphi}$. If one of those spaces is not empty then A has an invariant plane such that in an appropriate coordinate system the map acts in that plane as a rotation by the angle φ about the origin.

Every circle with center at the origin is invariant under the map. Thus in order to continue our analysis of linear maps we should first understand the behavior of iterates for a rotation of the circle. This is the first time in our survey when we will encounter the phenomenon of nontrivial recurrence, that is, iterates of a point coming back arbitrary close to the initial position without returning exactly to that position. A detailed study of rotations is our next task. The structure of general linear maps is discussed in Exercises 1.2.4 and 1.2.5.

Exercises

1.2.1. Prove (1.2.3).

1.2.2. Show that the eigenvalues of a linear map A depend continuously on A . Do they depend smoothly?

1.2.3. Show that hyperbolic linear maps are an open dense subset of the set of linear maps $\mathbb{R}^n \rightarrow \mathbb{R}^n$.

1.2.4. Suppose all eigenvalues of a linear map $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ have absolute value one. Then there exists an invariant subspace $C = C(A) \subset \mathbb{R}^n$ and a norm in \mathbb{R}^n such that A acts in C as an isometry and for every vector $v \in \mathbb{R}^n \setminus C$ the norm $\|A^n v\|$ grows polynomially as $|n| \rightarrow \infty$, that is, for some positive integer k and $c > 0$

$$\lim_{|n| \rightarrow \infty} \frac{\|A^n v\|}{\|v\| |n|^k} = c$$

(k and c may depend on v). Show how to determine the maximal value of k for a given map.

1.2.5. Use Proposition 1.2.8 and Exercise 1.2.4 to describe the asymptotic behavior of points for an arbitrary invertible linear map in terms of the decomposition

$$\mathbb{R}^n = E^+(A) \oplus E^-(A) \oplus E^0(A).$$

3. Rotations of the circle

We can use either multiplicative notation so that the circle is represented as the unit circle in the complex plane

$$S^1 = \{z \in \mathbb{C} \mid |z| = 1\} = \{e^{2\pi i \varphi} \mid \varphi \in \mathbb{R}\}$$

or additive notation

$$S^1 = \mathbb{R}/\mathbb{Z},$$

the factor group of the additive group of real numbers modulo the subgroup of integers. The logarithm map

$$e^{2\pi i \varphi} \mapsto \varphi$$

establishes an isomorphism between these representations. We will use the symbol R_α to denote the rotation by angle $2\pi\alpha$. In multiplicative notation

$$R_\alpha z = z_0 z \text{ with } z_0 = e^{2\pi i \alpha}.$$

Not surprisingly, in additive notation we have

$$R_\alpha x = x + \alpha \pmod{1},$$

where $\pmod{1}$ means that numbers which differ by an integer are identified. The iterates of the rotation are correspondingly

$$R_\alpha^n z = R_{n\alpha} z = z_0^n z \text{ or } R_\alpha^n x = x + n\alpha \pmod{1}. \quad (1.3.1)$$

A crucial distinction appears between the cases of rational and irrational α .

In the former case, write $\alpha = p/q$, where p, q are relatively prime integers. Then $R_\alpha^q x = x$ for all x so R_α^q is the identity map and after q iterates the transformation simply repeats itself.

The latter case is much more interesting. We begin with two general definitions which belong to topological dynamics.

Definition 1.3.1. A topological dynamical system $f: X \rightarrow X$ is called *topologically transitive* if there exists a point $x \in X$ such that its orbit $\mathcal{O}_f(x) := \{f^n(x)\}_{n \in \mathbb{Z}}$ is dense in X .

The definitions for noninvertible and continuous-time systems are similar.

Definition 1.3.2. A topological dynamical system $f: X \rightarrow X$ is called *minimal* if the orbit of every point $x \in X$ is dense in X , or, equivalently, if f has no proper closed invariant sets.

Proposition 1.3.3. *If α is irrational then the rotation R_α is minimal.*

Proof. Let $A \subset S^1$ be the closure of an orbit. If the orbit is not dense, the complement $S^1 \setminus A$ is a nonempty open invariant set which consists of disjoint intervals. Let I be the longest of those intervals (or one of the longest, if there are several of the same length). Since rotation preserves the length of any interval, the iterates $R_\alpha^n I$ do not overlap. Otherwise $S^1 \setminus A$ would contain an interval longer than I . Since α is irrational, no iterates of I can coincide; because then an endpoint x of an iterate of I would come back to itself and we would have $x + k\alpha = x \pmod{1}$ with $k\alpha = l$ an integer and $\alpha = l/k$ a rational number. Thus the intervals $R_\alpha^n I$ are all of equal length and all disjoint, but this is impossible because the circle has finite length and the sum of lengths of disjoint intervals can not exceed the length of the circle. \square

Irrational rotations serve as the starting point for a number of very fruitful generalizations. Let us discuss one of them. The circle is a compact abelian group and the rotation can be represented in group terms as the group multiplication or left translation

$$L_{g_0}: G \rightarrow G, \quad L_{g_0}g = g_0g. \quad (1.3.2)$$

The orbit of the unit element $e \in G$ is the cyclic subgroup $\{g_0^n\}_{n \in \mathbb{Z}}$ and it is easy to deduce from Proposition 1.3.3 that the circle does not have proper infinite closed subgroups.

Proposition 1.3.4. *If the translation L_{g_0} on a topological group G is topologically transitive then it is minimal.*

Proof. For $g, g' \in G$ denote by $A, A' \subset G$ the closures of the orbits of g and g' , respectively. Now $g_0^n g' = g_0^n g(g^{-1}g')$, so $A' = Ag^{-1}g'$ and $A' = G$ if and only if $A = G$. \square

Exercises

1.3.1. Prove that the decimal expansion of the number 2^n may begin with any finite combination of digits.

1.3.2. Let G be a metrizable compact topological group. Suppose for some $g_0 \in G$ the translation L_{g_0} is topologically transitive. Prove that G is abelian.

1.3.3. Define the following metric d_2 on the group \mathbb{Z} of all integers: $d_2(m, n) = \|m - n\|_2$ where

$$\|n\|_2 = 2^{-k} \quad \text{if } n = 2^k l \text{ with an odd number } l.$$

The completion of \mathbb{Z} with respect to that metric is called the group of 2-adic integers and is usually denoted by \mathbb{Z}_2 . It is a compact topological group. Let \mathbb{Z}_2^+ be the closure of the even integers with respect to the metric d_2 . \mathbb{Z}_2^+ is a subgroup of \mathbb{Z}_2 of index two.

Prove that for $g_0 \in \mathbb{Z}_2$ the translation $L_{g_0}: \mathbb{Z}_2 \rightarrow \mathbb{Z}_2$ is topologically transitive if and only if $g_0 \in \mathbb{Z}_2 \setminus \mathbb{Z}_2^+$.

This is an example of a class of systems called *adding machines*. We will encounter them again in Section 15.4.

4. Translations on the torus

This is a generalization of rotations and a special case of group translations. This example plays the central role in the theory of completely integrable Hamiltonian systems, which we will touch upon at the end of the next section. The phase space is the n -dimensional torus

$$\mathbb{T}^n = \underbrace{S^1 \times \cdots \times S^1}_{n \text{ times}} = \mathbb{R}^n / \mathbb{Z}^n = \underbrace{\mathbb{R}/\mathbb{Z} \times \cdots \times \mathbb{R}/\mathbb{Z}}_{n \text{ times}}.$$

A natural fundamental domain for $\mathbb{R}^n / \mathbb{Z}^n$ is the unit cube:

$$I^n = \{(x_1, \dots, x_n) \in \mathbb{R}^n \mid 0 \leq x_i \leq 1 \text{ for } i = 1, \dots, n\}.$$

In order to represent the torus, one should identify opposite faces of I^n so that the point $(x_1, \dots, x_{i-1}, 0, x_{i+1}, \dots, x_n)$ is identified with $(x_1, \dots, x_{i-1}, 1, x_{i+1}, \dots, x_n)$ because these two points represent the same element of the factor group.

Similar to the case of the circle there are two convenient coordinate systems on \mathbb{T}^n , namely,

- (1) multiplicative, where the elements of \mathbb{T}^n are represented as (z_1, \dots, z_n) with $z_i \in \mathbb{C}$ and $|z_i| = 1$ for $i = 1, \dots, n$, and
- (2) additive, when they are represented by n -vectors (x_1, \dots, x_n) , where each coordinate is defined mod 1.

The correspondence $(x_1, \dots, x_n) \mapsto (e^{2\pi i x_1}, \dots, e^{2\pi i x_n})$ establishes an isomorphism between these two representations. In additive notation let $\gamma = (\gamma_1, \dots, \gamma_n) \in \mathbb{T}^n$. The translation T_γ has the form

$$T_\gamma(x_1, \dots, x_n) = (x_1 + \gamma_1, \dots, x_n + \gamma_n) \pmod{1}.$$

If all coordinates of the vector γ are rational numbers, then T_γ is periodic. However, unlike the circle case, it is not true any more that minimality is the only alternative to periodicity. For example, if $n = 2$ and $\gamma = (\alpha, 0)$ where α is an irrational number then the torus \mathbb{T}^2 splits into a family of invariant circles $x_2 = \text{const.}$ and every orbit stays on one of these circles and fills it densely.

Proposition 1.4.1. *The translation T_γ is minimal if and only if the numbers $\gamma_1, \dots, \gamma_n$ and 1 are rationally independent, that is, if $\sum_{i=1}^n k_i \gamma_i$ is not an integer for any collection of integers k_1, \dots, k_n except for $k_1 = k_2 = \dots = k_n = 0$.*

Remark. One may give an algebraic proof of this proposition which involves the classification of all closed subgroups of \mathbb{T}^n and induction on dimension. We prefer an analytical approach which anticipates some of the most fruitful methods used for the study of smooth dynamical systems and will be developed further in Section 4.2.

Before we proceed to the proof of this proposition we need to establish some general criteria for topological transitivity.

Lemma 1.4.2. *Let $f: X \rightarrow X$ be a continuous map of a locally compact separable metric space X into itself. The map f is topologically transitive if and only if for any two nonempty open sets $U, V \subset X$ there exists an integer $N = N(U, V)$ such that $f^N(U) \cap V$ is nonempty.*

Proof. Let f be topologically transitive and suppose the orbit of $x \in X$ is dense. Then, in particular, this orbit intersects both U and V so $f^n(x) \in U$, $f^m(x) \in V$, where, say $m \geq n$. Consequently $f^{m-n}(U) \cap V$ is nonempty. (Recall that $f^{-1}(A) := \{x \in X \mid f(x) \in A\}$.)

Now let us assume that the intersection condition holds. Let U_1, U_2, \dots be a countable base of open subsets of X . This means that for any $x \in X$ and every open set with $x \in U \subset X$ there is an n such that $x \in U_n \subset U$. Let us furthermore choose U_1 in such a way that its closure \bar{U}_1 is compact. In order to prove topological transitivity, it is enough to construct an orbit which intersects every U_n . By assumption, there exists an integer N_1 such that $f^{N_1}(U_1) \cap U_2$ is nonempty. Let V_1 be a nonempty open set such that $\bar{V}_1 \subset U_1 \cap f^{-N_1}(U_2)$. Obviously \bar{V}_1 is compact. There exists an integer N_2 such that $f^{N_2}(V_1) \cap U_3$ is nonempty. Again, take an open set V_2 such that $\bar{V}_2 \subset V_1 \cap f^{-N_2}(U_3)$. By induction, we construct a nested sequence of open sets V_n such that $\bar{V}_{n+1} \subset V_n \cap f^{-N_{n+1}}(U_{n+2})$. The intersection $V = \bigcap_{n=1}^{\infty} \bar{V}_n = \bigcap_{n=1}^{\infty} V_n$ is nonempty because the \bar{V}_n are compact. If $x \in V$ then $f^{N_{n-1}}(x) \in U_n$ for every $n \in \mathbb{N}$. \square

Corollary 1.4.3. *A continuous open map f of a locally compact separable metric space is topologically transitive if and only if there are no two disjoint open nonempty f -invariant sets.*

Proof. If $U, V \subset X$ are open then the invariant open sets $\tilde{U} := \bigcup_{n \in \mathbb{Z}} f^n(U)$ and $\tilde{V} := \bigcup_{n \in \mathbb{Z}} f^n(V)$ are not disjoint, so that $f^n(U) \cap f^m(V) \neq \emptyset$ for some $n, m \in \mathbb{Z}$, and $f^{n-m}(U) \cap V \neq \emptyset$. \square

Corollary 1.4.4. *If $f: X \rightarrow X$ is topologically transitive then there is no f -invariant nonconstant continuous function $\varphi: X \rightarrow \mathbb{R}$.*

Proof. Let $\varphi: X \rightarrow \mathbb{R}$ be f -invariant, that is, $\varphi(f(x)) = \varphi(x)$ for all $x \in X$. Since it is not a constant, there exists $t \in \mathbb{R}$ such that both $\{x \in X \mid \varphi(x) > t\}$ and $\{x \in X \mid \varphi(x) < t\}$ are nonempty. Since φ is invariant, these sets are also invariant. Since φ is continuous, they are open. \square

Proof of Proposition 1.4.1. First let us show that if $\sum_{i=1}^n k_i \gamma_i = k$ and not all integers k_1, \dots, k_n are zero, then T_γ is not topologically transitive. We will construct a continuous T_γ -invariant function and then use Corollary 1.4.4. Our function is $\varphi(x) = \sin 2\pi (\sum k_i x_i)$. It is defined on \mathbb{T}^n by the periodicity of $\sin(x)$ and it is not constant by our assumption. On the other hand φ is invariant because

$$\begin{aligned} \varphi(T_\gamma x) &= \sin(2\pi \sum k_i (x_i + \gamma_i)) \\ &= \sin(2\pi \sum k_i x_i + 2\pi k) = \sin(2\pi \sum k_i x_i) = \varphi(x). \end{aligned}$$

To prove the converse it is enough to show that rational independence of $\gamma_1, \dots, \gamma_n, 1$ implies topological transitivity of T_γ . Since T_γ is a translation on a group, this will imply minimality by Proposition 1.3.4. We will use Corollary 1.4.3 and prove the contrapositive. Let U, V be two disjoint nonempty open T_γ -invariant sets. Let χ be the characteristic function of U . By invariance of U we have

$$\chi(T_\gamma x) = \chi(x).$$

Take the Fourier expansion

$$\chi(x_1, \dots, x_n) = \sum_{(k_1, \dots, k_n) \in \mathbb{Z}^n} \chi_{k_1, \dots, k_n} \exp\left(2\pi i \sum_{j=1}^n k_j x_j\right)$$

of χ . Then

$$\begin{aligned} \chi(T_\gamma x) &= \chi(x_1 + \gamma_1, \dots, x_n + \gamma_n) \\ &= \sum_{(k_1, \dots, k_n) \in \mathbb{Z}^n} \chi_{k_1, \dots, k_n} \exp\left(2\pi i \sum_{j=1}^n k_j (x_j + \gamma_j)\right) \\ &= \sum_{(k_1, \dots, k_n) \in \mathbb{Z}^n} \chi_{k_1, \dots, k_n} \exp\left(2\pi i \sum_{j=1}^n k_j \gamma_j\right) \exp\left(2\pi i \sum_{j=1}^n k_j x_j\right). \end{aligned}$$

Invariance of χ and uniqueness of the Fourier expansion imply that for every k_1, \dots, k_n we have $\chi_{k_1, \dots, k_n} = \chi_{k_1, \dots, k_n} \exp\left(2\pi i \sum_{j=1}^n k_j \gamma_j\right)$ or

$$\chi_{k_1, \dots, k_n} \left(1 - \exp 2\pi i \sum_{j=1}^n k_j \gamma_j\right) = 0,$$

which means that either $\chi_{k_1, \dots, k_n} = 0$ or $\exp(2\pi i \sum k_i \gamma_i) = 1$, that is, $\sum k_i \gamma_i$ is an integer. Since both U and its complement contain nonempty open sets, which have positive Lebesgue measure, χ is not constant almost everywhere. Therefore there is some $(k_1, \dots, k_n) \neq 0$ such that $\chi_{k_1, \dots, k_n} \neq 0$ and hence $\sum k_i \gamma_i$ is an integer. \square

Let us point out a relation between toral translations and linear maps. Let $A: \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$ be given by the matrix

$$\begin{pmatrix} R_{\varphi_1} & & & 0 \\ & \ddots & & \\ & & \ddots & \\ 0 & & & R_{\varphi_n} \end{pmatrix}.$$

Using a complex coordinate $z_k = x_{2k-1} + ix_{2k}$ in each two-plane $0 = x_1 = \dots = x_{2k-2} = x_{2k+1} = \dots = x_{2n}$ for $k = 1, \dots, n$ we can write the map A as $A(z_1, \dots, z_n) = (e^{i\varphi_1} z_1, \dots, e^{i\varphi_n} z_n)$. Let $\rho = (\rho_1, \dots, \rho_n)$ be a vector with nonnegative coordinates. The torus $\mathbb{T}_\rho^k = \{|z_1| = \rho_1, \dots, |z_n| = \rho_n\}$ is invariant with respect to A and its dimension k is equal to the number of nonzero coordinates in ρ . The restriction of A to that torus is obviously just the translation T_γ where γ is the k -dimensional vector composed of all φ_i 's for which $\rho_i \neq 0$.

Exercises

1.4.1. Prove that for any translation T_γ and any $x \in \mathbb{T}^n$ the closure $C(x)$ of the orbit of x is a finite union of tori of dimension k , $0 \leq k \leq n$, and that the restriction of T_γ to $C(x)$ is minimal.

1.4.2. Let X be a compact metrizable space which is perfect, that is, does not have isolated points. Prove that if a homeomorphism $f: X \rightarrow X$ is topologically transitive, that is, for some point $x \in X$ the whole orbit $\mathcal{O}(x) = \{f^n(x) \mid n \in \mathbb{Z}\}$ is dense, then there exists a point $y \in X$ whose positive semiorbit $\mathcal{O}^+(y) = \{f^n(y) \mid n = 0, 1, 2, \dots\}$ is dense.

1.4.3. Construct an example showing that the assertion of Exercise 1.4.2 is not true if X is an arbitrary compact metrizable space.

1.4.4. Prove that the map $A_\alpha: \mathbb{T}^2 \rightarrow \mathbb{T}^2$, $A_\alpha(x, y) = (x + \alpha, y + x) \pmod{1}$ is topologically transitive if and only if α is irrational.

5. Linear flow on the torus and completely integrable systems

In this section we consider continuous-time analogs of examples from the past two sections. We begin with the following system of differential equations on the 2-torus (we use additive notation)

$$\frac{dx_1}{dt} = \omega_1, \quad \frac{dx_2}{dt} = \omega_2. \quad (1.5.1)$$

This system of differential equations can be easily integrated explicitly. The resulting flow $\{T_\omega^t\}_{t \in \mathbb{R}}$ has the form

$$T_\omega^t(x_1, x_2) = (x_1 + \omega_1 t, x_2 + \omega_2 t) \pmod{1}. \quad (1.5.2)$$

We will present a geometric picture of this flow. As we already mentioned, the torus $\mathbb{T}^2 = \mathbb{R}^2/\mathbb{Z}^2$ can be represented as the unit square $I^2 = \{(x_1, x_2) \mid 0 \leq x_1 \leq 1, 0 \leq x_2 \leq 1\}$ with pairs of opposite sides identified: $(x, 0) \sim (x, 1)$ and $(0, x) \sim (1, x)$. In this representation the integral curves of the system (1.5.1) are pieces of straight lines with slope $\gamma = \omega_2/\omega_1$. The motion along the orbits is uniform with instantaneous “jumps” to the corresponding points when the orbit reaches the boundary of the square (compare with the suspension construction in Section 0.3). If we consider the successive moments when an orbit intersects the circle $C_1 = \{x_1 = 0\}$, the x_2 coordinate changes by exactly $\gamma \pmod{1}$ between two such returns. Thus by Proposition 1.3.3 if γ is irrational, the closure of every orbit contains the circle C_1 and since the images of this circle under the flow $\{T_\omega^t\}$ cover the whole torus, the flow is minimal in the sense similar to that of Definition 1.3.2, that is, every orbit is dense in \mathbb{T}^2 . If γ is rational, then every orbit is closed, as becomes immediately clear from (1.5.2).

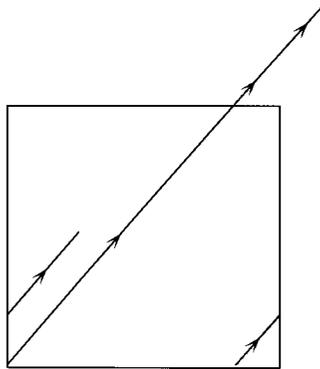


FIGURE 1.5.1. Linear flow on the torus

This example has a natural generalization to a torus of arbitrary dimension. Namely, let us consider the following system of differential equations on \mathbb{T}^n

$$\frac{dx_i}{dt} = \omega_i \text{ for } i = 1, \dots, n.$$

Again integration produces a one-parameter group of translations

$$T_\omega^t(x_1, \dots, x_n) = (x_1 + t\omega_1, \dots, x_n + t\omega_n) \pmod{1}. \quad (1.5.3)$$

Obviously the flow $\{T_\omega^t\}$ is minimal if for some t_0 the transformation $T_\omega^{t_0}$ is minimal. This remark together with Proposition 1.4.1 allows us to establish the criterion of minimality for this case.

Proposition 1.5.1. *The flow $\{T_\omega^t\}$ is minimal if and only if the numbers $\omega_1, \dots, \omega_n$ are rationally independent, that is, if $\sum_{i=1}^n k_i \omega_i \neq 0$ for any integers k_1, \dots, k_n unless $k_1 = \dots = k_n = 0$.*

Proof. Since $T_\omega^t = T_{t\omega}$ minimality follows from Proposition 1.4.1 once we show that for some real t and for any nonzero integer vector (k_1, \dots, k_n) the sum $\sum_{i=1}^n tk_i \omega_i$ is never an integer. But for any collection of integers k_1, \dots, k_n, k there is only one value of t such that

$$t \sum_{i=1}^n k_i \omega_i = k,$$

namely, $t = k / \sum k_i \omega_i$, unless $\sum_{i=1}^n k_i \omega_i = 0$, which cannot be the case by our assumption. The proof in one direction is finished because there are only countably many different integer vectors k_1, \dots, k_n, k and uncountably many values of t to take care of.

On the other hand, if $\sum_{i=1}^n k_i \omega_i = 0$ for some nonzero vector (k_1, \dots, k_n) , then the function $\sin 2\pi (\sum_{i=1}^n k_i x_i)$ is continuous, nonconstant, and invariant under the flow $\{T_\omega^t\}$. \square

Similarly to the discrete-time case, there is a close connection between linear flows on tori and solutions of certain linear systems of ordinary differential equations with constant coefficients. Let A be a $2n \times 2n$ real matrix with n pairs of distinct purely imaginary eigenvalues $\pm i\alpha_i$, $i = 1, \dots, n$. Consider the system of ordinary differential equations

$$\frac{dx}{dt} = Ax. \quad (1.5.4)$$

The solution of (1.5.4) is

$$x(t) = e^{tA}x(0).$$

By a coordinate change the matrix A can be transformed into

$$\begin{pmatrix} 0 & \alpha_1 & & & \\ -\alpha_1 & 0 & & & \\ & & \ddots & & \\ & & & 0 & \alpha_n \\ & & & -\alpha_n & 0 \end{pmatrix}$$

so that the solution for the time t is the linear transformation given by the matrix

$$\begin{pmatrix} R_{t\alpha_1} & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & R_{t\alpha_n} \end{pmatrix}.$$

Arguing exactly as at the end of the previous section we split \mathbb{R}^{2n} into invariant tori where the flow defined by (1.5.4) acts by translations.

A more important class of dynamical systems related to linear flows on the torus comes from Hamiltonian mechanics. Let us recall the most classical definition of a Hamiltonian system. Let H be a smooth function defined in an open subset U of Euclidean space \mathbb{R}^{2n} . The Hamiltonian equations for the Hamiltonian function H are

$$\begin{aligned} \frac{dx_i}{dt} &= \frac{\partial H}{\partial x_{i+n}}, & i &= 1, \dots, n, \\ \frac{dx_i}{dt} &= -\frac{\partial H}{\partial x_{i-n}}, & i &= n+1, \dots, 2n. \end{aligned} \tag{1.5.5}$$

The more general definition involves a $2n$ -dimensional smooth manifold M , a closed nondegenerate differential 2-form Ω on TM , that is, a form such that the exterior derivative $d\Omega = 0$ and the n -fold wedge product $\Omega^n \neq 0$, and a smooth function $H: M \rightarrow \mathbb{R}$. Then the Hamiltonian vector field V_H is defined by the condition that

$$\Omega(\xi, V_H(x)) = dH(\xi) \tag{1.5.6}$$

for $x \in M$, $\xi \in T_x M$. The Euclidean case (1.5.5) corresponds to $\Omega = \sum_{i=1}^n dx_i \wedge dx_{i+n}$. A more motivated and detailed description of Hamiltonian systems will be given in Section 5.5c.

It is not appropriate now to discuss a general notion of complete integrability, with its historical development and implications. It is enough for us to say that a Hamiltonian system is *completely integrable in the open set* $V \subset M$ if one can introduce coordinates $(I, \varphi) = (I_1, \dots, I_n; \varphi_1, \dots, \varphi_n)$ in V such that $\varphi_1, \dots, \varphi_n$ are defined mod 1, $I \in U \subset \mathbb{R}^n$, $\Omega = \sum_{i=1}^n d\varphi_i \wedge dI_i$ in these coordinates, and the Hamiltonian H depends only on I . Such coordinates are usually called *action-angle coordinates* for H . One immediately sees that in action-angle coordinates (1.5.6) implies that $V_H(I, \varphi) = (0, \dots, 0, -\partial H/\partial I_1, \dots, -\partial H/\partial I_n)$. Thus the action variables I_1, \dots, I_n are preserved by the Hamiltonian flow and on each torus $\mathbb{T}_n^c = \{I_i = c_i, i = 1, \dots, n\}$ the flow is linear. However, unlike the case of the linear ordinary differential equation, the *frequency vector* $\omega(I) = (\partial H/\partial I_1, \dots, \partial H/\partial I_n)$ is usually different for different values of $c = (c_1, \dots, c_n)$. Thus in general, a completely integrable system in V looks like a collection of invariant tori with linear flows whose frequency vector and hence the type of recurrence changes from one torus to another. We will return to completely integrable Hamiltonian systems in Section 5.5c. In particular the Liouville Theorem 5.5.21 explains why the above notion of complete integrability is natural.

Exercise

1.5.1. Consider a Lissajous figure on the plane

$$x(t) = A \sin(t + \varphi), \quad y(t) = B \sin(\omega t + \psi) \quad (t \in \mathbb{R}).$$

Prove that if ω is irrational then for any phases φ, ψ the set $\{(x(t), y(t))\}_{t \in \mathbb{R}}$ is dense in the rectangle $|x| \leq A, |y| \leq B$.

6. Gradient flows

Let $S^2 = \{(x, y, z) \mid x^2 + y^2 + z^2 = 1\}$ be the standard unit sphere in \mathbb{R}^3 . We will consider the flow that moves every point downward (or “southward”, if we think of S^2 as the surface of the globe and take the earth’s axis to be vertical) along a great circle (meridian) connecting the point $(0, 0, 1)$ (“the north pole”) and $(0, 0, -1)$ (“the south pole”). The speed of the motion is equal to the derivative of the vertical coordinate along the meridian. In other words, our flow is generated by integrating the following vector field v on the sphere:

$$v(x, y, z) = (xz, yz, -x^2 - y^2). \quad (1.6.1)$$

To see this note that the downward unit vector tangent to the sphere at (x, y, z) is given by $(xz, yz, -(x^2 + y^2))/\sqrt{x^2 + y^2}$. The absolute value of its z -coordinate $\sqrt{x^2 + y^2}$ gives the norm of the gradient vector, which is hence given by (1.6.1). The two poles are the only zeroes of this vector field and consequently they are fixed points for the flow. It is almost obvious that every point except for the north pole asymptotically approaches the south pole as time goes to plus infinity. In fact this convergence is exponential. Similarly, every point except for the south pole exponentially approaches the north pole as time goes to minus infinity.

To generalize this construction, let us consider a Riemannian metric on a compact smooth manifold M and a real-valued function F on M . At each point $x \in M$ that is not a critical point for F one can define the unique direction of fastest increase for F , that is, the unit tangent vector $\zeta(x) \in T_x M$ such that $\mathcal{L}_{\zeta(x)} F = \max_{\eta \in T_x M} \mathcal{L}_{\eta} F / \|\eta\|$, where $\mathcal{L}_{\eta} F$ denotes the Lie (directional) derivative of the function F along the vector η .

We define the gradient vector field ∇F by

$$\nabla F(x) = \begin{cases} \mathcal{L}_{\zeta(x)} F \cdot \zeta(x) & \text{if } x \text{ is noncritical,} \\ 0 & \text{if } x \text{ is critical.} \end{cases}$$

Suppose that in local coordinates (x_1, \dots, x_n) the Riemannian metric has the form $ds^2 = \sum g_{ij}(x_1, \dots, x_n) dx_i dx_j$. Then

$$\nabla F(x_1, \dots, x_n) = G^{-1}(x) \left(\frac{\partial F}{\partial x_1}, \dots, \frac{\partial F}{\partial x_n} \right),$$

where $G(x) = \{g_{ij}(x)\}$ and G^{-1} is the inverse matrix, so it is a smooth vector field on M . The flow generated by the gradient vector field ∇F is called the *gradient flow* of F .

From calculus we know that the gradient is orthogonal to level sets of the function. Via coordinate calculations this is still true in this setting:

Lemma 1.6.1. *The gradient vector field is orthogonal to the level sets.*

Our first example was the gradient flow on the two-sphere provided with the Riemannian metric induced from the standard Euclidean metric in \mathbb{R}^3 for the function $F(x, y, z) = -z$. Let us consider two more examples.

Let M be the two-dimensional torus embedded into \mathbb{R}^3 as a doughnut or bagel standing up, that is, in the position of a doughnut being dunked, and F as before be the negative of the height function z , $F(x, y, z) = -z$. The function F has four critical points on the torus, a maximum A , two saddles B and C , and a minimum D . All orbits of the gradient flow other than those fixed points and six special orbits described below approach the minimum D as time goes to $+\infty$ and the maximum A as it goes to $-\infty$. Two special orbits connect A with B , two more connect B with C , and the last two connect C with D .

Now let us tilt this torus a little bit, that is, change our embedding but keep the function F the same. Equivalently we can consider the same embedding but the function $F = -z + \epsilon x$ for small $\epsilon > 0$. Four critical points remain, as well as the special orbits connecting the maximum with the upper saddle and the lower saddle with the minimum. However, the orbits connecting the two saddles disappear. Instead of these two orbits we will have four: two connecting the maximum with the lower saddle and two connecting the upper saddle with the minimum.

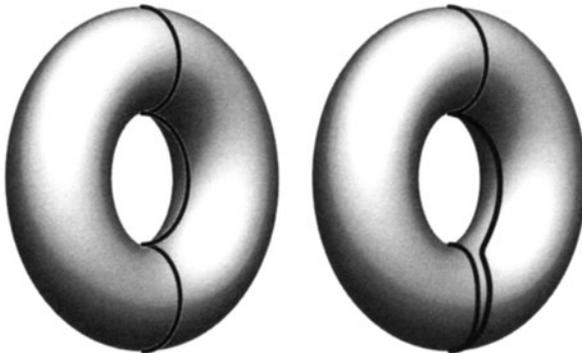


FIGURE 1.6.1. Gradient flows on the torus

Some features of the asymptotic behavior observed in our three examples remain true for general gradient flows. In order to describe those features we

need some general definitions from topological dynamics. We will consider a topological dynamical system defined on the phase space X with discrete or continuous time.

Definition 1.6.2. A point $y \in X$ is called an ω -limit point (correspondingly, an α -limit point) for a point $x \in X$ if there exists a sequence of moments of time going to $+\infty$ (correspondingly, to $-\infty$) such that the images of x converge to y .

The set of all ω -limit (α -limit) points for x is denoted by $\omega(x)$ (correspondingly, $\alpha(x)$) and is called its ω -limit (α -limit) set.

$\omega(x)$ and $\alpha(x)$ are obviously closed and invariant. It follows from the definition that for the dynamical system $\{\varphi^t\}$

$$\begin{aligned}\omega(x) &= \bigcap_{T=0}^{\infty} \overline{\bigcup_{t \geq T} \varphi^t x}, \\ \alpha(x) &= \bigcap_{T=0}^{-\infty} \overline{\bigcup_{t \leq T} \varphi^t x}.\end{aligned}\tag{1.6.2}$$

Thus if X is compact, the sets $\omega(x)$ and $\alpha(x)$ are nonempty and every point sooner or later comes to any given neighborhood of its ω -limit set and stays there.

Let us denote by $\omega_F(x)$ (correspondingly, $\alpha_F(x)$) the ω -limit (α -limit) set of the point $x \in M$ with respect to the gradient flow for the function F .

Proposition 1.6.3. *The sets $\omega_F(x)$ and $\alpha_F(x)$ consist of critical points of F , that is, fixed points of the gradient flow.*

Proof. Let $\{\varphi^t\}_{t \in \mathbb{R}}$ be the gradient flow of the function F . Note that $F \circ \varphi^t$ is nondecreasing (in t) and increases at noncritical points. Thus if $y \in X$ is a noncritical point for F then $F(\varphi^t(y)) > F(y)$ for any $t > 0$. Assume $y \in \omega_F(x)$. Fix $t_0 > 0$ and let $\delta_0 = F(\varphi^{t_0}(y)) - F(y)$. If $x_n \rightarrow y$ then by continuity of the gradient flow $F(\varphi^{t_0}(x_n)) \rightarrow F(y) + \delta_0$. In particular, if $y \in \omega(x)$ then there is a sequence $t_n \rightarrow \infty$ such that $\varphi^{t_n}(x) \rightarrow y$ and hence $F(\varphi^{t_0+t_n}(x)) > F(y) + \delta_0/2$ for sufficiently large n , and since F does not decrease along the orbits, $F(\varphi^t(x)) > F(y) + \delta_0/2$ for sufficiently large t . But this contradicts the convergence $\varphi^{t_n}(x) \rightarrow y$. \square

Proposition 1.6.4. *For any $x \in M$ and any F the set $\omega_F(x)$ is either a single point or an infinite set.*

Proof. Since M is compact, the set $\omega_F(x)$ is nonempty. Suppose $\omega_F(x)$ is finite and $y, z \in \omega_F(x)$, $y \neq z$. We have $y = \lim \varphi^{t_n}(x)$ and $z = \lim \varphi^{s_n}(x)$, where as before $\{\varphi^t\}$ is the gradient flow. Let B be a ball around y and S the boundary of B such that $(B \cup S) \cap \omega_f(x) = \{y\}$. Since the orbit of x enters and leaves B infinitely many times the intersection $\mathcal{O}(x) \cap S$ is an infinite set and by compactness of S it contains a limit point which must belong to $\omega(x)$. \square

Corollary 1.6.5. *If the function F has only isolated critical points then every orbit of the gradient flow of F converges to a critical point of F as $t \rightarrow +\infty$.*

We will see in Section 9.3 how this property of a gradient flow constructed in an auxiliary space helps in finding special orbits for certain dynamical systems.

From a certain formal point of view there is a duality between gradient flows and Hamiltonian dynamical systems. Let us mention it only for the most elementary Euclidean case. In \mathbb{R}^{2n} provided with the standard Euclidean metric the standard Hamiltonian form $\Omega = \sum_{i=1}^n dx_i \wedge dx_{i+n}$ can be expressed via the metric and the operator

$$I = \begin{pmatrix} 0 & \text{Id} \\ -\text{Id} & 0 \end{pmatrix}.$$

Namely, for any two tangent vectors $\xi, \eta \in T_x \mathbb{R}^{2n}$

$$\Omega(\xi, \eta) = \langle \xi, I\eta \rangle,$$

where $\langle \cdot, \cdot \rangle$ is the Euclidean scalar product. Accordingly the Hamiltonian vector fields $V_H = (\partial H/\partial x_{n+1}, \dots, \partial H/\partial x_{2n}, -\partial H/\partial x_1, \dots, -\partial H/\partial x_n)$ can be written as $V_H = I \nabla H$. It is curious to point out that the types of asymptotic behavior for a Hamiltonian vector field on a compact energy manifold $H = \text{const.}$ and for a gradient flow are in a sense diametrically opposite. Whereas for gradient vector fields the only recurrent behavior is represented by fixed points, for Hamiltonian systems nontrivial recurrence is a rule. We have seen this already for completely integrable Hamiltonian systems in the previous section. In general this fact follows from the Liouville Theorem (Proposition 5.5.12) and the Poincaré Recurrence Theorem (Theorem 4.1.19).

Exercises

1.6.1. *Prove that the ω -limit set of any point with respect to a gradient flow is connected.*

1.6.2. *Give an example of a C^∞ function F on a compact manifold and a point $x \in M$ such that $\omega_F(x)$ contains more than one point.*

1.6.3. *For every $g \geq 1$ construct a C^∞ function on the compact orientable surface of genus g with exactly 3 critical points. Describe the dynamics of the gradient flow for that function.*

7. Expanding maps

Consider the following noninvertible map E_2 of the circle: in multiplicative notation

$$E_2(z) = z^2, \quad |z| = 1,$$

or in additive notation

$$E_2(x) = 2x \pmod{1}. \quad (1.7.1)$$

Algebraically this map represents an endomorphism of the group $S^1 = \mathbb{R}/\mathbb{Z}$ onto itself. Geometrically it is a double cover of S^1 .

This is the first example where we will encounter simultaneously and in an essential way nontrivial recurrence as in Sections 1.3–1.5 and different asymptotic behavior for different orbits as in Sections 1.2. and 1.6. The combination of these two phenomena makes the orbit structure for this seemingly very simple transformation much more complicated than everything we have seen so far.

Definition 1.7.1. For a transformation $f: X \rightarrow X$ we let $P_n(f)$ be the number of periodic points of f with (not necessarily minimal) period n , that is, the number of fixed points for f^n .

The following proposition uncovers some of the features of the complicated orbit structure for E_2 .

Proposition 1.7.2. $P_n(E_2) = 2^n - 1$, periodic points for E_2 are dense in S^1 , and E_2 is topologically transitive.

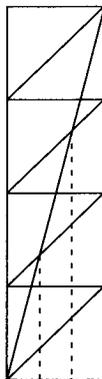


FIGURE 1.7.1. Periodic points for an expanding map

Proof. If $E_2^n(z) = z$, then $z^{2^n} = z$, and $z^{2^n-1} = 1$. Thus every root of unity of order $2^n - 1$ is a periodic point for E_2 of period n . There are exactly $2^n - 1$ such roots of unity. Furthermore, they are uniformly spread over the circle with equal intervals and when n becomes large these intervals become small.

To prove topological transitivity we will consider the binary intervals

$$\Delta_n^k = \left[\frac{k}{2^n}, \frac{k+1}{2^n} \right] \quad \text{for } n = 1, \dots \quad \text{and} \quad k = 0, 1, \dots, 2^n - 1.$$

Let $x = 0.x_1x_2\dots$ be the binary representation of the number $x \in [0, 1]$. Then $2x = x_1.x_2x_3\dots = 0.x_2x_3\dots \pmod{1}$. Thus

$$E_2(x) = 0.x_2x_3\dots \pmod{1}. \quad (1.7.2)$$

Let $k_1\dots k_n$ be the binary representation of the integer k , maybe with several zeroes at the beginning. Then $x \in \Delta_n^k$ if and only if $x_i = k_i$ for $i = 1, \dots, n$. Thus $E_2^n(\Delta_n^k) = S^1$ and since every interval $I \subset S^1$ contains a binary interval, $E_2^n(I) = S^1$ for some n . Thus for any nonempty open sets U, V there is an $n \in \mathbb{N}$ such that $E_2^n(U) \cap V$ is nonempty and by Lemma 1.4.2 E_2 is topologically transitive. \square

The maps

$$E_m: x \mapsto mx \pmod{1},$$

where m is an integer of absolute value greater than one, represent a straightforward generalization of the map E_2 . Not surprisingly these maps are also topologically transitive and have dense sets of periodic orbits. The proof of Proposition 1.7.2 holds verbatim with the replacement of the binary representation by the representation with base m for positive m and with a very minor modification for negative m .

Furthermore, besides periodic and dense orbits there are other types of asymptotic behavior for orbits of expanding maps. One can construct such orbits for E_2 (cf. Exercise 1.7.5) but the simplest and most elegant example appears for the map E_3 .

Proposition 1.7.3. *There exists a point $x \in S^1$ such that in additive notation its ω -limit set with respect to the map E_3 is the standard middle-third Cantor set K . In particular K is E_3 -invariant and contains a dense orbit.*

Proof. The middle-third Cantor set K can be described as the set of all points on the unit interval which have a representation in base 3 with only 0's and 2's as digits (see Exercise 1.7.4). Similarly to (1.7.2) in the base 3 representation the map E_3 acts as the shift of digits to the left. This implies that the set K is E_3 -invariant. It remains to show that E_3 has a dense orbit in K .

Every point in K has a unique representation in base three without ones. Let $x \in K$ and

$$0.x_1x_2x_3\dots \quad (1.7.3)$$

be such a representation. Let $h(x)$ be the number whose representation in base two is $0.\frac{x_1}{2}\frac{x_2}{2}\frac{x_3}{2}\dots$, that is, it is obtained from (1.7.3) by replacing twos by ones. Thus we have constructed a map $h: K \rightarrow [0, 1]$ which is continuous,

monotone (that is, $x > y$ implies $h(x) > h(y)$), and one-to-one except for the fact that binary rationals have two preimages each. Furthermore $h \circ E_3 = E_2 \circ h$. Let $D \subset [0, 1]$ be a dense set of points which does not contain binary rationals. Then $h^{-1}(D)$ is dense in K . This immediately follows from the fact that if Δ is an open interval such that $\Delta \cap K \neq \emptyset$ then $h(\Delta)$ is a nonempty interval, open, closed, or semiclosed. Now take any $x \in [0, 1]$ whose E_2 orbit is dense; the E_3 orbit of $h^{-1}(x) \in K$ is dense in K . \square

Let us emphasize again a crucial difference between all our earlier examples and the expanding maps. In most of the previous examples either the recurrent behavior was very simple, that is, only fixed points, as for contracting maps, hyperbolic linear maps, and gradient flows, or, if nontrivial recurrence was present, all recurrent orbits behaved similarly as for translations and linear flows on the torus. True, for general completely integrable systems different orbits behave differently and nontrivial recurrence takes place. But the phase space of such a system splits into invariant pieces (tori) and all orbits on each torus have the same structure. By contrast, for expanding maps orbits with different behavior (such as periodic, dense, or with a Cantor closure) are interspersed and cannot be separated. This makes the orbit structure very complicated and asymptotic behavior of an *individual* orbit very sensitive to the initial condition and unstable. Furthermore any two orbits will diverge from each other exponentially until they are separated some distance δ . Consequently it is impossible to predict the behavior of an orbit for a long time if the initial position is known only with limited accuracy. For example, iterating E_2 on a computer will clearly yield only as many useful iterations as there are significant binary digits in the initial data. Furthermore any increase in accuracy will only yield a very modest increase in the time over which one can make reasonable predictions: Although doubling the number of significant digits in the initial data and the computation is likely to double the time span of possible prediction, the required improvement in the measurement of the initial state is of astronomical (and illusory) magnitude. Likewise cutting the initial error in half yields only one more step of valid iteration.

Rather surprisingly we will see later in Section 2.4 that the orbit structure *as a whole* is remarkably stable in a certain sense.

There are also important examples of one-dimensional maps which are not expanding. Here is one that we will encounter in the exercises and many times later. For $\lambda \in \mathbb{R}$ let $f_\lambda: \mathbb{R} \rightarrow \mathbb{R}$, $f_\lambda(x) := \lambda x(1 - x)$. For $0 \leq \lambda \leq 4$ the f_λ map the unit interval $I = [0, 1]$ into itself. The family f_λ , $\lambda \in [0, 4]$, is referred to as the quadratic family. It is by far the most popular model in one-dimensional dynamics, both real and complex (in the latter case the maps are extended to \mathbb{C}).¹

Exercises

1.7.1. Calculate the number of periodic points of period n for the map E_m .

1.7.2. Consider the family f_λ .

- (1) Calculate $P_n(f_4)$.
- (2) Calculate $P_n(f_3)$,
- (3) Show that for $\lambda > 3$ there is a period-two orbit.
- (4) Show that for $\lambda \in (3, 1 + \sqrt{6}]$ there are no points of period higher than two.

1.7.3. Show that for any point $x \in S^1$ the set

$$P_x = \{y \in S^1 \mid \exists n \in \mathbb{N} \text{ such that } E_m^n(y) = x\}$$

is dense in S^1 .

1.7.4. Show that the middle-third Cantor set is the set of points on the unit interval which have a representation in base 3 with only 0's and 2's as digits.

1.7.5*. Consider the set T of all points on the unit interval which have a binary representation without two successive zeroes. Prove that T is a perfect nowhere dense set invariant with respect to the map E_2 . Prove that there is a point $x \in T$ whose orbit with respect to E_2 is dense in T .

1.7.6*. Find a point $x \in S^1$ such that $E_3^n(x) \in S^1 \setminus K$ for $n = 0, 1, \dots$ and the orbit closure consists of the orbit itself and the middle-third Cantor set K . In other words, the points of the orbit are isolated in the orbit closure but accumulate to the perfect set K .

8. Hyperbolic toral automorphisms

Hyperbolic toral automorphisms are an invertible analog of the expanding maps E_m . They have very similar properties and analyzing them will give us a chance to preview some of the methods used in the theory of hyperbolic dynamical systems.

We consider the following linear map of \mathbb{R}^2 :

$$L(x, y) = (2x + y, x + y).$$

If two vectors (x, y) and (x', y') represent the same element of \mathbb{T}^2 , that is, if $(x - x', y - y') \in \mathbb{Z}^2$, then $L(x, y)$ and $L(x', y')$ also represent the same element of \mathbb{T}^2 . Thus L defines a map $F_L: \mathbb{T}^2 \rightarrow \mathbb{T}^2$:

$$F_L(x, y) = (2x + y, x + y) \pmod{1}.$$

The map F_L is invertible because the matrix $\begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}$ has determinant one, so L^{-1} also has integer entries. Moreover F_L is an automorphism of the abelian

group $\mathbb{T}^2 = \mathbb{R}^2/\mathbb{Z}^2$. Before we proceed to the study of the map F_L let us discuss some properties of the linear map L .

First, the eigenvalues of L are

$$\lambda_1 = \frac{3 + \sqrt{5}}{2} > 1 \text{ and } \lambda_1^{-1} = \lambda_2 = \frac{3 - \sqrt{5}}{2} < 1.$$

Since the matrix L is symmetric, the eigenvectors are orthogonal. The eigenvectors corresponding to the first eigenvalue belong to the line $y = \frac{\sqrt{5} - 1}{2}x$. The family of lines parallel to it is invariant under L and L uniformly expands distances on those lines by a factor λ_1 . Similarly, there is an invariant family of contracting lines $y = \frac{-\sqrt{5} - 1}{2}x + \text{const.}$

Figure 1.8.1 gives an idea of the action of F_L on the fundamental square $I = \{(x, y) \mid 0 \leq x < 1, 0 \leq y \leq 1\}$. The lines with arrows show the eigendirections.

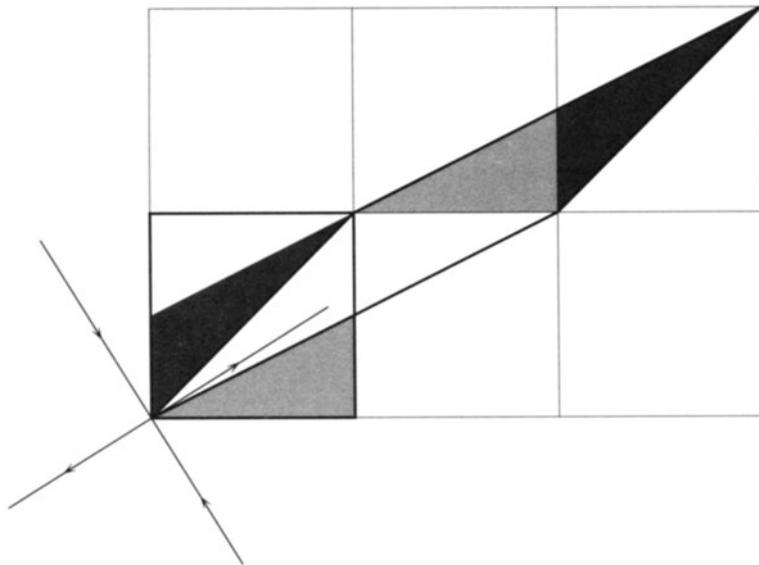


FIGURE 1.8.1. The hyperbolic toral map

Proposition 1.8.1. *Periodic points of F_L are dense, F_L is topologically transitive, and $P_n(F_L) = \lambda_1^n + \lambda_1^{-n} - 2$, where $P_n(F_L)$ is as in Definition 1.7.1.*

Proof. First, let us show that points with rational coordinates are periodic points of F_L . Let $x = s/q, y = t/q$, where s, t, q are integers. Then $F_L \begin{pmatrix} s & t \\ q & q \end{pmatrix} = \begin{pmatrix} 2s + t & s + t \\ q & q \end{pmatrix}$, that is, it is a rational point whose coordinates also have

denominator q . But there are only q^2 different points on \mathbb{T}^2 whose coordinates can be represented as rational numbers with denominator q and all iterates $F_L^n(s/q, t/q)$, $n = 0, 1, 2, \dots$, belong to that finite set. Thus they must repeat, that is, $F_L^n(s/q, t/q) = F_L^m(s/q, t/q)$ for some integers n, m . But since F_L is invertible, $F_L^{n-m}(s/q, t/q) = (s/q, t/q)$. Thus we have proved the density of periodic orbits. Before we proceed further, let us show that points with rational coordinates are the only periodic points for F_L .

Assume $F_L^n(x, y) = (x, y)$. But $F_L^n(x, y) = (ax + by, cx + dy) \pmod{1}$ where a, b, c, d are integers. Thus we have

$$\begin{aligned} ax + by &= x + k, \\ cx + dy &= y + l \end{aligned}$$

for some integers k, l . Since 1 is not an eigenvalue for L^n we determine (x, y) uniquely from a, b, c, d, k, l .

$$x = \frac{(d-1)k - bl}{(a-1)(d-1) - cb}, \quad y = \frac{(a-1)l - ck}{(a-1)(d-1) - cb}.$$

Thus x, y are rational numbers.

The L -invariant families of lines

$$y = \frac{\sqrt{5}-1}{2}x + \text{const.} \quad \text{and} \quad y = \frac{-\sqrt{5}-1}{2}x + \text{const.} \tag{1.8.1}$$

project to \mathbb{T}^2 as F_L -invariant families of orbits of linear flows with irrational slopes. Thus the projection of each line is everywhere dense on the torus.

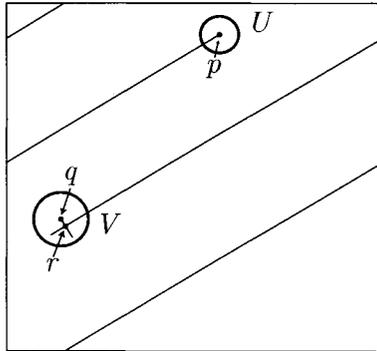


FIGURE 1.8.2. Topological transitivity

Now we are ready to prove that F_L is topologically transitive. Take arbitrary nonempty open subsets U, V of \mathbb{T}^2 . Let $p \in U$, $q \in V$ be two periodic points and let n be their common period. The line from the first family passing through

the point p is invariant under F_L^n and F_L^n expands it with coefficient $\lambda_1^n > 1$. Similarly, the line from the second family passing through the point q is F_L^n invariant and contracting. Let r be a point of intersection of these two lines. $F_L^{kn}(r) \rightarrow q$ as $k \rightarrow +\infty$ and $F_L^{kn}(r) \rightarrow p$ as $k \rightarrow -\infty$. Thus if k is large and positive then $F_L^{-kn}(r) \in U$, $F_L^{kn}(r) \in V$ and $F_L^{2kn}(U) \cap V$ is nonempty.

Finally, let us calculate $P_n(F_L)$. As before, if $F_L^n(x, y) = (x, y)$ then $(a - 1)x + by$ and $cx + (d - 1)y$ are integers. The map $G = F_L^n - \text{Id}: (x, y) \mapsto ((a - 1)x + by, cx + (d - 1)y) \pmod{1}$ is a well-defined noninvertible map of the torus onto itself. Every point of the torus, including the integer point $(0, 0)$, has the same number of preimages which is $|\det(L^n - \text{Id})| = |(\lambda_1^n - 1)(\lambda_1^{-n} - 1)| = \lambda_1^n + \lambda_1^{-n} - 2$. (We have used the fact that the number of points of \mathbb{Z}^2 in $G([0, 1] \times [0, 1])$ is the area of $G([0, 1] \times [0, 1])$.) This is exactly the number of different points $(x, y) \in \mathbb{T}^2$ for which both numbers $(a - 1)x + by$ and $cx + (d - 1)y$ are integers, that is, the number of fixed points for F_L^n . \square

Let us compare the asymptotic properties of the map F_L with those of the toral translations discussed in Section 1.4.

According to Proposition 1.4.1 if the coordinates of the vector γ are rationally independent from 1, then the translation T_γ is topologically transitive; as we have just proved, the automorphism F_L possesses the same property. On the other hand, all orbits of T_γ are dense whereas for F_L dense orbits coexist with a dense set of periodic orbits, each of the latter obviously not being dense. Thus the recurrence of orbits represented by their density is uniform with respect to initial conditions for a translation and is highly sensitive to initial conditions for F_L .

Another aspect of asymptotic behavior is related to regularity of recurrence with respect to time. Topological transitivity implies that iterates of any open set from time to time intersect any other open set. A stronger version of recurrence is reflected by the following property.

Definition 1.8.2. A topological dynamical system $f: X \rightarrow X$ is called *topologically mixing* if for any two open nonempty sets $U, V \subset X$ there exists a positive integer $N = N(U, V)$ such that for every $n > N$ the intersection $f^n(U) \cap V$ is nonempty.

By Lemma 1.4.2 every topologically mixing map is topologically transitive. On the other hand, no translation T_γ is topologically mixing. This follows from the fact that translations preserve the natural metric on the torus induced by the standard Euclidean metric on \mathbb{R}^n and from the following general criterion.

Lemma 1.8.3. *If a topological dynamical system $f: X \rightarrow X$ preserves a metric on X which generates the topology of X , then f is not topologically mixing.*

Proof. Fix an invariant metric on X , take three different points $x, y, z \in X$, and let δ be one tenth of the minimum of the pairwise distances between those points. Let U, V_1, V_2 be δ -balls around x, y, z correspondingly. Since f preserves the diameter of any set, the diameter of $f^n(U)$ does not exceed 2δ whereas the

distance between any two points $p \in V_1$, $q \in V_2$ is greater than 7δ . Thus for each n either $f^n(U) \cap V_1$ or $f^n(U) \cap V_2$ is empty. \square

By contrast, we have the following statement.

Proposition 1.8.4. *The automorphism F_L is topologically mixing.*

Proof. The expanding lines $y = \frac{\sqrt{5}-1}{2}x + \text{const.}$ are orbits of the linear flow T_ω^t , where $\omega = \left(1, \frac{\sqrt{5}-1}{2}\right)$. By Proposition 1.5.1 this flow is minimal. Any open set U contains a piece J of an expanding line. Let us fix $\epsilon > 0$. Then there exists $T = T(\epsilon)$ such that every segment of an expanding line of length T intersects any ϵ -ball on the torus. The existence of at least one such segment follows from the topological transitivity of the flow T_ω^t . But since all segments of a given length are translations of one another this property holds for all segments. Thus for any fixed open set V for any sufficiently large n , $f^n(J)$ intersects V and so does $f^n(U)$. \square

Similarly to the expanding maps of the previous section the behavior of orbits of the map F_L depends on the initial point in a very sensitive way. Any two orbits will diverge from each other exponentially in the future or in the past and often in both past and future (always up to a certain distance). Again this presents difficulties for numerical calculations, for example, if an orbit is computed numerically when the initial condition is known to three decimal places, for example, to within $1/2000$. Then, since $(\frac{3+\sqrt{5}}{2})^8 > 2000$, the error may be of order one already after eight iterations, rendering the calculation useless. Again, cutting the initial error in half barely yields one more step of valid iteration.

The construction of the example from this section can be generalized.

Let $L: \mathbb{R}^m \rightarrow \mathbb{R}^m$ be an $m \times m$ matrix with integer entries, with determinant $+1$ or -1 , and without eigenvalues of absolute value 1, that is, a hyperbolic matrix. Then $L\mathbb{Z}^m = \mathbb{Z}^m$ and L is invertible on \mathbb{Z}^m , so L determines an invertible map of the m -torus $\mathbb{R}^m/\mathbb{Z}^m$ which has properties very similar to those of the map F_L discussed earlier. We will call such a map a *hyperbolic toral automorphism*. Furthermore if one drops the determinant condition, the resulting map still can be defined on the torus although it ceases to be invertible. Those maps are called *hyperbolic toral endomorphisms*. For $m = 1$ these are the expanding maps of the circle.

Exercises

1.8.1. Calculate the number $P_n(H)$ of periodic points of period n for a general hyperbolic automorphism H of the torus. Show that $\lim_{n \rightarrow \infty} P_n(H)/\lambda^n$ exists for some $\lambda > 0$.

1.8.2. Every integer matrix L with determinant ± 1 defines a map of the torus. Show that the resulting map has finitely many periodic points of each period if and only if no eigenvalue of L is a root of unity.

1.8.3*. Show that periodic orbits of any hyperbolic toral endomorphism F_L are dense.

9. Symbolic dynamical systems

a. Sequence spaces. We now introduce a class of topological dynamical systems of particular importance for the theory of smooth dynamical systems. One reason is that in many respects symbolic systems serve as models for smooth ones; it is often easier to see many properties in the symbolic case first and then to carry them over to the smooth case. Second, restrictions of some smooth dynamical systems to various (often important) invariant sets look very much like symbolic systems. Furthermore, symbolic systems can be used to “code” some smooth systems.

Let us consider for each natural number $N \geq 2$ the space

$$\Omega_N = \{\omega = (\dots, \omega_{-1}, \omega_0, \omega_1, \dots) \mid \omega_i \in \{0, 1, \dots, N-1\} \text{ for } i \in \mathbb{Z}\}$$

of two-sided sequences of N symbols and a similar one-sided space

$$\Omega_N^R = \{\omega = (\omega_0, \omega_1, \omega_2, \dots) \mid \omega_i \in \{0, 1, \dots, N-1\} \text{ for } i \in \mathbb{N}_0\}.$$

We can define a topology by noting that Ω_N is the direct product of \mathbb{Z} copies of the finite set $\{0, 1, \dots, N-1\}$, each with the discrete topology, and using the product topology.

Notice that if we consider the finite set $\{0, 1, \dots, N-1\}$ as the finite group $\mathbb{Z}/n\mathbb{Z}$ then this product is naturally a compact abelian topological group.

Fix integers $n_1 < n_2 < \dots < n_k$ and numbers $\alpha_1, \dots, \alpha_k \in \{0, 1, \dots, N-1\}$ and call the subset

$$C_{\alpha_1, \dots, \alpha_k}^{n_1, \dots, n_k} = \{\omega \in \Omega_N \mid \omega_{n_i} = \alpha_i \text{ for } i = 1, \dots, k\} \quad (1.9.1)$$

a *cylinder* and the number k of fixed digits the *rank* of that cylinder. Cylinders in the space Ω_N^R are defined similarly.

An alternative way to define the topology in the space Ω_N (and similarly in Ω_N^R) is by declaring that all cylinders are open sets and that they form a base for the topology. Then every cylinder is also closed because the complement to a cylinder is a finite union of cylinders. The most general open set is a countable union of cylinders.

One more way is to introduce a metric

$$d_\lambda(\omega, \omega') = \sum_{n=-\infty}^{\infty} \frac{|\omega_n - \omega'_n|}{\lambda^{|n|}}$$

for any fixed $\lambda > 1$. These metrics are particularly convenient for large λ , say for $\lambda = 10N$, because any symmetric cylinder $C_{\alpha_{-n}, \dots, \alpha_n}^{-n, \dots, n}$ of rank $2n + 1$ is a ball with respect to such a metric. Ω_N is a perfect, totally disconnected compact space so it is homeomorphic to a Cantor set.

The different metrics d_λ not only define the same topology on Ω_N (although they are not equivalent as metrics) but also determine a Hölder structure. This means that the notion of Hölder-continuous function with respect to the metric d_λ does not depend on λ . That class of Hölder-continuous functions plays an extremely important role in applications to differentiable dynamics (see Chapters 19 and 20) and can be described as follows.

Let φ be a continuous complex-valued function defined on Ω_N or on a closed subset and write $\omega = (\dots, \omega_{-1}, \omega_0, \omega_1, \dots)$ and $\omega' = (\dots, \omega'_{-1}, \omega'_0, \omega'_1, \dots)$. Then for $n = 0, 1, \dots$ let

$$V_n(\varphi) := \max\{|\varphi(\omega) - \varphi(\omega')| \mid \omega_k = \omega'_k \text{ for } |k| \leq n\}.$$

Since Ω_N is compact, φ is uniformly continuous and $V_n(\varphi) \rightarrow 0$ as $n \rightarrow \infty$. We will say that φ has *exponential type* if for some $a, c > 0$

$$V_n(\varphi) \leq ce^{-an}.$$

It is not difficult to see that φ has exponential type if and only if it is Hölder continuous with respect to some (and hence any) metric d_λ . (See Exercise 1.9.1.)

An equivalent way to express independence of the class of Hölder-continuous functions of λ is to point out that for any λ, μ the identity map $\text{Id}: (\Omega_N, d_\lambda) \rightarrow (\Omega_N, d_\mu)$ is Hölder continuous, that is, there exist $a, c > 0$ such that for any $\omega, \omega' \in \Omega_N$ we have

$$d_\mu(\omega, \omega') < cd_\lambda(\omega, \omega')^a. \quad (1.9.2)$$

All of the above discussion translates with obvious changes to the spaces Ω_N^R .

b. The shift transformation. Let us consider the left shift in Ω_N

$$\sigma_N: \Omega_N \rightarrow \Omega_N, \quad \sigma_N(\omega) = \omega' = (\dots, \omega'_0, \omega'_1, \dots), \quad (1.9.3)$$

where $\omega'_n = \omega_{n+1}$.

σ_N is a one-to-one map and takes cylinders into cylinders. Thus it is a homeomorphism of Ω_N . Sometimes the shift σ_N is called a *topological Bernoulli shift*.

Similarly let us define the *one-sided N-shift* $\sigma_N^R: \Omega_N^R \rightarrow \Omega_N^R$ by

$$\sigma_N^R(\omega_0, \omega_1, \omega_2, \dots) = (\omega_1, \omega_2, \omega_3, \dots).$$

This is a continuous noninvertible transformation of the space Ω_N^R onto itself.

The shifts σ_N^R and σ_N possess a number of properties already familiar to us from Sections 1.7 and 1.8.

Proposition 1.9.1. *Periodic points for the shifts σ_N and σ_N^R are dense in Ω_N and Ω_N^R , correspondingly, $P_n(\sigma_N) = P_n(\sigma_N^R) = N^n$, and both transformations σ_N and σ_N^R are topologically mixing.*

Proof. Periodic orbits for a shift are periodic sequences, that is, $(\sigma_N)^m \omega = \omega$ if and only if $\omega_{n+m} = \omega_n$ for all $n \in \mathbb{Z}$ and similarly for σ_N^R . In order to prove the density of periodic points it is enough to find a periodic point in every cylinder. Since any cylinder in Ω_N contains a symmetric cylinder of rank $2m + 1$ for some m such as

$$C_{\alpha_{-m}, \dots, \alpha_m}^{-m, \dots, m} =: C_\alpha^m,$$

where $\alpha = \alpha_{-m}, \dots, \alpha_m$, it is enough to consider only such cylinders. But the sequence obtained by simply repeating the finite sequence $\alpha_{-m}, \dots, \alpha_m$, that is, ω where $\omega_n = \alpha_{n'}$ for $|n'| \leq m$, $n' = n \pmod{2m + 1}$, obviously lies in our cylinder and has period $2m + 1$.

Every periodic sequence ω of period n is uniquely determined by its coordinates $\omega_0, \dots, \omega_{n-1}$. There are N^n different finite sequences $(\omega_0, \dots, \omega_{n-1})$.

Finally, in order to prove topological mixing, it is enough to show that for any $\alpha = \alpha_{-m}, \dots, \alpha_m$ and $\beta = \beta_{-m}, \dots, \beta_m$ and n sufficiently large $\sigma_N^n(C_\alpha^m)$ intersects C_β^m . We take $n > 2m + 1$, say $n = 2m + k + 1$ with $k > 0$. Consider any sequence ω such that

$$\omega_i = \alpha_i \text{ for } |i| \leq m, \quad \omega_i = \beta_{i-n} \text{ for } i = m + k + 1, \dots, 3m + k + 1.$$

Obviously, $\omega \in C_\alpha^m$ and $\sigma_N^n(\omega) \in C_\beta^m$.

The arguments for the one-sided shift are completely similar. \square

Remark. The map $\pi: \Omega_2^R \rightarrow K$, $\pi(\omega_0, \omega_1, \dots) = 0.\beta(\omega_0)\beta(\omega_1)\dots$, where K is the middle-third Cantor set and $\beta(0) = 0$, $\beta(1) = 2$, is a homeomorphism and obviously $\pi \circ \sigma_2^R = E_3 \circ \pi$. Thus Proposition 1.7.3 implies topological transitivity of σ_2^R . This is the simplest example of the situation where a restriction of a smooth system to an invariant set looks like a shift. Accordingly, the correspondence h , described in the proof of Proposition 1.7.3, is the simplest example of coding. We will discuss this topic in greater detail in Sections 2.4 and 2.5.

Definition 1.9.2. The restriction of the shifts σ_N or σ_N^R to any closed, shift-invariant subset of Ω_N or Ω_N^R , respectively, is called a *symbolic dynamical system*.

Properties of symbolic dynamical systems vary widely. These systems provide a rich source of examples and counterexamples for topological dynamics and ergodic theory.

c. Topological Markov chains. Here we will consider only one special (although probably the most important) class of symbolic dynamical systems.

Let $A = (a_{ij})_{i,j=0}^{N-1}$ be an $N \times N$ matrix whose entries a_{ij} are either zeroes or ones. (We call such a matrix a 0-1 matrix.) Let

$$\Omega_A := \{\omega \in \Omega_N \mid a_{\omega_n \omega_{n+1}} = 1 \text{ for } n \in \mathbb{Z}\}. \quad (1.9.4)$$

In other words, the matrix A determines all admissible transitions between the symbols $0, 1, \dots, N-1$. The set Ω_A is obviously shift invariant.

Definition 1.9.3. The restriction

$$\sigma_N \upharpoonright_{\Omega_A} =: \sigma_A$$

is called the *topological Markov chain*¹ determined by the matrix A . Sometimes σ_A is also called a *subshift of finite type*.

There is a useful geometric representation for topological Markov chains. Let us identify the symbols $0, 1, \dots, N-1$ with points x_0, \dots, x_{N-1} and connect x_i with x_j by an arrow if $a_{ij} = 1$. This way we obtain a graph G_A with N vertices and a certain number of oriented edges. We will call a finite or infinite sequence of vertices of G_A an *admissible path* or *admissible sequence* if any two consecutive vertices in the sequence are connected by an oriented arrow. A point of Ω_A corresponds to a doubly infinite path in G_A with a marked origin; the topological Markov chain σ_A corresponds to moving the origin to the next vertex.

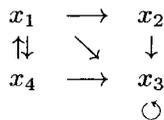


FIGURE 1.9.1. A Markov graph

The following simple combinatorial lemma is a key to the study of topological Markov chains:

Lemma 1.9.4. *For every $i, j \in \{0, 1, \dots, N-1\}$, the number N_{ij}^m of admissible paths of length $m+1$ that begin at x_i and end at x_j is equal to the entry a_{ij}^m of the matrix A^m .*

Proof. We will use induction on m . First, it follows immediately from the definition of the graph G_A that $N_{ij}^1 = a_{ij}$. Let us show that

$$N_{ij}^{m+1} = \sum_{k=0}^{N-1} N_{ik}^m a_{kj}. \quad (1.9.5)$$

For every $k \in \{0, \dots, N-1\}$ every admissible path of length $m+1$ connecting x_i and x_k produces exactly one admissible path of length $m+2$ connecting x_i and x_j by adding x_j to it, if and only if $a_{kj} = 1$. This proves (1.9.5). Now, assuming by induction that $N_{ij}^m = a_{ij}^m$ for all ij , we obtain from (1.9.5) that $N_{ij}^{m+1} = a_{ij}^{m+1}$. \square

Every admissible closed path of length $m + 1$ with marked origin, that is, a path that begins and ends at the same vertex of G_A , produces exactly one periodic point of σ_A of period m . Thus we have

Corollary 1.9.5. $P_n(\sigma_A) = \text{tr } A^n$.

Topological Markov chains can be classified according to recurrence properties of various orbits they contain. Some principal elements of this classification are given in Exercises 1.9.4–1.9.9. Now we will concentrate on the most interesting special class of topological Markov chains that possess the strongest recurrence properties.

Definition 1.9.6. A 0-1 matrix A is called *transitive* if for some positive m all entries of the matrix A^m are positive numbers. We will call a topological Markov chain σ_A *transitive* if A is a transitive matrix.

Lemma 1.9.7. *If all entries of A^m are positive then for any $n \geq m$ all entries of A^n are positive too.*

Proof. First notice that if $a_{ij}^n > 0$ for all i, j , then for each j there is a k such that $a_{kj} = 1$. Otherwise $a_{ij}^n = 0$ for every n and i . Now use induction. Assume that $a_{ij}^n > 0$ for all i, j ; then $a_{ij}^{n+1} = \sum_{k=0}^{N-1} a_{ik}^n a_{kj} > 0$ because $a_{kj} = 1$ for at least one k . \square

Lemma 1.9.8. *If A is transitive and $\alpha = (\alpha_{-k}, \dots, \alpha_k)$ is admissible, that is, $a_{\alpha_i, \alpha_{i+1}} = 1$ for $i = -k, \dots, k-1$, then the intersection $\Omega_A \cap C_\alpha^k =: C_{\alpha, A}^k$ is nonempty and contains a periodic point.*

Proof. Take m such that $a_{\alpha_k, \alpha_{-k}}^m > 0$. Then one can extend the sequence α to an admissible sequence of length $2k + m + 1$ which begins and ends with α_{-k} . Repeating this sequence periodically we obtain a periodic point in $C_{\alpha, A}^k$. \square

Proposition 1.9.9. *If A is a transitive matrix then the topological Markov chain σ_A is topologically mixing and its periodic orbits are dense in Ω_A .*

Proof. The density of periodic orbits follows immediately from Lemma 1.9.8.

In order to establish topological mixing it is enough to show that if for two sequences $\alpha = (\alpha_{-k}, \dots, \alpha_k)$ and $\beta = (\beta_{-k}, \dots, \beta_k)$ the cylinders $C_{\alpha, A}^k$ and $C_{\beta, A}^k$ are nonempty, then for any sufficiently large n the set $\sigma_A^n(C_{\alpha, A}^k) \cap C_{\beta, A}^k$ is also nonempty. Take $n \geq 2k + 1 + m$, where m is taken from Definition 1.9.6, say $n = 2k + 1 + m + l$ with $l \geq 0$. Then $a_{\alpha_k \beta_{-k}}^{m+l} > 0$ by Lemma 1.9.7, so one can construct an admissible sequence of length $4k + 2 + m + l$ whose first $2k + 1$ symbols are identical to α and the last $2k + 1$ symbols to β . By Lemma 1.9.8 this sequence can be extended to a periodic element of Ω_A which obviously belongs to $\sigma_A^n(C_{\alpha, A}^k) \cap C_{\beta, A}^k$. \square

There is a natural class of symbolic systems more general than Markov chains.

Definition 1.9.10. Let $A: \{1, \dots, N\}^{n+1} \rightarrow \{0, 1\}$ and $\Omega_A := \{\omega \in \Omega_N \mid A(\omega_m, \dots, \omega_{m+n}) = 1 \text{ for } m \in \mathbb{Z}\}$. Then the restriction σ_A of σ_N to Ω_A is called an n -step topological Markov chain.

From the point of view of their intrinsic dynamics n -step topological Markov chains are the same as topological Markov chains, since they can be described as topological Markov chains over the alphabet $\{1, \dots, N\}^n$ by taking the matrix A given by $A_{(i_1, \dots, i_n), (j_1, \dots, j_n)} = 1$ if $j_k = i_{k+1}$ for $k = 1, \dots, n-1$ and $A_{(i_1, \dots, i_n), (j_n)} = 1$.

Let $\lambda_1, \dots, \lambda_N$ be the eigenvalues of the matrix A taken with their multiplicities and ordered in decreasing order of their absolute values. By Corollary 1.9.5 we have $P_n(\sigma_A) = \sum_{i=1}^n \lambda_i^n$. For a transitive matrix A a very precise asymptotic of the last expression can be found. It is based on some results about positive matrices which we will present here both for the sake of completeness and with an eye on future uses.

d. The Perron–Frobenius theorem for positive matrices.

Theorem 1.9.11. (Perron–Frobenius Theorem)² Let L be an $N \times N$ matrix with nonnegative entries such that for a certain power L^n all entries are positive. Then L has one (up to a scalar) eigenvector e with positive coordinates and no other eigenvectors with nonnegative coordinates. Moreover, the eigenvalue corresponding to e is simple, positive, and greater than the absolute values of all other eigenvalues.

Corollary 1.9.12. $P_n(\sigma_A) = \lambda_{\max}^n + \mu_n$ for a transitive 0-1 matrix A , where $\lambda_{\max} > 1$ is the eigenvalue corresponding to the positive eigenvector and $|\mu_n| < C\lambda^n$ for some $C > 0$ and $\lambda < \lambda_{\max}$.

Proof of Corollary 1.9.12. All statements except for $\lambda_{\max} > 1$ follow immediately from Corollary 1.9.5 and Theorem 1.9.11. Let $x = (x_0, \dots, x_{N-1})$, $x_i > 0$ for $i = 0, \dots, N-1$, and $Ax = \lambda_{\max}x$. Then $A^n x = \lambda_{\max}^n x$, that is, $\lambda_{\max}^n x_i = \sum_{j=0}^{N-1} a_{ij}^n x_j$ with $a_{ij}^n \geq 1$. Thus $\lambda_{\max}^n x_i \geq \sum_{j=0}^{N-1} x_j > x_i$; hence $\lambda_{\max}^n > 1$ and $\lambda_{\max} > 1$. \square

Proof of Theorem 1.9.11. Let us denote by P the set of all vectors in \mathbb{R}^N with nonnegative coordinates and by σ the unit simplex in P , that is, $\sigma = \{(x_1, \dots, x_N) \mid x_i \geq 0, \sum_{i=1}^N x_i = 1\}$. By assumption $LP \subset P$. Thus for every $x \in \sigma$ there exists a unique vector $Tx \in \sigma$ proportional to Lx . Thus we have defined a map $T: \sigma \rightarrow \sigma$. Obviously for each convex subset $S \subset \sigma$ both the image TS and the preimage are convex. By assumption $L^n P \subset \text{Int } P$; hence $T^n \sigma \subset \text{Int } \sigma$.

For the map T the extreme points of the image of any closed convex set are among images of its extreme points (see Definition A.2.8).

The set $\sigma_0 = \bigcap_{n=0}^{\infty} T^n \sigma \subset \text{Int } \sigma$ is closed, convex, and strictly T -invariant (that is, $T\sigma_0 = \sigma_0$).

Let us show that σ_0 has no more than N extreme points. Let $x \in \sigma_0 \subset T^n \sigma$. Then x is a convex linear combination of extreme points of $T^n \sigma$, but as we pointed out all extreme points of $T^n \sigma$ are among the images of the vertices e_1, \dots, e_N of σ . Thus $x = \sum_{i=1}^N \lambda_i^{(n)} T^n e_i$, where $\lambda_i^{(n)} \geq 0$ and $\sum_{i=1}^N \lambda_i^{(n)} = 1$. One can find a subsequence $n_k \rightarrow \infty$ such that $T^{n_k} e_i$ and $\lambda_i^{(n_k)}$ converge for all $i = 1, \dots, N$. Let $\lim T^{n_k} e_i = p_i$, $\lim \lambda_i^{(n_k)} = \lambda_i$. We have $x = \sum_{i=1}^N \lambda_i p_i$. If x is different from p_i , $i = 1, \dots, N$, it is not an extreme point of σ_0 .

The set of extreme points of σ_0 is thus finite and T -invariant, so all those points are fixed points for T^m for some m . Each of these points corresponds to an eigenvector of L^m with positive coordinates. We are going to show that L^m may have only one (up to a scalar multiple) such eigenvector. Assuming that there are at least two, we consider two cases.

Case 1: All eigenvectors have the same eigenvalue. Thus we have $e, f \in \text{Int } P$, $L^m e = \lambda e$, $L^m f = \lambda f$. It is actually enough to assume only that $f \in \text{Int } P$, because then there exists a positive number α such that the vector $e - \alpha f$ belongs to the boundary of P . Since $T^m(e - \alpha f) = e - \alpha f$, this contradicts the assumption that for large n , $L^n P \subset \text{Int } P$.

Case 2: There are two eigenvectors with different eigenvalues

$$e, f \in P, \quad L^m e = \lambda e, \quad L^m f = \mu f.$$

Obviously λ and μ are positive numbers, so we can assume that $\lambda > \mu$. Consider the plane generated by e and f . The lines containing e and f divide it into four sectors.

Consider the sector S bounded by the half-lines containing f and $-e$. If $x \in S$ then $x = \alpha f - \beta e$, $\alpha, \beta > 0$, so $L^{nm} x = \alpha \mu^n f - \beta \lambda^n e$, that is, the direction of that vector approaches that of $-e$ as $n \rightarrow \infty$. In particular, for a large n , $L^{nm} x$ does not belong to P . However, since $f \in \text{Int } P$, $f - \epsilon e \in P$ for small $\epsilon > 0$. Thus we have found a vector in P which eventually leaves P , contradicting our assumption. This completes the proof of uniqueness of an eigenvector for L because our argument implies that σ_0 consists of a single point which is then fixed for L and thus generates an eigenvector $e \in P$ for L .

Let us show that if $Le = \lambda e$ and μ is another eigenvalue of L , then $\lambda > |\mu|$.

If μ is real, consider an eigenvector f with eigenvalue μ and the plane generated by e and f . We already proved that no other vectors with eigenvalues $\pm \lambda$ exist. Assume $|\mu| > \lambda$. Then as before for $\alpha, \beta > 0$ the direction of the vector $L^n(\alpha e + \beta f)$ approaches the direction of f as $n \rightarrow \infty$ and hence for large n these vectors are outside P . But if ϵ is small then $e + \epsilon f \in P$, a contradiction. Similarly if μ is complex, say $\mu = \rho \cdot e^{i2\pi\varphi}$, one finds an invariant two-plane where L acts as a multiplication by ρ and rotation by φ . If $\rho > \lambda$ we obtain a similar contradiction by considering the action of L on the 3-dimensional space generated by e and that plane.

If $\rho = \lambda$ we take a vector in that 3-space which lies on the boundary of P . This vector either eventually returns to ∂P (if φ is rational) or comes arbitrary

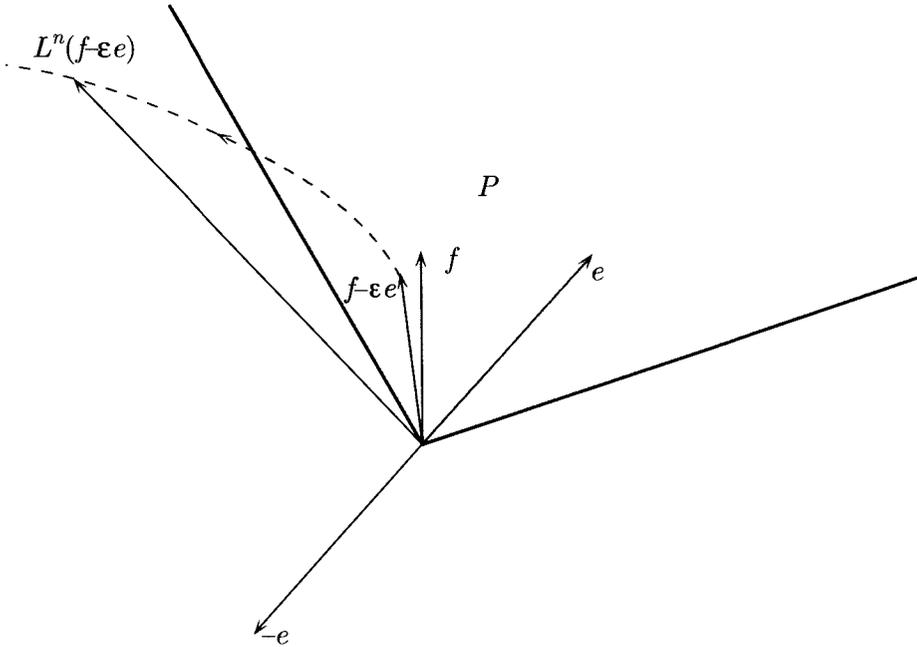


FIGURE 1.9.2. The case of two eigenvalues

close to it (if φ is irrational) contradicting the fact that $L^m P$ lies strictly inside P for large m .

It remains to prove that λ is a simple eigenvalue. We already proved that the space of eigenvectors with eigenvalue λ is one-dimensional. The remaining possibility is that the root space corresponding to λ is more than one-dimensional. But then there exists $f \in \mathbb{R}^N$ such that $Lf = \lambda(f + e)$. Then for small ϵ

$$e - \epsilon f \in P \quad \text{and} \quad L^n(e - \epsilon f) = \lambda^n(-\epsilon f + (1 - \epsilon n)e).$$

For large positive n the direction of the latter vector approaches that of $-e$, that is, it leaves P , which is impossible. \square

Exercises

1.9.1. Justify the statements made at the end of subsection a. Namely:

- (1) Show that a function $\varphi: \Omega_N \rightarrow \mathbb{C}$ is of exponential type if and only if it is Hölder continuous with respect to the metric d_λ for some $\lambda > 1$.
- (2) Prove (1.9.2).

1.9.2. Let d_* be the following metric on Ω_N

$$d_*(\omega, \omega') = \sum_{n=-\infty}^{\infty} \frac{|\omega_n - \omega'_n|}{n^2}.$$

Show that it determines the same topology as any of the metrics d_λ . Show that every function of exponential type is Hölder continuous with respect to d_* but there are also some functions not of exponential type that are Hölder continuous with respect to that metric.

1.9.3. Consider the natural homeomorphism $H: \Omega_2^{\mathbb{R}} \rightarrow K$, where K is the middle-third Cantor set:

$$H(\omega_0, \omega_1, \omega_2, \dots) := 0.\alpha(\omega_0)\alpha(\omega_1)\dots,$$

where $\alpha(0) = 0$, $\alpha(1) = 2$, and the number on the right-hand side is written in base 3.

Prove that the class of functions of exponential type on $\Omega_2^{\mathbb{R}}$ becomes the class of functions which are Hölder continuous with respect to the metric on K induced by the Euclidean metric on $[0, 1]$.

Let us assume from now on that A is a 0-1 $m \times m$ matrix which has at least one 1 in each row and each column.

1.9.4. Prove that for every $i \in \{0, \dots, m-1\}$ the set $\Omega_{A,i} = \{\omega \in \Omega_A \mid \omega_0 = i\}$ is nonempty.

1.9.5. Prove that if there is an element $\omega \in \Omega_A$ that contains the symbol i at least twice then there is a periodic element $\omega' \in \Omega_A$ such that $\omega'_0 = i$.

1.9.6. Let us call symbols i satisfying the condition of the previous problem essential. Prove that any ω -limit point (see Definition 1.6.2) of any element of Ω_A contains only essential symbols.

1.9.7. Let us call two essential symbols i and j equivalent if there exist $\omega, \omega' \in \Omega_A$, $k_1 < k_2$, $l_1 < l_2$ such that

$$\omega_{k_1} = \omega'_{l_2} = i, \quad \omega_{k_2} = \omega'_{l_1} = j.$$

Prove that the set of all essential symbols splits into disjoint subsets of mutually equivalent symbols (that is, this is an equivalence relation).

1.9.8. Prove that σ_A has a dense positive semiorbit if and only if all symbols are essential and equivalent.

1.9.9. Show that under the condition of the previous problem there exists a positive integer N and a decomposition of Ω_A into closed disjoint subsets $\Lambda_1, \dots, \Lambda_N = \Lambda_0$ such that $\sigma_A \Lambda_i = \Lambda_{i+1}$ for $i = 0, 1, \dots, N-1$ and the restriction of $(\sigma_A)^N$ to each Λ_i is topologically mixing. Moreover the decomposition

of Ω_A into Λ_i 's corresponds to a decomposition of the set $\{1, \dots, m\}$ into N equal groups such that every element $\omega \in \Omega_A$ has only symbols from one group in positions equal modulo N .

This is called the *spectral decomposition*.³

1.9.10. Let $B_k = \{\omega \in \Omega_2 \mid \forall m, n \in \mathbb{Z}, m > n, \quad |\sum_{i=n}^m (-1)^{\omega_i}| \leq k\}$. Prove that B_k is a closed σ_2 -invariant subset of Ω_2 . Denote $S_k = \sigma_2|_{B_k}$. Prove that S_k is topologically transitive but not topologically mixing.

1.9.11. Prove that there exists a 0-1 matrix A and a continuous map $H: \Omega_A \rightarrow B_2$ such that the diagram

$$\begin{array}{ccc} \Omega_A & \xrightarrow{\sigma_A} & \Omega_A \\ H \downarrow & & \downarrow H \\ B_2 & \xrightarrow{S_2} & B_2 \end{array} \quad (1.9.6)$$

is commutative and all but two points in B_2 have exactly one preimage.

1.9.12*. Prove that there is no homeomorphism satisfying (1.9.6), that is, the map S_2 is not C^0 equivalent or topologically conjugate (see Definition 2.1.1 and Definition 2.3.1 in the sequel) to any topological Markov chain.

Equivalence, classification, and invariants

In the previous chapter in the process of studying various examples we encountered a number of useful concepts related to the asymptotic behavior of dynamical systems. Our list includes so far the growth of the number of periodic orbits, their spatial distribution (for example, density), topological transitivity, minimality, α - and ω -limit sets, and topological mixing. This list will be considerably extended and systematized in Chapters 3 and 4. Before doing that, we are going to look into the problem of studying the asymptotic behavior of smooth dynamical systems from a different angle.

1. Smooth conjugacy and moduli for maps

a. Equivalence and moduli. We can consider the properties in question as some features of the global orbit structure independent of a particular choice of coordinates. From the global point of view a coordinate change is given by a diffeomorphism (in the case of a smooth structure) or a homeomorphism (for the topological situation) between the phase spaces. Thus we can introduce natural equivalence relations between dynamical systems associated with various classes of coordinate changes and interpret the problem of the description of the orbit structure as the classification of dynamical systems with respect to those equivalence relations.

We begin our discussion with the discrete-time case.

Definition 2.1.1. Two C^r maps $f: M \rightarrow M$ and $g: N \rightarrow N$ are said to be C^m *equivalent* or C^m *conjugate* ($m \leq r$) if there exists a C^m diffeomorphism $h: M \rightarrow N$ such that $f = h^{-1} \circ g \circ h$. h , or the existence of such h , is referred to as a (*smooth*) *conjugacy*.

In other words, C^m equivalence means that f differs from g by a C^m coordinate change. This certainly looks like a natural equivalence relation in differentiable dynamics both from the general structural point of view and with regard to applications.

Definition 2.1.2. Let U be an open subset in the space $C^r(M, M)$ of C^r maps of M into itself with the C^r topology. A continuous functional $F: U \rightarrow \mathbb{R}$ is called a C^m modulus if there exists a $\delta > 0$ such that $F(f) = F(g)$ for any two maps $f, g \in U$ that are C^m equivalent via a diffeomorphism h with $\text{dist}_{C^m}(h, id) < \delta$.

The condition of closeness to the identity becomes particularly important for continuous-time systems. We will discuss appropriate examples in the next section. Right now we will show that at least in many interesting cases with nontrivial recurrence such as those described in Sections 1.7 and 1.8, there are many C^1 moduli.

Let p be a periodic point of period n for f . Obviously for any g that is C^m equivalent to f we have $g^n h(p) = h f^n(p) = h(p)$ so $q = h(p)$ is a periodic point for g of the same period. Thus, $P_n(f) = P_n(g)$ for $n \in \mathbb{N}$. Furthermore, if $m \geq 1$, one has for every (not necessarily periodic) point x and for every n

$$Df_x^n = Dh_{g^n h x}^{-1} \circ Dg_{h x}^n \circ Dh_x.$$

In particular if $f^n p = p$ then

$$Df_p^n = (Dh_p)^{-1} Dg_{h p}^n Dh_p$$

because in this case $g^n h(p) = h(p)$ and $(Dh^{-1})_{h p} = (Dh_p)^{-1}$. Thus the linear operators Df_p^n and $Dg_{h p}^n$ are conjugate and in particular they have the same eigenvalues. Let us call the set of eigenvalues of Df_p^n the *spectrum* of the periodic point p .

Invoking Proposition 1.1.4 one sees that every periodic point p of f , whose spectrum does not contain one, determines several C^1 moduli. If we assume for simplicity that the eigenvalues of Df_p^n are simple, those eigenvalues can serve as the moduli. Since such periodic orbits are isolated, their spectra can be perturbed separately, at least for any finite collection of points. Thus, the moduli obtained from different periodic orbits are in a certain sense independent.

b. Local analytic linearization. On the other hand, at least in some cases *locally* the spectrum is a complete invariant of smooth conjugacy. Right now we will give a very simple example of such a situation which also sheds some light on certain methods used in more analytical aspects of the theory of smooth dynamical systems. Different approaches to the problem of local smooth conjugacy will be discussed in Sections 2.8 and 6.6.

Proposition 2.1.3. Let $I = [-\delta, \delta]$, $f: I \rightarrow I$ be a real analytic contracting map, $f(0) = 0$, and $0 \neq \mu := f'(0)$. Then there are an interval $J_1 \subset I$ containing 0 and a real-analytic diffeomorphism $h: J_1 \rightarrow J_2 \subset \mathbb{R}$ preserving the origin and conjugating f with the linear map $x \mapsto \mu x$.

There are also versions of Proposition 2.1.3 for C^∞ and C^r maps. The C^∞ version is contained in Theorem 6.6.6.

Proof. We will show how to find a formal power series for the conjugacy and then prove convergence of this series. This is a very simple example of the majorization method which is widely used in many local and semilocal problems concerning conjugacy of dynamical systems.

By a formal power series we mean an expression $u = \sum_{i=0}^{\infty} u_i x^i$ which we do not assume to converge anywhere except for $x = 0$. Given two formal power series $u = \sum_{i=0}^{\infty} u_i x^i$ and $v = \sum_{i=0}^{\infty} v_i x^i$ we say that v majorizes u if $|u_i| \leq v_i$ for all i . We write this as $u \prec v$ and call v a majorant for u . We use the same notation when u and v are analytic functions. Thus, for example, $1 \prec (1-x)^{-1}$ and $1 + 2x \prec (1-x)^{-2}$.

Write $f(x) = \sum_{i=1}^{\infty} f_i x^i$ and let $\epsilon = \min_{n \geq 2} |f_n|^{-1/(n-1)} > 0$ (since f is defined on a neighborhood of 0). Then

$$f'(x) := f(\epsilon x)/\epsilon = \lambda x + \sum_{i=2}^{\infty} f_i \epsilon^{i-1} x^i \prec |\lambda|x + \frac{x^2}{1-x} =: |\lambda|x + F(x).$$

If $h'(\lambda x) = f'(h'(x))$ and $h := h'/\epsilon$ then

$$f(h(x)) = \epsilon f'(h(x)/\epsilon) = \epsilon f'(h'(x)) = \epsilon h'(\lambda x) = h(\lambda x),$$

so we may assume $f \prec |\lambda|x + F(x)$. If $h(x) = x + \tilde{h}(x)$ and $f(x) = \lambda x + \tilde{f}(x)$ then $h(\lambda x) = f(h(x))$ becomes $\lambda x + \tilde{h}(\lambda x) = \lambda h(x) + \tilde{f}(h(x))$ or $\tilde{h}(\lambda x) - \lambda \tilde{h}(x) = \tilde{f}(h(x))$. For the coefficients of $h(x) = x + \sum_{i=2}^{\infty} h_n x^n$ this means

$$(\lambda^k - \lambda)h_k = (\tilde{f} \circ h)_k. \quad (2.1.1)$$

The right hand side involves only coefficients of h of order lower than k because \tilde{f} begins with quadratic terms, so this determines the coefficients of h uniquely by recursion starting from $h_0 = 0$ and $h_1 = 1$. It remains to show that this power series converges. To that end let

$$c = \max_{k \geq 2} \frac{1}{|\lambda^k - \lambda|} = \frac{1}{|\lambda|} \max_{k \geq 2} \frac{1}{|1 - \lambda^{k-1}|} = \frac{1}{|\lambda|} \frac{1}{1 - |\lambda|}$$

and \bar{h} the power series with coefficients $\bar{h}_k = |h_k|$. Then for $k \geq 2$

$$\bar{h}_k = |h_k| = \frac{1}{|\lambda^k - \lambda|} (\tilde{f} \circ h)_k \leq \frac{1}{|\lambda^k - \lambda|} (F \circ \bar{h})_k \leq c(F \circ \bar{h})_k.$$

Here we used that if $u \prec U$ and $v \prec V$ with $V(0) = 0$ then $u \circ v \prec U \circ V$. To see this, note that the k th order coefficient on either side of $u \circ v \prec U \circ V$ is a polynomial with positive coefficients in the coefficients of u and v on the one hand and U and V on the other hand. The polynomials are the same on either side and the arguments on the right side are larger in absolute value, giving larger values.

Thus $\bar{h} \prec x + cF \circ \bar{h} = x + c \frac{\bar{h}^2}{1 - \bar{h}} \prec \frac{x + c\bar{h}^2}{1 - \bar{h}}$, because $x \prec x(1 + \bar{h} + \bar{h}^2 + \dots) = x/1 - \bar{h}$. Note that in the expression $\frac{x + cu^2}{1 - u} = (x + cu^2)(1 + u + u^2 + \dots)$, where u is any formal power series, the k th order coefficient is given by a polynomial $P_k(u_0, \dots, u_{k-1})$ in lower order coefficients of u all of whose coefficients are positive. Thus

$$H = \frac{x + cH^2}{1 - H} \quad (2.1.2)$$

defines a formal power series H by $H_k = P_k(H_0, \dots, H_{k-1})$ and inductively

$$|h_k| = \bar{h}_k \leq P_k(h_0, \dots, h_{k-1}) \leq P_k(H_0, \dots, H_{k-1}) = H_k, \quad (2.1.3)$$

or $h \prec H$. H converges because $(c+1)H^2 - H + x = 0$ or $a^2H^2 - 2aH + 2ax = 0$, where $a = 2(c+1)$, so $(1 - aH)^2 = 1 - 2aH + a^2H^2 = 1 - 2ax$ and

$$1 + 2aH \prec 1 + 2aH + 2a^2H^2 + \dots = (1 - aH)^{-2} = \frac{1}{1 - 2ax},$$

hence $2aH \prec \frac{1}{1 - 2ax} - 1 = \frac{2ax}{1 - 2ax}$ and $H \prec \frac{x}{1 - 2ax}$ converges for $|2ax| < 1$. \square

Proposition 2.1.3 is a particular case of Theorem 2.8.2. (Although Theorem 2.8.2 deals with transformations in a complex domain one easily sees from the proof that if f has real coefficients then so does the conjugating map h .) The proof of Theorem 2.8.2 will serve as an illustration of the fast-converging iteration method, sometimes called the Newton method. The method, described in a general way in Section 2.7, is one of the most powerful and versatile tools in the theory of smooth dynamical systems, especially for problems related to smooth equivalence. Its special importance is due to the fact that it is applicable to situations where, unlike in our case, no hyperbolicity is present.

c. Various types of moduli. Let us come back to the discussion of the global orbit structure. Our construction of independent moduli associated with periodic orbits shows that in the case of infinitely many periodic orbits as for the expanding maps E_m (Section 1.7) and the hyperbolic toral automorphism F_L (Section 1.8) there are infinitely many invariants of local C^1 equivalence. It turns out that for both those cases which represent the simplest examples of hyperbolic systems (cf. Section 6.4 and Part 4 of this book), the spectra of periodic orbits form a complete system of invariants for C^1 and even C^∞ equivalence in a neighborhood of the maps E_m and F , correspondingly. For C^1 conjugacy of toral maps this is contained in Theorem 20.4.3. A reasonable description of the set of possible values for eigenvalues of *all* periodic points remains an open problem.

A modulus of a different kind gives substantial although not complete information about smooth equivalence in a neighborhood of the circle rotation

R_α . The rotation number (see Definition 11.1.2) is a C^0 modulus and for some irrational values of α its levels determine the smooth equivalence class (see Theorem 12.3.1).

On the other hand, in many situations the set of all C^r equivalence classes for $r \geq 1$ is both too big and does not admit any reasonable structure. The case of $r = 0$, that is, the *topological* equivalence of *smooth* dynamical systems, is strikingly different and will be discussed in Section 2.3.

Now we will give an idea of how C^r classification might look for systems with a very simple recurrence pattern. For that purpose we consider a monotone analytic map φ of the unit interval $I = [0, 1]$ to itself, which fixes the endpoints, has no other fixed points, and is such that $\varphi'(0) > 1$, $\varphi'(1) < 1$, for example,

$$\varphi(x) = -\frac{x^2}{2} + \frac{3}{2}x. \quad (2.1.4)$$

Thus $x = 0$ and $x = 1$ are attracting fixed points for negative and positive iterates of the map φ , respectively, that is, any point in between tends to 0 and 1 for negative and positive iterates accordingly (cf. Proposition 1.1.6).

First we will show that such a map φ intrinsically defines two smooth affine structures on the open interval $(0, 1)$.

Lemma 2.1.4. *Any C^1 map defined in a neighborhood of the origin on the real line and commuting with a linear contraction $\Lambda: x \rightarrow \lambda x$, $|\lambda| < 1$, is linear.*

Proof. Let $f: [-\epsilon, \epsilon] \rightarrow \mathbb{R}$ be such a map. First, f must preserve the origin because $f(0)$ is a fixed point for Λ which is unique. Furthermore, the commutation condition implies $f(\lambda x) = \lambda(f(x))$ and inductively

$$f(\lambda^n x) = \lambda^n f(x). \quad (2.1.5)$$

Since f is differentiable at 0, $f(\lambda^n x)/\lambda^n x$ has a limit as $n \rightarrow \infty$ that is independent of x and will be denoted by a . By (2.1.5)

$$a = \lim_{n \rightarrow \infty} \frac{f(\lambda^n x)}{\lambda^n x} = \frac{f(x)}{x},$$

that is, $f(x) = ax$. □

Remark. The differentiability condition is important. By contrast there are many nonlinear Lipschitz maps commuting with Λ .

Corollary 2.1.5. *Let h_1, h_2 be two diffeomorphisms satisfying the assertion of Proposition 2.1.3. Then there exists a real number μ such that $h_2(x) = h_1(\mu x)$.*

In other words, each real-analytic contracting map preserves a uniquely defined smooth affine structure. For the map φ we have been discussing, the two structures, defined near the endpoints of the interval, meet in the middle. The transition function between the two structures at any fundamental domain $J = [a, \varphi(a)]$ provides an infinite-dimensional space of moduli for φ . Let us